

K-Means Clustering.

Supervised Machine Learning {

- Linear Regression → Regression (Continuous)
- Logistic Regression → Classification (Categories)

~~H.W.~~ 

→ Why Logistic Regression has regression in its name, but it actually helps us to perform Classification Task.

Unsupervised Machine Learning (e.g. Algorithm (K-Means))

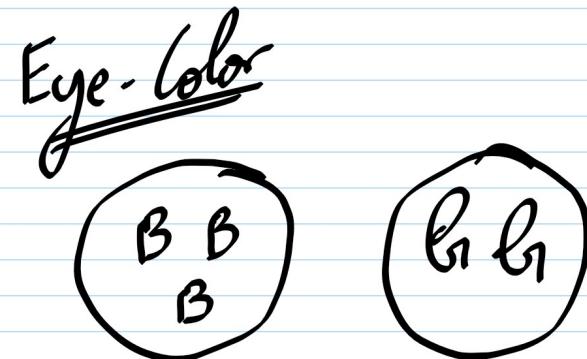
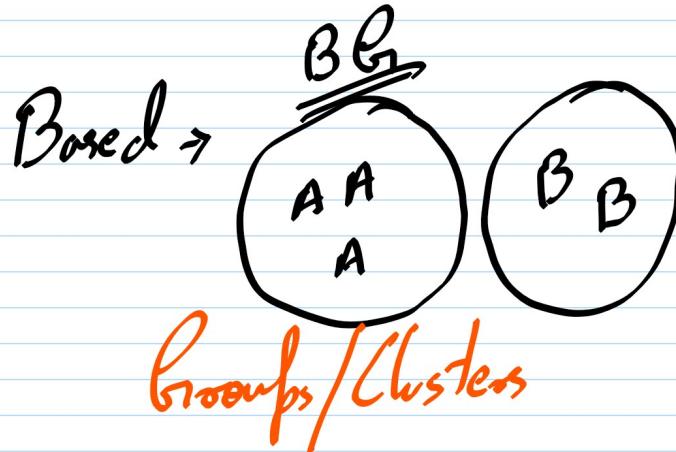
- DBScan, Hierarchical Clustering.

K-Means Clustering

→ It is an unsupervised ML model used for cluster or group the data together.

→ e.g.

	<u>Age</u>	<u>Blood Group</u>	<u>Gender</u>	<u>Eye-color</u>	<u>label</u>
→ -	-	A	M	B	↓ Absent
→ -	-	B	M	B	
→ -	-	A	F	B	
⋮	⋮	A	F	B	
⋮	⋮	B	M	B	

