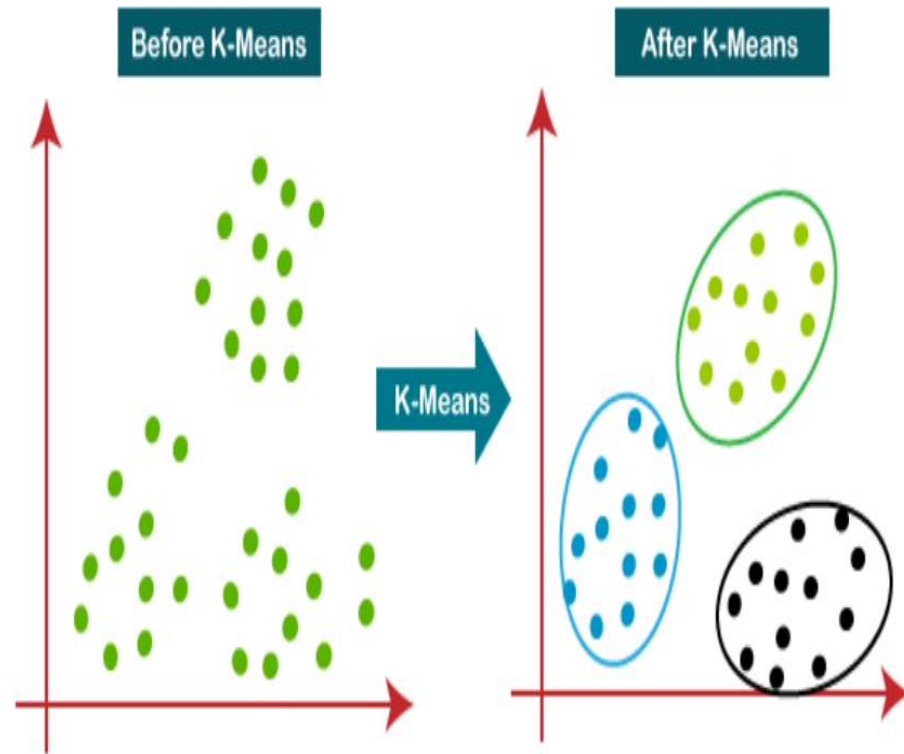# K Means Clustering

# What is K-Means Algorithm?

- K-Means Clustering is an Unsupervised Learning algorithm, which groups the <span style="color:red">unlabeled dataset into different clusters.</span>

- Here <span style="color:red">K defines the number of pre-defined clusters</span> 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, and so on.

- It allows us to cluster the <span style="color:red">data into different groups</span> and a convenient way to <span style="color:red">discover the categories of groups</span> in the unlabeled dataset on its own without the need for any training.

- 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 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.
- Hence each cluster has datapoints with some commonalities, and it is away from other clusters.
- The below diagram explains the working of the K-means Clustering Algorithm:

Before K-Means / After K-Means

How does the K-Means Algorithm Work?

The working of the K-Means algorithm is explained in the below steps:

**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).

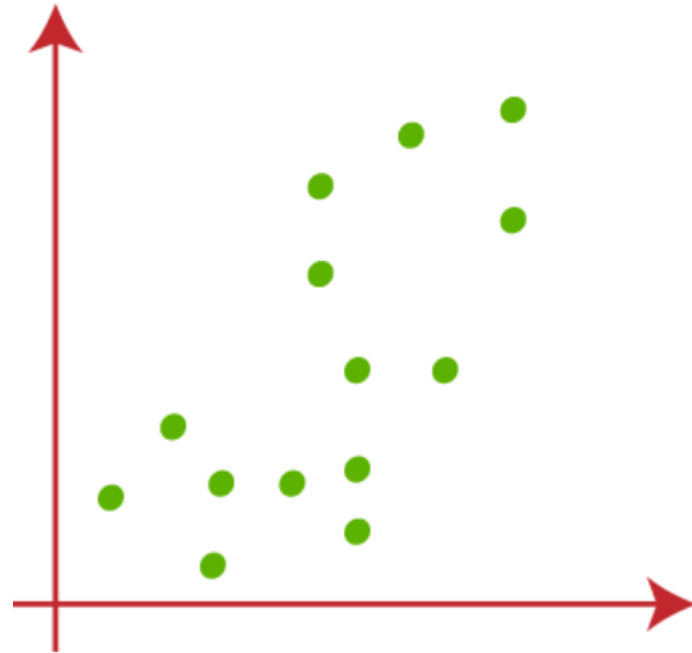**Step-3:** Assign each data point to their closest centroid, which will form the predefined K clusters.

