

Write a program of cluster analysis using simple k-means algorithm python programming language.

Cluster Analysis:

Cluster Analysis is a statistical method for processing data. It works by organizing items into groups, or clusters on the basis of how closely associated they are.

K-means algorithm:

K-means algorithm is a simple two steps clustering process. The first step is cluster assignment and the second one is the move centroid step. However, this unsupervised algorithm can easily create, implement and handle massive datasets.

Steps involved in k-means Algorithm:

- 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 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.

K-means Algorithm using python programming:

```
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

dataset = pd.read_csv('content/sample_data/Mall-
customers.csv')
X = dataset.iloc[:, [3,4]].values

from sklearn.cluster import KMeans
wcss = list()

for i in range(1,11):
    kmeans = KMeans(n_clusters=i, init =
        'k-means++', random_state=42)
    kmeans.fit(X)
    wcss_list.append(kmeans.inertia_)

plt.plot(range(1,11), wcss_list)
plt.title('The Elbow Method Graph')
plt.xlabel('Number of clusters(k)')
plt.ylabel('wcss_list')
plt.show()

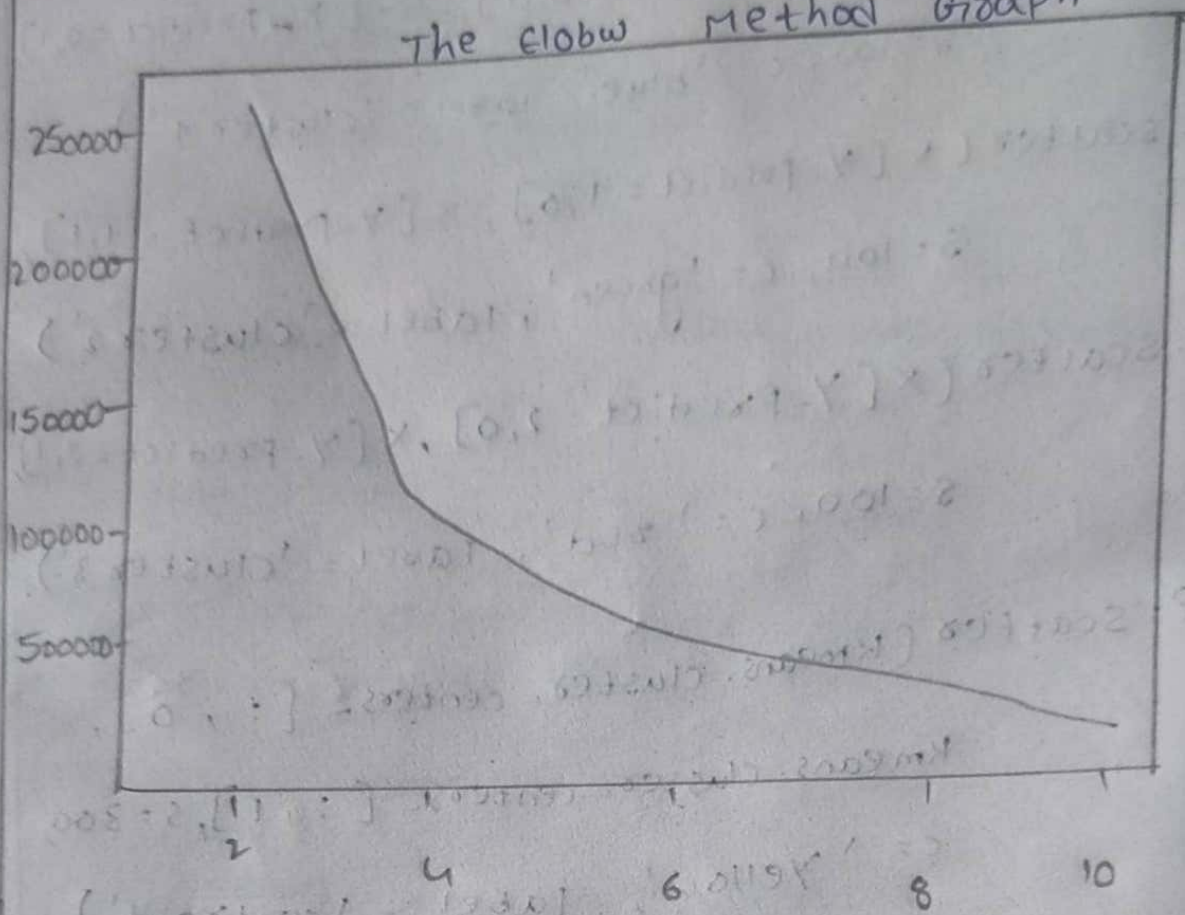
kmeans = KMeans(n_clusters=5, init='k-means++',
    random_state=42)
```

```
Y_predict = kmeans.fit_predict(X)
plt.scatter(X[Y_predict == 0,0], X[Y_predict == 0,1],
    s=100, c='blue', label='cluster 1')
plt.scatter(X[Y_predict == 1,0], X[Y_predict == 1,1],
    s=100, c='green', label='cluster 2')
plt.scatter(X[Y_predict == 2,0], X[Y_predict == 2,1],
    s=100, c='red', label='cluster 3')
plt.scatter(kmeans.cluster_centers_[:,0],
    kmeans.cluster_centers_[:,1], s=300,
    c='yellow', label='centroid')

plt.title('clusters of customers')
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending score (1-100)')
plt.legend()
plt.show()
```


Output

The Elbow Method Graph



Number of clusters (K)

clusters of customers

