

K-Means

OPTIMIZATION

Mobile service providers need to setup their network – where to set up the towers so that all its users receive the maximum signal strength

Stringent law enforcement by the government – station the patrol vans- area of high crime rates are in the vicinity of patrol vans

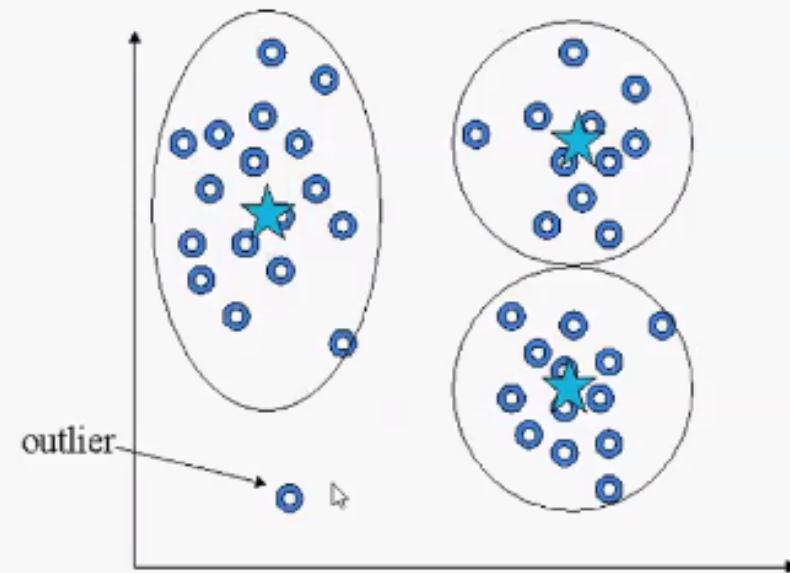
A hospital care chain wants to open series of emergency care wards- max accident prone areas

