# 1. What is clustering? Explain in detail.
# 2. Explain K Means clustering.

K-Means Clustering is an unsupervised learning algorithm that is used to solve the clustering problems in machine learning or data science.
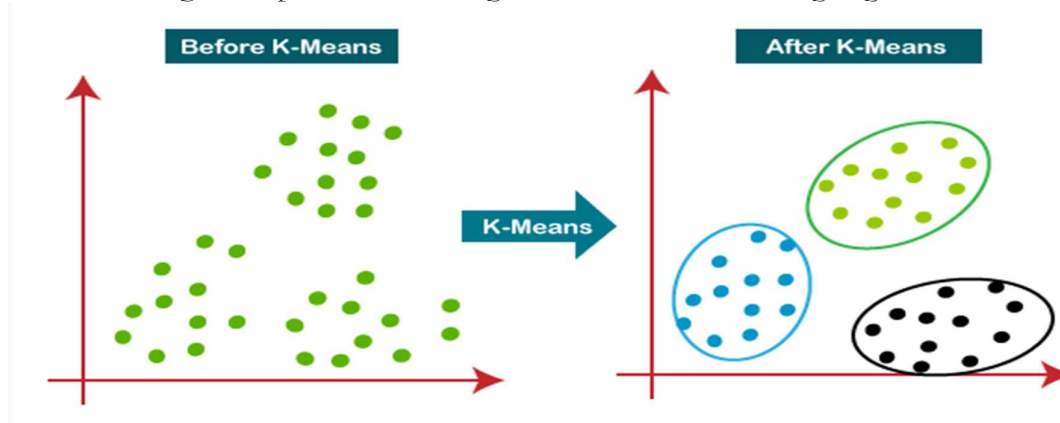
K-Means Clustering is an Unsupervised Learning algorithm, which groups the unlabeled dataset into different clusters. Here K defines the number of pre-defined clusters that need to be created in the process, as if K=2, there will be two clusters, and for K=3, there will be three clusters, and so on.

- 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 allows us to cluster the data into different groups and a convenient way to discover the categories of groups 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:

o Determines the best value for K center points or centroids by an iterative process.

o Assigns each data point to its closest k-center. Those data points which are near to the particular k-center, create a cluster.

The below diagram explains the working of the K-means Clustering Algorithm:



Working of K-Means Algorithm:

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).

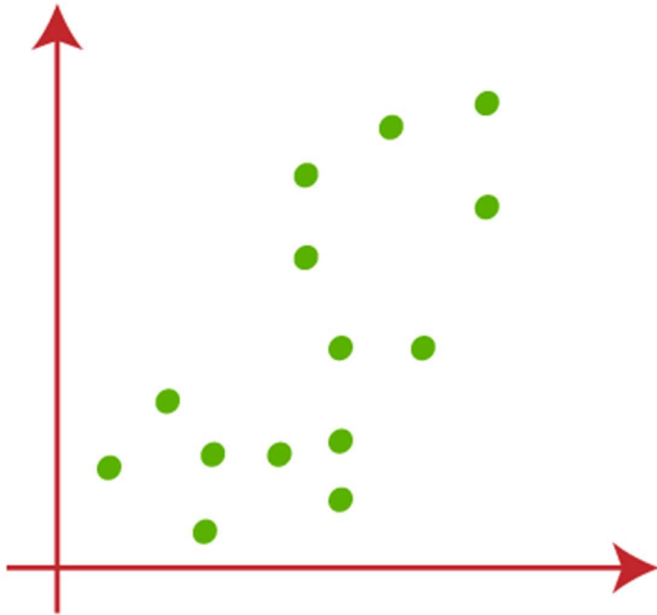**Step-3:** Assign each data point to their closest centroid, which will form the predefined K clusters.

