

# Unsupervised Machine learning

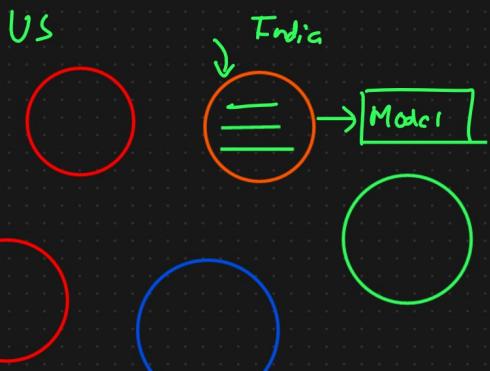
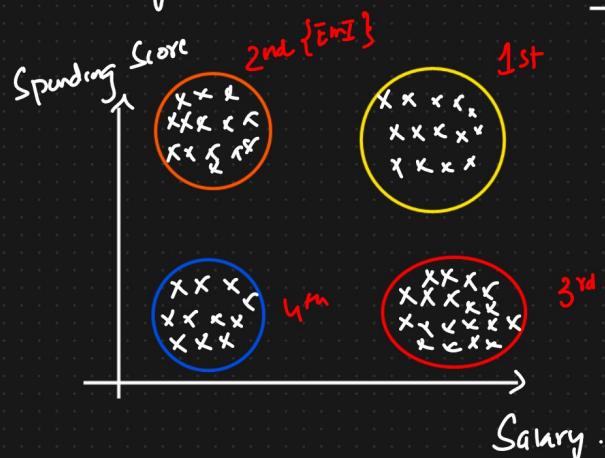
## Clustering [grouping]

- ① K Mean Clustering
- ② Hierarchical clustering
- ③ DBSCAN clustering

Performance Metrics Validate clustering algorithm.