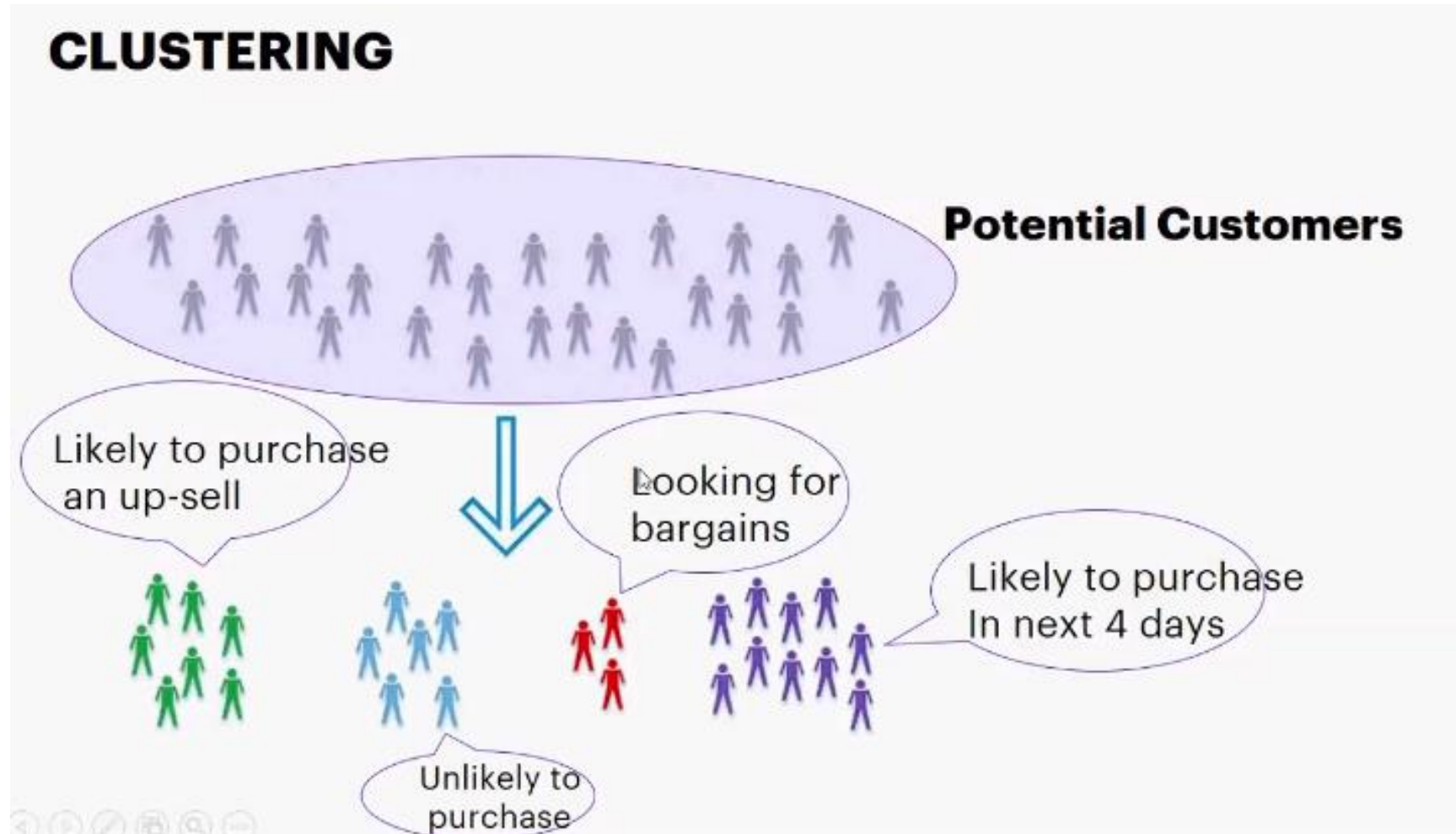
IDENTIFYING ANOMALY- OUTLIERS

Fraud detection- Communication service providers

Fraud transaction- Banking/Finance

Medicine – abnormal cells – diagnosis

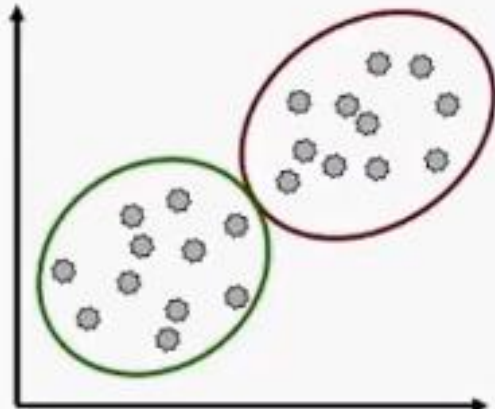




DIFFERENCE BETWEEN CLUSTERING & CLASSIFICATION

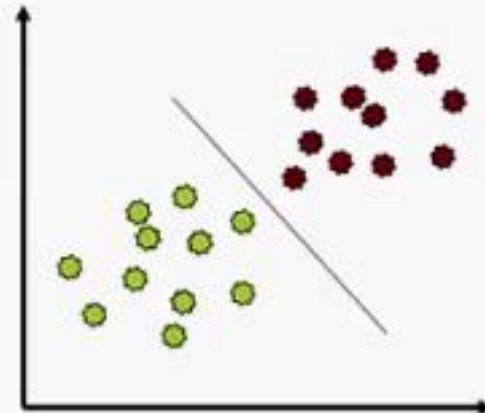
CLUSTERING

- Data is not labeled
- Group points that are "close" to each other
- Identify structure or patterns in data
- Unsupervised learning

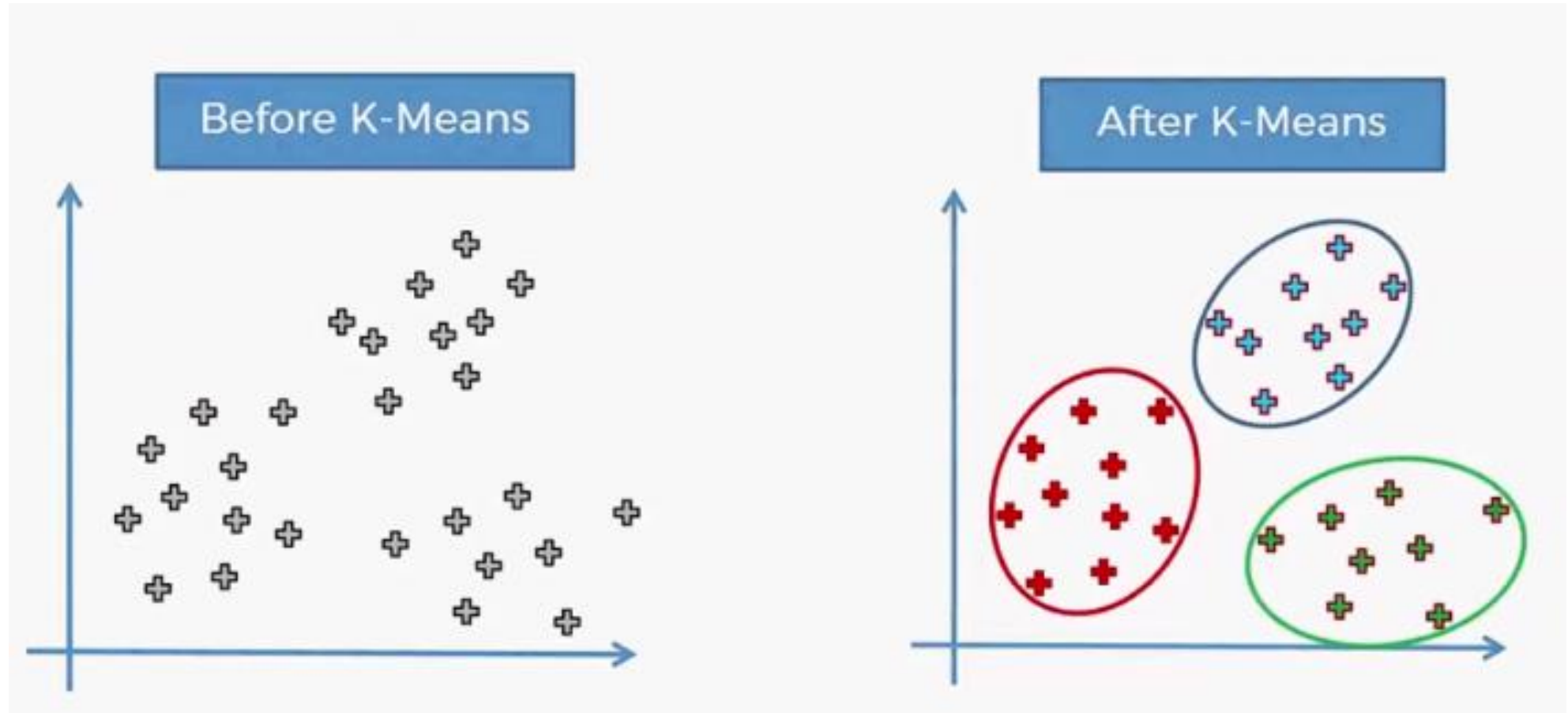


CLASSIFICATION

- Labeled data points
- Want a "rule" that assigns labels to new points
- Supervised learning



