

## K means Clustering – Introduction

K-Means Clustering is an Unsupervised Machine Learning algorithm, which groups the unlabeled dataset into different clusters.

## K means Clustering

[Unsupervised Machine Learning](#) is the process of teaching a computer to use unlabeled, unclassified data and enabling the algorithm to operate on that data without supervision. Without any previous data training, the machine's job in this case is to organize unsorted data according to parallels, patterns, and variations.

The goal of [clustering](#) is to divide the population or set of data points into a number of groups so that the data points within each group are more comparable to one another and different from the data points within the other groups. It is essentially a grouping of things based on how similar and different they are to one another.

We are given a data set of items, with certain features, and values for these features (like a vector). The task is to categorize those items into groups. To achieve this, we will use the K-means algorithm; an unsupervised learning algorithm. 'K' in the name of the algorithm represents the number of groups/clusters we want to classify our items into.

**Open In App**

(It will help if you think of items as points in an n-dimensional space). The algorithm will categorize the items into k groups or clusters of similarity. To calculate that similarity, we will use the Euclidean distance as a measurement.

The algorithm works as follows:

1. First, we randomly initialize k points, called means or cluster centroids.
2. We categorize each item to its closest mean and we update the mean's coordinates, which are the averages of the items categorized in that cluster so far.
3. We repeat the process for a given number of iterations and at the end, we have our clusters.

The “points” mentioned above are called means because they are the mean values of the items categorized in them. To initialize these means, we have a lot of options. An intuitive method is to initialize the means at random items in the data set. Another method is to initialize the means at random values between the boundaries of the data set (if for a feature  $x$ , the items have values in  $[0,3]$ , we will initialize the means with values for  $x$  at  $[0,3]$ ).

The above algorithm in pseudocode is as follows:

```
Initialize k means with random values
```

```
--> For a given number of iterations:
```

```
    --> Iterate through items:
```

```
        --> Find the mean closest to the item by calculating  
            the euclidean distance of the item with each of the means
```

```
        --> Assign item to mean
```

```
        --> Update mean by shifting it to the average of the items in that cluster
```

## Import the necessary Libraries:

We are importing Numpy for statistical computations, Matplotlib to plot the graph, and `make_blobs` from `sklearn.datasets`.

Python3

```
import numpy as np  
import matplotlib.pyplot as plt  
from sklearn.datasets import make_blobs
```

## Create the custom dataset with make\_blobs and plot it

Python3

```
X,y = make_blobs(n_samples = 500,n_features = 2,centers = 3,  
fig = plt.figure(0)  
plt.grid(True)  
plt.scatter(X[:,0],X[:,1])  
plt.show()
```

**Output:**

Clustering dataset

---

## Initialize the random centroids

## Python3

```
k = 3

clusters = {}
np.random.seed(23)

for idx in range(k):
    center = 2*(2*np.random.random((X.shape[1],))-1)
    points = []
    cluster = {
        'center' : center,
        'points' : []
    }

    clusters[idx] = cluster

clusters
```

### Output:

```
{0: {'center': array([0.06919154, 1.78785042]), 'points': []},
 1: {'center': array([ 1.06183904, -0.87041662]), 'points': []},
 2: {'center': array([-1.11581855,  0.74488834]), 'points': []}}
```

### Plot the random initialize center with data points

## Python3

```
plt.scatter(X[:,0],X[:,1])
plt.grid(True)
for i in clusters:
    center = clusters[i]['center']
    plt.scatter(center[0],center[1],marker = '*',c = 'red')
plt.show()
```

**Output:**

Data points with random center

---

**Define euclidean distance**

## Python3

```
def distance(p1,p2):  
    return np.sqrt(np.sum((p1-p2)**2))
```

**Create the function to Assign and Update the cluster center**

## Python3

#Implementing E step

```
def assign_clusters(X, clusters):  
    for idx in range(X.shape[0]):  
        dist = []  
  
        curr_x = X[idx]  
  
        for i in range(k):  
            dis = distance(curr_x, clusters[i]['center'])  
            dist.append(dis)  
            curr_cluster = np.argmin(dist)  
            clusters[curr_cluster]['points'].append(curr_x)  
    return clusters
```

#Implementing the M-Step

```
def update_clusters(X, clusters):  
    for i in range(k):  
        points = np.array(clusters[i]['points'])  
        if points.shape[0] > 0:  
            new_center = points.mean(axis = 0)  
            clusters[i]['center'] = new_center  
  
            clusters[i]['points'] = []  
    return clusters
```