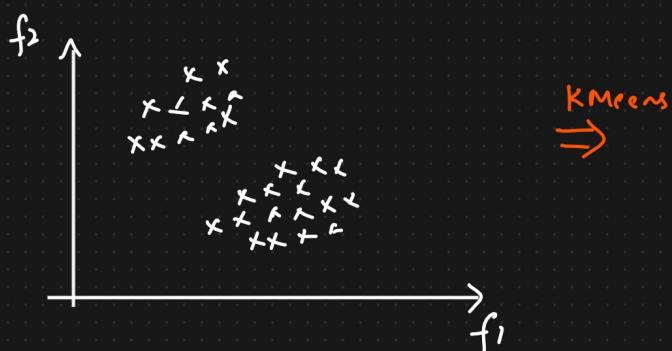


### Silhouette score

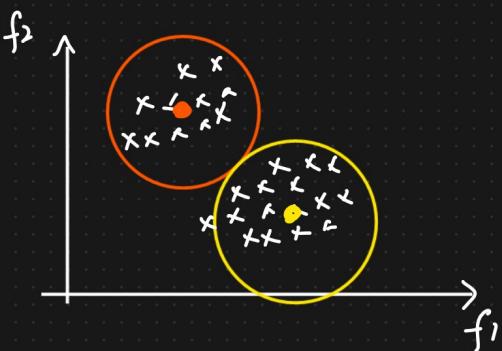


## ① K Means Clustering

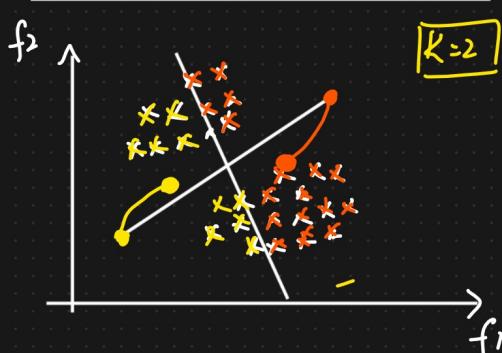
### Geometric Intuition



KMeans  
⇒



### K Means Mathematical Intuition

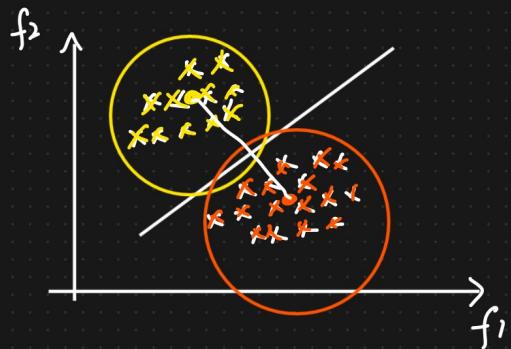
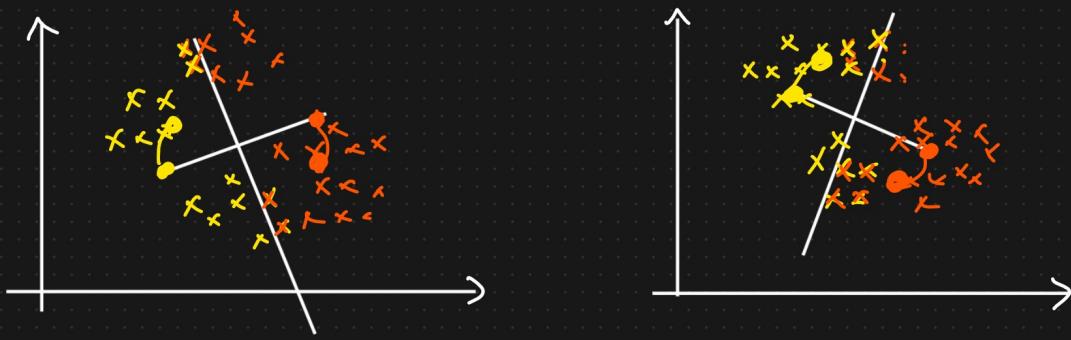


|K=2|

### Steps.

K centroids

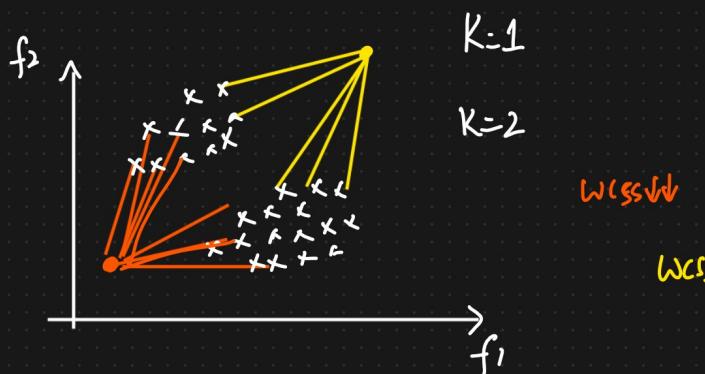
- ① Initialize some K value
- ② Label all the points based on distance to the Nearest Centroid
- ③ Move the Centroid



How do we select the K value Elbow method  $\Rightarrow$  Knee locator

WCSS = Within Cluster Sum of Square.

