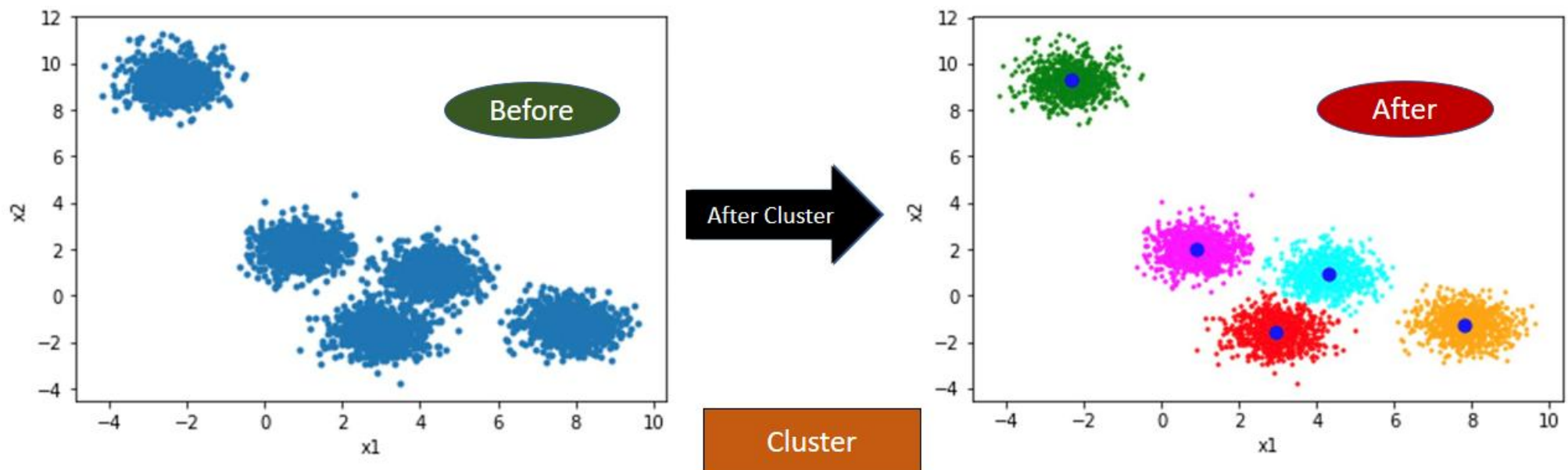


Clustering is the task of dividing the population or data points into a number of groups such that data points in the same groups are more similar to other data points in the same group and dissimilar to the data points in other groups.



Applications of Clustering: Real-World Scenarios

Clustering is a widely used technique in the industry. It is actually being used in almost every domain, ranging from banking to recommendation engines, document clustering to image segmentation.

- Customer Segmentation
- Document Clustering
- Image Segmentation
- Recommendation Engines

Perform K-Means Cluster



Tasks:

The k-means cluster algorithm mainly performs two important tasks:

- Determines the best value for K center points or centroids by an iterative process.
- Assigns each data point to its closest k-center (also called centroid). Those data points which are near to the particular k-center, create a cluster.

Perform K-Means Cluster



Steps:

How does the K-Means Algorithm Work?

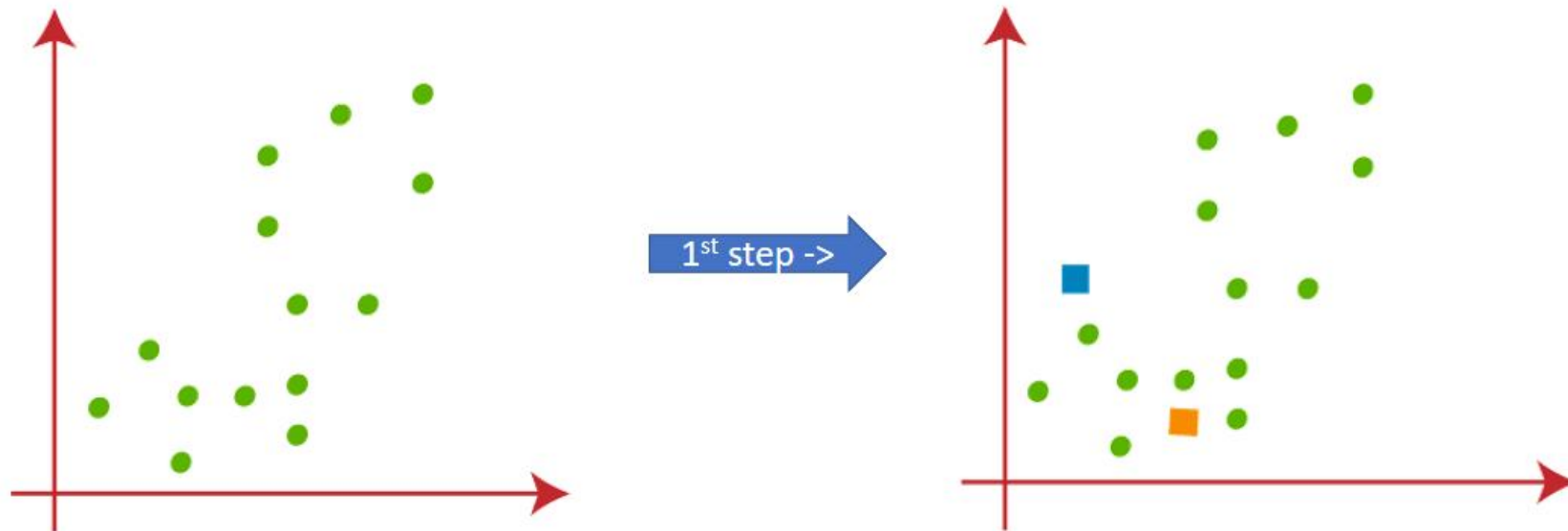
The working of the K-Means algorithm is explained in the below steps:

- Step-1: Select the number K to decide the number of clusters.
- Step-2: Select random K points or centroids.
- 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 data point 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.

Perform K-Means Cluster

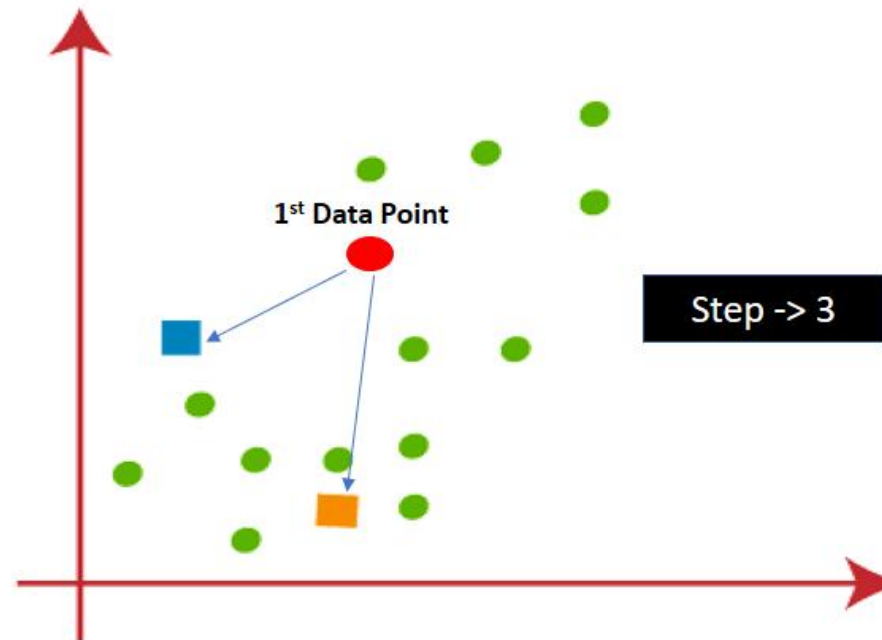


Assume, we have two variables M1 and M2. The X & Y axis scatter plot of these 2 variables is given below:

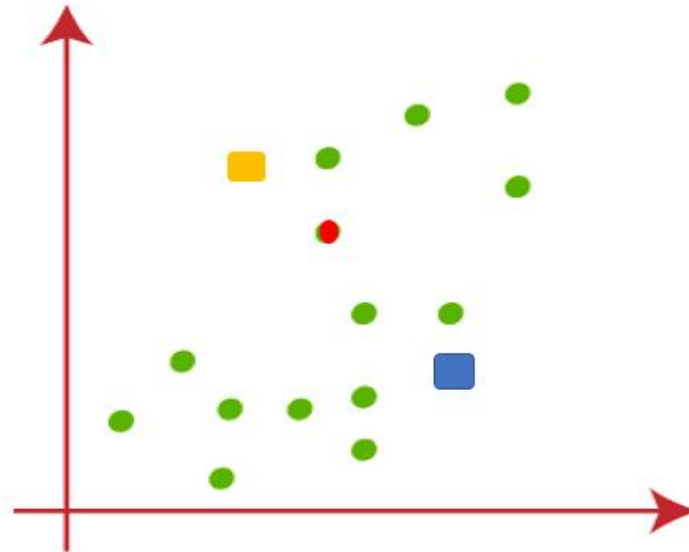


Note: We need to choose some random k points or centroid to form the cluster. These points can be either the points from the dataset or any other point.

