

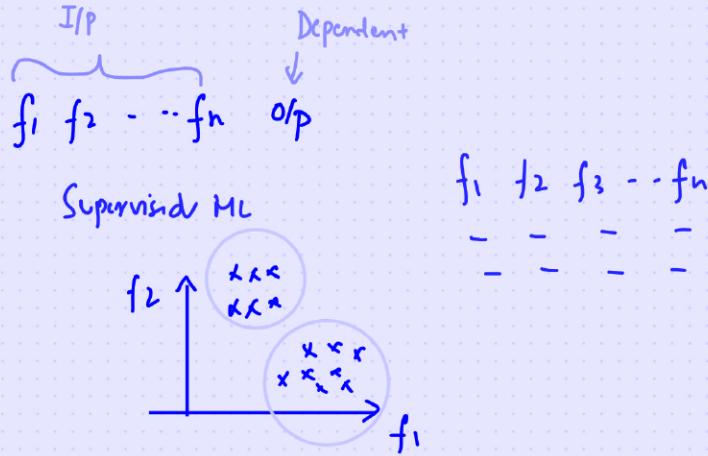
# Unsupervised Machine Learning

## Clustering Algorithms

① K Means Clustering

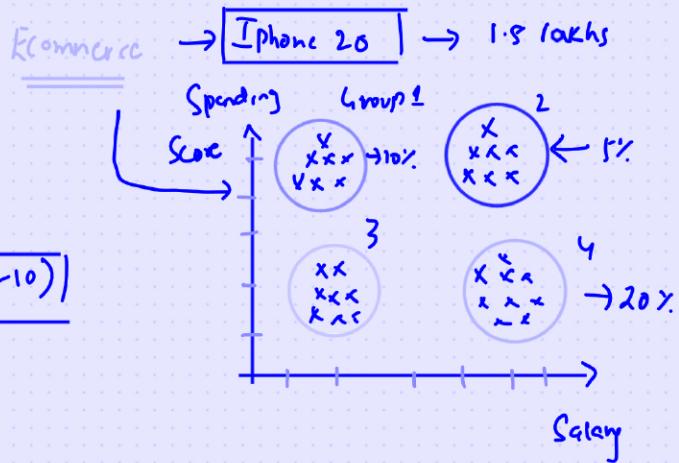
② Hierarchical Clustering.

③ DBScan Clustering.



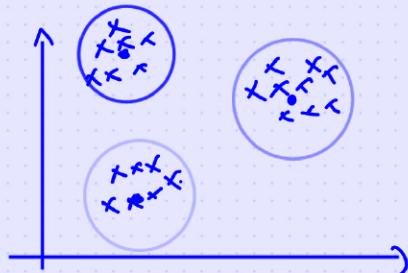
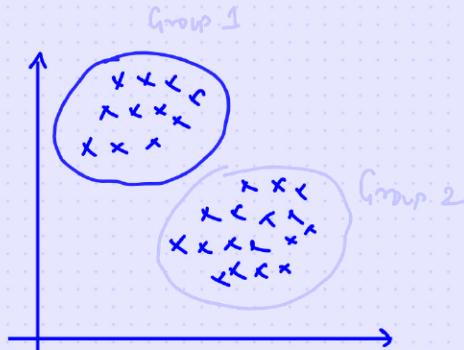
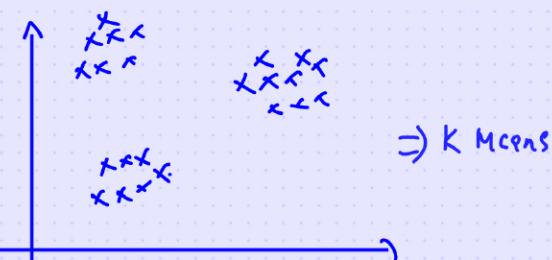
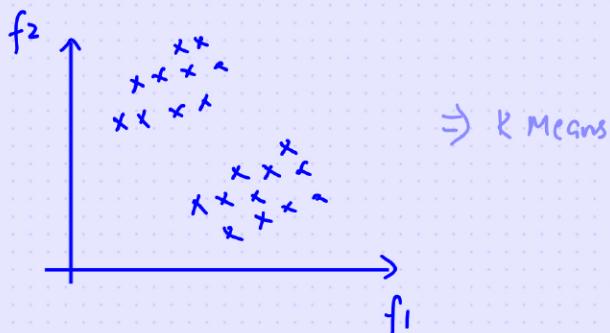
## Dataset