K-Means



STEPS

Choose the number K of clusters

Select at random K points, the centroids

Assign each data point to the closest centroid

Compute and place the new centroid of each cluster

Reassign each data point to the new closest centroid

if any reassignment took place, go to previous step

else finish

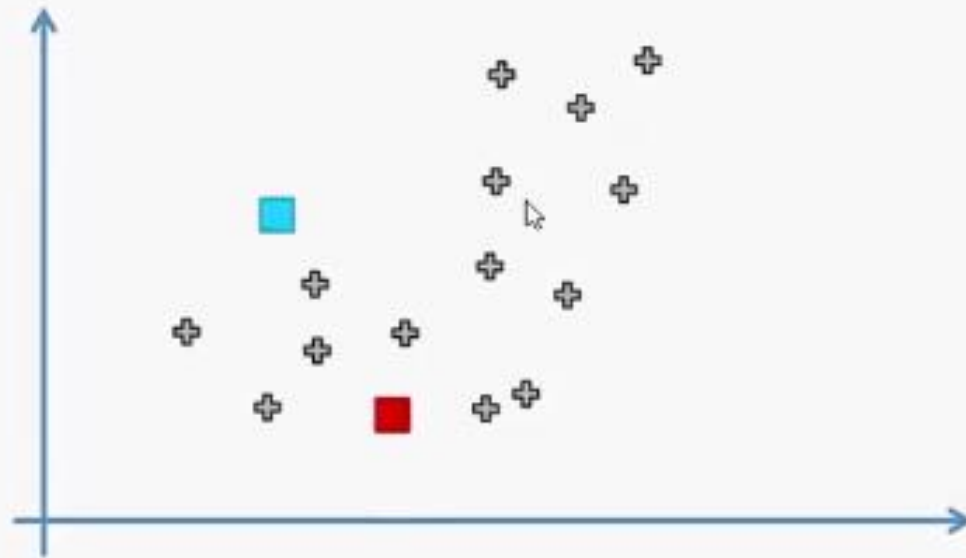
K MEANS ALGORITHM

STEP 1: Choose the number K of clusters: $K = 2$



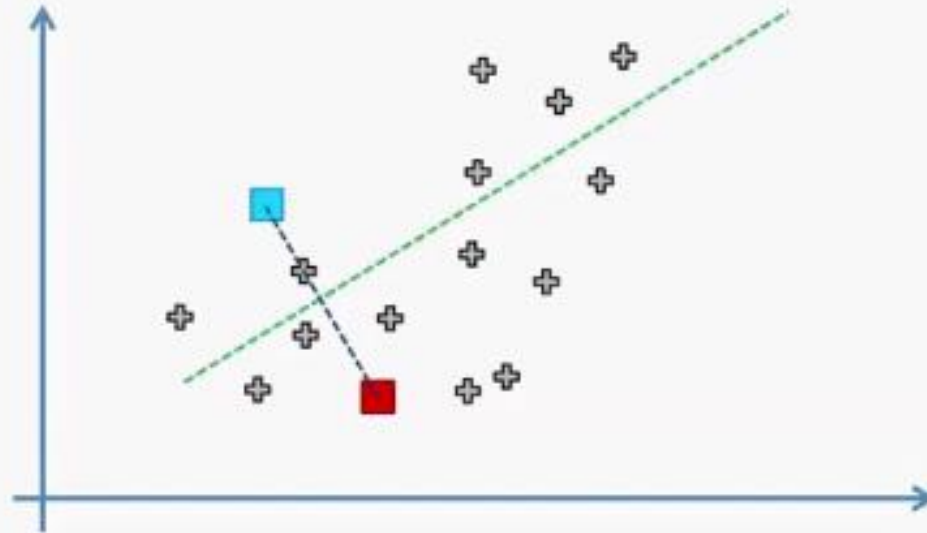
K-Means

STEP 2: Select at random K points, the centroids (not necessarily from your dataset)



