



clustering

Clustering

What is Clustering?

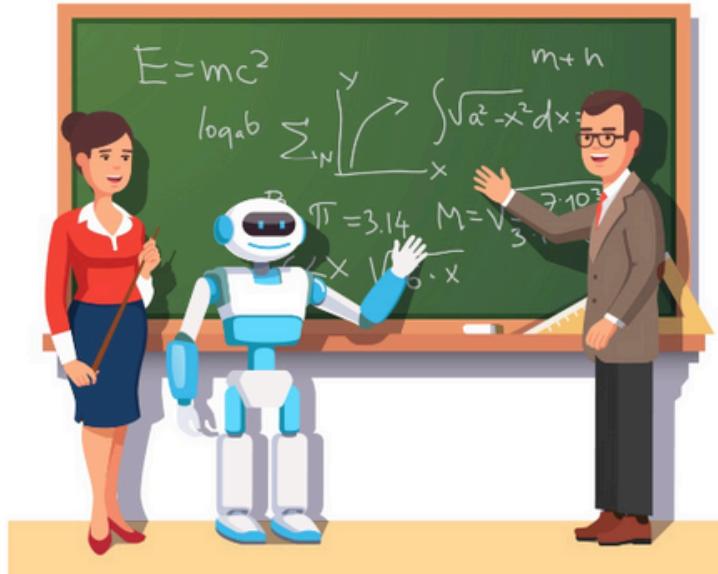
Clustering - grouping
Unlabelled Data

Clustering

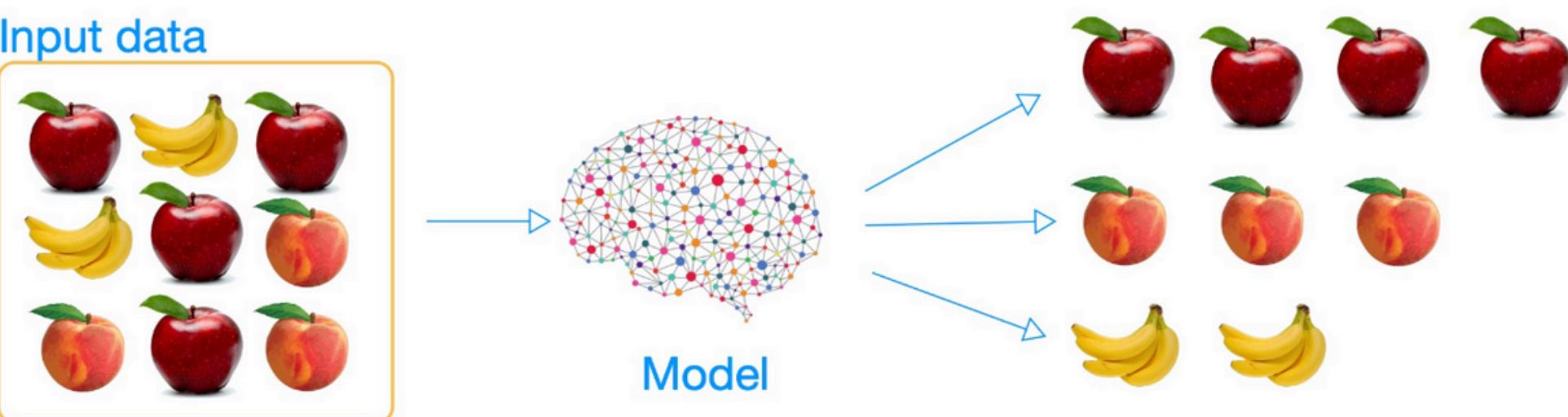
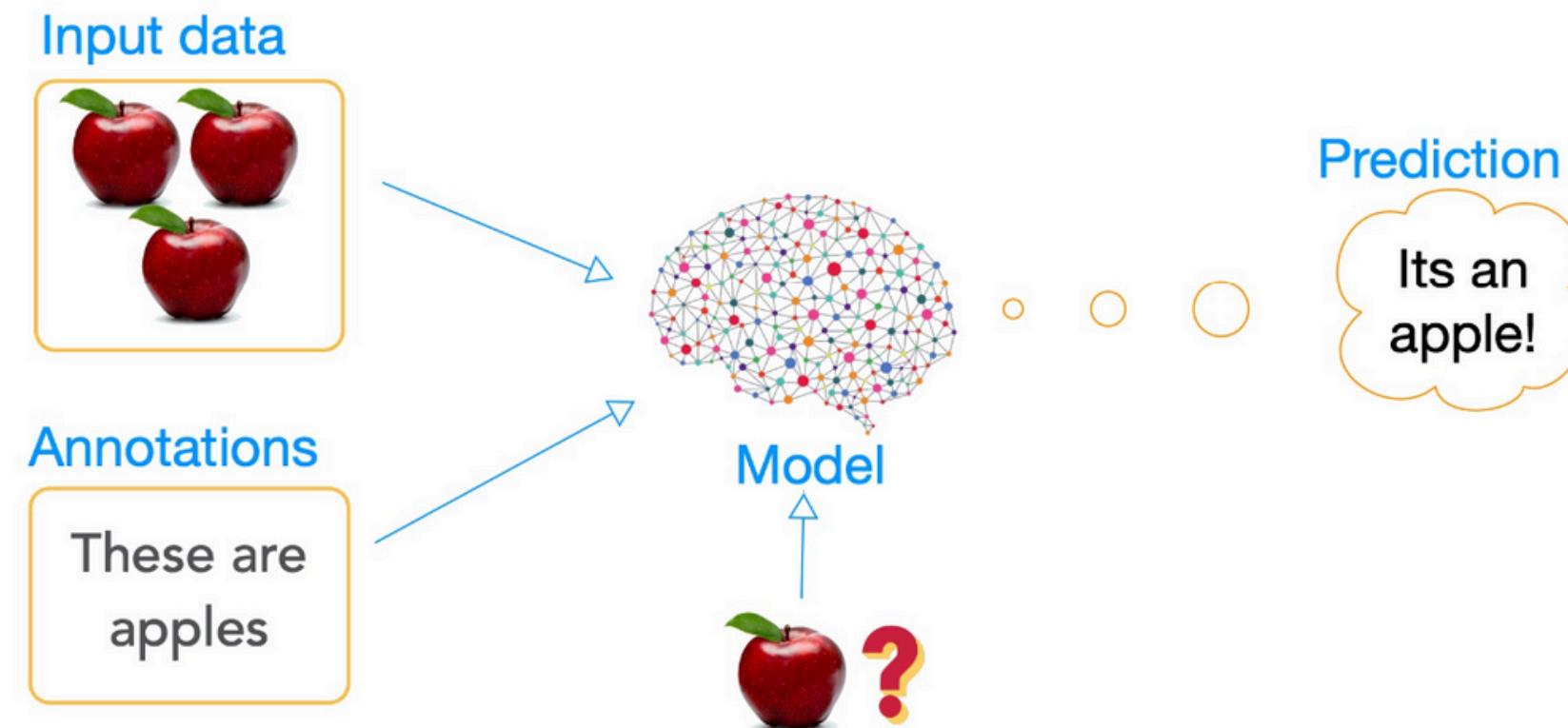
What is Clustering?