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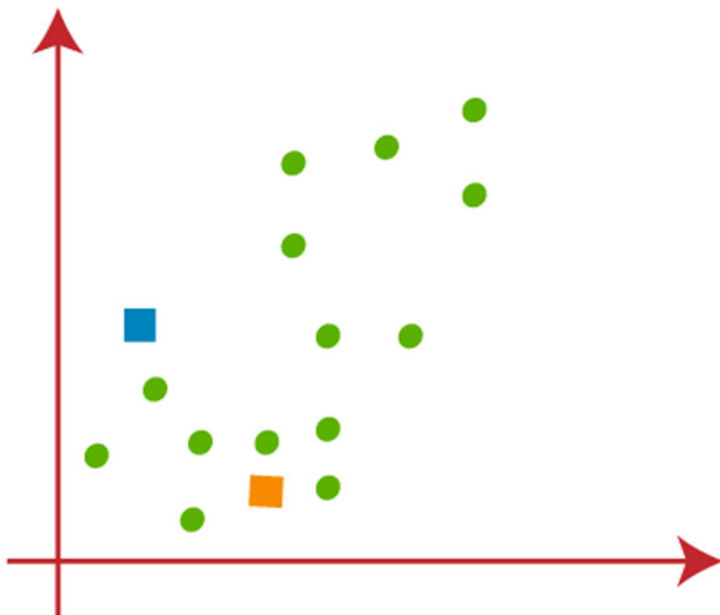
**Step-6:** If any reassignment occurs, then go to step-4 else go to FINISH.

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

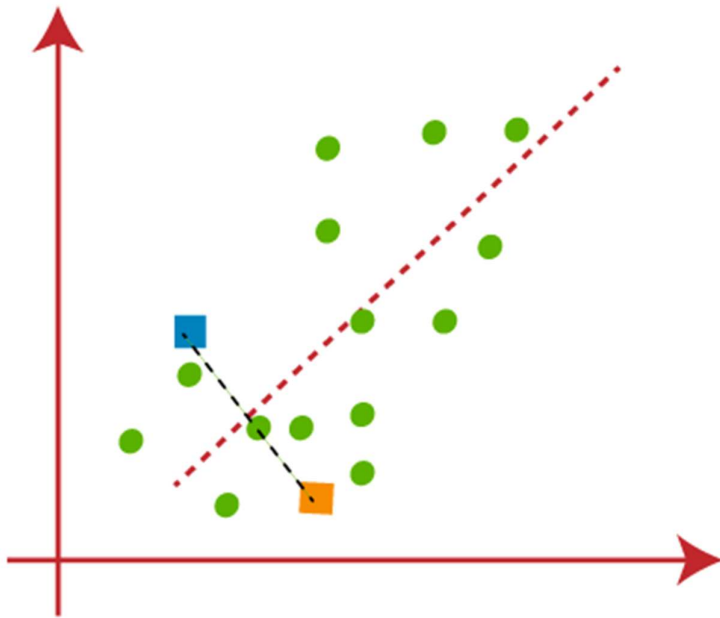Suppose we have two variables M1 and M2. The x-y axis scatter plot of these two variables is given below:



- o 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.
- o 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:
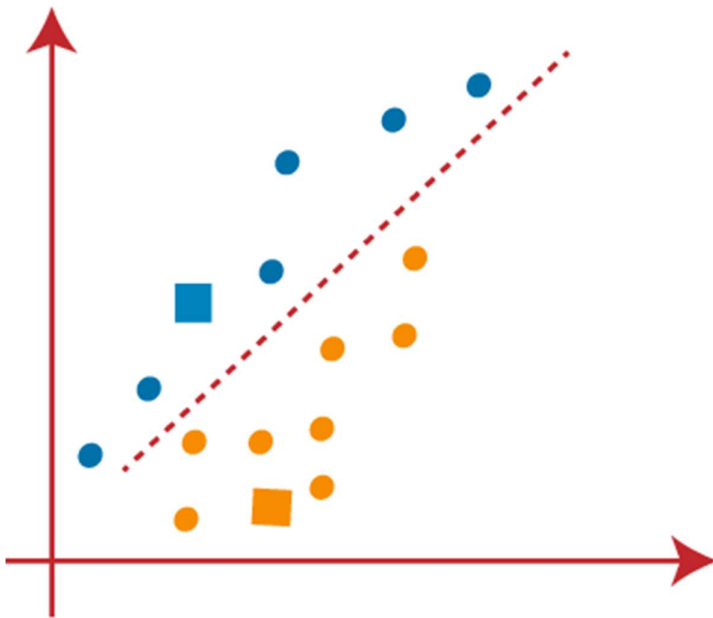