Salary    Spending Score (1-10)



① K Means clustering

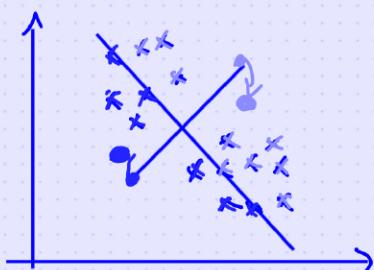
## Geometric Intuition



## K-Means Mathematical Intuition

Steps :

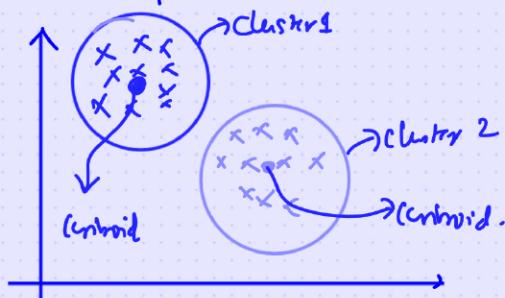
K=2



① Initialize some  $K \rightarrow$  centroids

② Points that are nearest to the centroid  $\rightarrow$  Group

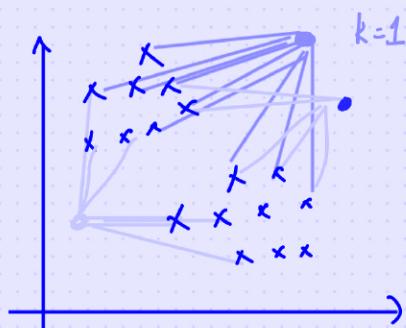
③ Move the centroids  $\rightarrow$  Mean



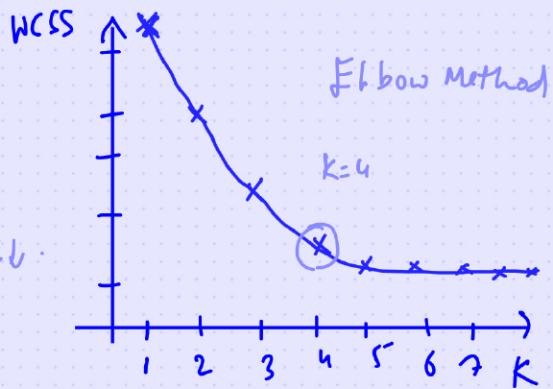
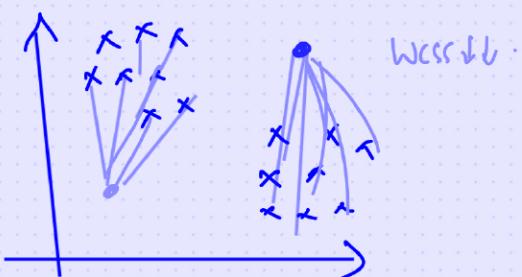
## How do we select the K Value?

WCSS = Within Cluster Sum of Squares

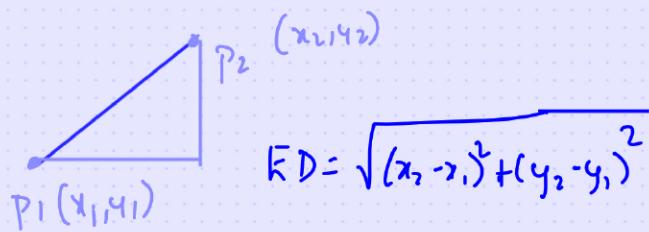
Initialize  $K=1$  to 20



$$WCSS = \sum_{i=1}^K (\text{Distance between point to the nearest centroid})^2$$



