

# Implementing K-Means Clustering

Devopriya Tirtho

16.02.04.033

Department of Computer Science and Engineering  
Ahsanullah University of Science and Technology  
Dhaka, Bangladesh

**Abstract**—‘Machine Learning’ is the kind of learning which tries to train a device according to learning to work smarter. In Machine Learning, there are two types of learning. One is Supervised Learning and the other one is Unsupervised Learning. In first intuition, we may think without any labelling how we can train the device to learn something. There are some advanced algorithms which are based on Unsupervised Learning. K-Means Clustering is one of them. In this experiment, we will implement the K-Means Clustering algorithm.

**Index Terms**—Machine Learning, K-Means Cluster, Unsupervised Learning, Supervised Learning.

## I. INTRODUCTION

**K-Means Clustering** is a form of **Unsupervised Learning**. We have already known about the **Supervised Learning**, now it is the time to work with **Unsupervised Learning**. In **Unsupervised Learning**, there is no labelled data. We have to predict the classes according to their happenings. In **K-Means Clustering**, we are given **K** number of clusters to divide our datapoints. The clusters are generated according to their **mean** value, that is why this model is called **K-Means Clustering**. The distance measurement we consider here is the **Euclidean Distance** from one point to another. The word **cluster** means there are some datapoints together and the main goal is to keep the distance low of in-cluster points and keep the distance high of the clusters.

## II. TASK

As this classifier is a form of **unsupervised learning**, there is no labelled data. So, we are given some datapoints at random. These data have to be clustered. Here, there are some examples of our given datapoints:

$W = \{(-7.87157, -4.86573), (5.86288, 0.99790), \dots\}$

For implementing the **K-Means Clustering**, there are some gradually incremented tasks which we have to solve on the way of implementing. These tasks are:

- Firstly, we have to take all the input datapoints from the file and plot the points.
- Secondly, the implementation part comes. We have to implement the **K-Means Clustering** model with the cluster value of **k** which is taken from user.
- Lastly, we have to color the corresponding datapoints according to clusters.

## III. EXPERIMENTAL DESIGN

- For the first task, we have to plot the datapoints which are taken from the input file and visualize the occurrences.

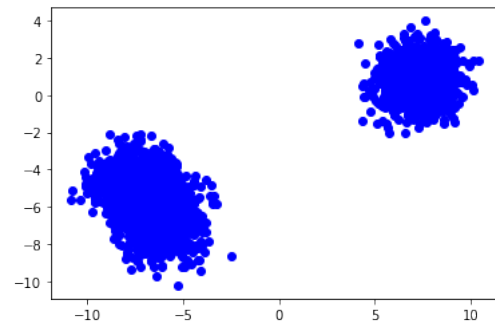


Fig. 1. Visualization of the Datapoints

- Now, it is the time for the implementation of **K-Means Clustering** model. Before the implementation, we need to understand how the algorithm works.

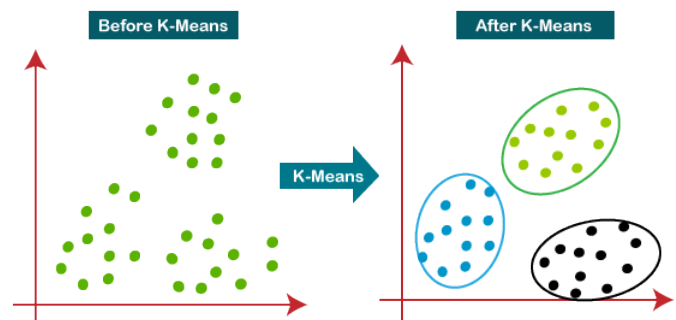


Fig. 2. K-Means Clustering

- Firstly, we need to take the input of cluster number **K**
- Then, we have to find the centroids by shuffling the datapoints and take random datapoints as centroids.
- After that, we have to iterate over all the datapoints and keep changing the centroid according to the newest mean value. For finding the distance, we will follow the **Euclidean Distance's** value.

