

→ ~~to~~ unsupervised ML

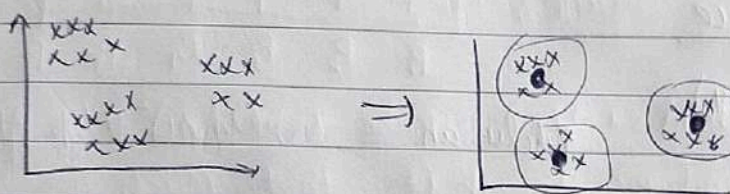
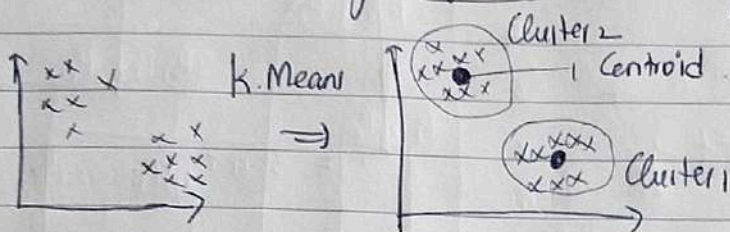
||

Group (or) Cluster

## unsupervised ML

- ① K-Means Algorithm
  - ② Hierarchical Clustering
  - ③ DBScan - Clustering
  - ④ Silhouette - Scoring
- } Validate

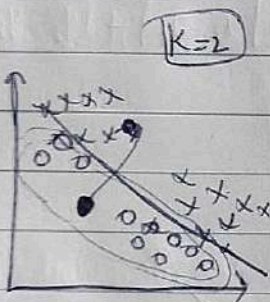
(It have centroid to each cluster then we draw perpendicular line the which  
K-Means Clustering Algorithm is near combine and make clusters the  
 Centroids change. This steps goes on  
 till all clusters are correct



euclidean Distance

(or)

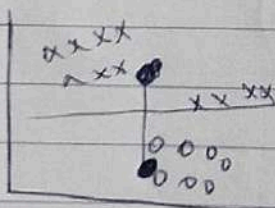
Manhattan Distance



① Initize some  $K \rightarrow$  centroid

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

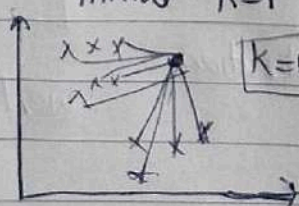
③ move the centroids  $\rightarrow$  Average



How do we select the  $K$  value?

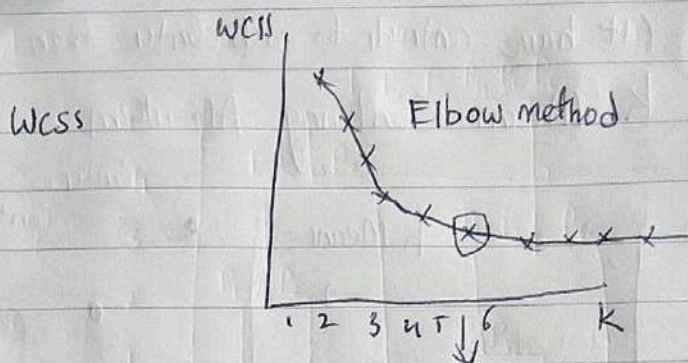
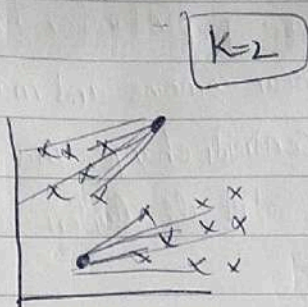
Wcss = within cluster sum of squares

Initial  $K=1$  to 20

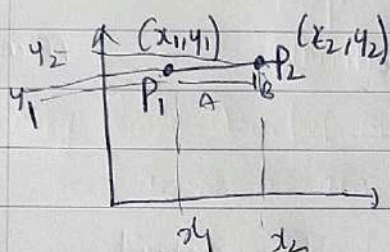


$$K=1 \quad Wcss = \sum_{i=1}^n (\text{distance b/w points to nearest centroid})^2$$





### \* Euclidean Distance

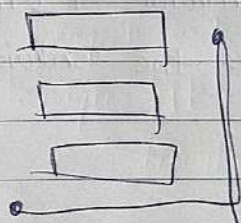


$$\text{Euclidean} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

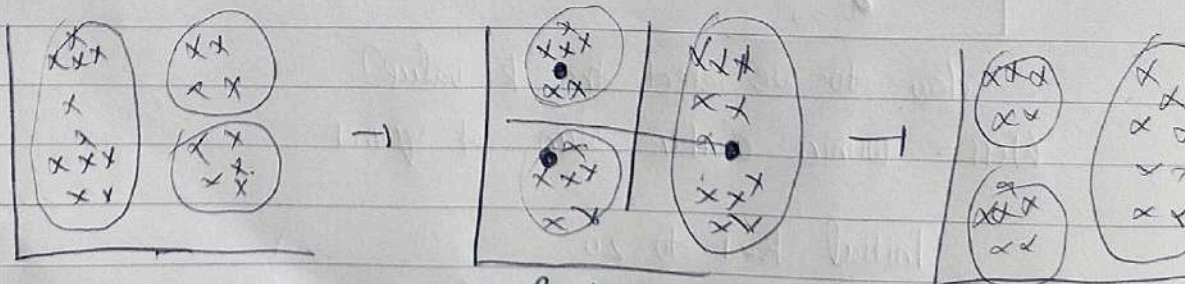
$$\text{Manhattan dist} = |x_2 - x_1| + |y_2 - y_1|$$