Supervised Learning

(e.g. Regression, Classification)

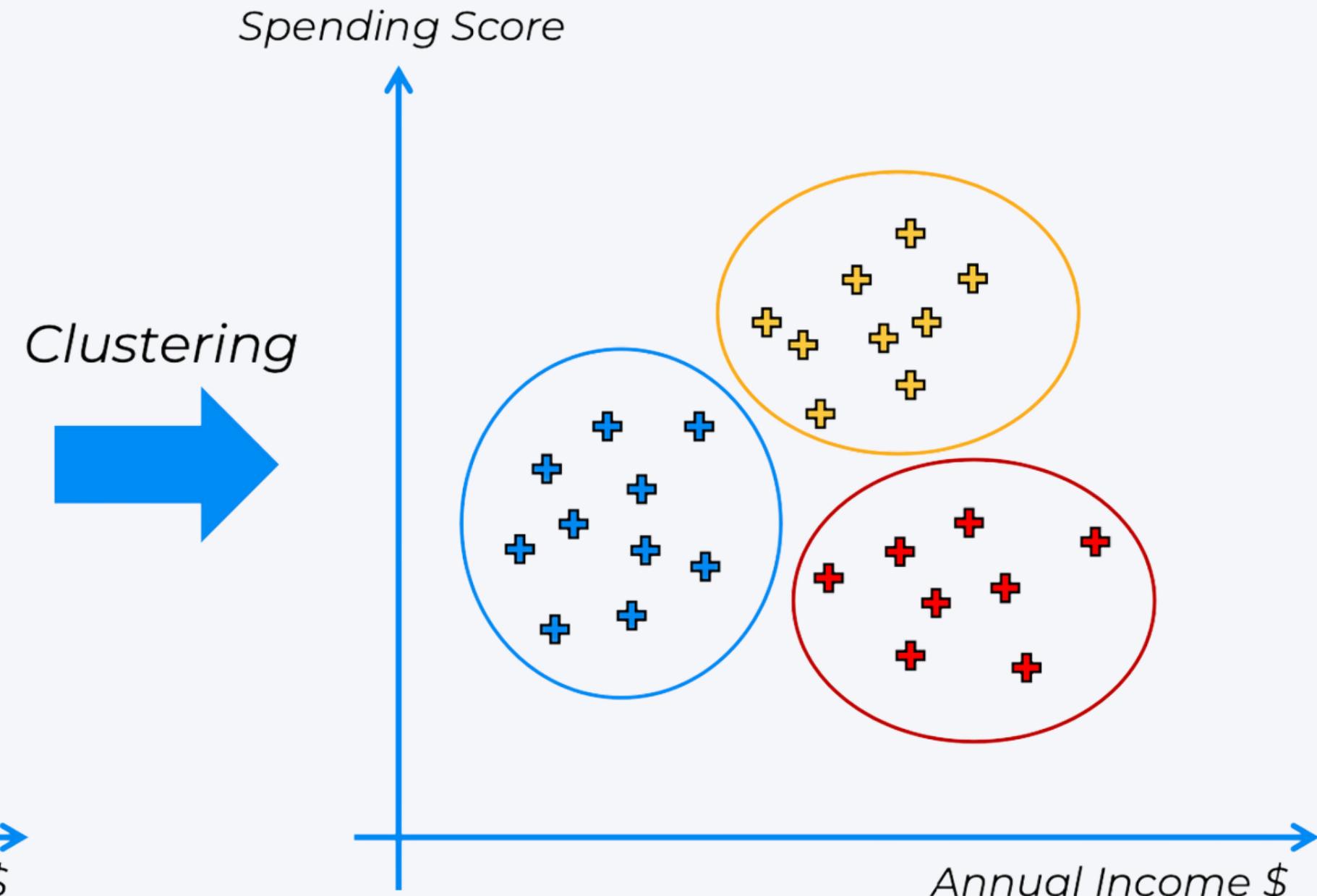


Unsupervised Learning (e.g. Clustering)



Clustering

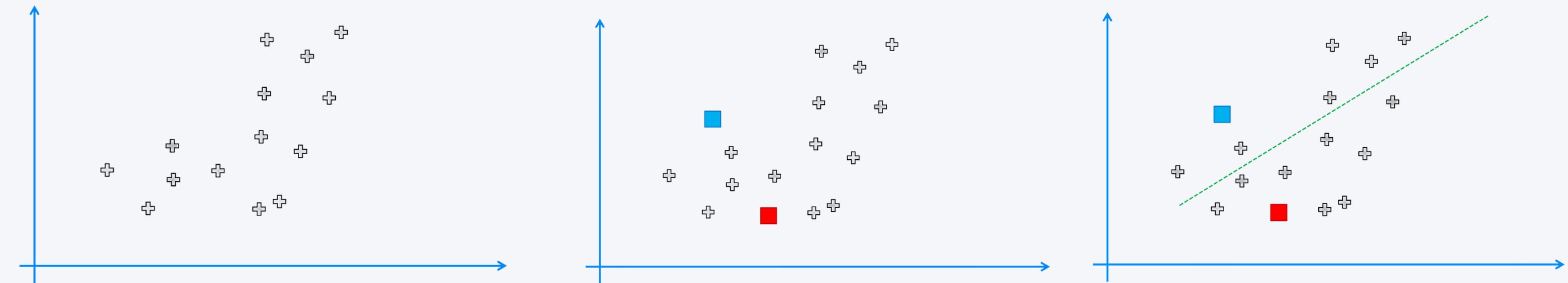
What is Clustering?



K-Means clustering

K-Means Clustering

The Intuition Behind K-Means Clustering



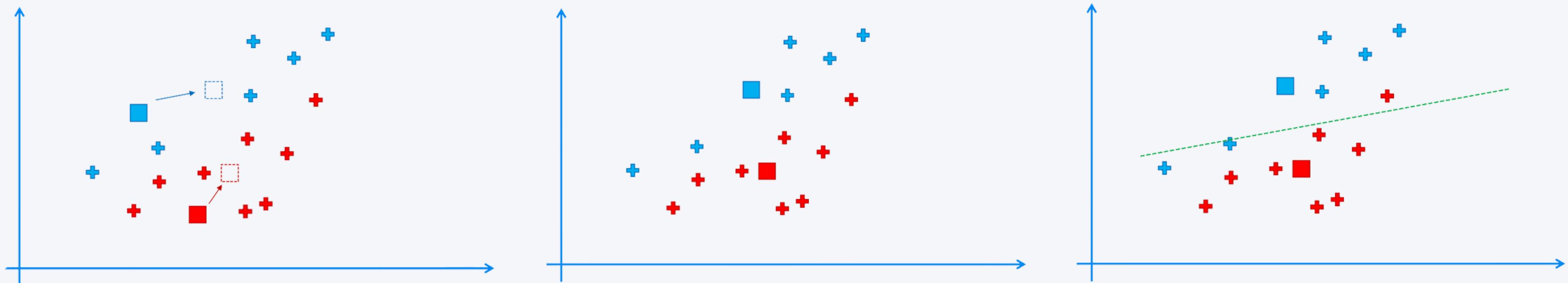
No labels data
How many clusters ?
Let consider two

Replace them randomly ,
any point (centroid)

K-means assign distance
line to the closest these
centroid , easiest way like
this

K-Means Clustering

The Intuition Behind K-Means Clustering



Calculate mass or gravity of each clusters ..