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \quad (1)$$

- For stopping, we may define a fixed iteration number or we have to keep track of the changing datapoints and centroids. When there is no change in fixing the centroid, the algorithm stops. Our main task is to plot the datapoints close to the centroids.
- Finally, we need to color the corresponding datapoints of each cluster differently.

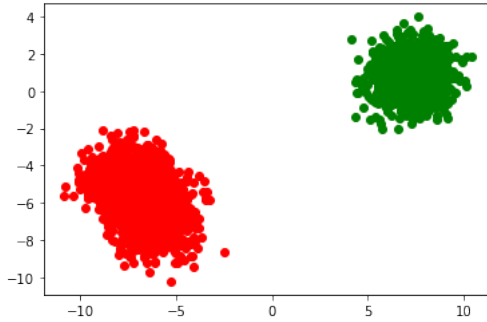


Fig. 3. Visualization of the Clusters

#### IV. RESULT ANALYSIS

After implementing the algorithm, we have seen different colored datapoint at different clusters. From the picture, we can see that, datapoints which are very close to each other are colored in same and datapoints which are far, they are of different color. The in-cluster distance of the datapoints is less while, the distance between different clusters are high.

#### V. PYTHON CODE

```
1 # -*- coding: utf-8 -*-
2 """160204033_A2_05.ipynb
3
4 Automatically generated by Colaboratory.
5
6 Original file is located at
7 https://colab.research.google.com/drive/1
8 OKyARfCG4HyAg-Ff-MnHkmlsT_gSiIYP
9 """
10 import io
11 import random
12 import pandas as pd
13 import numpy as np
14 import matplotlib.pyplot as plt
15 data = pd.read_csv('data_k_mean.txt', sep=" ", header
16 = None)
17 x=[i for i in data[0]]
18 y=[i for i in data[1]]
19 plt.plot(x ,y ,marker="o",linestyle = 'None',color="
20 blue")
21
22 k = int(input("Enter the Value of Clusters: "))
23 datalen=len(data[0])
24 print(datalen)
25 centroids=[]
26 for i in range(k):
27     index=random.randint(0,datalen-1)
28     centroids.append([data[0][index],data[1][index]])
```

```
27 print(centroids)
28
29 import math
30 lst=[]
31 converge=False
32 for m in range (0,200):
33
34     if(m>0):
35         centroidsTmp=[]
36         for j in range(k):
37             classwisedata1=[i[0] for i in lst if i[2]==j]
38             classwisedata2=[i[1] for i in lst if i[2]==j]
39
40             centroidsTmp.append([sum(classwisedata1)/len(
41 classwisedata1),sum(classwisedata2)/len(
42 classwisedata2)])
43
44         if(centroids==centroidsTmp):
45             converge=True
46         else:
47             centroids=centroidsTmp
48
49     if(converge):
50         break
51     lst=[]
52     print("iteration: ",m)
53     for i in range(0,datalen):
54         distance=[]
55         for j in range(k):
56             distance.append(math.sqrt(pow(centroids[j][0]-
57 data[0][i],2)+pow(centroids[j][1]-data[1][i],2))
58 )
59
60     lst.append([data[0][i],data[1][i],distance.index
61 (min(distance))])
62
63 print(lst)
64
65 import matplotlib.pyplot as plt
66 st=['red','green','blue','yellow']
67 for i in range(k):
68     x1=[j[0] for j in lst if j[2]==i ]
69     y1=[j[1] for j in lst if j[2]==i ]
70     colIndex=i%4
71     plt.plot(x1,y1 ,marker="o",linestyle = 'None',
72 color=st[colIndex])
```

#### VI. CONCLUSION

**K-Means Classifier** is a classifier which is based on **Unsupervised Learning**. There are several classifiers, but **K-Means Classifier** works better and the algorithm is simple as well. The knowledge of implementation will help us to construct other models in future.