Now we will assign each data point of the scatter plot to its closest K-point or centroid. We will compute it by applying some mathematics that we have studied to calculate the distance between two points. So, we

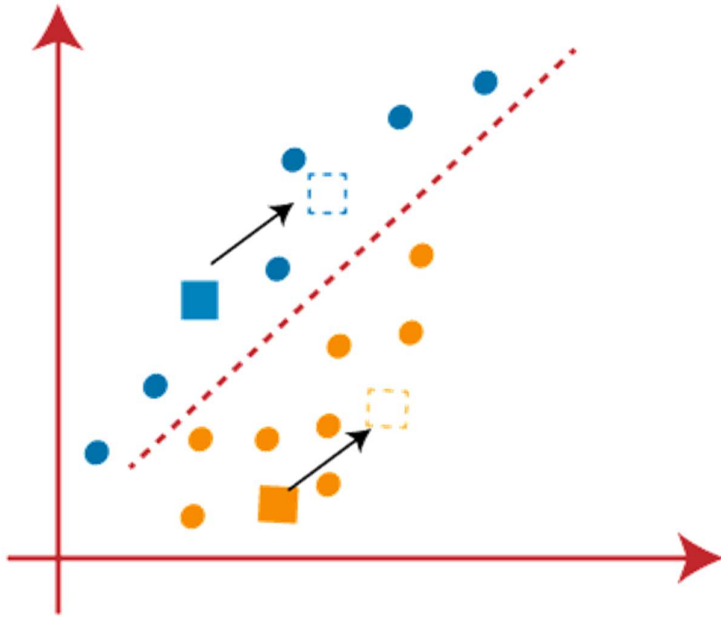will draw a median between both the centroids. Consider the below image:



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.
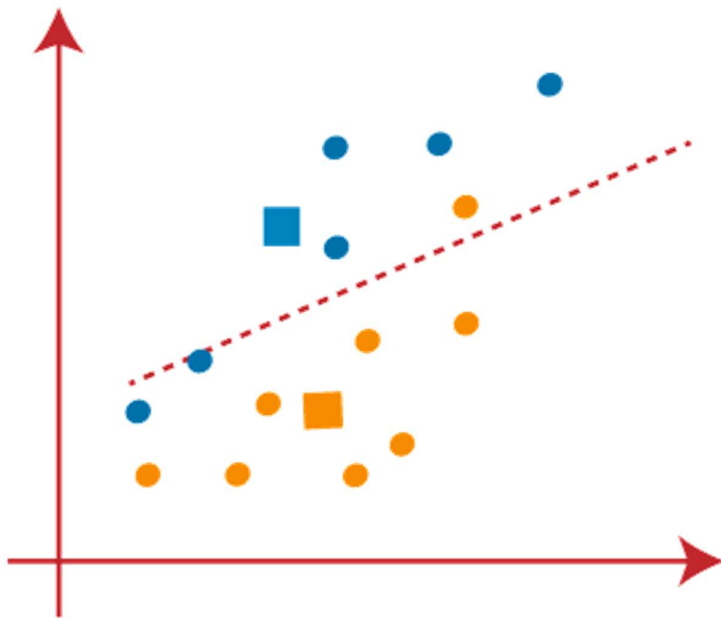