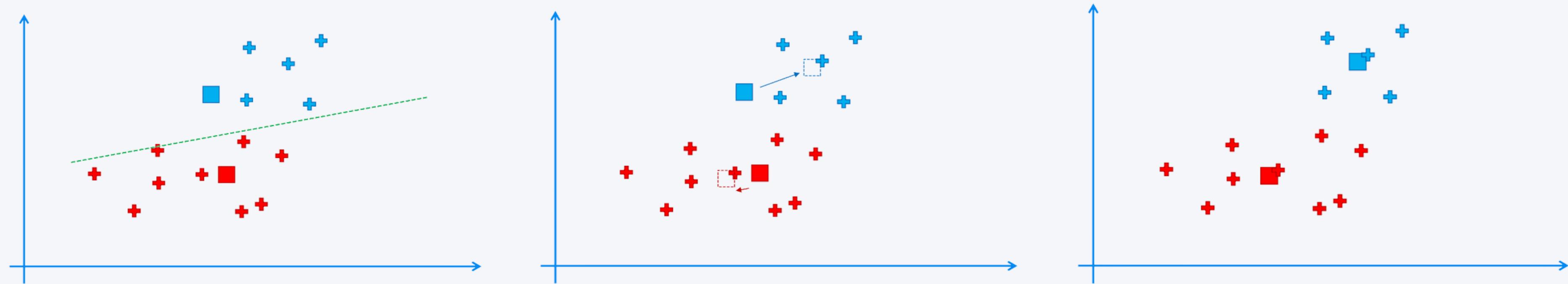
For take all x and get average and for blue take all y and get the average so you will find the cluster mass

move the centriod to new positions and repeat the process

Place the line to the closest centroid

K-Means Clustering

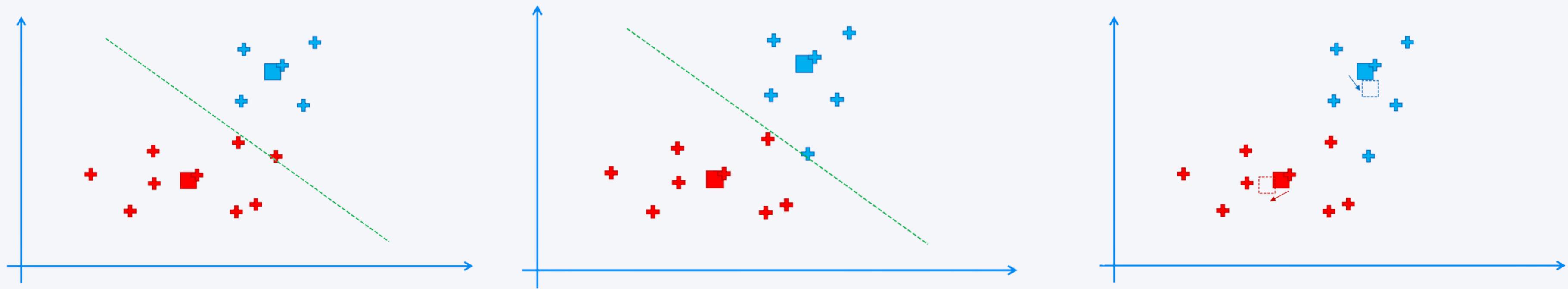
The Intuition Behind K-Means Clustering



Repeat the process ..

K-Means Clustering

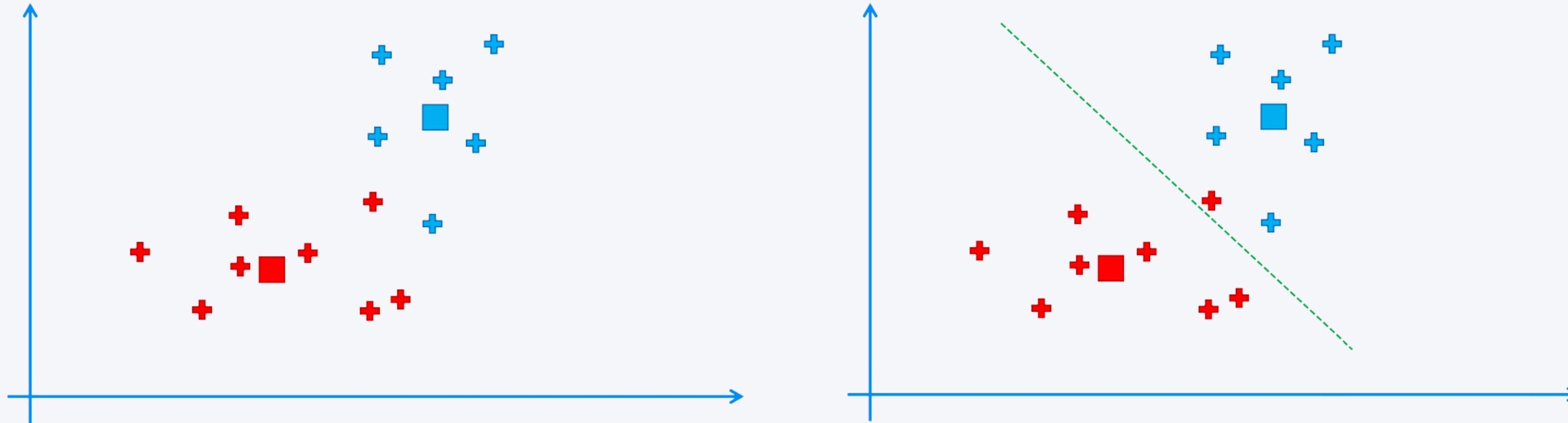
The Intuition Behind K-Means Clustering



Repeat the process ..