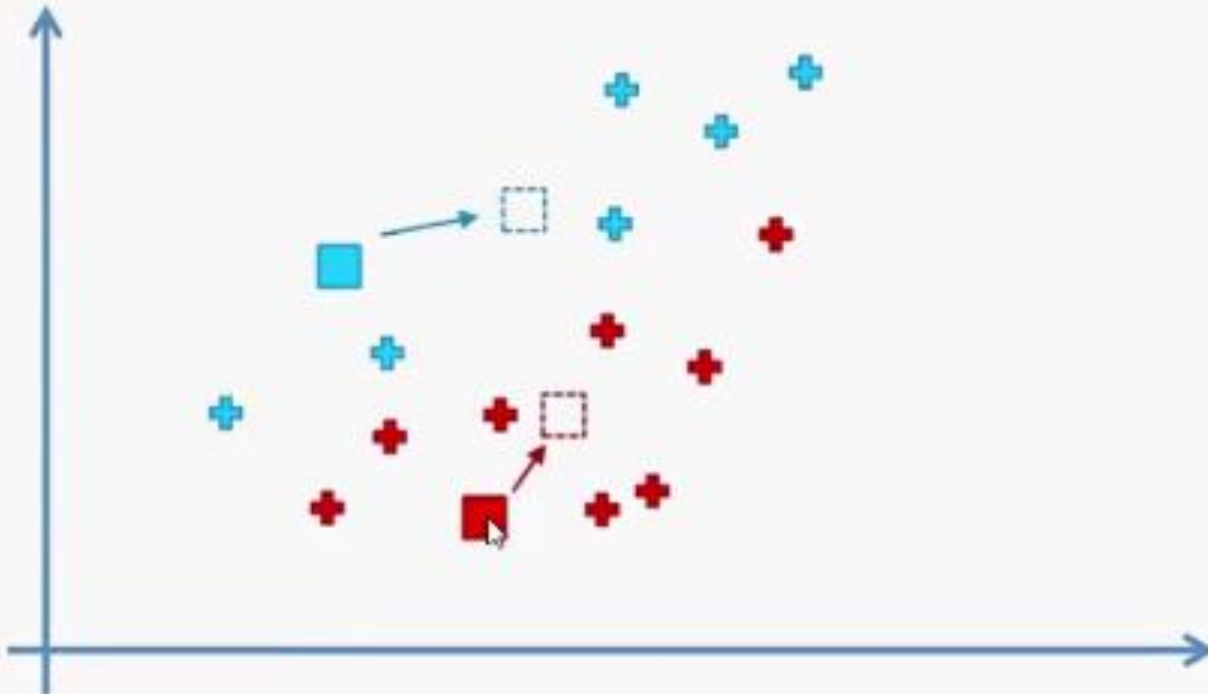
K-Means

STEP 3: Assign each data point to the closest centroid → That forms K clusters

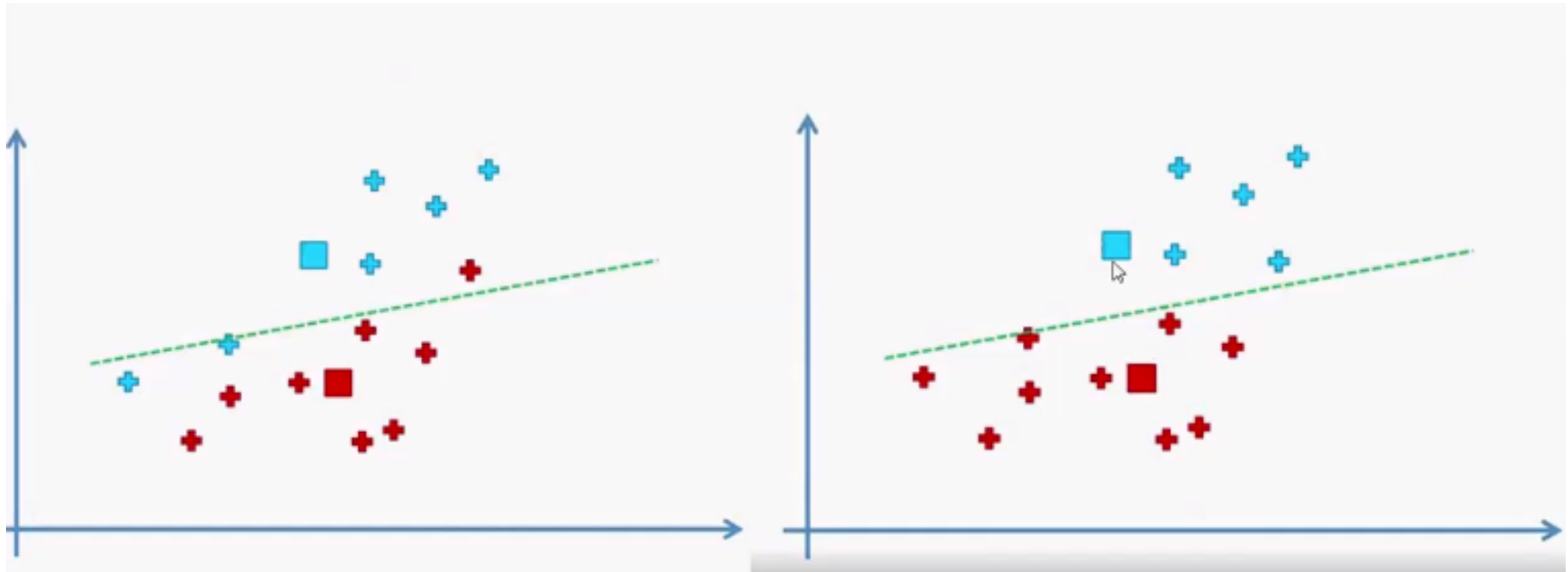


K-Means

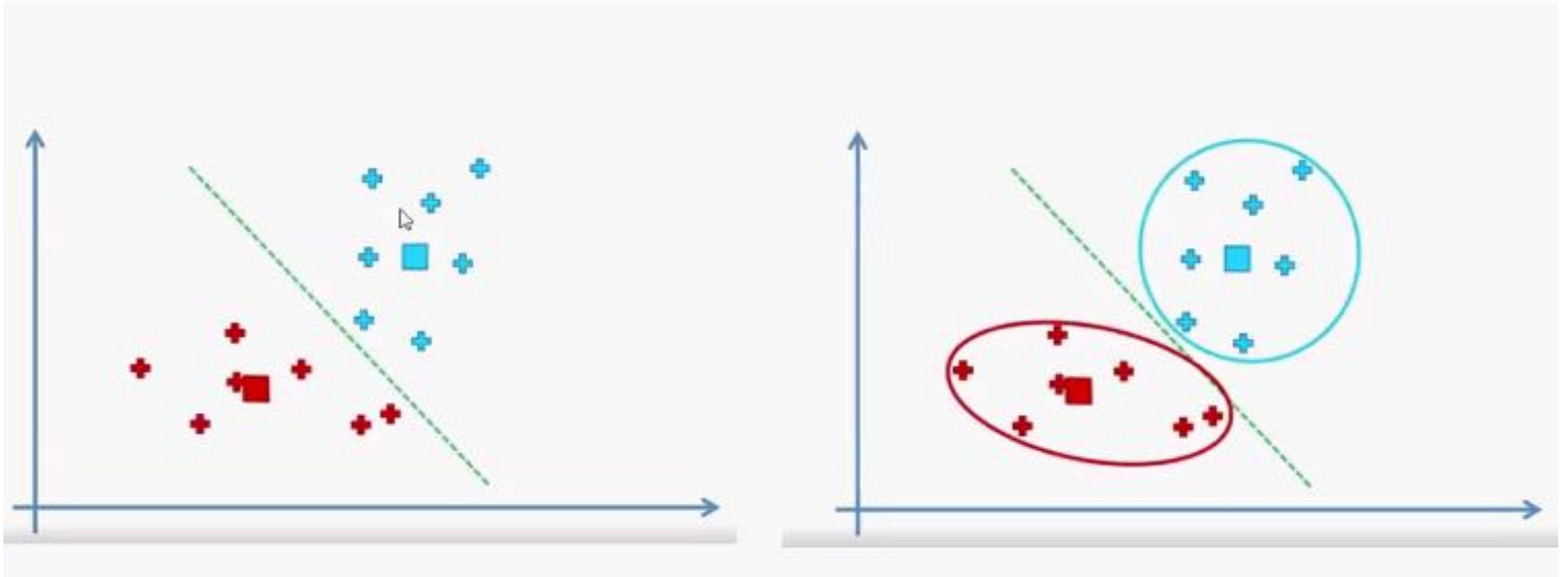
STEP 4: Compute and place the new centroid of each cluster