As we need to find the closest cluster, so we will repeat the process by choosing **a new centroid**. To choose the new centroids, we will compute the center of gravity of these centroids, and will find new centroids as
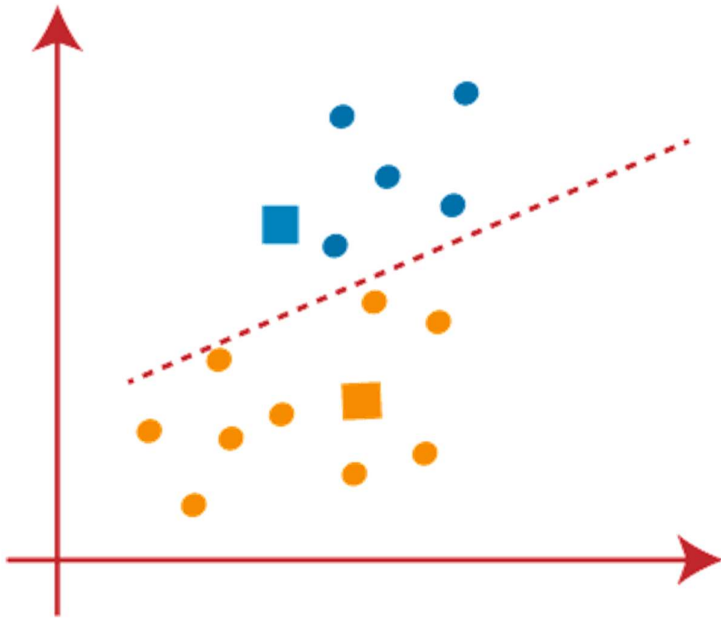
below:



- o  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.
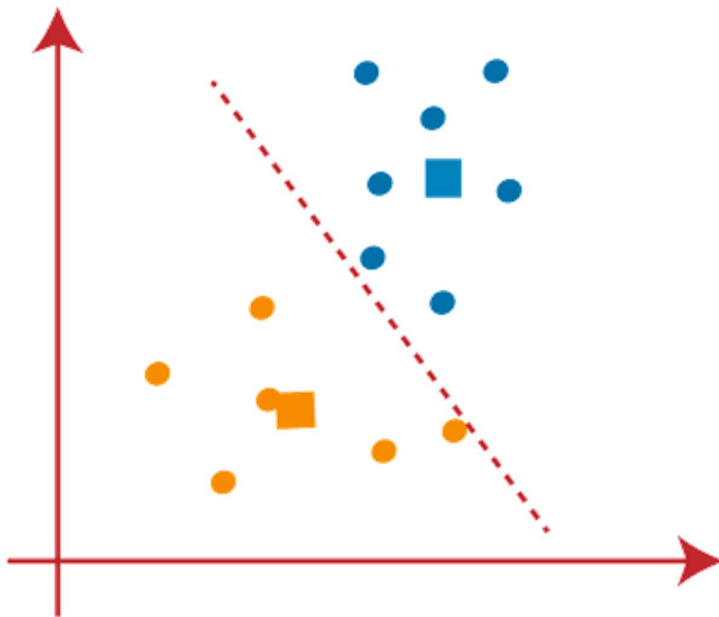
As reassignment has taken place, so we will again go to the step-4, which is finding new centroids or K-points.
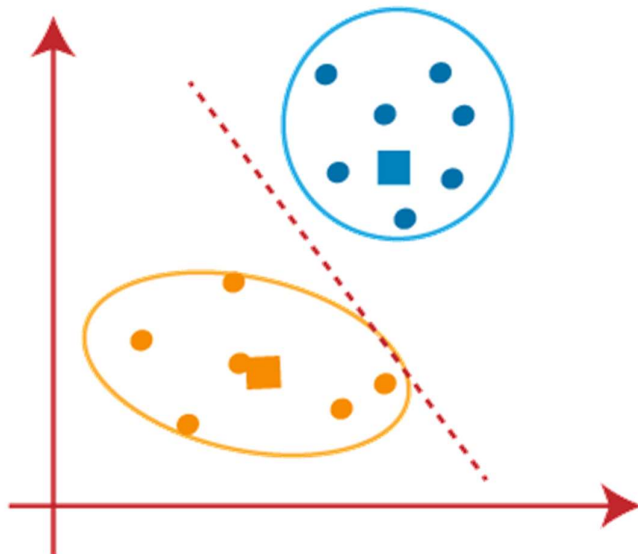
- o We will repeat the process by finding the center of gravity of centroids, so the new centroids will be as shown in the below image:
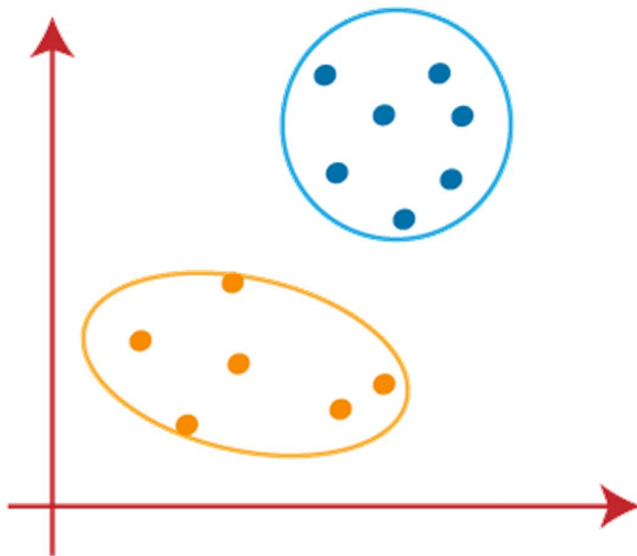
o   As we got the new centroids so again will draw the median line and reassign the data points. So, the image will be:



o   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 below image:

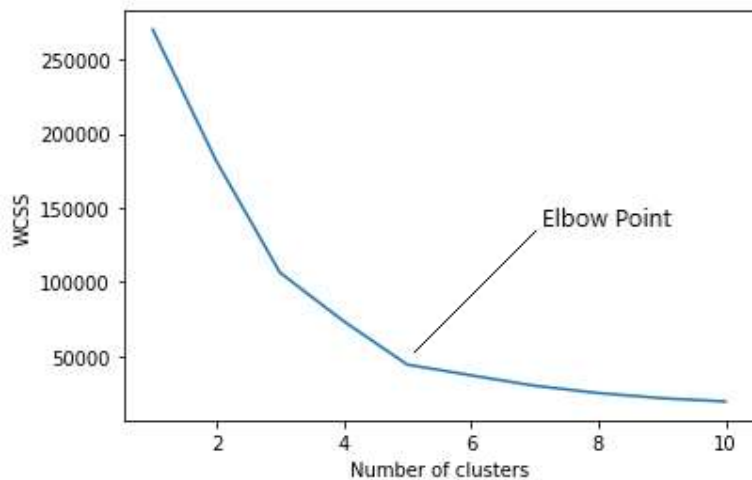### 3. Explain Elbow Method in K Means clustering

The Elbow Method is a technique used in data analysis and machine learning for determining the optimal number of clusters in a dataset. It involves plotting the variance explained by different numbers of clusters and identifying the "elbow" point, where the rate of variance decreases sharply levels off, suggesting an appropriate cluster count for analysis or model training.

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:

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

### K Means Clustering Using the Elbow Method

In the Elbow method, we are actually varying the number of clusters (K) from 1 – 10. For each value of K, we are calculating WCSS (Within-Cluster Sum of Square). WCSS is the sum of the squared distance between each point and the centroid in a cluster. When we plot the WCSS with the K value, the plot looks like an Elbow. As the number of clusters increases, the WCSS value will start to decrease. WCSS value is largest when K = 1. When we analyze the graph, we can see that the graph will rapidly change at a point and thus creating an elbow shape. From this point, the graph moves almost parallel to the X-axis. The K value corresponding to this point is the optimal value of K or an optimal number of clusters.

**Pros of Elbow Method:**

1. Simplicity: It's easy to understand and implement. You just plot the WCSS vs number of clusters (k) and look for the elbow.
2. Visualization: The graph provides a visual aid for choosing k, making it easier to understand the relationship between cluster number and data fit.
3. Efficiency: It's computationally cheap compared to other methods for choosing k.