K-Means Clustering

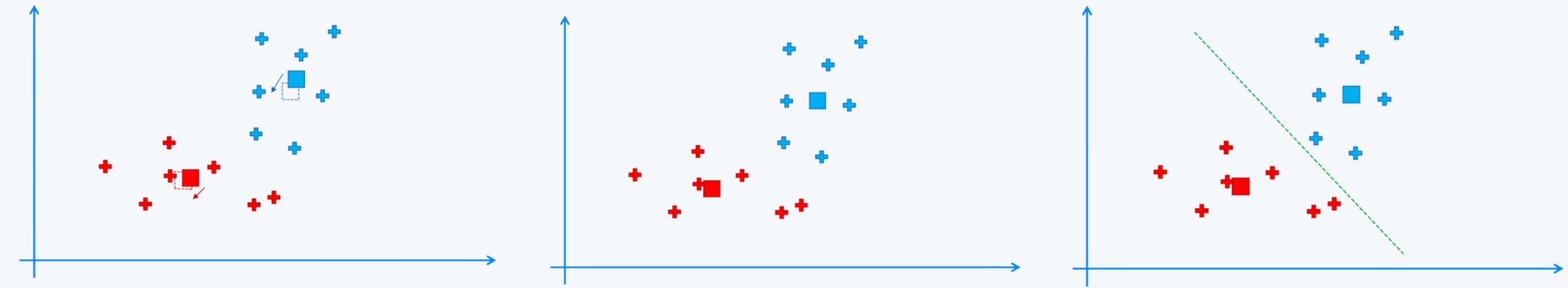
The Intuition Behind K-Means Clustering



Repeat the process

K-Means Clustering

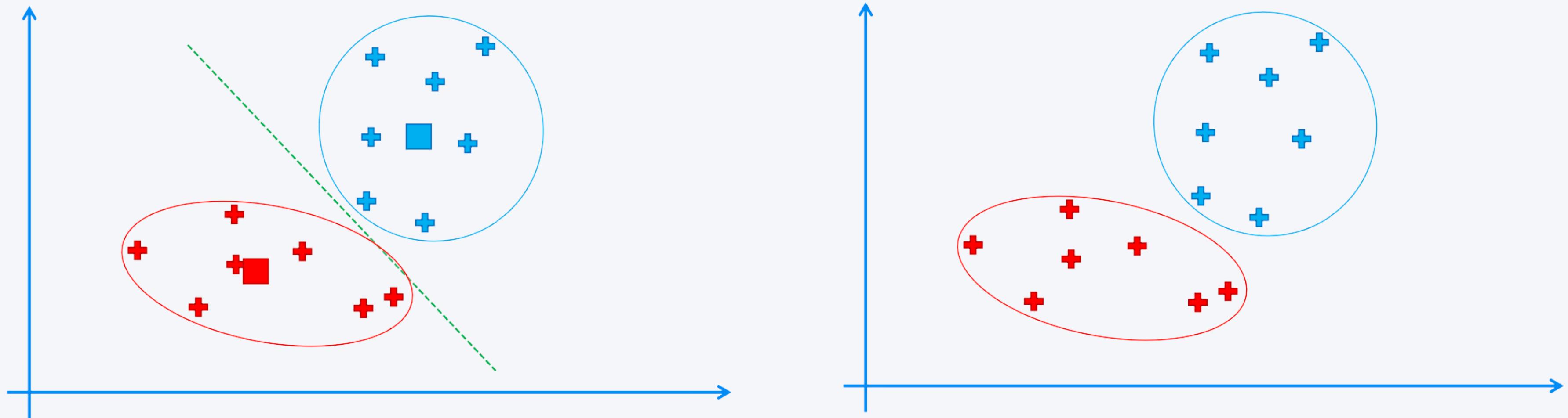
The Intuition Behind K-Means Clustering



Repeat until the cluster wont change

K-Means Clustering

The Intuition Behind K-Means Clustering



So Here is our final clusters

The question is how many clusters should be/should we have !

- Sometimes from the knowledge domain
- Sometimes we need to know how many we should have

K-Means Clustering

The Elbow Method

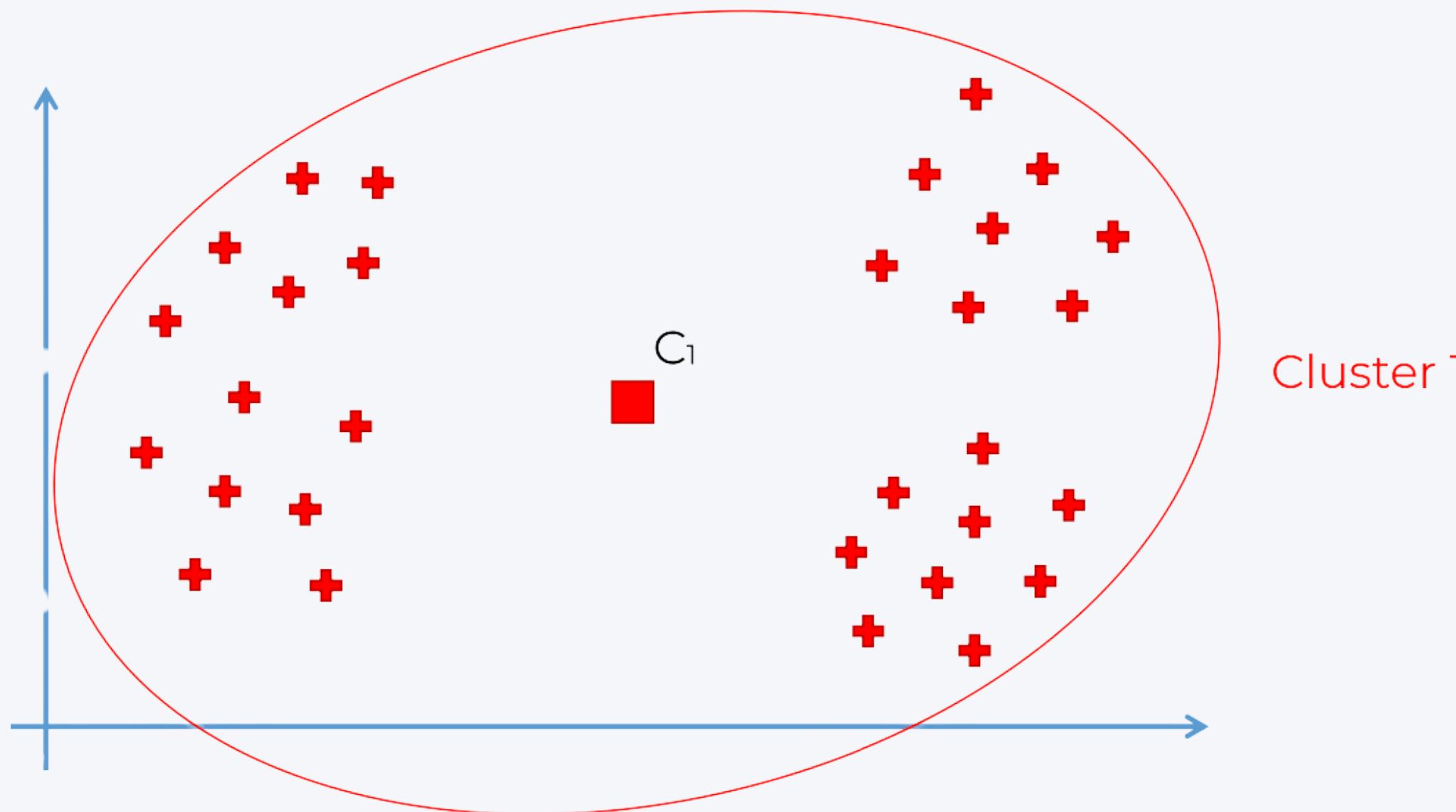
Within Cluster Sum of Squares:

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

It basically looks at the distance between each point and the centroid and square it