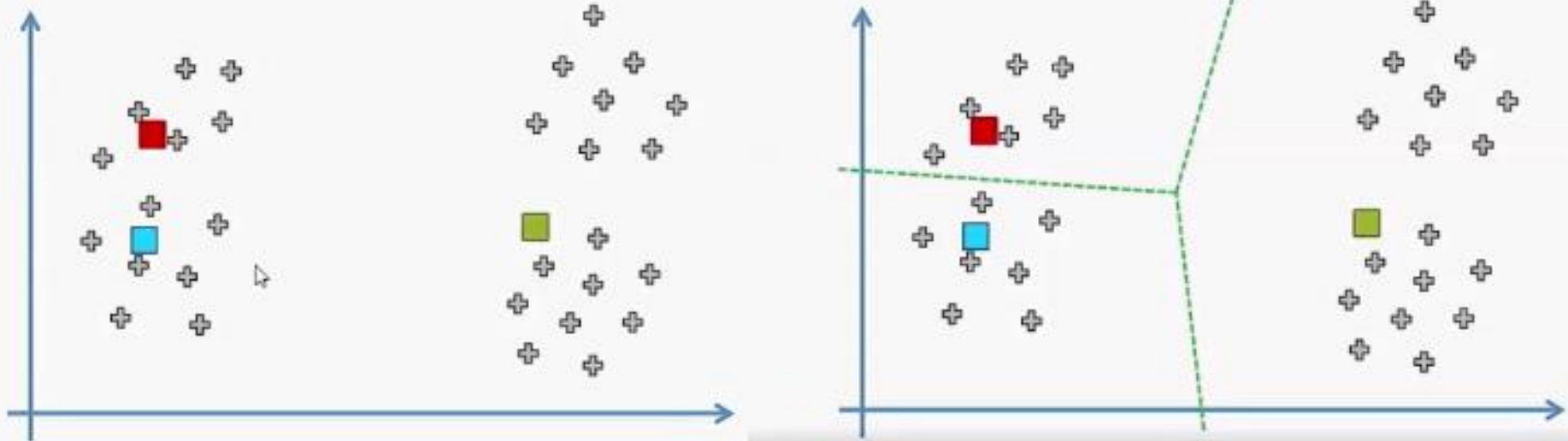
K-Means



K-Means

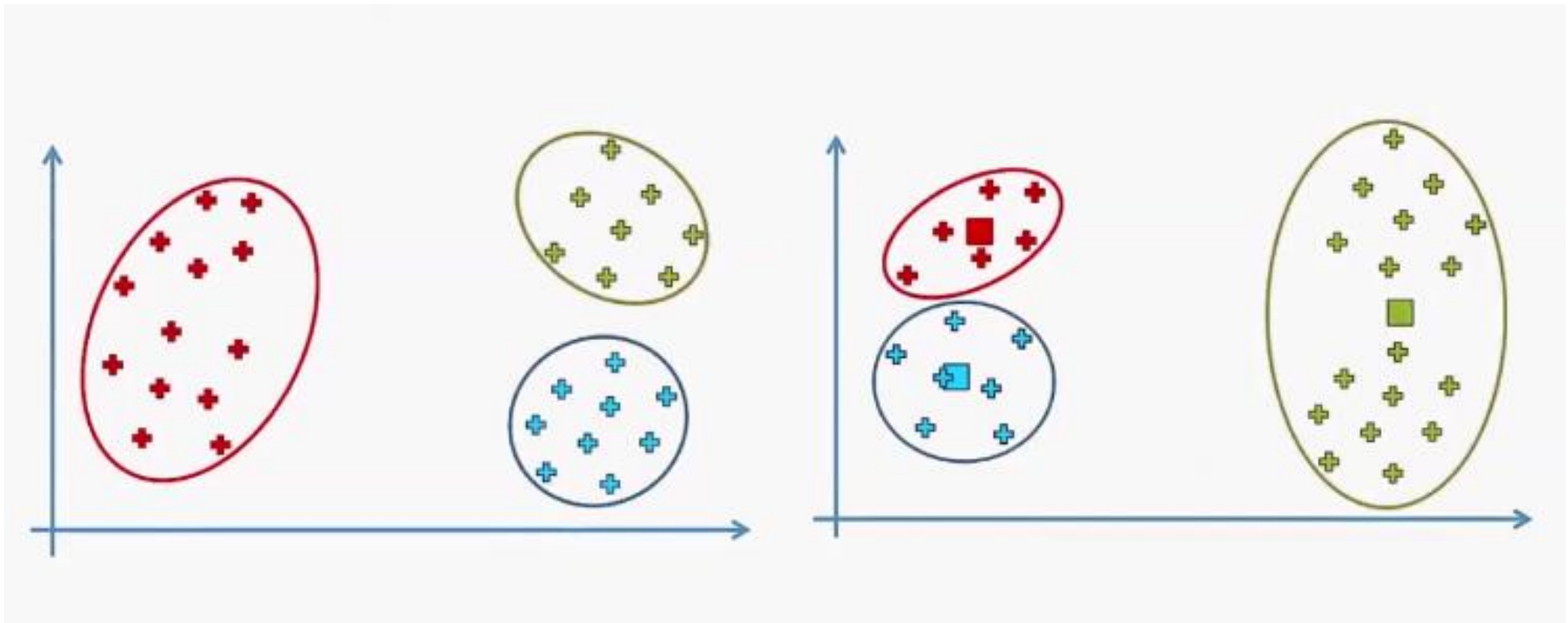


BAD RANDOM POINT INITIALIZATION

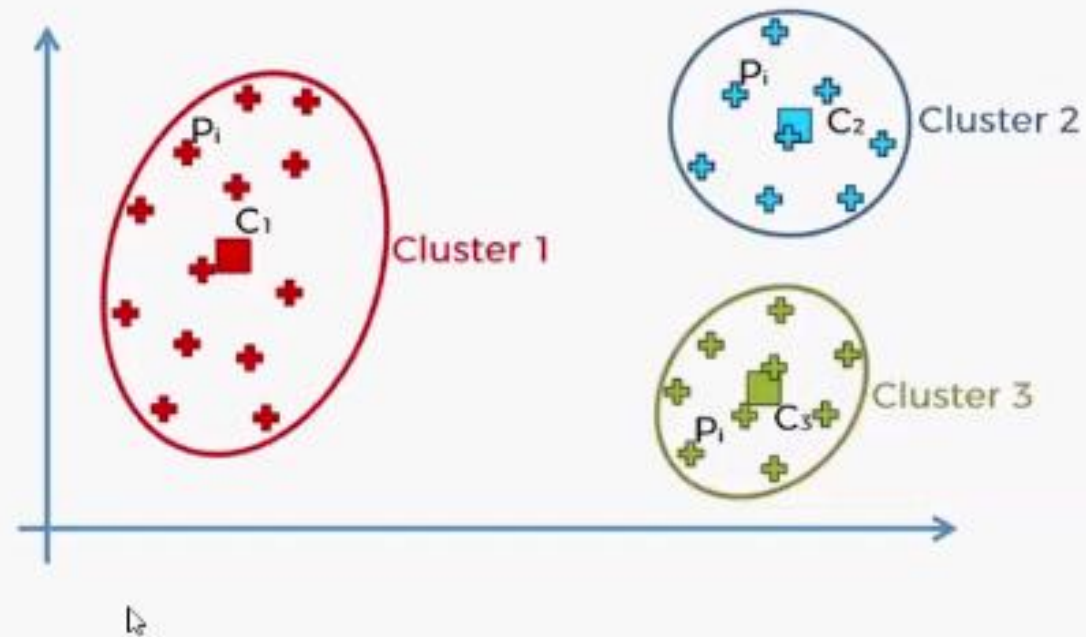


K-Means ++

Using K-means++ we can achieve optimum initial centroids as shown in Fig 1. Where as Fig 2 shows the incorrect initial centroids assignment.



