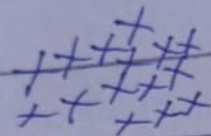
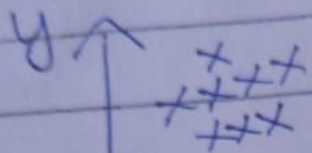
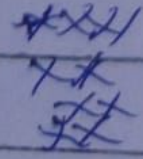
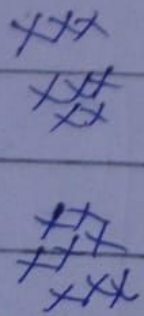
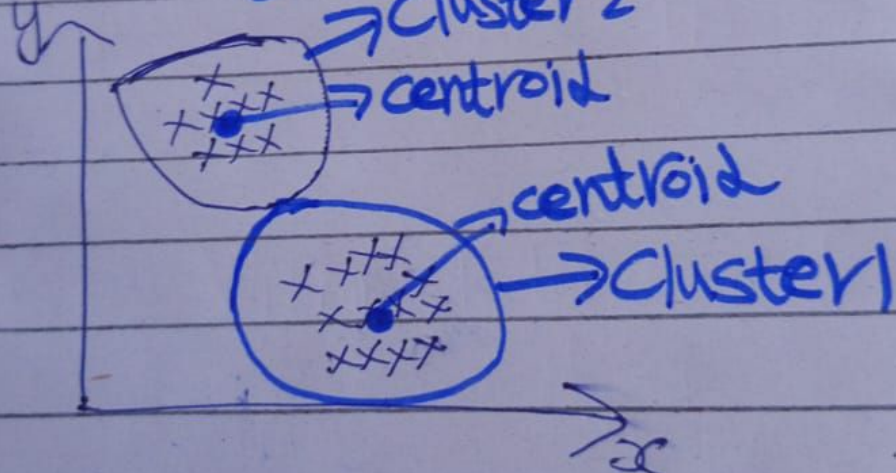


# ⇒ K Means Clustering Algorithm:

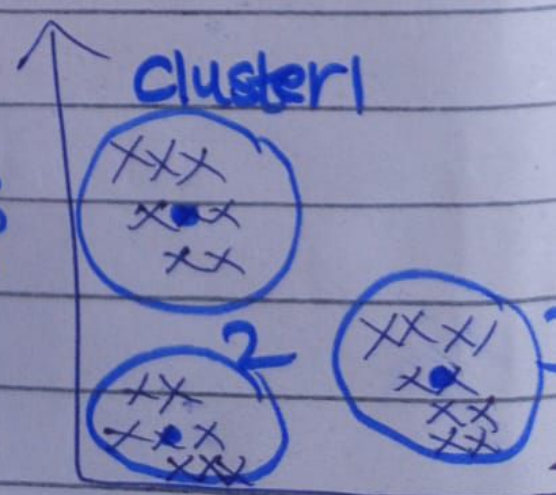
• Geometric Intuition:



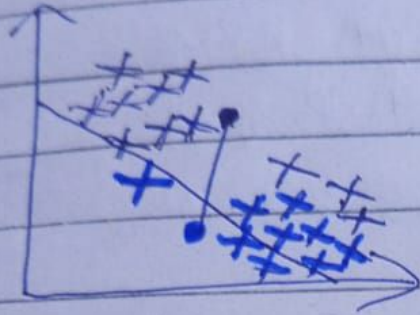
↓ After applying K Means



K Means



Let's suppose we have datapoints



- Centroid 1
- Centroid 2

**Step 01: Initialize some  $K$  values**  
↓  
centroids

Let  $K=2$  and we just select 2 random centroids

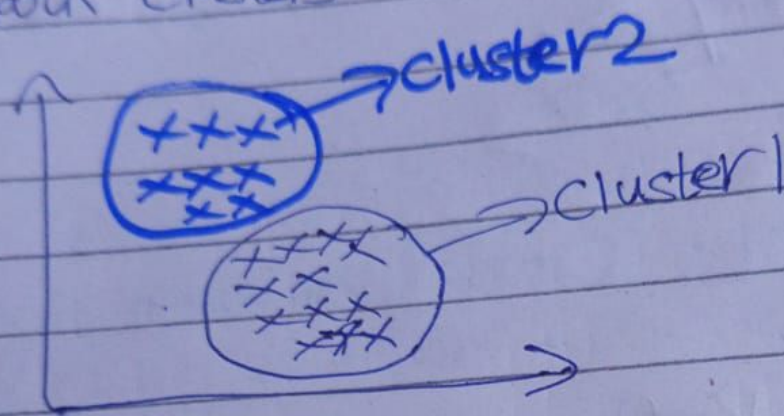
**Step 02:** Now we will find distance of all points from centroid 1 and centroid 2. Points that are nearest to centroid 1 mark them blue in that group