K-Means Clustering

The Elbow Method



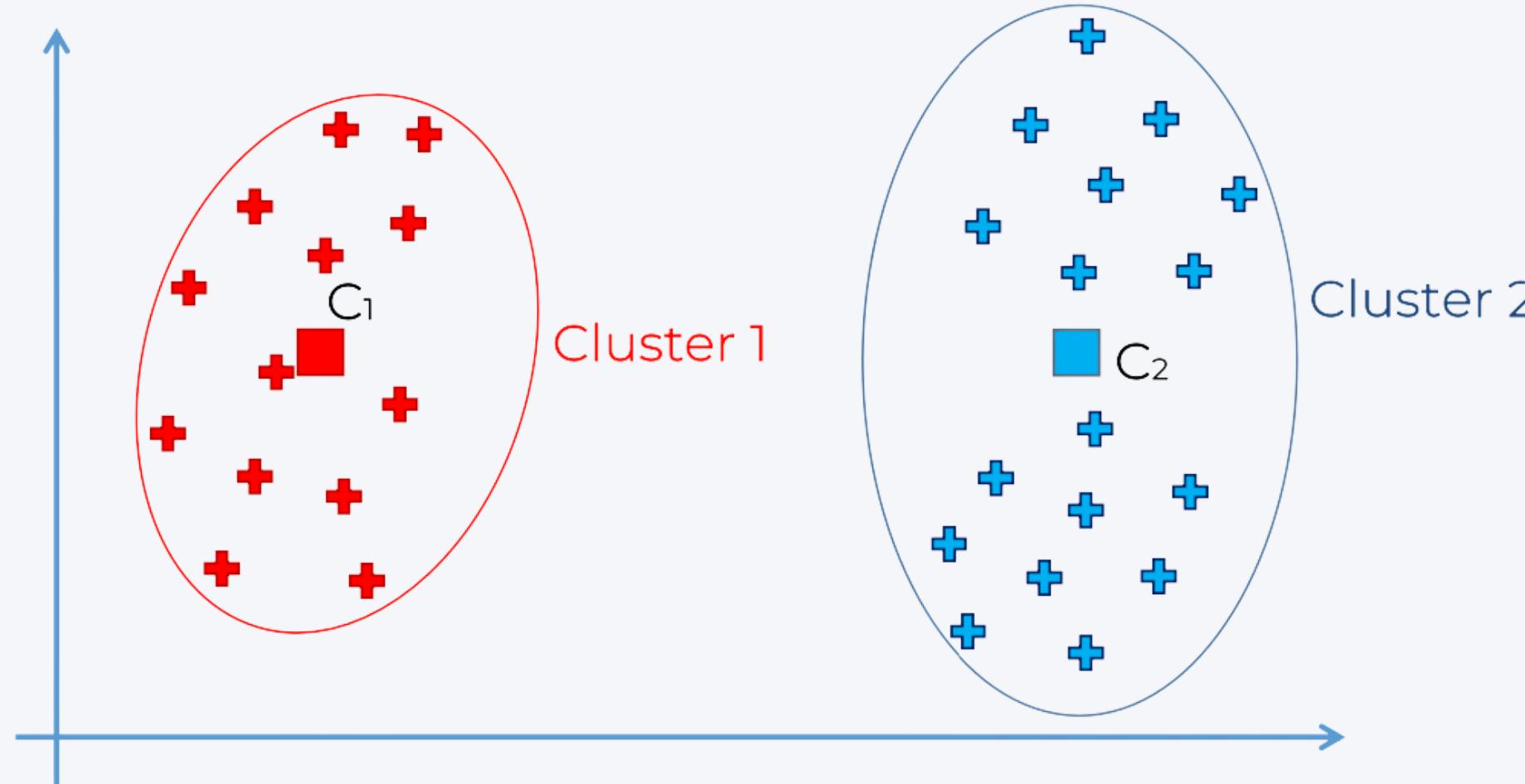
Within Cluster Sum of Squares:

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

In this example, measure each points with the centroid and square it to find the WCSS

K-Means Clustering

The Elbow Method



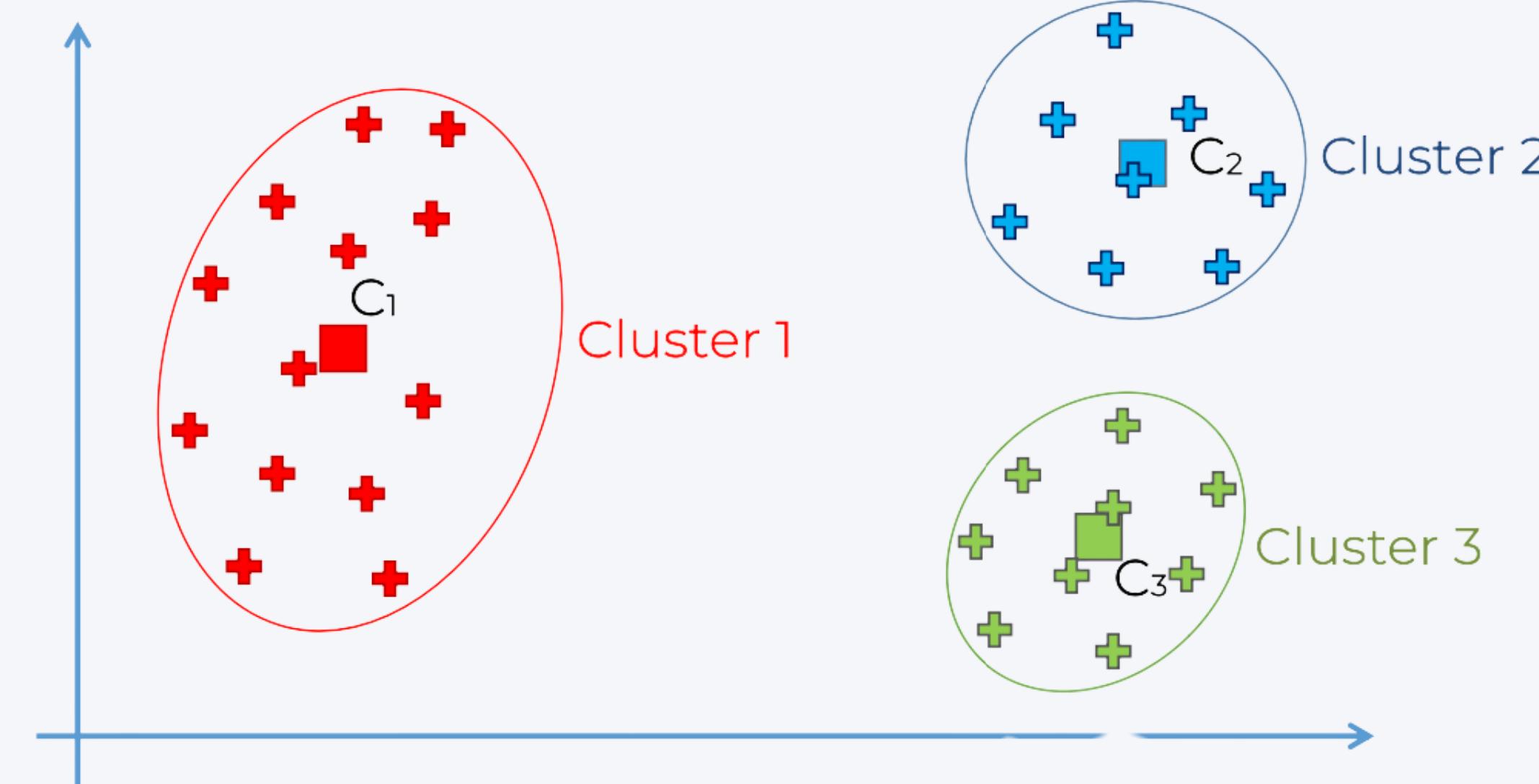
Case if we have two clusters

Within Cluster Sum of Squares:

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

K-Means Clustering

The Elbow Method



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

Case if we have three clusters and so on ...

