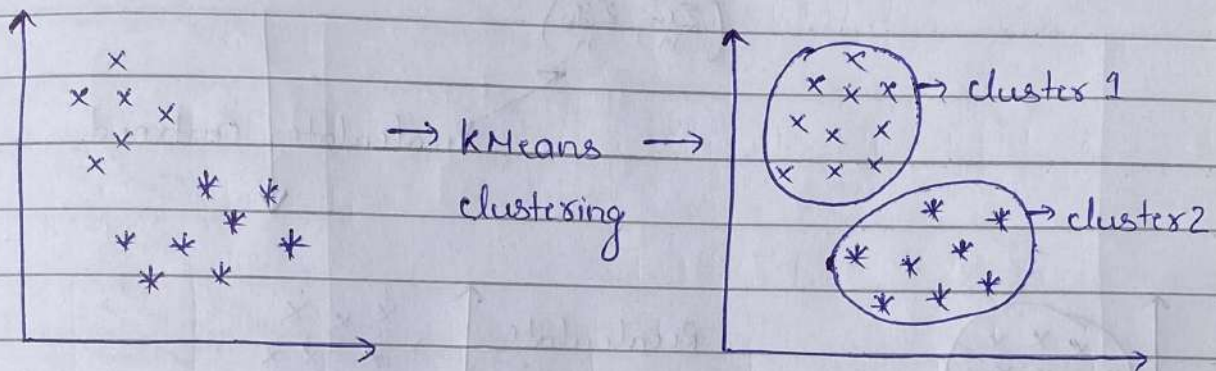# K Means Clustering

→ It is an unsupervised machine learning algorithm which is used to solve clustering problem by grouping the unlabeled dataset into different clusters.
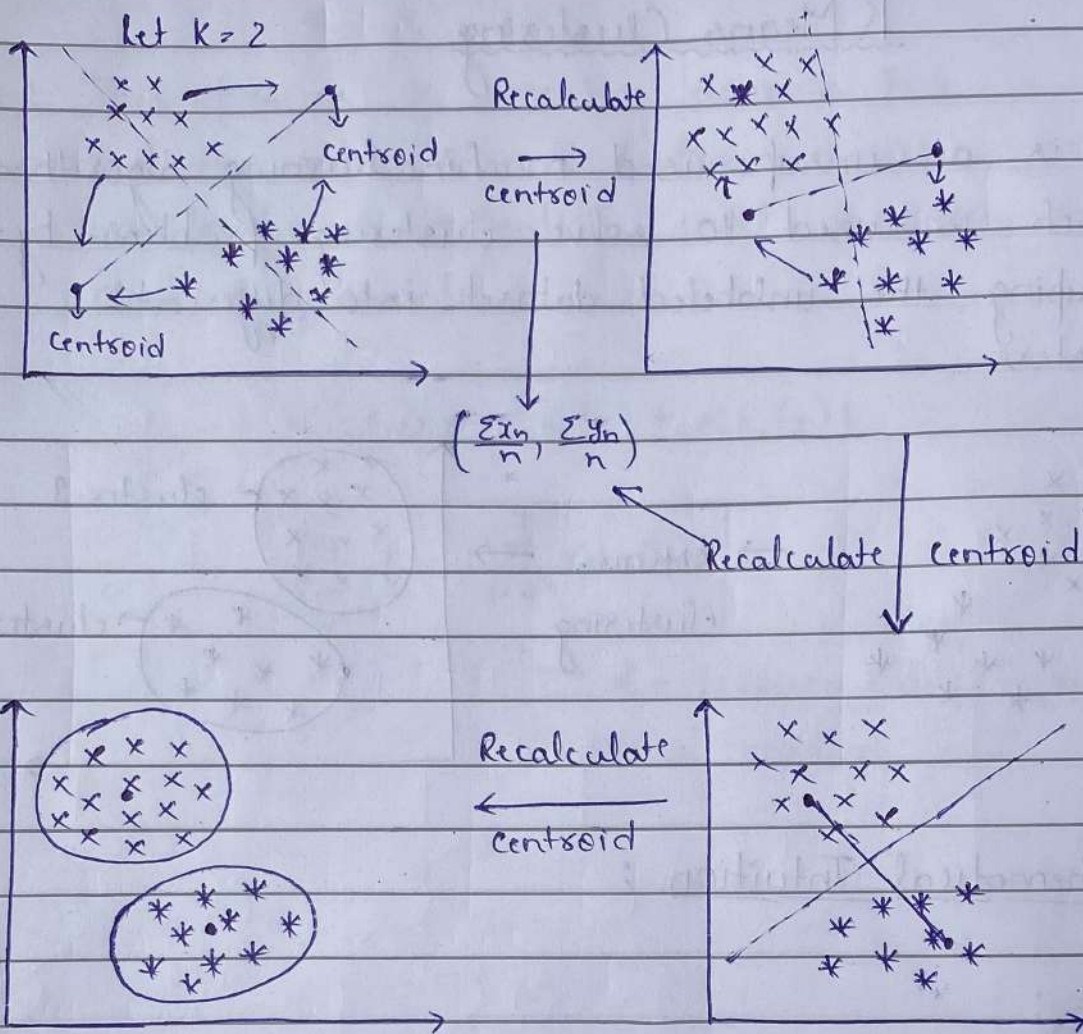


## * Mathematical Intuition :

### Steps :

(i) Initialize some k value.