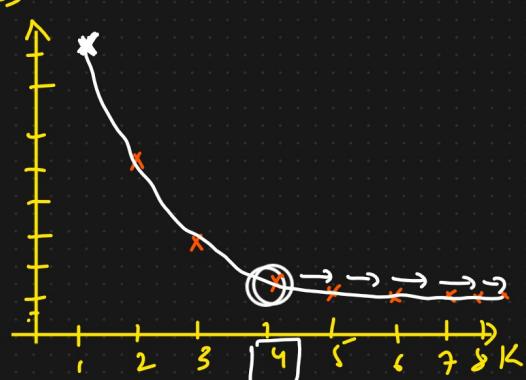
Initialize K: 1 to 10



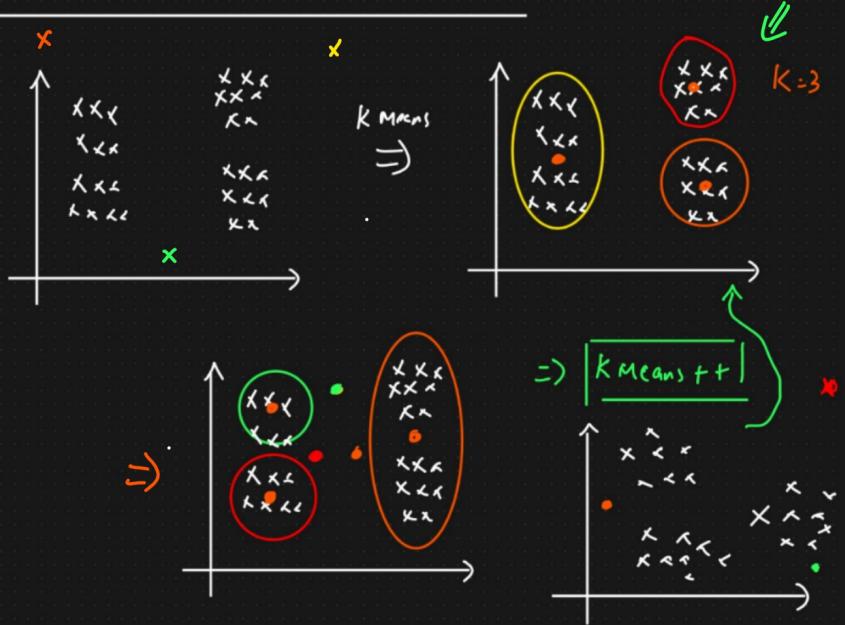
$WCSS = \sum_{i=1}^K (\text{Distance between point to the nearest centroid})^2$

Eucleidian Distance

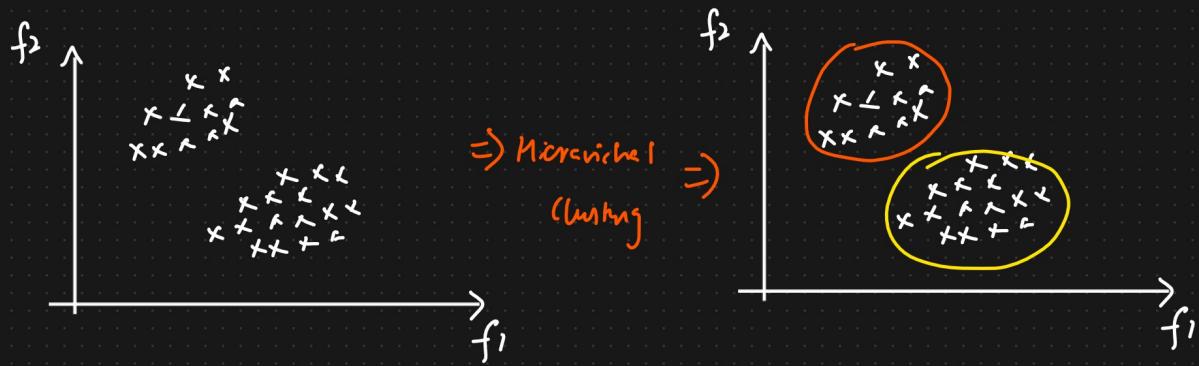
K = Number of centroids:



## Random Initialization Trap ( $K$ Means + r) $\Rightarrow$

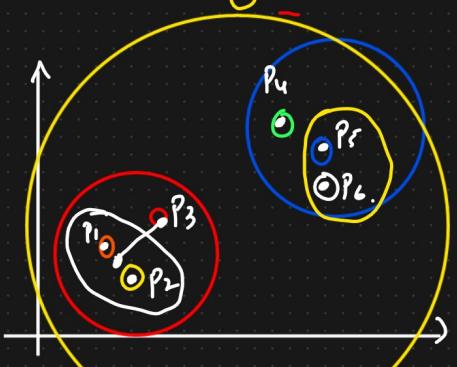


## ② Hierarchical Clustering [Agglomerative Clustering]



### HC

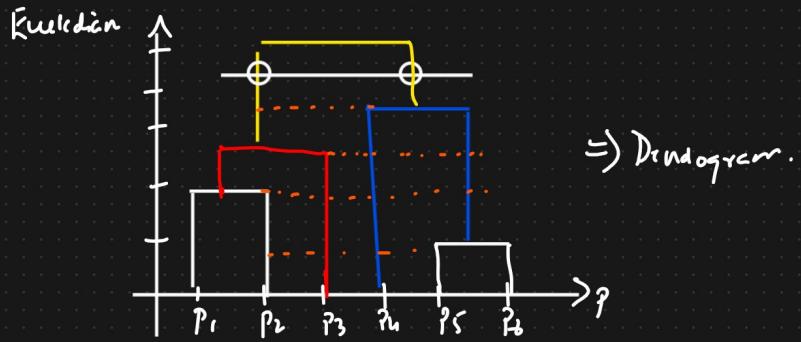
- ① Agglomerative
  - ② Divisive
- $\Rightarrow$  Geometric Intuition