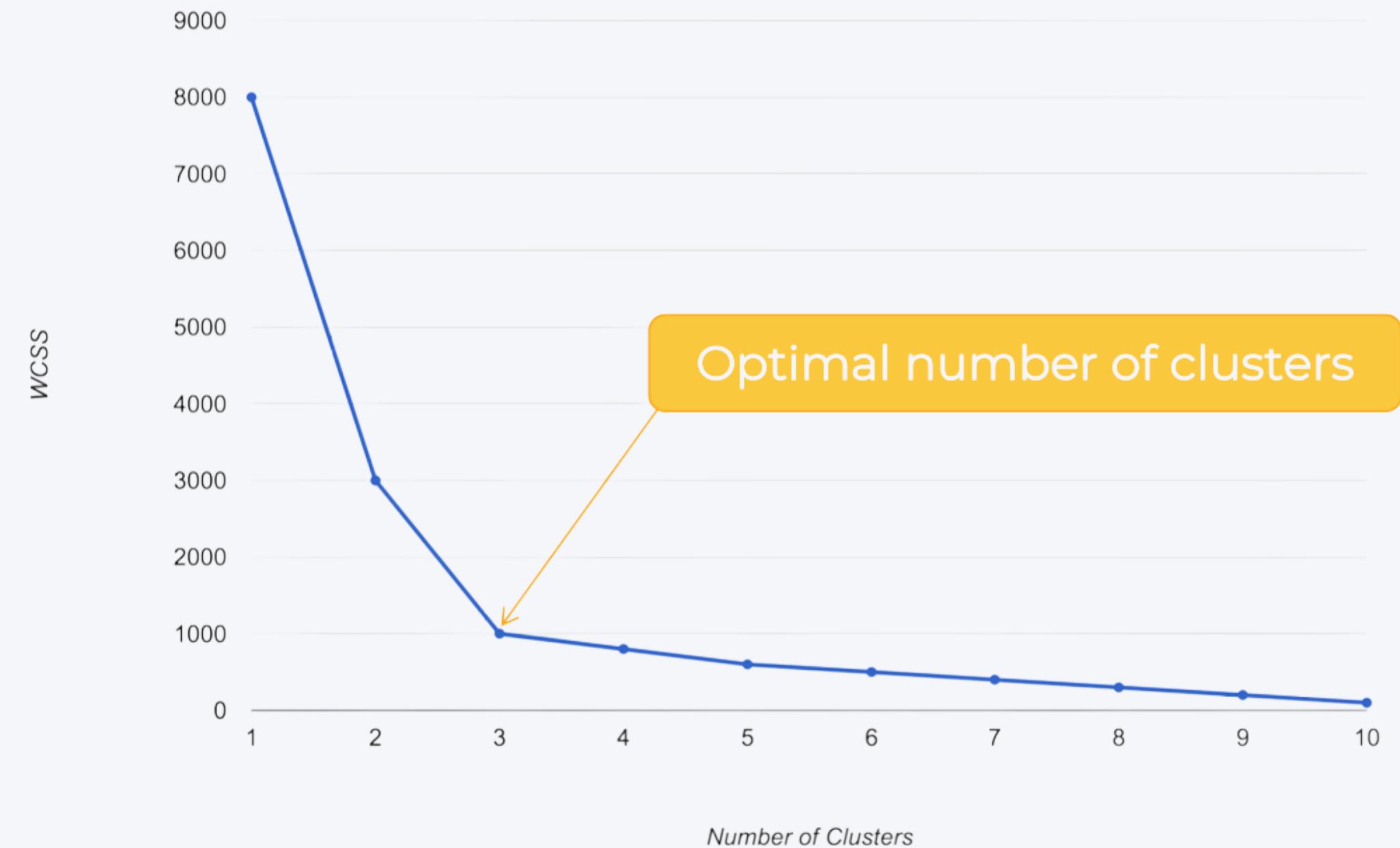
K-Means Clustering

The Elbow Method

- If you notice, we need the clusters to find WCSS .. so we need to run it first
- It is **backward !!**
- The more clusters we have, the smaller WCSS we have .. So we can continue increasing numbers of clusters until we get max number of clusters = number of points so WCSS will become 0
- The optimal number of clusters is three because it is the elbow , where WCSS stop dropping as rapidly.

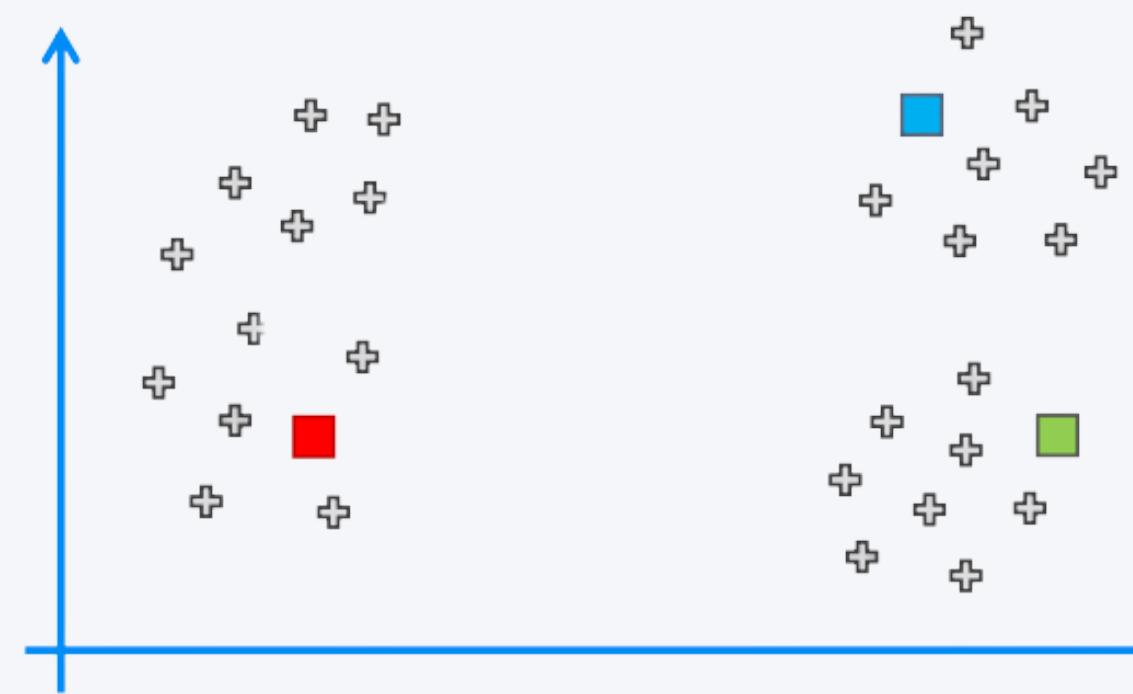
Judgment call !! sometimes not clear, more than candidate .. so

