# Dr. SNS RAJALAKSHMI COLLEGE OF ARTS AND SCIENCE

## (AUTONOMOUS)

Accredited by NAAC (Cycle- III) with 'A+' Grade

## DEPARTMENT OF B.SC CS (GCD)

## 22UDA501 – INTRODUCTION TO DATA ANALYTICS
## UNIT- III

Dr.SNSRCAS  B.Sc CS(GCD)

# K-Means

K-Means Clustering is an [Unsupervised Learning algorithm](#), which groups the unlabeled dataset into different clusters.

Here K defines the number of pre-defined clusters that need to be created in the process, as if K=2, there will be two clusters, and for K=3, there will be three clusters

It is an iterative algorithm that divides the unlabeled dataset into k different clusters in such a way that each dataset belongs only one group that has similar properties.
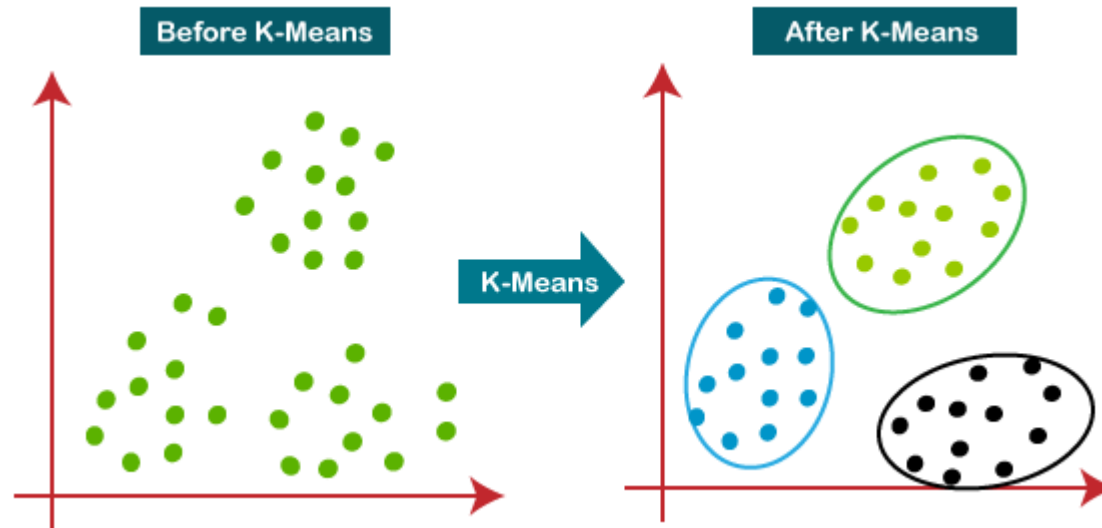
- It is a centroid-based algorithm, where each cluster is associated with a centroid.
- The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters.

- The algorithm takes the unlabeled dataset as input, divides the dataset into k-number of clusters, and repeats the process until it does not find the best clusters.
- The value of k should be predetermined in this algorithm.