### Steps

- ① For each point we will consider it as a separate cluster
- ② Find the nearest point and create a new cluster

Kullback Distance

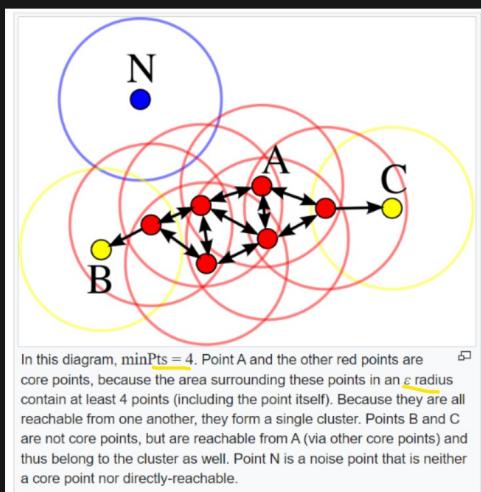


## k Means Vs hierarchical Clustering

## Scalability and Flexibility

① Dataset size  $\xrightarrow{\text{large}} \text{K Means}$   
 $\xrightarrow{\text{Small}} \text{Hierarchical clustering}$

### ③ DB Scan Clustering



- $\longrightarrow$  Core points
  - $\longrightarrow$  Border points
  - $\longrightarrow$  Noise/outliers

} Non linear  
Clustering

## Hyperparameters

$$\textcircled{1} \boxed{\minpls = 4} \quad \textcircled{2} \quad E = \text{radius}$$

① No. of points within the  $E \geq \min$  is

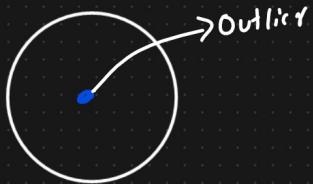


## ② Border points

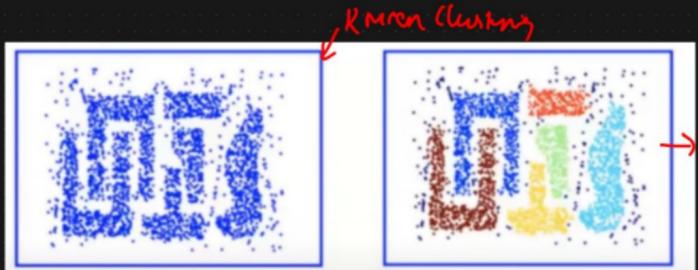
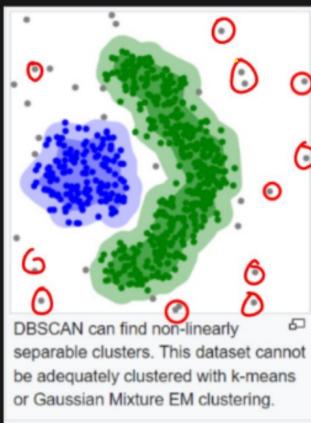
No. of data points within the radius is less than minpts=4



③ Outliers [DBSCAN is robust to outliers].



Some Examples after we apply DBScan Clustering



The left image depicts a more traditional clustering method that does not account for multi-dimensionality. Whereas the right image shows how DBSCAN can convert the data into different shapes and dimensions in order to find similar clusters.

Silhouette Score [Validate the clustering model].

① First step:

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

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

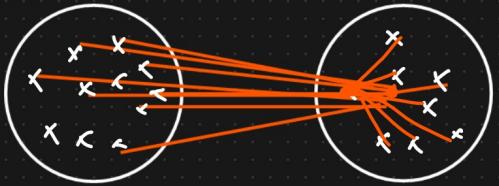
be the mean distance between i and all other data points in the same cluster, where  $|C_I|$  is the number of points belonging to cluster  $C_I$ , and  $d(i, j)$  is the distance between data points i and j in the cluster  $C_I$  (we divide by  $|C_I| - 1$  because we do not include the distance  $d(i, i)$  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).

a(i)



b(i) > a(i)  $\Rightarrow$  clustering is good.

b(i)



②

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

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

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:

$$s(i) = \begin{cases} 1 - a(i)/b(i), & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1, & \text{if } a(i) > b(i) \end{cases}$$

From the above definition it is clear that

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