**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.

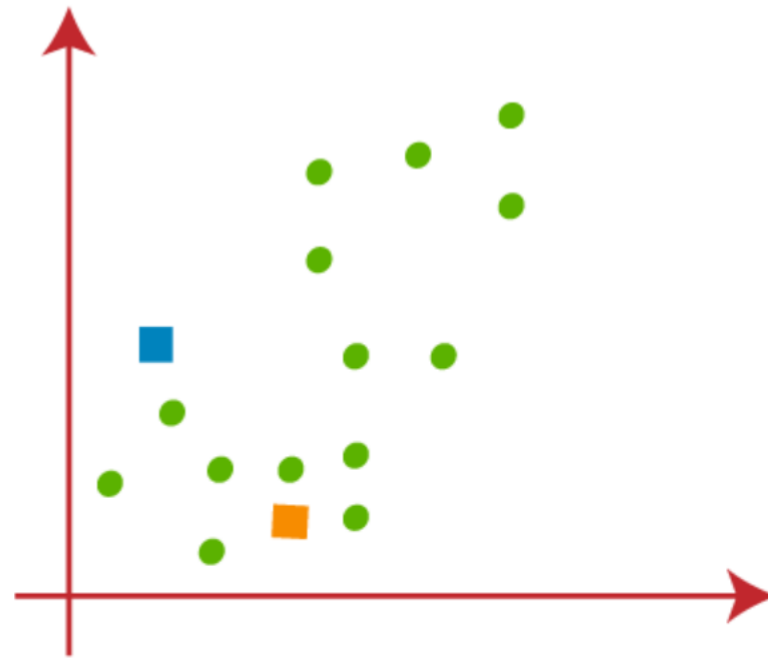**Step-6:** If any reassignment occurs, then go to step-4 else go to FINISH.

**Step-7**: The model is ready.

- Let's understand the above steps by considering the visual plots:
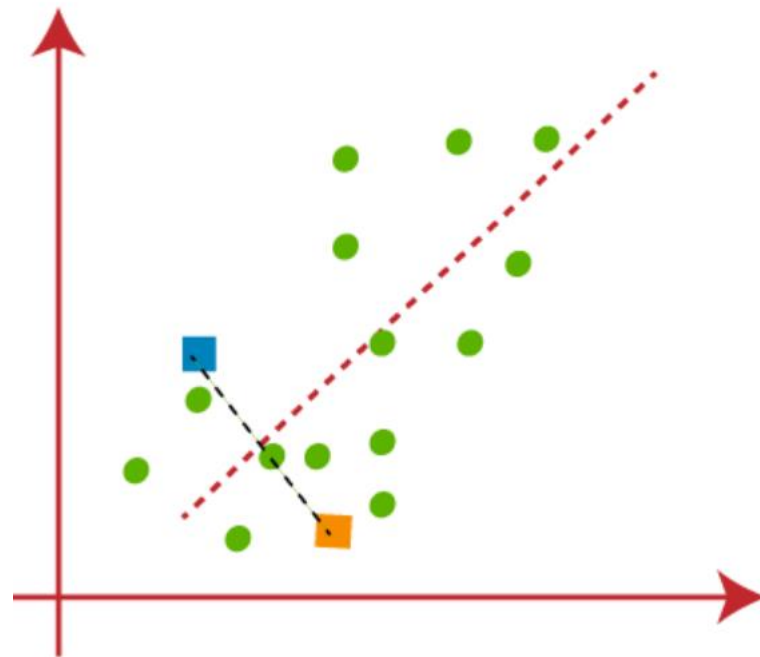- Suppose we have two variables M1 and M2. The x-y axis scatter plot of these two variables is given below:

- Let's take number k of clusters, i.e., K=2, to identify the dataset and to put them into different clusters. It means here we will try to group these datasets into two different clusters.
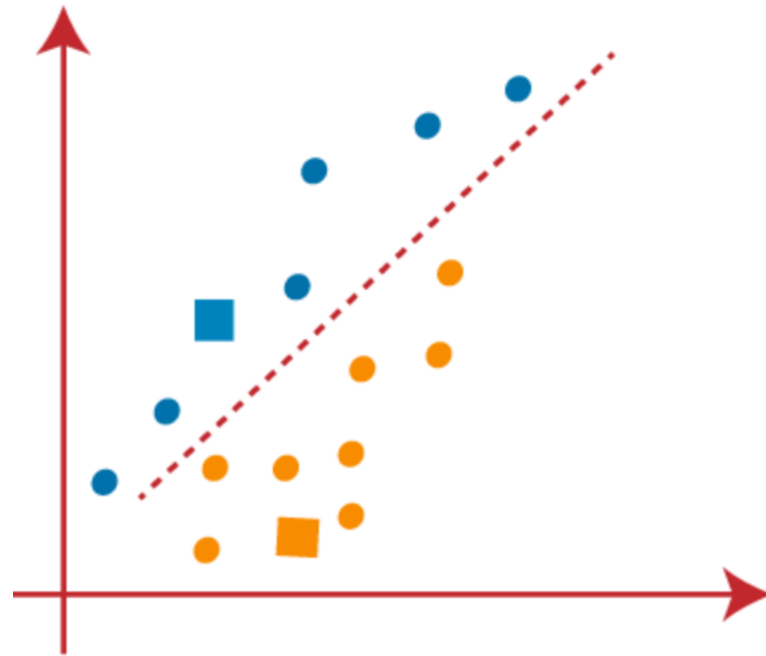
- We need to choose some random k points or centroid to form the cluster.
- These points can be either the points from the dataset or any other point. So, here we are selecting the below two points as k points, which are not the part of our dataset.
- Consider the below image:

- Now 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.

- Consider the below image:

- From the above image, 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.

- As we need 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 (average position) of these centroids, and will find new centroids as below:

$$\text{New Centroid} = \frac{1}{n}\sum_{i=1}^{n} x_i$$

where $n$ is the number of points in the cluster, and $x_i$ are the data points in the cluster.

- Next, 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:

- From the above 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.

- As reassignment has taken place, so we will again go to the step-4, which is finding new centroids or K-points.

- 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:

- As we got the new centroids so again will draw the median line and reassign the data points.
- So, the image will be:

- We can see in the above image; there are no dissimilar data points on either side of the line, which means our model is formed. Consider the below image:



- As our model is ready, so we can now remove the assumed centroids, and the two final clusters will be as shown in the above image

- How to choose the value of "K number of clusters" in K-means Clustering?

- Here are some different ways to find the optimal number of clusters, but here we are discussing the most appropriate method to find the number of clusters or value of K. The method is given below:

**Elbow Method**

- The Elbow method is one of the most popular ways to find the optimal number of clusters. This method uses the concept of WCSS value. WCSS stands for Within Cluster Sum of Squares, which defines the total variations within a cluster. The formula to calculate the value of WCSS (for 3 clusters) is given below:

$$WCSS = \sum_{P_i \text{ in Cluster1}} distance(P_i \ C_1)^2 + \sum_{P_i \text{ in Cluster2}} distance(P_i \ C_2)^2 + \sum_{P_i \text{ in CLuster3}} distance(P_i \ C_3)^2$$

- In the above formula of WCSS,
- $\sum P_{i\ in\ Cluster1}$ distance$(P_i\ C_1)^2$: It is the sum of the square of the distances between each data point and its centroid within a cluster1 and the same for the other two terms.
- To measure the distance between data points and centroid, we can use any method such as Euclidean distance or Manhattan distance.
- To find the optimal value of clusters, the elbow method follows the below steps:
  - It executes the K-means clustering on a given dataset for different K values (ranges from 1-10).
  - For each value of K, calculates the WCSS value.
  - Plots a curve between calculated WCSS values and the number of clusters K.
  - The sharp point of bend or a point of the plot looks like an arm, then that point is considered as the best value of K.

- Since the graph shows the sharp bend, which looks like an elbow, hence it is known as the elbow method. The graph for the elbow method looks like the below image:

-

# KModes clustering

- **KModes clustering** is one of the unsupervised Machine Learning algorithms that is used to cluster **categorical variables.**

- KMeans uses mathematical measures (distance) to cluster continuous data. The lesser the distance, the more similar our data points are. Centroids are updated by Means.

- But for categorical data points, we cannot calculate the distance. So we go for KModes algorithm. It uses the dissimilarities(total mismatches) between the data points.