The Elbow Method

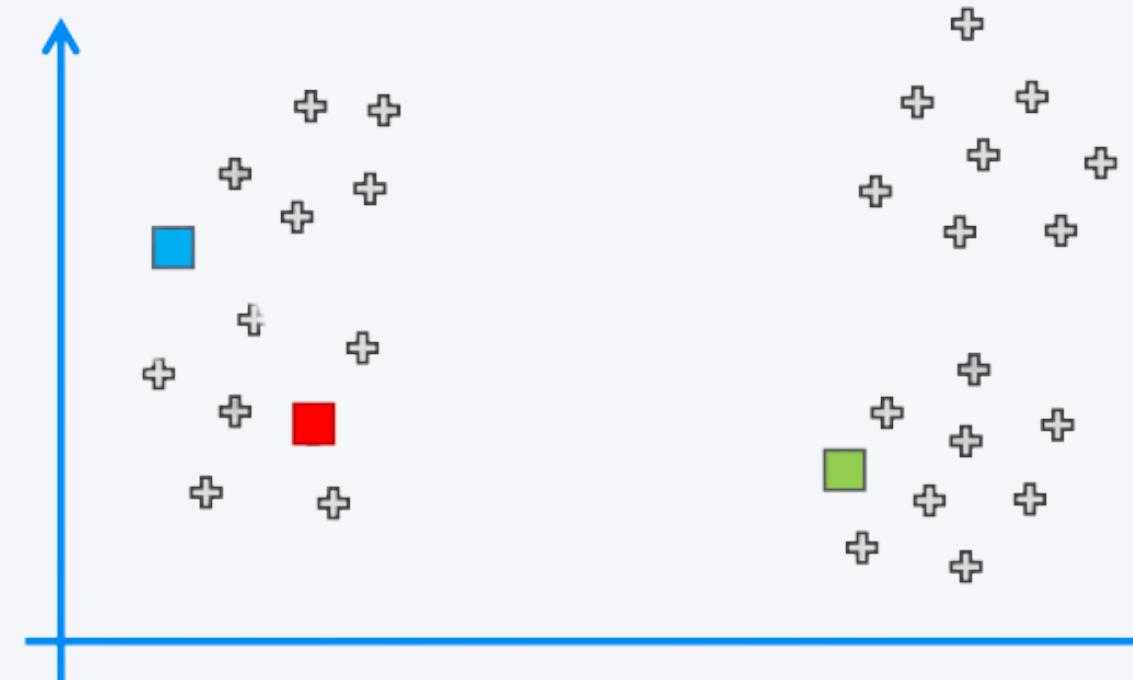
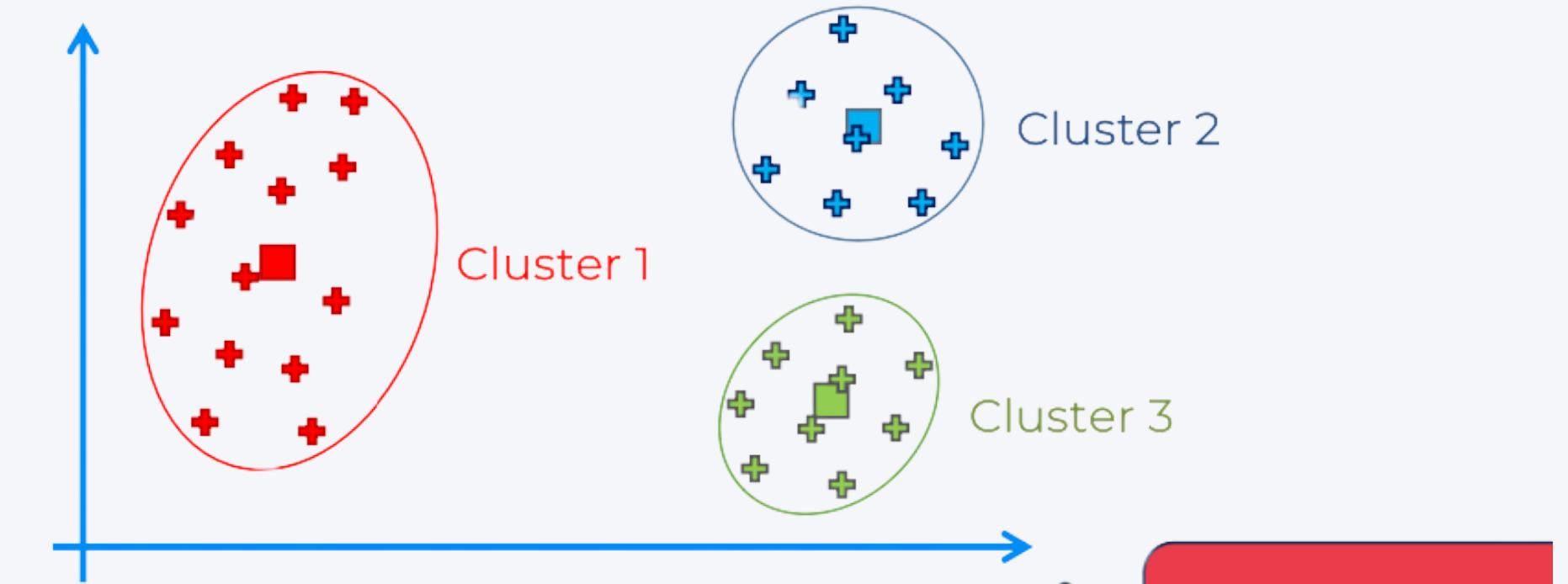