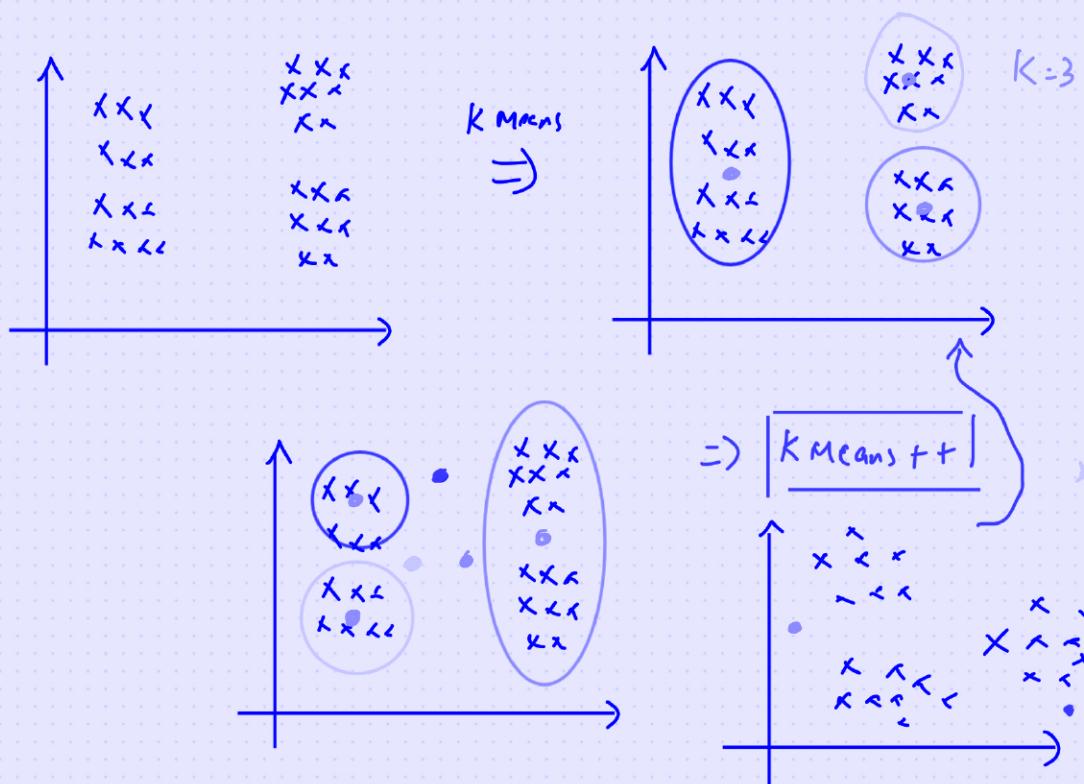
Euclidean Distance or Manhattan Distance