Perform K-Means Cluster



Perform K-Means Cluster

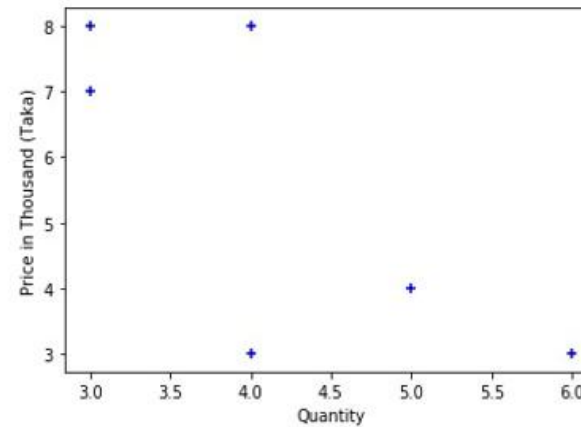


Let's see an EXAMPLE

Products	Quantity	Price(K)
FaceWash	3	7
Cream	5	4
Shoes	4	3
Bags	4	8
Jacket	6	3
Shirt	3	8

```
In [6]: plt.xlabel('Quantity')
plt.ylabel('Price in Thousand (Taka)')
plt.scatter(dataframe['Quantity'], dataframe['Price(K)'], marker='+', color='blue')
```

```
Out[6]: <matplotlib.collections.PathCollection at 0x2a9e524a048>
```



Products	Quantity	Price(K)
FaceWash	3	7
Cream	5	4
Shoes	4	3
Bags	4	8
Jacket	6	3
Shirt	3	8



$$c_1 = (3, 7) \text{ and } c_2 = (5, 4)$$

* For First data point (3, 7) • Facewash :

Distance from $c_1 = 0$ * (1)

$$\text{Distance from } c_2 = \sqrt{(5-3)^2 + (4-7)^2}$$

$$= \sqrt{9+9}$$

$$= 4.24$$

Products	Quantity	Price(K)
FaceWash	3	7
Cream	5	4
Shoes	4	3
Bags	4	8
Jacket	6	3
Shirt	3	8



$$c_1 = (3, 7) \text{ and } c_2 = (5, 4)$$

* For First data point (3, 7) • Facewash :

Distance from $c_1 = 0$ * (C₁)

$$\text{Distance from } c_2 = \sqrt{(5-3)^2 + (4-7)^2}$$

$$= \sqrt{9+9}$$

$$= 4.24$$

* For second data point (5, 4) • Cream :

$$\text{Distance from } c_1 = \sqrt{(5-3)^2 + (4-7)^2}$$

$$= \sqrt{4+9}$$

$$= \sqrt{13} = 3.60$$

Distance from $c_2 = 0$ * (C₂)

Class - 09

Products	Quantity	Price(K)
FaceWash	3	7
Cream	5	4
Shoes	4	3
Bags	4	8
Jacket	6	3
Shirt	3	8



$$c_1 = (3, 7) \text{ and } c_2 = (5, 4)$$

* For First data point (3, 7) • Facewash :

Distance from $c_1 = 0$ * (c_1)

$$\begin{aligned} \text{Distance from } c_2 &= \sqrt{(5-3)^2 + (4-7)^2} \\ &= \sqrt{4+9} \\ &= 4.24 \end{aligned}$$

* For second data point (5, 4) • Cream :

$$\begin{aligned} \text{Distance from } c_1 &= \sqrt{(5-3)^2 + (4-7)^2} \\ &= \sqrt{4+9} \\ &= \sqrt{13} = 3.60 \end{aligned}$$

Distance from $c_2 = 0$ * (c_2)

* For third data point (4, 3) Shoes :

$$\begin{aligned} \text{Distance from } c_1 &= \sqrt{(4-3)^2 + (3-7)^2} \\ &= \sqrt{1+16} \\ &= 4.123 \end{aligned}$$

$$\begin{aligned} \text{Distance from } c_2 &= \sqrt{(4-5)^2 + (3-4)^2} \\ &= \sqrt{1+1} \\ &= 1.41 * (c_2) \end{aligned}$$

$$\begin{aligned} \text{So new centroid} &= \left(\frac{5+4}{2}, \frac{4+3}{2} \right) \\ c_2 &= (4.5, 3.5) \end{aligned}$$