# The k-means [clustering](clustering) algorithm mainly performs two tasks:

- Determines the best value for K center points or centroids by an iterative process.
- Assigns each data point to its closest k-center. Those data points which are near to the particular k-center, create a cluster.
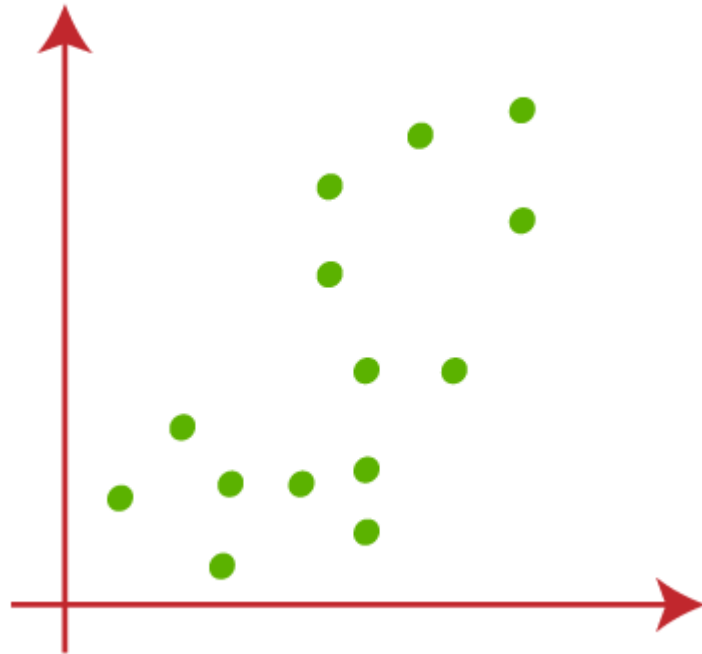
# K-means Clustering Algorithm:
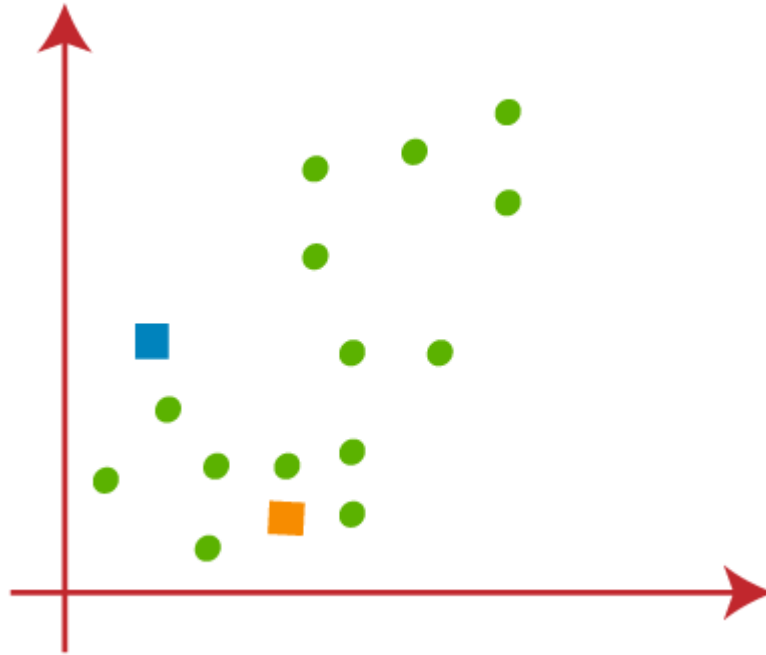
# How does the K-Means Algorithm Work?

- **Step-1:** Select the number K to decide the number of clusters.
- **Step-2:** Select random K points or centroids. (It can be other from the input dataset).
- **Step-3:** Assign each data point to their closest centroid, which will form the predefined K clusters.
- **Step-4:** Calculate the variance and place a new centroid of each cluster.
- **Step-5:** Repeat the third steps, which means reassign each datapoint to the new closest centroid of each cluster.
- **Step-6:** If any reassignment occurs, then go to step-4 else go to FINISH.
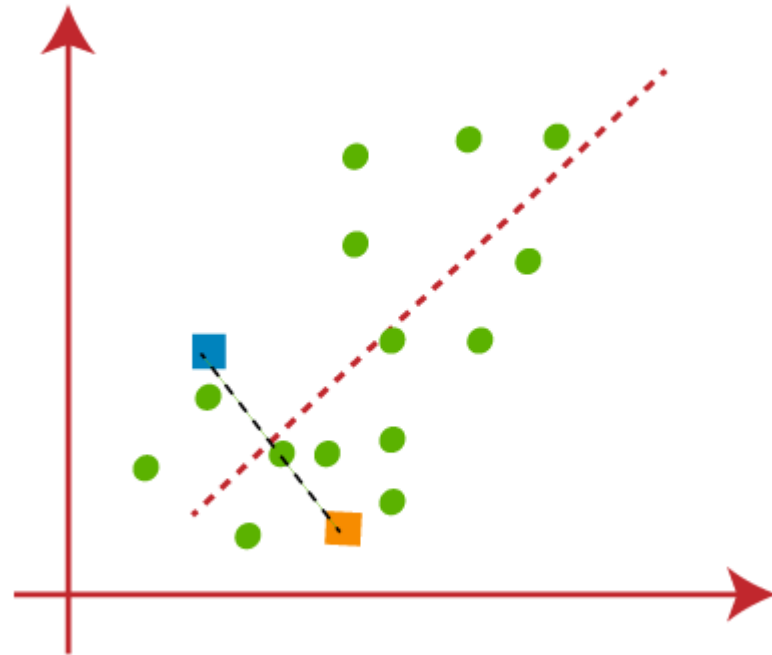- **Step-7**: The model is ready.