**Elbow Method Drawbacks:**

The elbow method, while a useful tool for determining the optimal number of clusters in K-means clustering, has some drawbacks:

1. Subjectivity: The choice of the "elbow point" can be subjective and might vary between individuals analyzing the same data.
2. Non-Gaussian Data: It assumes that clusters are spherical and equally sized, which may not hold for complex datasets with irregularly shaped or differently sized clusters.
3. Sensitivity to Initialization: K-means itself is sensitive to initial cluster centroids, which can affect the WCSS values and, consequently, the choice of the optimal K.
4. Inefficient for Large Datasets: For large datasets, calculating WCSS for a range of K values can be computationally expensive and time-consuming.
5. Unsuitable for All Distributions: The elbow method is not suitable for all data distributions, especially when clusters have varying densities or are non-convex.
6. Limited to K-means: It specifically applies to K-means clustering and may not be suitable for other clustering algorithms with different objectives.

4. What is Agglomerative Hierarchical clustering?

The agglomerative hierarchical clustering algorithm is a popular example of HCA. To group the datasets into clusters, it follows the **bottom-up approach**. It means, this algorithm considers each dataset as a single cluster
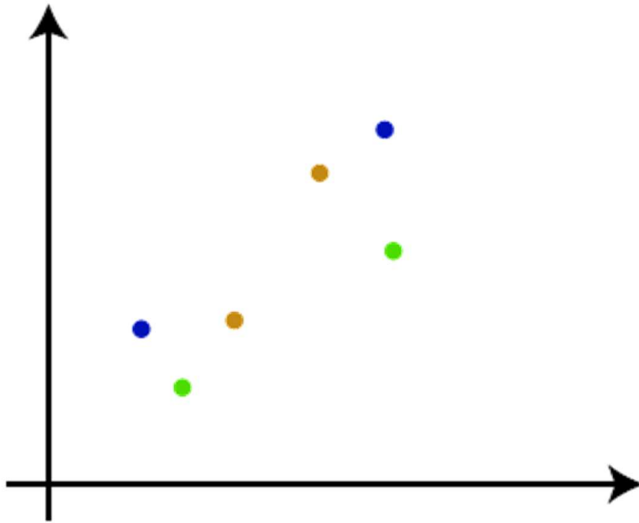
at the beginning, and then start combining the closest pair of clusters together. It does this until all the clusters are merged into a single cluster that contains all the datasets.

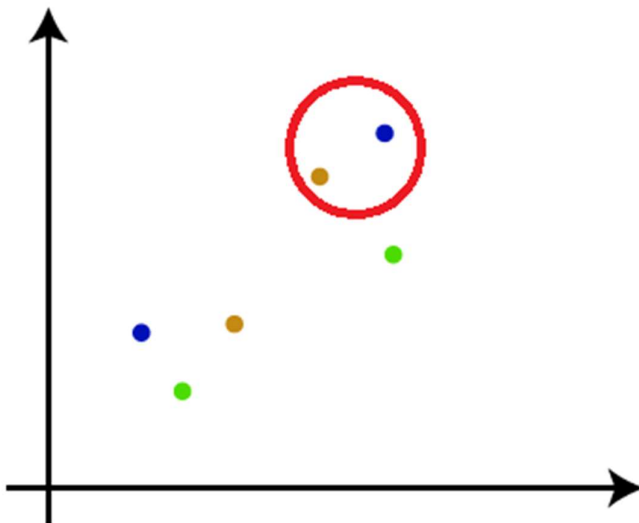This hierarchy of clusters is represented in the form of the dendrogram.

How the Agglomerative Hierarchical clustering Work?