Products	Quantity	Price(K)
FaceWash	3	7
Cream	5	4
Shoes	4	3
Bags	4	8
Jacket	6	3
Shirt	3	8



$$C_1 = (3, 7) \text{ and } C_2 (4.5, 3.5)$$

For 4th data point (4, 8) bags :

$$\begin{aligned} \text{Distance from } C_1 &= \sqrt{(4-3)^2 + (8-7)^2} \\ &= \sqrt{1+1} \\ &= 1.41 \text{ (*) } (C_1) \end{aligned}$$

$$\begin{aligned} \text{Distance from } C_2 &= \sqrt{(4-4.5)^2 + (8-3.5)^2} \\ &= 0.25 + 20.25 \\ &= 20.50 \end{aligned}$$

$$\begin{aligned} \therefore \text{New centroid} &= \left(\frac{3+4}{2}, \frac{7+8}{2} \right) \\ C_1 &= (3.5, 7.5) \end{aligned}$$

Products	Quantity	Price(K)
FaceWash	3	7
Cream	5	4
Shoes	4	3
Bags	4	8
Jacket	6	3
Shirt	3	8



For 5th data point (6,3) Jacket :

$$\begin{aligned}
 \text{Distance from } c_1 &= \sqrt{(6-3.5)^2 + (3-7.5)^2} \\
 &= 6.25 + 20.25 \\
 &= 26.5
 \end{aligned}$$

$$\begin{aligned}
 \text{Distance from } c_2 &= \sqrt{(4.5-6)^2 + (3-3.5)^2} \\
 &= (2.25 + 0.25) \\
 &= 2.50 \quad \text{✱} \quad (c_2)
 \end{aligned}$$

$$\begin{aligned}
 \text{New Centroid} &= \left(\frac{5+4+6}{3}, \frac{4+3+3}{3} \right) \\
 c_2 &= (5, 3.33)
 \end{aligned}$$

Products	Quantity	Price(K)
FaceWash	3	7
Cream	5	4
Shoes	4	3
Bags	4	8
Jacket	6	3
Shirt	3	8



$$c_1 = (3, 7) \text{ and } c_2 = (4.5, 3.5)$$

For 4th data point (4, 8) bags :

$$\begin{aligned} \text{Distance from } c_1 &= \sqrt{(4-3)^2 + (8-7)^2} \\ &= \sqrt{1+1} \\ &= 1.41 \quad (*) \quad (c_1) \end{aligned}$$

$$\begin{aligned} \text{Distance from } c_2 &= \sqrt{(4-4.5)^2 + (8-3.5)^2} \\ &= 0.25 + 20.25 \\ &= 20.50 \end{aligned}$$

$$\begin{aligned} \therefore \text{New centroid} &= \left(\frac{3+4}{2}, \frac{7+8}{2} \right) \\ c_1 &= (3.5, 7.5) \end{aligned}$$

For 5th data point (6, 3) Jacket :

$$\begin{aligned} \text{Distance from } c_1 &= \sqrt{(6-3.5)^2 + (3-7.5)^2} \\ &= 6.25 + 20.25 \\ &= 26.5 \end{aligned}$$

$$\begin{aligned} \text{Distance from } c_2 &= \sqrt{(4.5-6)^2 + (3-3.5)^2} \\ &= (2.25 + 0.25) \\ &= 2.50 \quad (*) \quad (c_2) \end{aligned}$$

$$\begin{aligned} \text{New centroid} &= \left(\frac{5+4+6}{3}, \frac{4+3+3}{3} \right) \\ c_2 &= (5, 3.33) \end{aligned}$$

Products	Quantity	Price(K)
FaceWash	3	7
Cream	5	4
Shoes	4	3
Bags	4	8
Jacket	6	3
Shirt	3	8



$$c_1 = (3.5, 7.5) \text{ and } c_2 = (5, 3.33)$$

For 6th data point (3, 8) shirt :

$$\begin{aligned}
 \text{Distance from } c_1 &= \sqrt{(3-3.5)^2 + (8-7.5)^2} \\
 &= \sqrt{.25 + .25} \\
 &= 0.70 \quad * \quad (c_1)
 \end{aligned}$$

$$\begin{aligned}
 \text{Distance from } c_2 &= \sqrt{(3-5)^2 + (8-3.33)^2} \\
 &= \sqrt{4 + 2.16} \\
 &= 2.48
 \end{aligned}$$

$$\text{New centroid} = \left(\frac{3+4+3}{3}, \frac{7+8+8}{3} \right)$$

$$c_1 = (3.33, 7.67)$$

$$c_2 = (5, 3.33)$$



Let's Do it with PYTHON