- The lesser the dissimilarities the more similar our data points are. It uses Modes instead of means.

How does the KModes algorithm work?

- Unlike Hierarchical clustering methods, we need to upfront specify the K.
- Pick K observations at random and use them as leaders/clusters.
- Calculate the dissimilarities and assign each observation to its closest cluster.
- Define new modes for the clusters.
- Repeat 2–3 steps until there are is no re-assignment required
- I hope you got the basic idea of the KModes algorithm by now. So let us quickly take an example to illustrate the working step by step.

**Example:** Imagine we have a dataset that has the information about hair color, eye color, and skin color of persons. We aim to group them based on the available information(maybe we want to suggest some styling ideas)

Hair color, eye color, and skin color are all categorical variables. Below 👇 is how our dataset looks like.

| person | hair color | eye color | skin color |
|--------|-----------|-----------|-----------|
| P1 | blonde | amber | fair |
| P2 | brunette | gray | brown |
| P3 | red | green | brown |
| P4 | black | hazel | brown |
| P5 | brunette | amber | fair |
| P6 | black | gray | brown |
| P7 | red | green | fair |
| P8 | black | hazel | fair |

*Image of our data*

Alright, we have the sample data now. Let us proceed by defining the number of clusters(K)=3

# Step 1: Pick K observations at random and use them as leaders/clusters

I am choosing P1, P7, P8 as leaders/clusters

| Leaders | | | |
|---|---|---|---|
| **P1** | blonde | amber | fair |
| **P7** | red | green | fair |
| **P8** | black | hazel | fair |

| person | hair color | eye color | skin color |
|---|---|---|---|
| **P1** | blonde | amber | fair |
| **P2** | brunette | gray | brown |
| **P3** | red | green | brown |
| **P4** | black | hazel | brown |
| **P5** | brunette | amber | fair |
| **P6** | black | gray | brown |
| **P7** | red | green | fair |
| **P8** | black | hazel | fair |

*Leaders and Observations*

# Step 2: Calculate the dissimilarities(no. of mismatches) and assign each observation to its closest cluster

Iteratively compare the cluster data points to each of the observations. Similar data points give 0, dissimilar data points give 1.

| Leaders | | | |
|---|---|---|---|
| P1 | blonde | amber | fair |
| P7 | red | green | fair |
| P8 | black | hazel | fair |

| person | hair color | eye color | skin color |
|---|---|---|---|
| P1 | blonde | amber | fair |
| P2 | brunette | gray | brown |
| P3 | red | green | brown |
| P4 | black | hazel | brown |
| P5 | brunette | amber | fair |
| P6 | black | gray | brown |
| P7 | red | green | fair |
| P8 | black | hazel | fair |

Comparing leader/Cluster P1 to the observation P1 gives 0 dissimilarities.

| Leaders | | | |
|---|---|---|---|
| **P1** | blonde | amber | fair |
| **P7** | red | green | fair |
| **P8** | black | hazel | fair |

| person | hair color | eye color | skin color |
|---|---|---|---|
| **P1** | blonde | amber | fair |
| **P2** | brunette | gray | brown |
| **P3** | red | green | brown |
| **P4** | black | hazel | brown |
| **P5** | brunette | amber | fair |
| **P6** | black | gray | brown |
| **P7** | red | green | fair |
| **P8** | black | hazel | fair |

Comparing leader/cluster P1 to the observation P2 gives 3(1+1+1) dissimilarities.

Likewise, calculate all the dissimilarities and put them in a matrix as shown below and assign the observations to their closest cluster(cluster that has the least dissimilarity)

| | Cluster 1 (P1) | Cluster 2 (P7) | Cluster 3 (P8) | Cluster |
|---|---|---|---|---|
| P1 | 0 ✓ | 2 | 2 | Cluster 1 |
| P2 | 3 ✓ | 3 | 3 | Cluster 1 |
| P3 | 3 | 1 ✓ | 3 | Cluster 2 |
| P4 | 3 | 3 | 1 ✓ | Cluster 3 |
| P5 | 1 ✓ | 2 | 2 | Cluster 1 |
| P6 | 3 | 3 | 2 ✓ | Cluster 3 |
| P7 | 2 | 0 ✓ | 2 | Cluster 2 |
| P8 | 2 | 2 | 0 ✓ | Cluster 3 |

*Dissimilarity matrix (Image by Author)*

After step 2, the observations P1, P2, P5 are assigned to cluster 1; P3, P7 are assigned to Cluster 2; and P4, P6, P8 are assigned to cluster 3.