$$MD = |x_2 - x_1| + |y_2 - y_1|$$

## Random Initialization Trap (K Means ++)

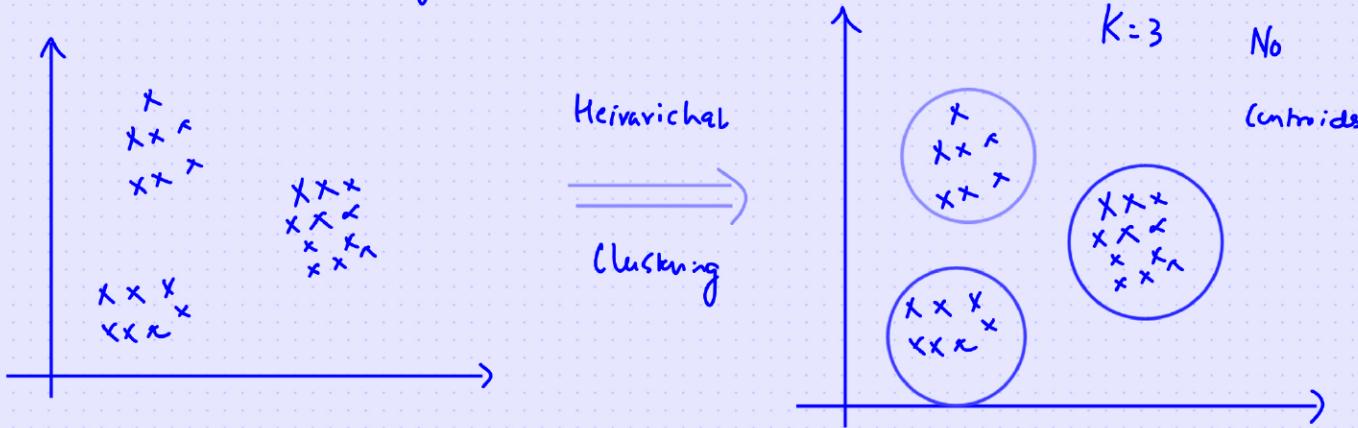
---



## K Means ++ Initialization Technique

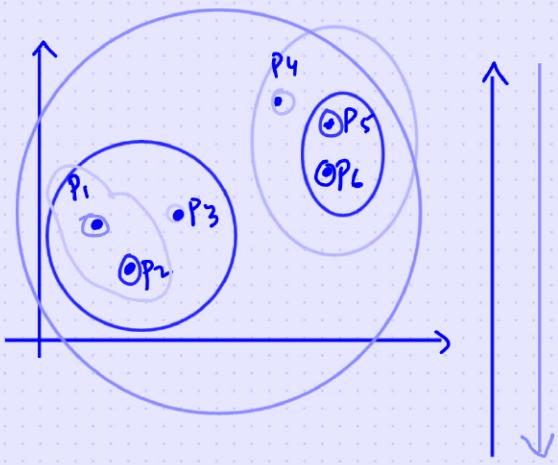
---

## Hierarchical Clustering



## HC

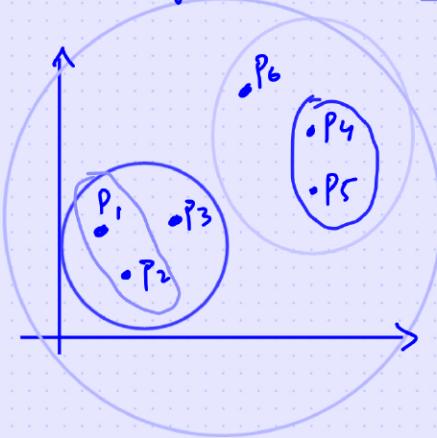
- ① Agglomerative
  - ② Divisive
- }  $\Rightarrow$  Geometric Intuition



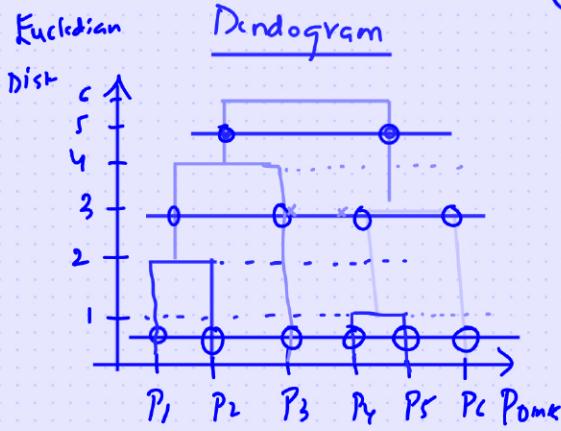
## Steps

- ① For each point initially we will consider it as a separate cluster
- ② Find the nearest point and create a new cluster
- ③ Keep on doing the same process [Step 2] until we get a single cluster.

How many clusters? K=2



## Threshold



- ④ Select the longest vertical line such that no horizontal line passes through it.