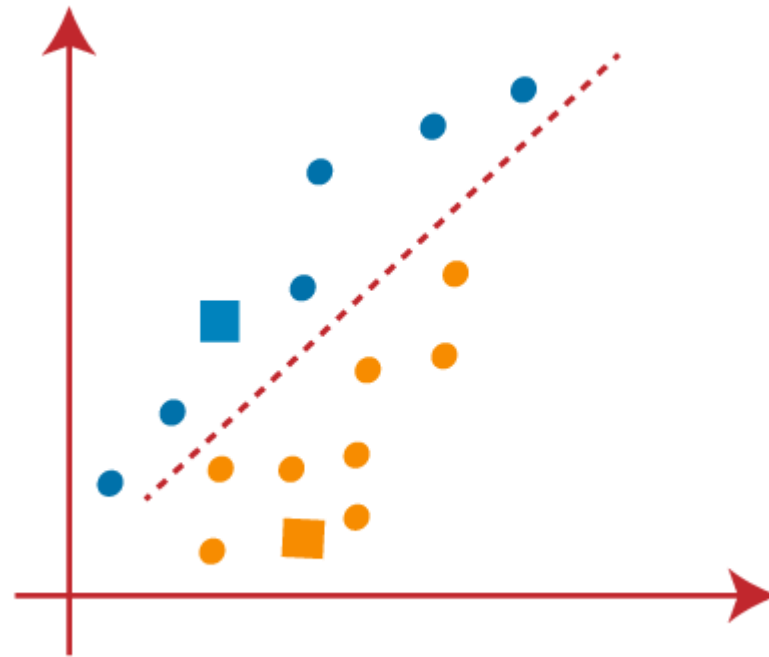we have two variables M1 and M2. The x-y axis scatter plot of these two variables