The working of the AHC algorithm can be explained using the below steps:
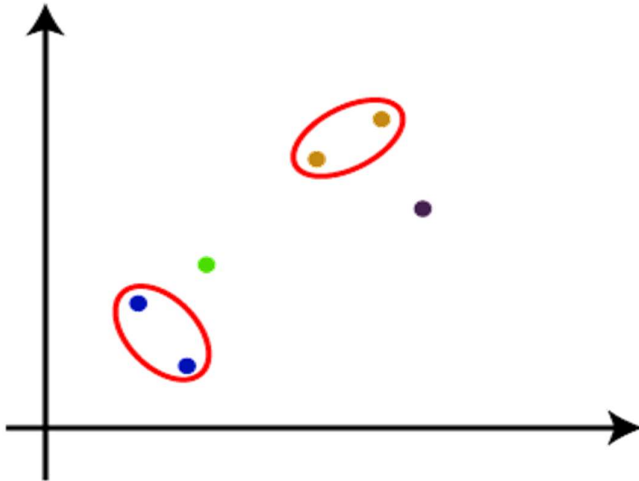
- o  **Step-1:** Create each data point as a single cluster. Let's say there are N data points, so the number of clusters will also be N.
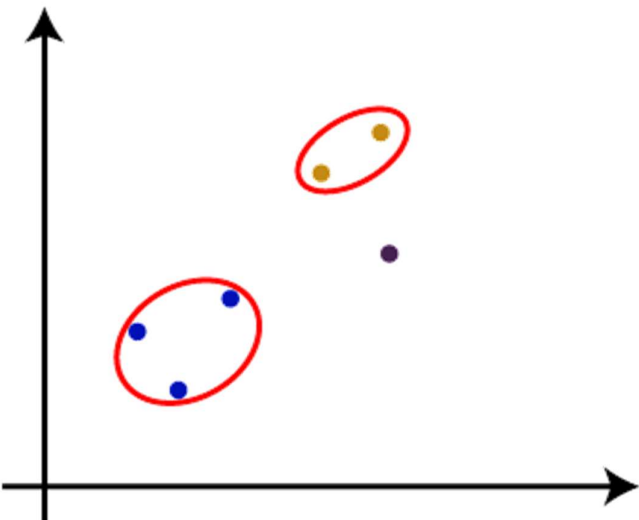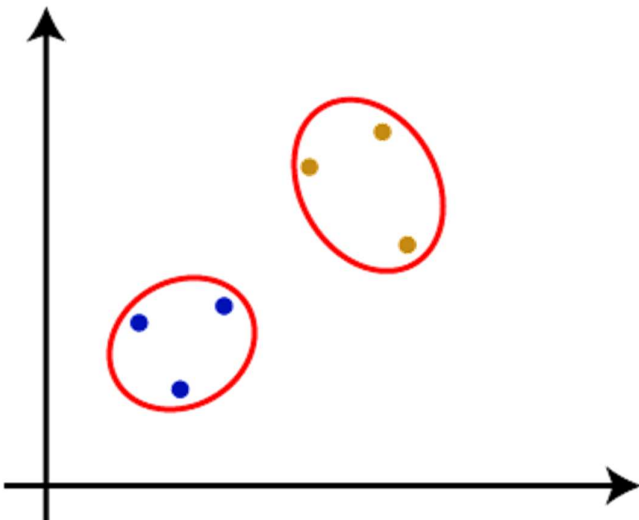


- o  **Step-2:** Take two closest data points or clusters and merge them to form one cluster. So, there will now be N-1 clusters.
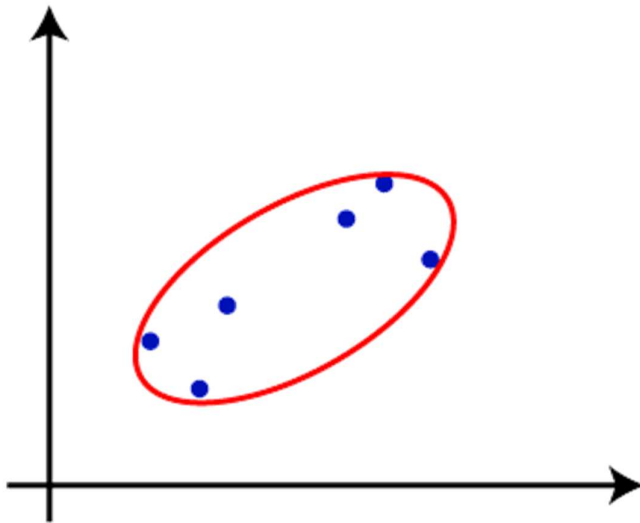
- **Step-3**: Again, take the two closest clusters and merge them together to form one cluster. There will be N-2 clusters.