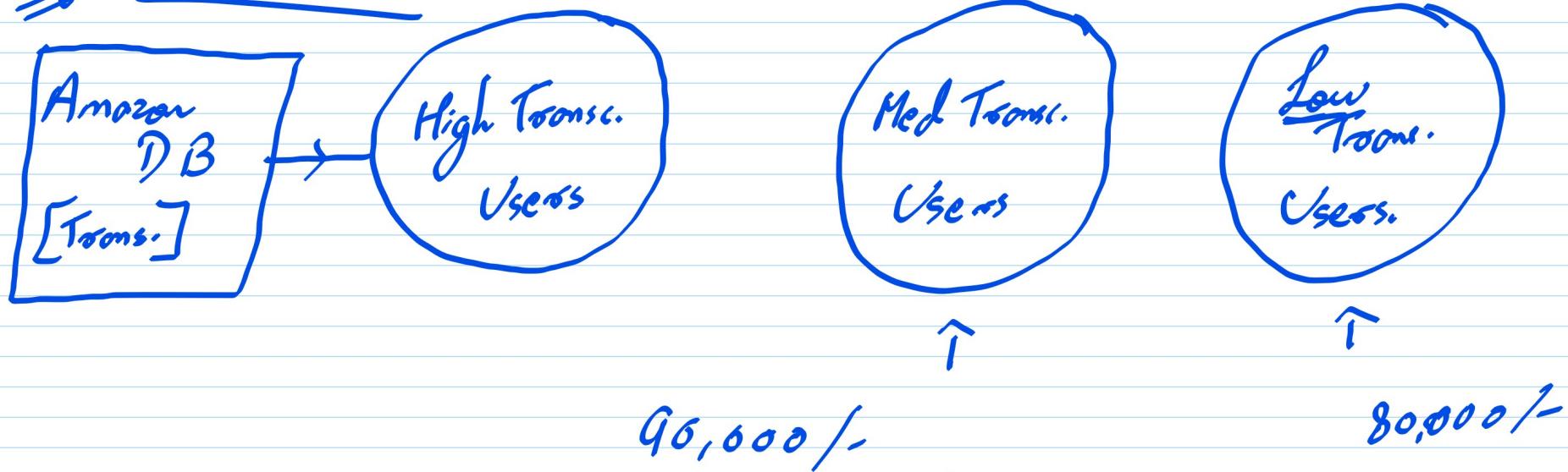


→ K-Means groups the data point together on the basis of certain common attributes.

Real World Application :-

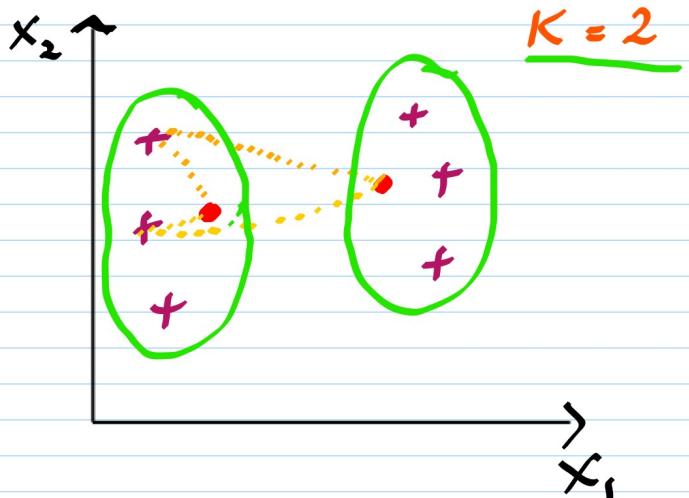
\Rightarrow Customer Profiling \Rightarrow Grouping similar kind of customers.

e.g. E-commerce.



If targets particular group of customers.

Q) How K-Means Clustering Work?



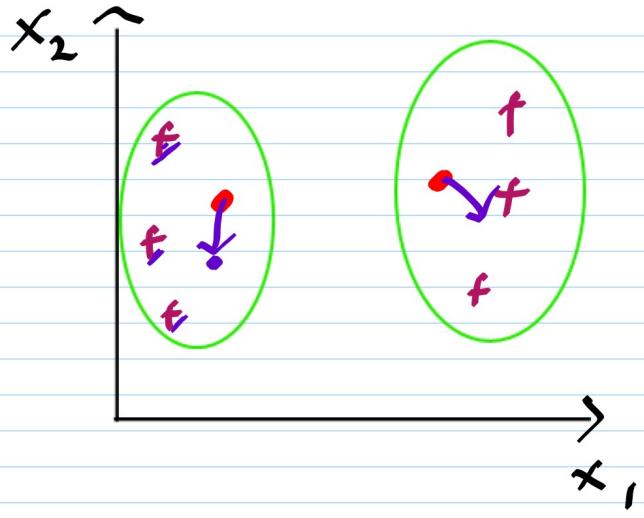
K-means Clustering
↳ Number of desired Clusters.

→ Centroid (It will take two center points).

① Randomly chooses the centroid and assigning each data point to its nearest centroid to form a cluster.

Distance is calculated based on Euclidean Distance .

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



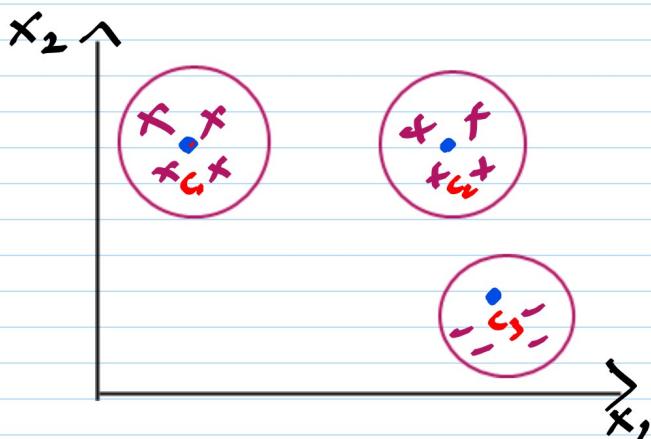
- ② Take the average value of the clusters and move the centroid to the average value.
- ③ Again calculate the distance of all the data points from the centroid and assigns those point to nearest centroid.
 → Perform the steps again and again until the movement of the centroids stops.