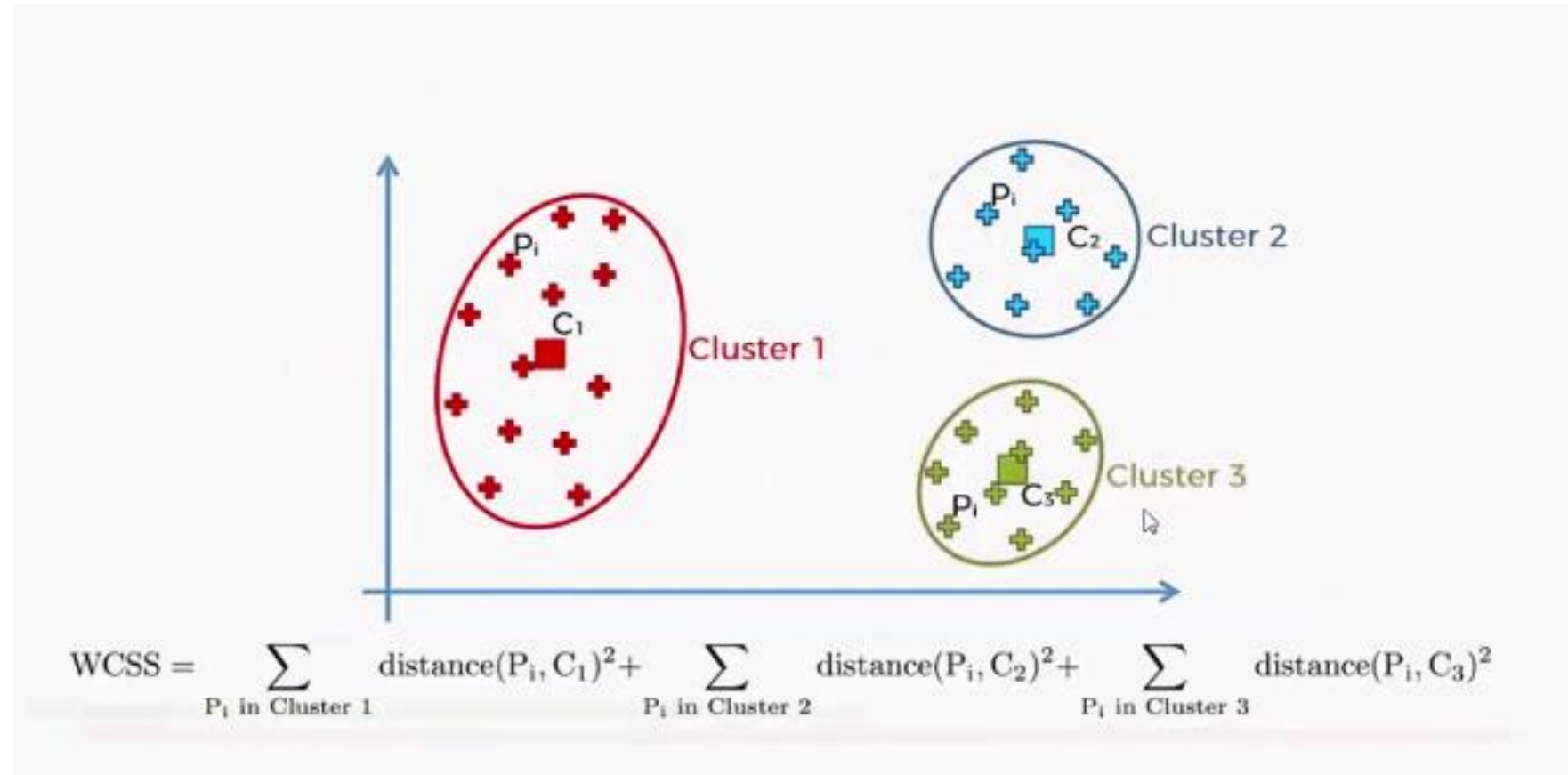
SELECTING THE NUMBER OF CLUSTERS



WITHIN CLUSTER SUM OF SQUARES

$$WCSS = \sum_{P_i \text{ in Cluster 1}} \text{distance}(P_i, C_1)^2 + \sum_{P_i \text{ in Cluster 2}} \text{distance}(P_i, C_2)^2 + \sum_{P_i \text{ in Cluster 3}} \text{distance}(P_i, C_3)^2$$

K-Means