- **Step-4:** Repeat Step 3 until only one cluster left. So, we will get the following clusters. Consider the below images:
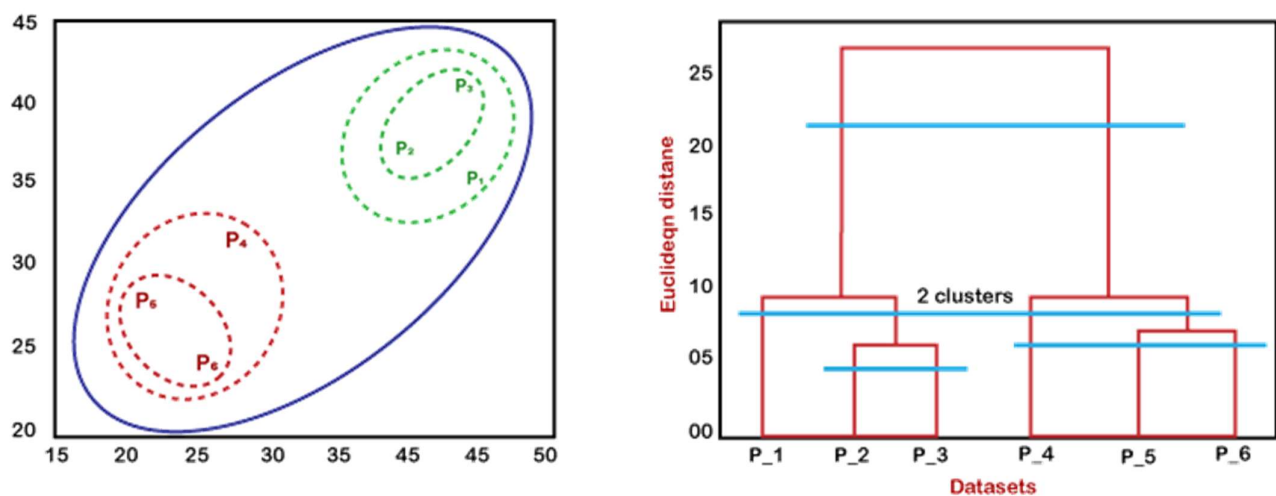
- o **Step-5:** Once all the clusters are combined into one big cluster, develop the dendrogram to divide the clusters as per the problem.

**Woking of Dendrogram in Hierarchical clustering**

The dendrogram is a tree-like structure that is mainly used to store each step as a memory that the HC algorithm performs. In the dendrogram plot, the Y-axis shows the Euclidean distances between the data points, and the x-axis shows all the data points of the given dataset.

The working of the dendrogram can be explained using the below diagram:



In the above diagram, the left part is showing how clusters are created in agglomerative clustering, and the right part is showing the corresponding dendrogram.

- o As we have discussed above, firstly, the datapoints P2 and P3 combine together and form a cluster, correspondingly a dendrogram is created, which connects P2 and P3 with a rectangular shape. The hight is decided according to the Euclidean distance between the data points.
- o In the next step, P5 and P6 form a cluster, and the corresponding dendrogram is created. It is higher than of previous, as the Euclidean distance between P5 and P6 is a little bit greater than the P2 and P3.
- o Again, two new dendrograms are created that combine P1, P2, and P3 in one dendrogram, and P4, P5, and P6, in another dendrogram.
- o At last, the final dendrogram is created that combines all the data points together.

We can cut the dendrogram tree structure at any level as per our requirement.


4. Explain Hierarchical clustering.
6. Explain working of dendrogram in Hierarchical clustering.
7. Explain Association Rule mining.
8. Explain apriori algorithm.
9. Explain Eclat algorithm
10. Explain F-P Growth algorithm