Q) What is the optimum number of clusters should we take?

$$K = 2, 3, 4, 5, 6, 7, \dots \sim [1, 10]$$

\downarrow
Clusters

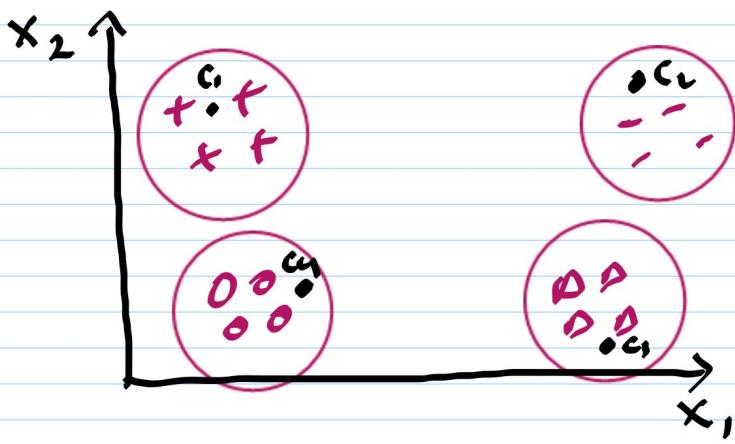
WCSS / Inertia (Within Cluster Sum of Square).



c = Centroid
 x_i = Individual data-point.

$$\text{WCSS} = \sum \sum (x_i - c_i)^2$$

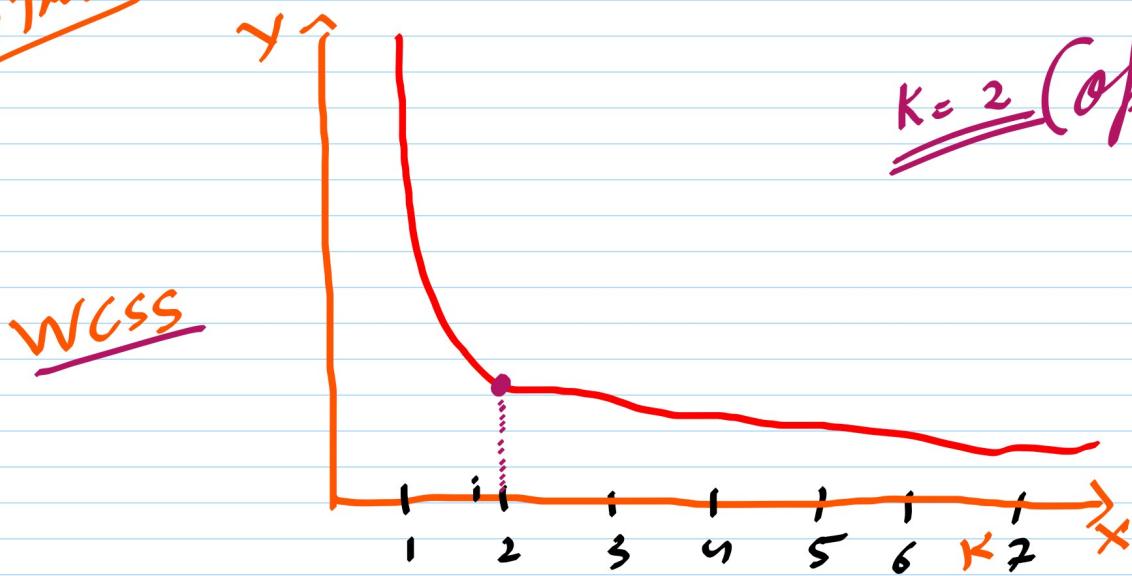
$$\underline{\text{WCSS}} = \sum (x_i - c_1)^2 + \sum (x_i - c_2)^2 + \sum (x_i - c_3)^2$$



$$\Rightarrow \underline{\text{WCSS}} = \sum (x_i - c_1)^2 + \sum (x_i - c_2)^2 + \sum (x_i - c_3)^2 + \sum (x_i - c_4)^2$$

Increase the number of Clusters, WCSS ↓

Elbow Method



$k=2$ (optimal k-value)