**Step 03:** Move the centroids to the average of the points in their group

And then again repeat step 2 and 3



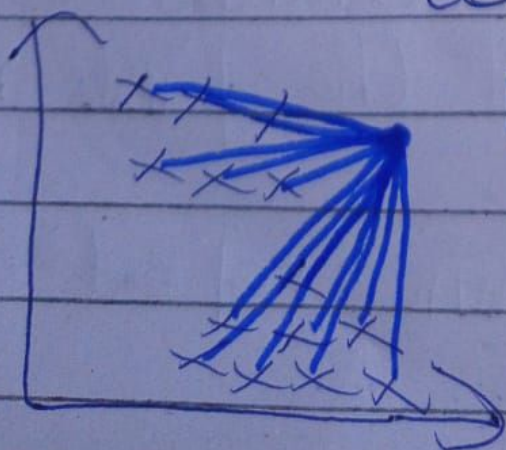
We will do these steps repeatedly until we not get 2 separate groups or clusters and then we just mark clusters which is goal.



⇒ How to find the K-value:

WCSS ⇒ Within cluster sum of squares

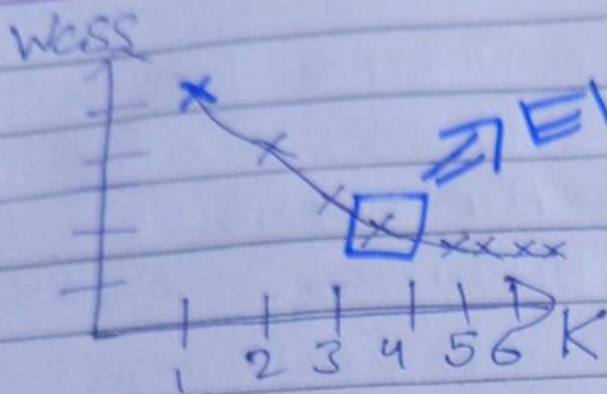
WCSS ⇒  $\sum_{i=1}^n (\text{distance blw points to nearest centroid})^2$



K=1 (only 1 centroid)



WCSS will be high for  $K=1$  as distance is more

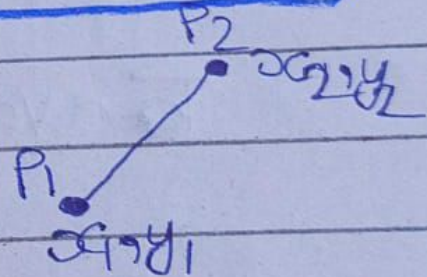


⇒ Elbow Method  
(We select  $K$ -value where WCSS abruptly decrease and then stabilize)

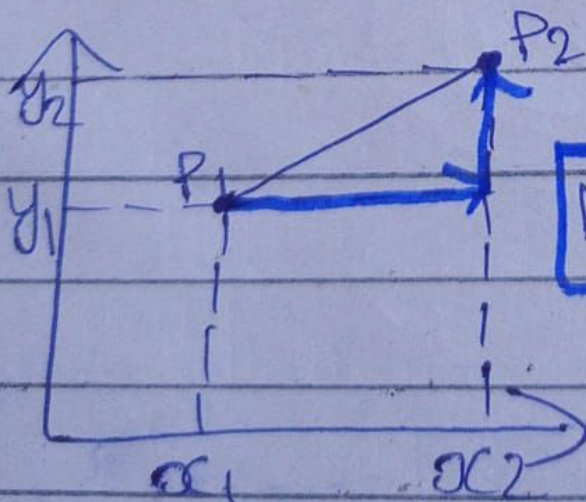
For  $K=2$  it WCSS decreases  
To find distance between points and centroid we use

### ① Euclidean Distance:

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



### ② Manhattan Distance



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



- We use Manhattan Distance for road things as there are buildings
- In case of Air traffic we use Euclidean Distance

### ⇒ Random Initialization Trap (KMeans++):

As for now we select centroids randomly but selection goes wrong i.e. centroids are very near to each other. This is

### called Random Initialization Trap

- To avoid this we use KMeans++ Initialization it ensures that centroids are at maximum distance from each other.