

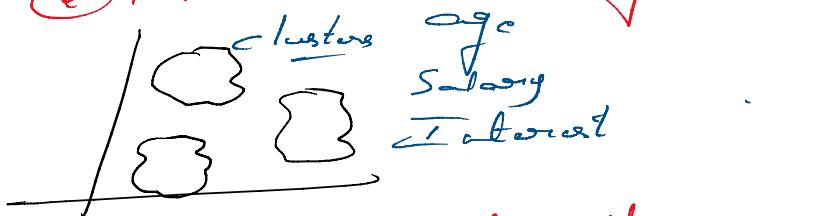
Unsupervised Learning

Training machine learning model with only input variables

→ unsupervised learning is clustering algorithm which clusters entire data into different groups based on similarity.

Eg. ① Market Basket Analysis

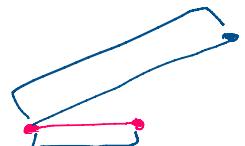
② Mall Customer Segmentation



→ K-Means Algorithm

- * K-Mean algorithm comes under unsupervised learning and also called as clustering algorithm.
- * K-Mean is a clustering algorithm which is used to classify unlabelled data $[x]$ into groups/clusters based on similarity.

Eg. - Mall Customer Segmentation



$\sqrt{K} \rightarrow$ No of clusters

Similarity → Nearest distance b/w observations

Distance Measures.

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

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

How K-Mean Algorithm works?

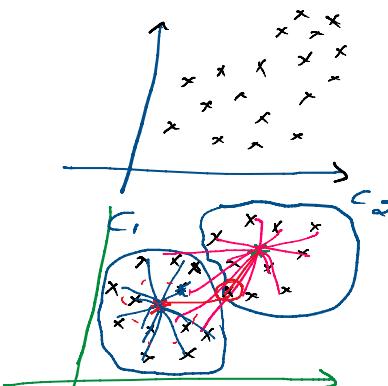
① Plot data

② Define no of clusters
 $K=2 \rightarrow$ Create two clusters

③ It will randomly Create two clusters.

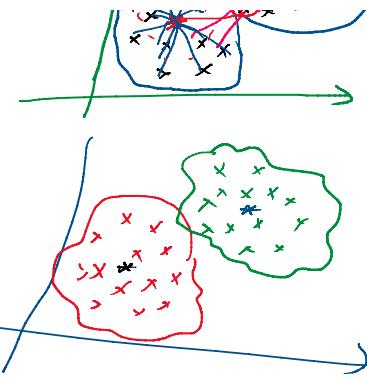
④ Initialize Centroids for each Cluster.

Centroids → Cluster centers



④ Initialize centroids in cluster.

Centroids \rightarrow Cluster Centers
Centroids are found by taking average of all the points in each cluster.



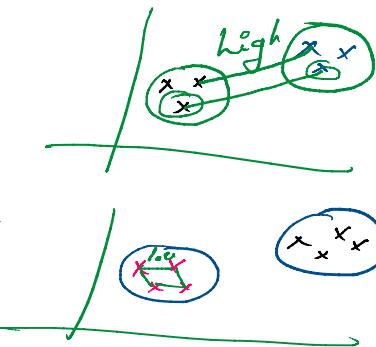
⑤ Assign each observation to the nearest cluster based on distance.
 \rightarrow find distance between data and centroid, if it's nearer to first cluster then it belongs to first cluster.

⑥ Reinitialize the centroids

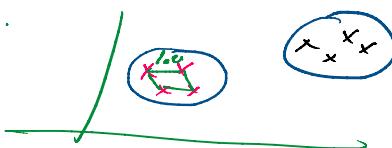
⑦ Repeat 5th step until you get clearer clusters.

Aim of K-Mean Algorithm

* Intercluster distance should be high.
 \rightarrow Distance b/w observations in two clusters should be high.



* Intradistance distance should be low.
 \rightarrow Distance of observations within the cluster should be very low.



How to Evaluate K-Mean Model?

$$\text{Silhouette Score} = \frac{b - a}{\max(a, b)}$$

$a \rightarrow$ distance within the cluster / intradistance

$b \rightarrow$ distance between the clusters / interdistance

Range of Silhouette = $[-1, 1]$.

value is near to 1 \rightarrow Clearer clusters

value is near to -1 \rightarrow Clusters are not clear / overlapping.