*Note: If all the clusters have the same dissimilarity with an observation, assign to any cluster randomly. In our case, the observation P2 has 3 dissimilarities with all the leaders. I randomly assigned it to Cluster 1.*

# Step 3: Define new modes for the clusters

Mode is simply the **most observed value**.

Mark the observations according to the cluster they belong to. Observations of Cluster 1 are marked in Yellow, Cluster 2 are marked in Brick red, and Cluster 3 are marked in Purple.

| person | hair color | eye color | skin color |
|--------|-----------|-----------|-----------|
| P1 | blonde | amber | fair |
| P2 | brunette | gray | brown |
| P3 | red | green | brown |
| P4 | black | hazel | brown |
| P5 | brunette | amber | fair |
| P6 | black | gray | brown |
| P7 | red | green | fair |
| P8 | black | hazel | fair |

*Looking for Modes (Image by author)*

Considering one cluster at a time, for each feature, look for the Mode and update the new leaders.

**Explanation:** Cluster 1 observations(P1, P2, P5) has brunette as the most observed hair color, amber as the most observed eye color, and fair as the most observed skin color.

*Note: If you observe the same occurrence of values, take the mode randomly. In our case, the observations of Cluster 3(P3, P7) have one occurrence of brown, fair skin color. I randomly chose brown as the mode.*

Below are our new leaders after the update.

| New Leaders | | | |
|---|---|---|---|
| | hair color | eye color | skin color |
| **Cluster 1** | brunette | amber | fair |
| **Cluster 2** | red | green | fair |
| **Cluster 3** | black | hazel | brown |

*Obtained new leaders*

**Repeat steps 2-4**

After obtaining the new leaders, again calculate the dissimilarities between the observations and the newly obtained leaders.

| New Leaders | | | |
|---|---|---|---|
| | hair color | eye color | skin color |
| **Cluster 1** | brunette | amber | fair |
| **Cluster 2** | red | green | fair |
| **Cluster 3** | black | hazel | brown |

| person | hair color | eye color | skin color |
|---|---|---|---|
| **P1** | blonde | amber | fair |
| **P2** | brunette | gray | brown |
| **P3** | red | green | brown |
| **P4** | black | hazel | brown |
| **P5** | brunette | amber | fair |
| **P6** | black | gray | brown |
| **P7** | red | green | fair |
| **P8** | black | hazel | fair |

Comparing Cluster 1 to the observation P1 gives 1 dissimilarity.

| New Leaders | | | |
|---|---|---|---|
| | **hair color** | **eye color** | **skin color** |
| **Cluster 1** | brunette | amber | fair |
| **Cluster 2** | red | green | fair |
| **Cluster 3** | black | hazel | brown |
| | | | |
| **person** | **hair color** | **eye color** | **skin color** |
| **P1** | blonde | amber | fair |
| **P2** | brunette | gray | brown |
| **P3** | red | green | brown |
| **P4** | black | hazel | brown |
| **P5** | brunette | amber | fair |
| **P6** | black | gray | brown |
| **P7** | red | green | fair |
| **P8** | black | hazel | fair |

Comparing Cluster 1 to the observation P2 gives 2 dissimilarities.

Likewise, calculate all the dissimilarities and put them in a matrix. Assign each observation to its closest cluster.

| | **Cluster 1** | **Cluster 2** | **Cluster 3** | **Cluster** |
|---|---|---|---|---|
| **P1** | 1 ✓ | 2 | 3 | Cluster 1 |
| **P2** | 2 ✓ | 3 | 2 | Cluster 1 |
| **P3** | 3 | 1 ✓ | 2 | Cluster 2 |
| **P4** | 3 | 3 | 0 ✓ | Cluster 3 |
| **P5** | 0 ✓ | 2 | 3 | Cluster 1 |
| **P6** | 3 | 3 | 1 ✓ | Cluster 3 |
| **P7** | 2 | 0 ✓ | 3 | Cluster 2 |
| **P8** | 2 | 2 | 1 ✓ | Cluster 3 |

The observations P1, P2, P5 are assigned to Cluster 1; P3, P7 are assigned to Cluster 2; and P4, P6, P8 are assigned to Cluster 3.