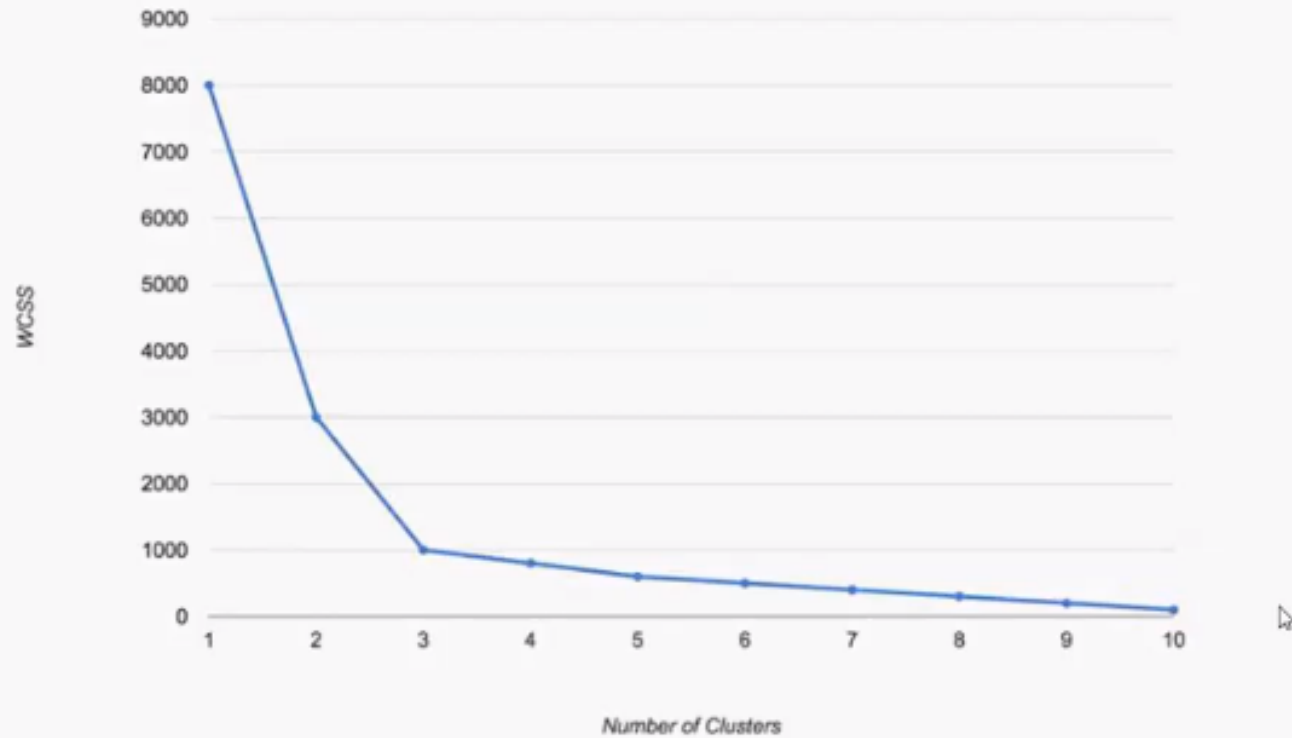


If the number of clusters = to the number of points

What is WCSS value ?



CHOOSING OPTIMUM NUMBER OF CLUSTERS- ELBOW METHOD



HIERARCHICAL CLUSTERING

Agglomerative
Divisive

AGGLOMERATIVE CLUSTERING

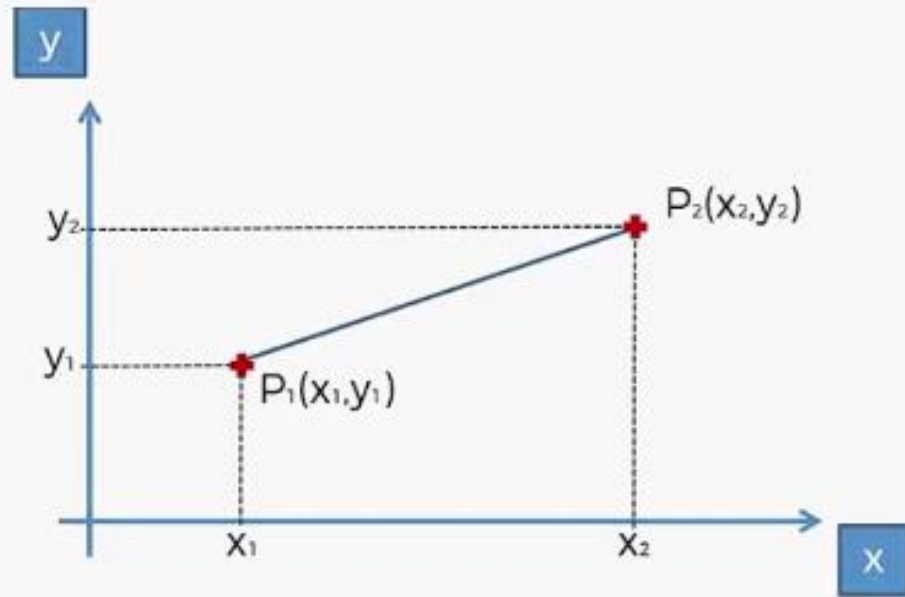
Step 1: Make each data point a single point cluster