To choose some random k points or centroid to form the cluster. These points can be either the points from the dataset or any other point. So, here we are selecting the below two points as k points, which are not the part of our dataset.
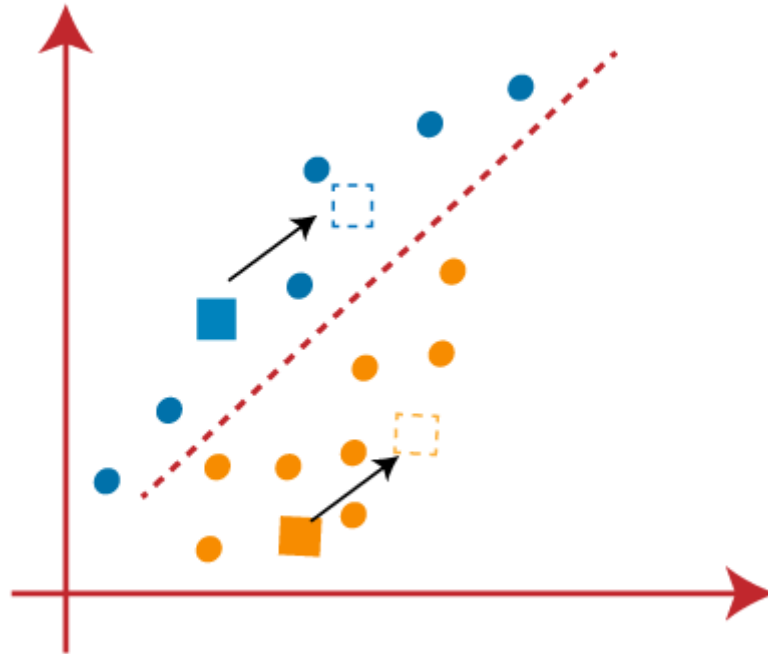
we will assign each data point of the scatter plot to its closest K-point or centroid. We will compute it by applying some mathematics that we have studied to calculate the distance between two points. So, we will draw a median between both the centroids.
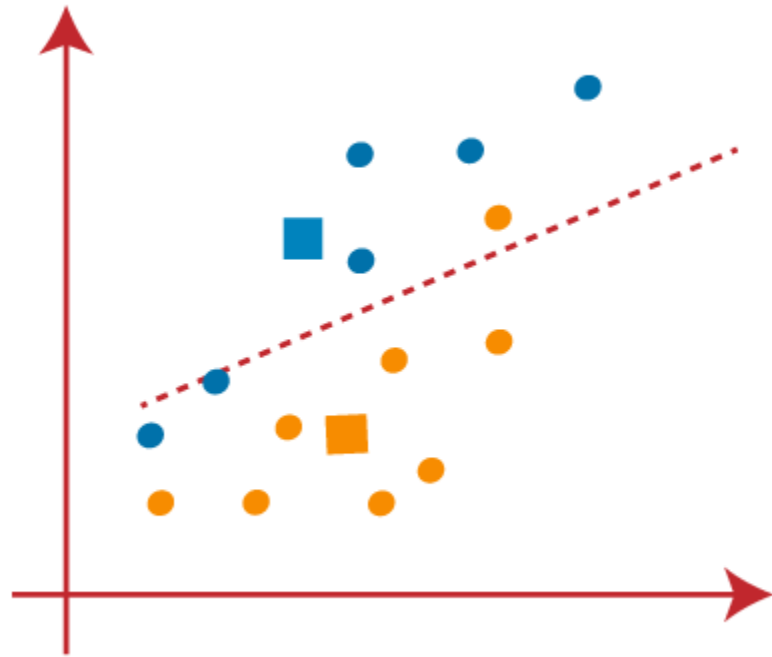
it is clear that points left side of the line is near to the K1 or blue centroid, and points to the right of the line are close to the yellow centroid. Let's color them as blue and yellow for clear visualization.
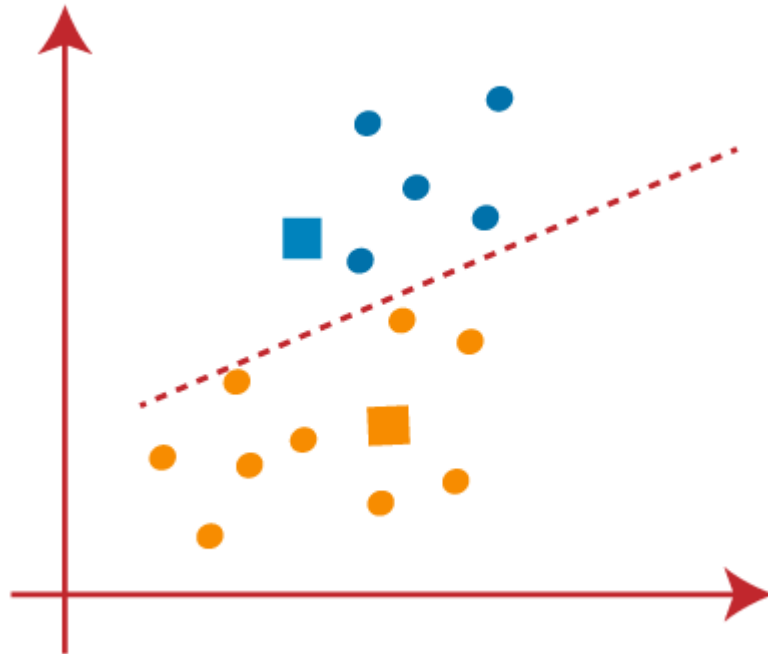
To find the closest cluster, so we will repeat the process by choosing **a new centroid**. To choose the new centroids, we will compute the center of gravity of these centroids, and centroids as below:

we will reassign each datapoint to the new centroid. For this, we will repeat the same process of finding a median line. The median will be like below image:

we can see, one yellow point is on the left side of the line, and two blue points are right to the line. So, these three points will be assigned to new centroids
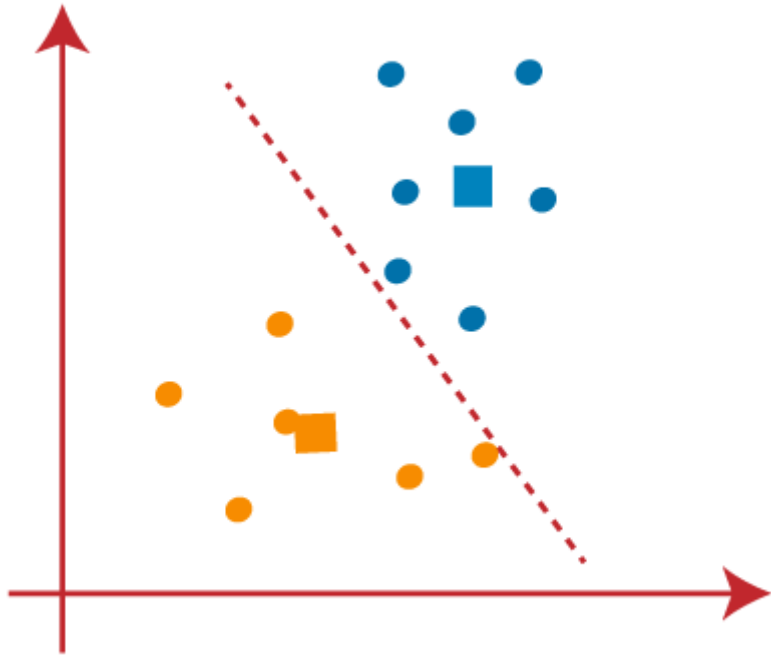
We will repeat the process by finding the center of gravity of centroids, so the new centroids will be as shown in the below image:
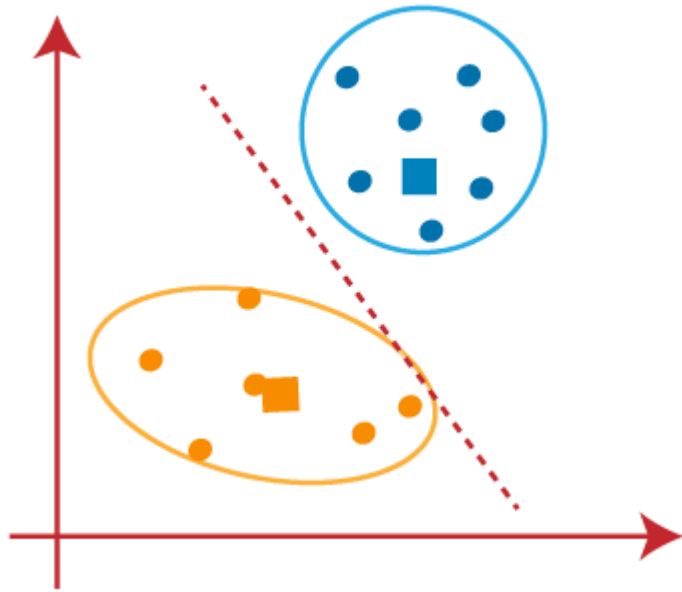
we got the new centroids so again will draw the median line and reassign the data points. So, the image will be:

There are no dissimilar data points on either side of the line, which means our model is formed. Consider the below image:

We can now remove the assumed centroids, and the two final clusters will be as shown in the below image