

# Unsupervised Learning

- It requires only input variable ( $x$ ).

## K-means

- K-means comes, whole unsupervised learning which makes use of unlabelled data.
- K-means is a clustering algorithm which clusters unlabelled data into different cluster based on similarity.

Ex:- Tinder, customer segmentation, fraud detection.

Labelled data  $\rightarrow$  Both  $x \& y$

Unlabelled data  $\rightarrow$  only  $x$ .

similarity  $\rightarrow$  nearest distance.

$K \rightarrow$  no. of cluster.

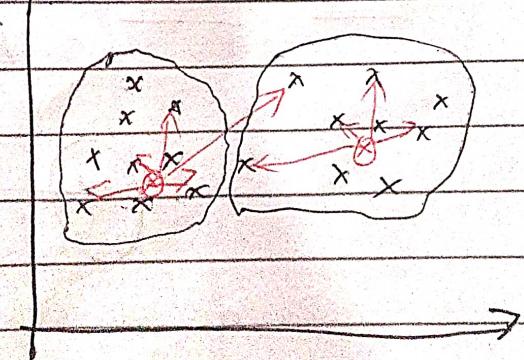
- K-mean is distance based algorithm

- It will make use of distance measures like euclidean, manhattan, & minkowski to find distance b/w any two observations

How K-mean works:-

① plot data

② Define no. of clusters,  $K = 2$   
(group)



③ It will create a random cluster.

④ It will initialise the centroids

centroids :- cluster centers, centroids  
are calculated by taking average of all  
the observation in the cluster.

$$\text{no of cluster} = \frac{\text{no of observations}}{\text{no of centroids}}$$

⑤ find distance b/w centroid and all other  
observation.

⑥ Assign each observation to the nearest  
centroid based on distance.

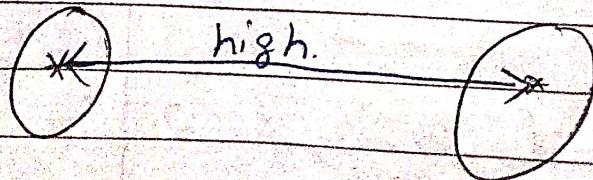
⑦ It will repeat steps from 4<sup>th</sup> step  
till it gets clearer cluster / no overlapping

### Goal / Aim:

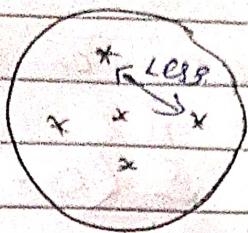
- Inter cluster distance should be high.

- Intra cluster distance should be less

→ Inter cluster :- distance b/w clusters should  
be high.



→ Intra cluster :- distance within cluster should be less



→ How to evaluate k-mean?

$$\text{silhouette score} = \frac{b - a}{\max(a, b)}$$

a → Intra cluster distance

b → Inter cluster distance.

$$\text{Range} = [-1, 1]$$

→ near 1 → clearer clusters

→ near 0 @ 1 → clusters are overlapping

→ Imp:-  
• scaling  
• handling outliers.

→ How to choose optimal value for K?

- Elbow method

- It is an iterative process which works iteratively for different k-values and gives the optimal k.

$$K = [2, 3, 4, \dots, 11]$$

- Take different k-values.

- Implement K-mean.

- calculate NCSS

[within cluster, sum of squared]

$$NCSS = \sum_{i=1}^n (c_p - x_i)^2$$

- NCSS is distance b/w centroid & observation.

- plot graph of K vs NCSS

- choose the point after which you see breakdown in NCSS