**Create the function to Predict the cluster for the datapoints**

### Python3

```
def pred_cluster(X, clusters):  
    pred = []  
    for i in range(X.shape[0]):  
        dist = []  
        for j in range(k):  
            dist.append(distance(X[i], clusters[j]['center']))  
        pred.append(np.argmin(dist))  
    return pred
```

### Assign, Update, and predict the cluster center

### Python3

```
clusters = assign_clusters(X, clusters)  
clusters = update_clusters(X, clusters)  
pred = pred_cluster(X, clusters)
```

### Plot the data points with their predicted cluster center

### Python3

```
plt.scatter(X[:,0], X[:,1], c = pred)  
for i in clusters:  
    center = clusters[i]['center']  
    plt.scatter(center[0], center[1], marker = '^', c = 'red')  
plt.show()
```

### Output:



## K-means Clustering

---

### Example 2:

Import the necessary libraries

Python3

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib.cm as cm
from sklearn.datasets import load_iris
from sklearn.cluster import KMeans
```

**Load the Dataset**

## Python3

```
X, y = load_iris(return_X_y=True)
```

### Elbow Method

Finding the ideal number of groups to divide the data into is a basic stage in any unsupervised algorithm. One of the most common techniques for figuring out this ideal value of k is the elbow approach.

## Python3

```
#Find optimum number of cluster
sse = [] #SUM OF SQUARED ERROR
for k in range(1,11):
    km = KMeans(n_clusters=k, random_state=2)
    km.fit(X)
    sse.append(km.inertia_)
```

### Plot the Elbow graph to find the optimum number of cluster

## Python3

```
sns.set_style("whitegrid")
g=sns.lineplot(x=range(1,11), y=sse)
g.set(xlabel = "Number of cluster (k)",
      ylabel = "Sum Squared Error",
      title = 'Elbow Method')
plt.show()
```

Output:

## Elbow Method

---

From the above graph, we can observe that at  $k=2$  and  $k=3$  elbow-like situation. So, we are considering  $K=3$

### Build the Kmeans clustering model

Python3

```
kmeans = KMeans(n_clusters = 3, random_state = 2)
kmeans.fit(X)
```

Output:

KMeans

KMeans(n\_clusters=3, random\_state=2)

### Find the cluster center

### Python3

```
kmeans.cluster_centers_
```

Output:

```
array([[5.006      , 3.428      , 1.462      , 0.246      ],
       [5.9016129 , 2.7483871 , 4.39354839, 1.43387097],
       [6.85      , 3.07368421, 5.74210526, 2.07105263]])
```

Predict the cluster group:

### Python3

```
pred = kmeans.fit_predict(X)
pred
```

Output:

```
array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 1, 1, 2, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 1, 2, 2, 2, 2, 1, 2, 2, 2,
       2, 2, 2, 1, 1, 2, 2, 2, 2, 1, 2, 1, 2, 1, 2, 2, 1, 1, 2, 2, 2, 2,
       2, 1, 2, 2, 2, 2, 1, 2, 2, 2, 1, 2, 2, 2, 1], dtype=int32)
```

**Plot the cluster center with data points**

## Python3

```
plt.figure(figsize=(12,5))
plt.subplot(1,2,1)
plt.scatter(X[:,0],X[:,1],c = pred, cmap=cm.Accent)
plt.grid(True)
for center in kmeans.cluster_centers_:
    center = center[:2]
    plt.scatter(center[0],center[1],marker = '^',c = 'red')
plt.xlabel("petal length (cm)")
plt.ylabel("petal width (cm)")