K-Means Clustering

K-Means++



K-Means



K-Means



K-Means Clustering

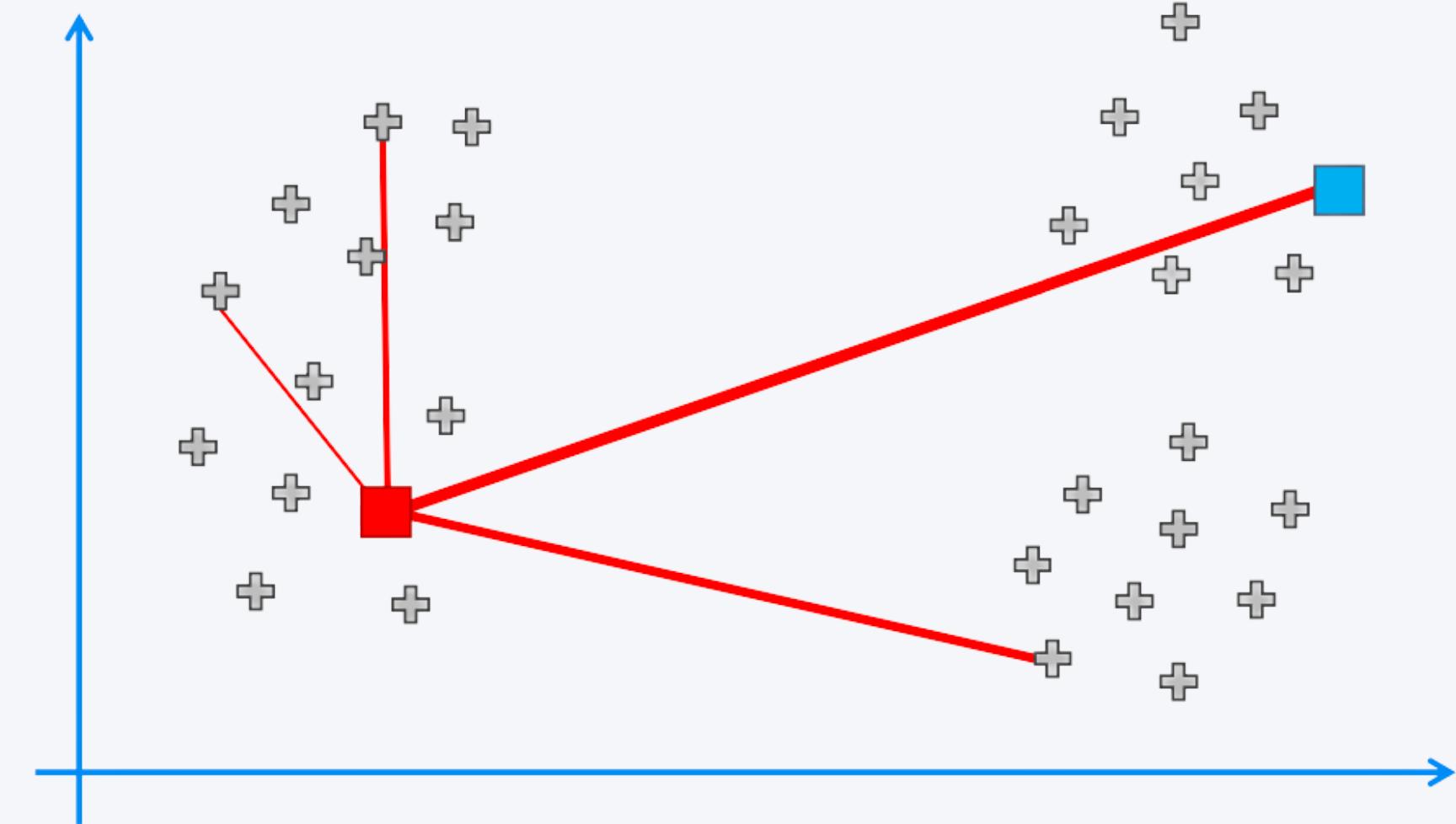
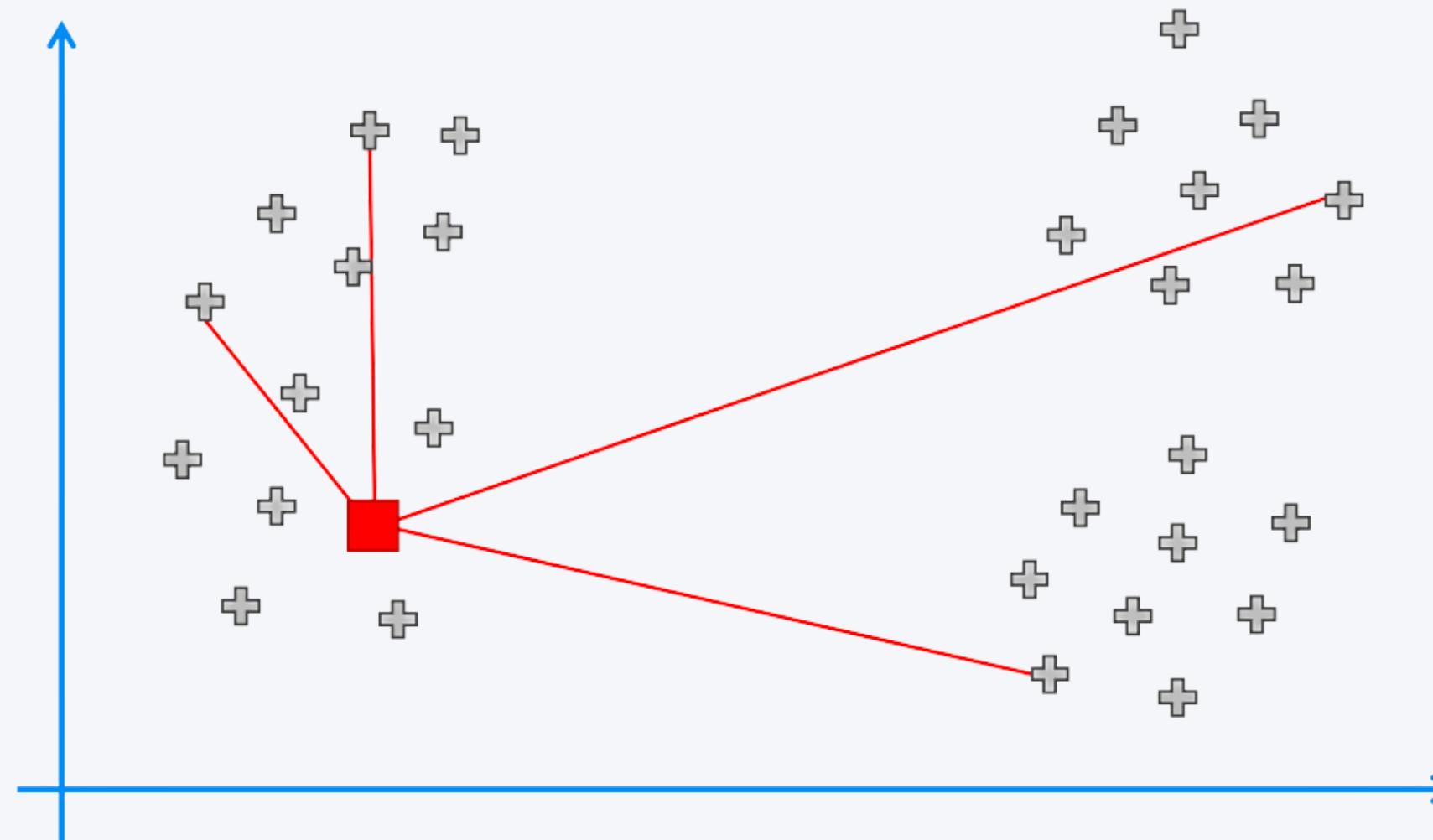
K-Means++

K-Means++ Initialization Algorithm:

- **Step 1:** Choose first centroid at random among data point
- **Step 2:** For each of the remaining data points compute the distance (D) to the nearest out of already selected centroids
- **Step 3:** Choose next centroid among remaining data points using weighted random selection – weighted by D^2
- **Step 4:** Repeat Steps 2 and 3 until all k centroids have been selected
- **Step 5:** Proceed with standard k-means clustering

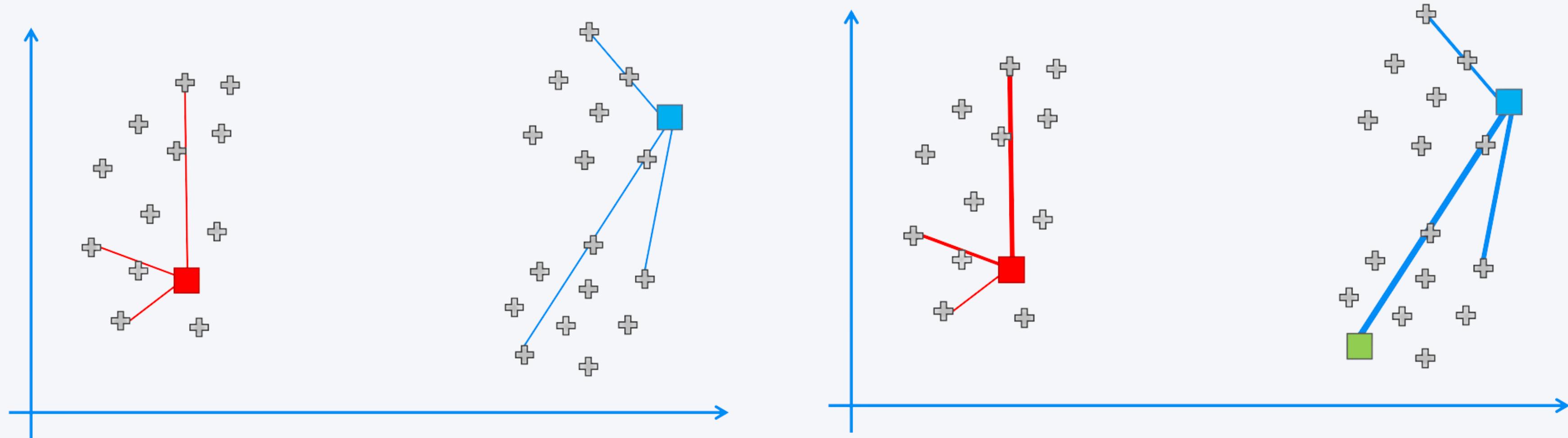
K-Means Clustering

K-Means++



K-Means Clustering

K-Means++



K-Means Clustering

K-Means++