IRON MAN → U.S

Air traffic



In this scenario, K-means<sub>1</sub> went wrong :-



Random initialization  
centroid

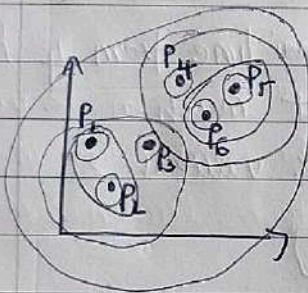
→ went wrong



k mean ++ - Initialization Technique - far far distance with each other

## Hierarichal Clustering Intuition

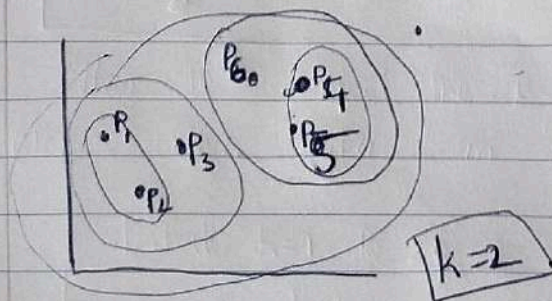
- ① Agglomerative } geometric
- ② Divisive



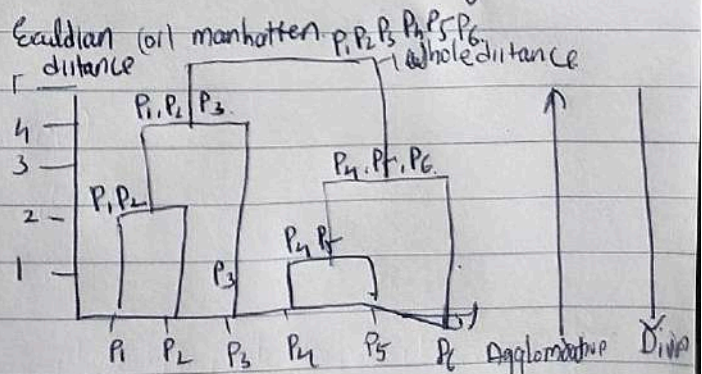
Steps:-

- ① For each data point initially will consider it as a separate ~~class~~ cluster. ( $P_1, P_2, P_3, P_4, P_5, P_6$ )
- ② We find the nearest point and create a new cluster.
- ③ Keep on doing the same process until get a single cluster.

Dendrogram (we combine all cluster we have know how many clusters we have) to take



Threshold - is 4

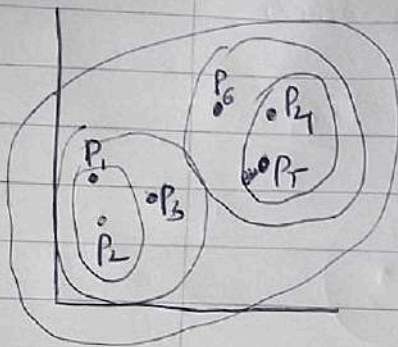


Dendrogram

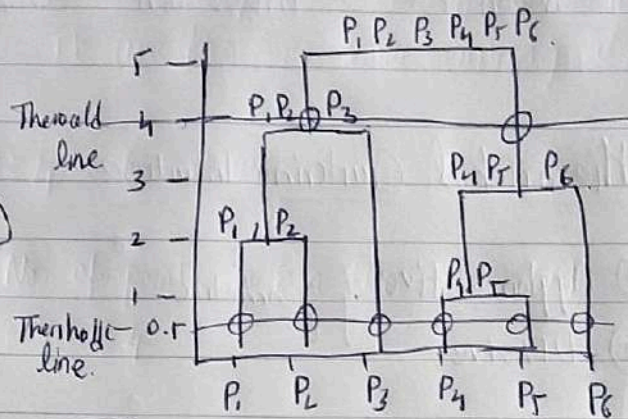


$K = \text{clusters}$  if  $K=2$  we are having 2 clusters

cosine similarity or euclidean distance



$K=2$



Keep Threshold  $\rightarrow 4$  make ~~an~~ line from 4. how many lines touches or goes through it that much is  $K$  value. see above 2 are touches so our ~~to~~  $K$  value is 2.

if Threshold is  $0.5$   $K$  value is 6  $K=6$

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

Threshold (such as  $y$  axis)



## K-mean Vs Hierarchical clustering

### Scalability & Flexibility

① Dataset size  $\rightarrow$  Huge  $\rightarrow$  K-Means

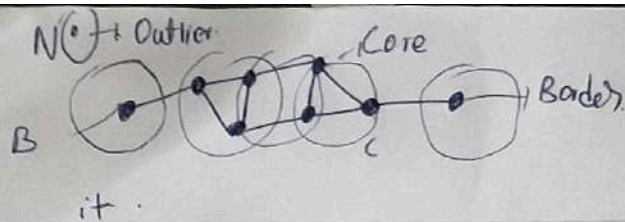
Small  $\rightarrow$  Hierarchical cluster (Dendrogram not see clear)

② K-mean  $\rightarrow$  Numerical data.

Hierarchical clustering  $\rightarrow$  Numerical and other also } Important  
(Variety of data)

③ Centroid  $\rightarrow$  Elbow method  $\rightarrow$  No. of centroids

Core point  
Border point  
Outlier

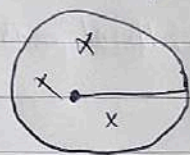


DBSCAN - Clustering - (Data having noise it can be handle well)  
it use for non linear clustering algs

Core point - minimum = 4  $\epsilon$  = radius

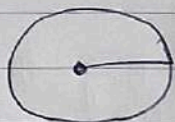
① - No. of points within the radius ( $\epsilon$ ) should be greater  $\geq 4$ .

Border point - minimum = 4



- No. of data points within in this radius will be less than minimum points.

Outlier



- no other datapoints are exists