(ii) Randomly initialize k centroids

(iii) Assign data points to nearest centroid

(iv) Recalculate the centroid value by mean of datapoints.

(v) Repeat 3rd and 4th step until centroid of current iteration become same of its previous iteration.

• Centroid calculation :

$$C = \left( \frac{\Sigma x_n}{n} , \frac{\Sigma y_n}{n} \right)$$
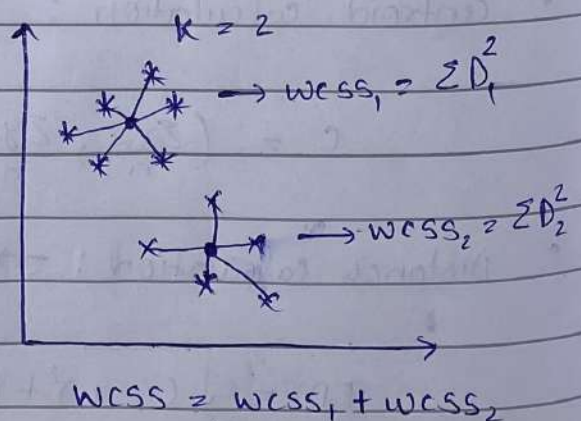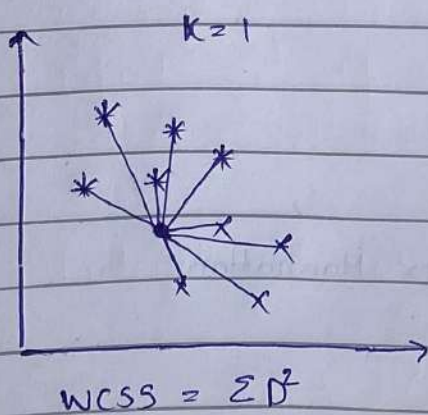
• Distance calculation : → Eulidean or Manhattan
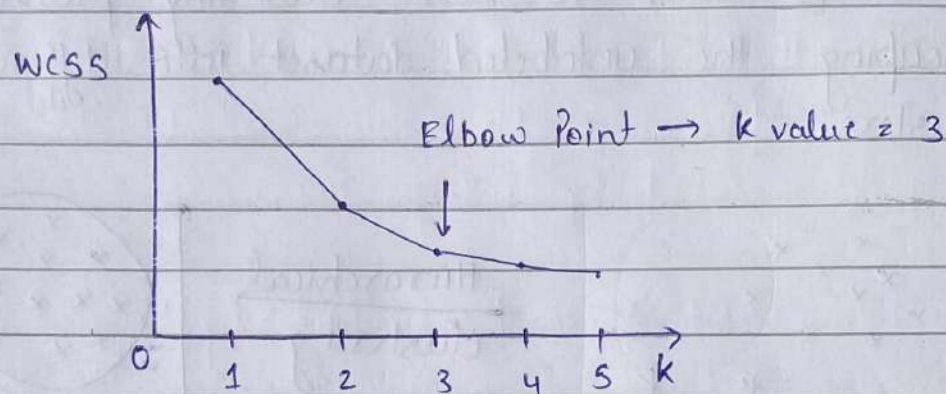
$$ED = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Let k=2

Recalculate
centroid →
centroid

centroid

centroid

$\left(\dfrac{\Sigma x_n}{n}, \dfrac{\Sigma y_n}{n}\right)$

Recalculate | centroid

Recalculate
←
centroid

Q How we can select the k value ?

- WCSS → within cluster sum of square.

$$WCSS = \sum_{j=1}^{K} \left[ \begin{array}{c} \text{Distance between points to nearest} \\ \text{centroid} \end{array} \right]^2$$

k=1

k=2

$\longrightarrow wess_1 = \Sigma D_1^2$

$\longrightarrow wess_2 = \Sigma D_2^2$

WCSS = $\Sigma D^2$

WCSS = $wcss_1 + wcss_2$

Note: As the value of k increases, the value of wcss decreases.



Elbow Point → k value = 3

* **Random Initialization Trap :**



Idially

But sometimes leads to this when centroid initializes very close to each other



Prevent this → KMeans++

Note: KMeans++ is an initialization technique which ensures that the centroids should initialize far from each other.