# K Means Vs Hierarchical Clustering

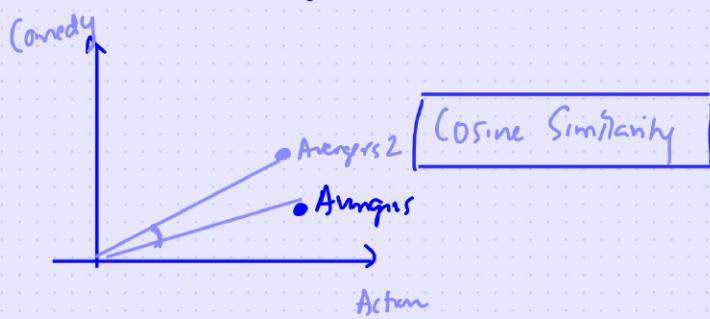
## Scalability And Flexibility

① Dataset size → Huge ⇒ K Means

→ Small ⇒ Hierarchical clustering

② Types of Data → Numerical data → K Means or Hierarchical

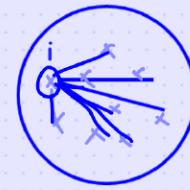
→ Variety of data → Hierarchical



# Silhouette Clustering

(1)

$$a(i)$$



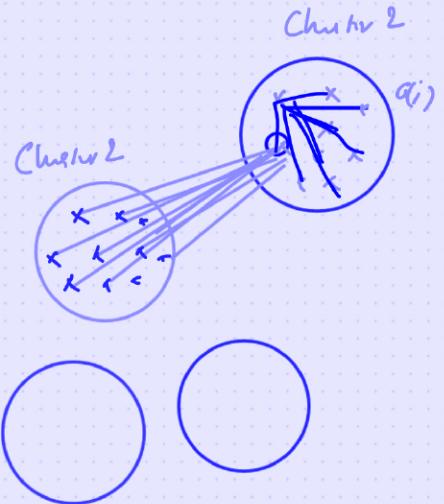
For data point  $i \in C_I$  (data point  $i$  in the cluster  $C_I$ ), let

$$a(i) = \frac{1}{|C_I| - 1} \left[ \sum_{j \in C_I, i \neq j} d(i, j) \right]$$

be the mean distance between  $i$  and all other data points in the same cluster, where  $|C_I|$  is the number of points belonging to cluster  $i$ , and  $d(i, j)$  is the distance between data points  $i$  and  $j$  in the cluster  $C_I$  (we divide by  $|C_I| - 1$  because we do not include the distance  $d(i, i)$  in the sum). We can interpret  $a(i)$  as a measure of how well  $i$  is assigned to its cluster (the smaller the value, the better the assignment).

(2)

$$b(i)$$



We then define the mean dissimilarity of point  $i$  to some cluster  $C_J$  as the mean of the distance from  $i$  to all points in  $C_J$  (where  $C_J \neq C_I$ ).

For each data point  $i \in C_I$ , we now define

$$b(i) = \min_{J \neq I} \frac{1}{|C_J|} \left[ \sum_{j \in C_J} d(i, j) \right]$$

to be the *smallest* (hence the **min** operator in the formula) mean distance of  $i$  to all points in any other cluster (i.e., in any cluster of which  $i$  is not a member). The cluster with this smallest mean dissimilarity is said to be the "neighboring cluster" of  $i$  because it is the next best fit cluster for point  $i$ .

### ③ Silhouette Score

Small  $\left| \frac{a_i}{b_i} \right| \Rightarrow \text{Good cluster}$

We now define a *silhouette* (value) of one data point  $i$ :

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}, \text{ if } |C_I| > 1$$

and

$$s(i) = 0, \text{ if } |C_I| = 1$$

Which can be also written as:

$$s(i) = \begin{cases} \left| \frac{b(i) - a(i)}{b(i)} \right| & \text{if } a(i) < b(i) \\ 0 & \text{if } a(i) = b(i) \\ \left| \frac{b(i) - a(i)}{a(i)} \right| & \text{if } a(i) > b(i) \end{cases}$$

From the above definition it is clear that

$$\{-1 \leq s(i) \leq 1\}$$

$\left| \frac{a_i}{b_i} \right| \Rightarrow \text{bad cluster}$