Step 2: Take the two closest data points and make them one cluster

Step3: Take the two closest clusters and make them one cluster

Step4: Repeat Step3 until there is only one cluster

EUCLIDEAN DISTANCE



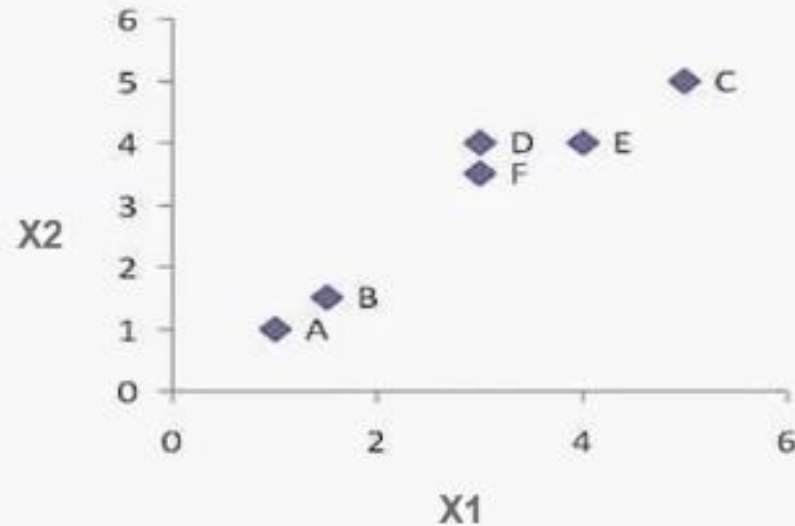
$$\text{Euclidean Distance between } P_1 \text{ and } P_2 = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

HIERARCHICAL CLUSTERING

Simple Example:

Assume that there are 6 objects namely A, B, C, D, E and F and each object has two measured features X1 and X2.

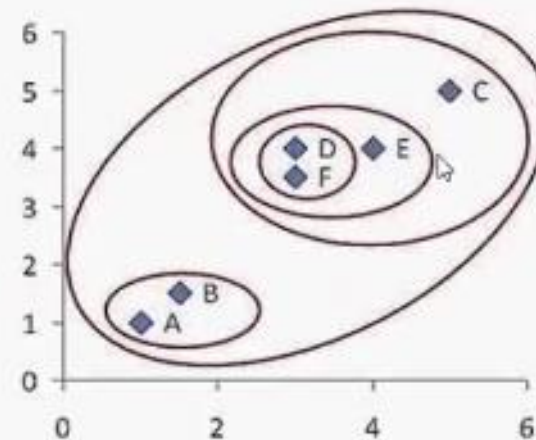
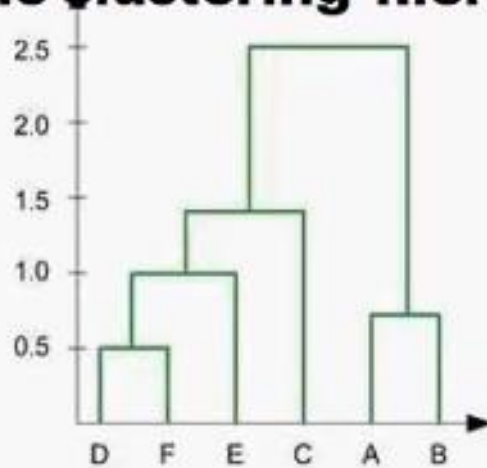
	X1	X2
A	1	1
B	1.5	1.5
C	5	5
D	3	4
E	4	4
F	3	3.5



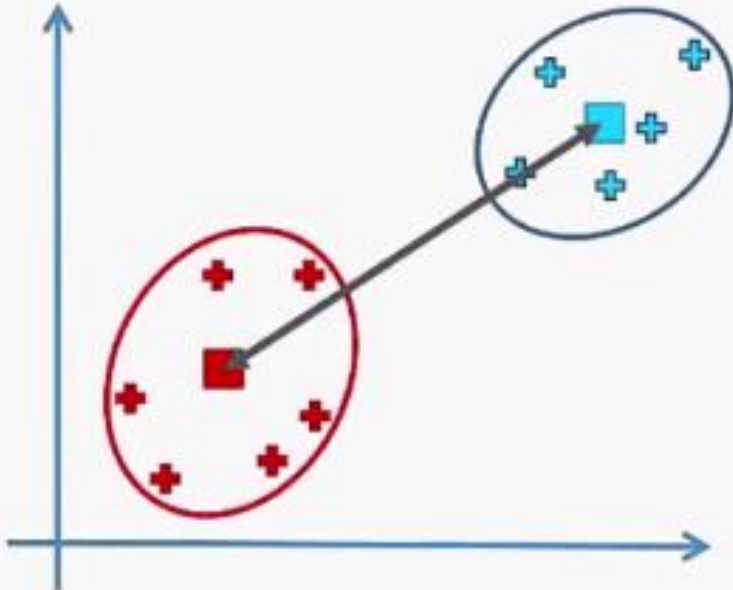
SINGLE LINKAGE HIERARCHICAL CLUSTERING

The dendrogram is drawn based on the distances to merge the clusters above

The hierarchy is given as $((D, F), E), C), (A, B)$. We can also plot the clustering hierarchy into XY space



DISTANCE BETWEEN TWO CLUSTERS



1. closest points
2. Furthest points
3. Average distance
4. Distance between centroids

Hierarchical Clustering