plt.subplot(1,2,2)
plt.scatter(X[:,2],X[:,3],c = pred, cmap=cm.Accent)
plt.grid(True)
for center in kmeans.cluster_centers_:
    center = center[2:4]
    plt.scatter(center[0],center[1],marker = '^',c = 'red')
plt.xlabel("sepal length (cm)")
plt.ylabel("sepal width (cm)")
plt.show()
```

Output:

**Article Tags :** Machine Learning Python

### Recommended Articles

1. [Analysis of test data using K-Means Clustering in Python](#)
2. [ML I Determine the optimal value of K in K-Means Clustering](#)
3. [ML I Mini Batch K-means clustering algorithm](#)
4. [Image compression using K-means clustering](#)
5. [K-Means Clustering in R Programming](#)
6. [Difference between K means and Hierarchical Clustering](#)
7. [Image Segmentation using K Means Clustering](#)
8. [K- means clustering with SciPy](#)
9. [K means clustering using Weka](#)
10. [Clustering Text Documents using K-Means in Scikit Learn](#)
11. [K-Means clustering on the handwritten digits data using Scikit Learn in Python](#)
12. [Analyzing Decision Tree and K-means Clustering using Iris dataset](#)
13. [K-Means Clustering using PySpark Python](#)
14. [DBSCAN Clustering in ML I Density based clustering](#)
15. [Difference between CURE Clustering and DBSCAN Clustering](#)
16. [Difference Between Agglomerative clustering and Divisive clustering](#)
17. [ML I Random Initialization Trap in K-Means](#)
18. [Understanding Types of Means I Set 1](#)
19. [ML I K-means++ Algorithm](#)
20. [Color Quantization using K-Means in Scikit Learn](#)
21. [Clustering in R Programming](#)
22. [Clustering in Machine Learning](#)
23. [Different Types of Clustering Algorithm](#)
24. [ML I Unsupervised Face Clustering Pipeline](#)
25. [Hierarchical Clustering in Machine Learning](#)

[Read Full Article](#)

---



A-143, 9th Floor, Sovereign Corporate  
Tower, Sector-136, Noida, Uttar Pradesh -  
201305



[feedback@geeksforgeeks.org](mailto:feedback@geeksforgeeks.org)



## Company

[About Us](#)

[Legal](#)

[Careers](#)

[In Media](#)

[Contact Us](#)

[Advertise with us](#)

[GFG Corporate Solution](#)

[Placement Training Program](#)

[Apply for Mentor](#)

## Explore

[Job-A-Thon Hiring Challenge](#)

[Hack-A-Thon](#)

[GfG Weekly Contest](#)

---

Offline Classes (Delhi/NCR)

DSA in JAVA/C++

Master System Design

Master CP

GeeksforGeeks Videos

## Languages

Python

Java

C++

PHP

GoLang

SQL

R Language

Android Tutorial

## DSA Concepts

Data Structures

Arrays

Strings

Linked List

Algorithms

Searching

Sorting

Mathematical

Dynamic Programming

## DSA Roadmaps

DSA for Beginners



Basic DSA Coding Problems  
DSA Roadmap by Sandeep Jain  
DSA with JavaScript  
Top 100 DSA Interview Problems  
All Cheat Sheets

## **Web Development**

HTML  
CSS  
JavaScript  
Bootstrap  
ReactJS  
AngularJS  
NodeJS  
Express.js  
Lodash

## **Computer Science**

GATE CS Notes  
Operating Systems  
Computer Network  
Database Management System  
Software Engineering  
Digital Logic Design  
Engineering Maths

## **Python**

Python Programming Examples  
Django Tutorial

Python Projects

Python Tkinter

OpenCV Python Tutorial

Python Interview Question

## **Data Science & ML**

Data Science With Python

Data Science For Beginner

Machine Learning Tutorial

Maths For Machine Learning

Pandas Tutorial

NumPy Tutorial

NLP Tutorial

Deep Learning Tutorial

## **DevOps**

Git

AWS

Docker

Kubernetes

Azure

GCP

## **Competitive Programming**

Top DSA for CP

Top 50 Tree Problems

Top 50 Graph Problems

Top 50 Array Problems

Top 50 String Problems

Top 50 DP Problems

Top 15 Websites for CP

## **System Design**

What is System Design

Monolithic and Distributed SD

Scalability in SD

Databases in SD

High Level Design or HLD

Low Level Design or LLD

Crack System Design Round

System Design Interview Questions

## **Interview Corner**

Company Wise Preparation

Preparation for SDE

Experienced Interviews

Internship Interviews

Competitive Programming

Aptitude Preparation

## **GfG School**

CBSE Notes for Class 8

CBSE Notes for Class 9

CBSE Notes for Class 10

CBSE Notes for Class 11

CBSE Notes for Class 12

English Grammar

## **Commerce**

Accountancy  
Business Studies  
Economics  
Human Resource Management (HRM)  
Management  
Income Tax  
Finance  
Statistics for Economics

## **UPSC**

Polity Notes  
Geography Notes  
History Notes  
Science and Technology Notes  
Economics Notes  
Important Topics in Ethics  
UPSC Previous Year Papers

## **SSC/ BANKING**

SSC CGL Syllabus  
SBI PO Syllabus  
SBI Clerk Syllabus  
IBPS PO Syllabus  
IBPS Clerk Syllabus  
Aptitude Questions  
SSC CGL Practice Papers

## **Write & Earn**

Write an Article