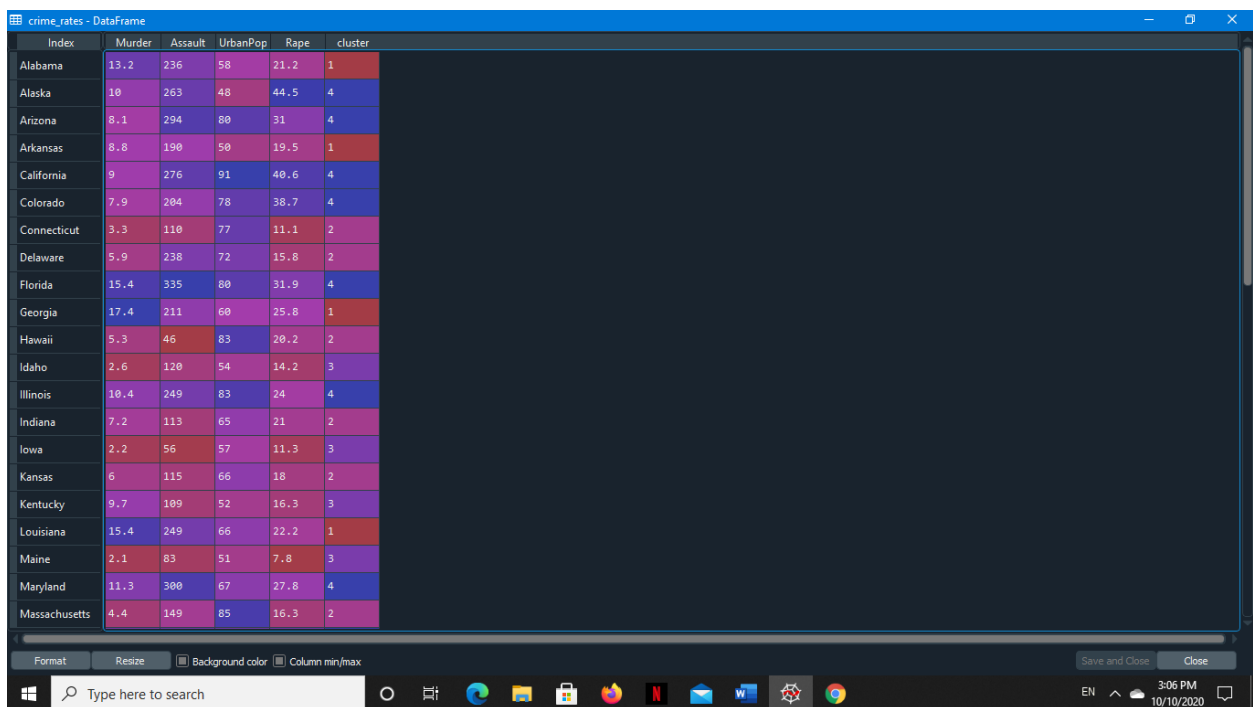2.) Perform K-Means Clustering for the crime data and identify the number of clusters formed and draw inferences.

In this Assignment first we have to load the data set using some library and the data set is in the form of excel and load that dataset using Common Separated values (CSV).

import numpy as np # linear algebra

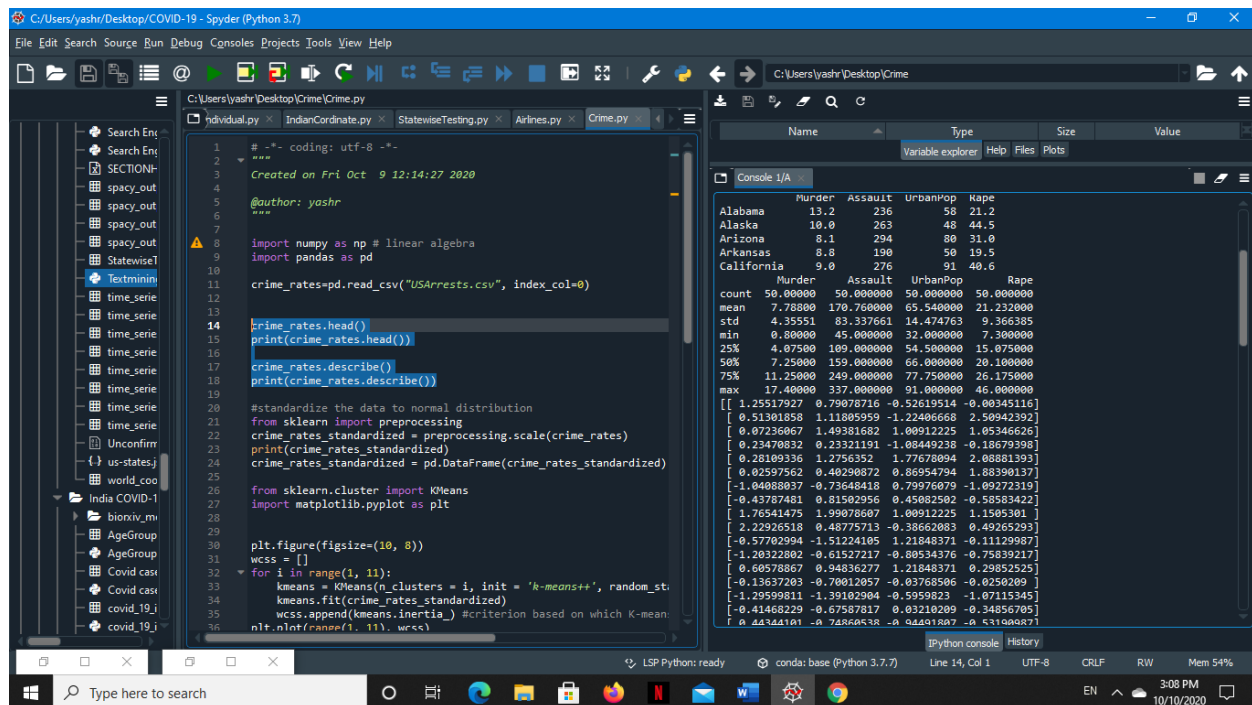import pandas as pd

crime_rates=pd.read_csv("USArrests.csv", index_col=0)



Now we have to perform some data preprocessing technique to access and clean the data set.

crime_rates.head()

print(crime_rates.head())

crime_rates.describe()

print(crime_rates.describe())



#standardize the data to normal distribution

from sklearn import preprocessing

crime_rates_standardized = preprocessing.scale(crime_rates)

print(crime_rates_standardized)

crime_rates_standardized = pd.DataFrame(crime_rates_standardized)

##### Now we are going to perform Clustering

```python
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt


plt.figure(figsize=(10, 8))
wcss = []
for i in range(1, 11):
    kmeans = KMeans(n_clusters = i, init = 'k-means++', random_state = 42)
    kmeans.fit(crime_rates_standardized)
    wcss.append(kmeans.inertia_) #criterion based on which K-means clustering works
plt.plot(range(1, 11), wcss)
plt.title('The Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.show()
```

| Index | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| 0 | 1.25518 | 0.790787 | -0.526195 | -0.00345116 |
| 1 | 0.513019 | 1.11806 | -1.22407 | 2.50942 |
| 2 | 0.0723607 | 1.49382 | 1.00912 | 1.05347 |
| 3 | 0.234708 | 0.233212 | -1.08449 | -0.186794 |
| 4 | 0.281093 | 1.27564 | 1.77678 | 2.08881 |
| 5 | 0.0259756 | 0.402909 | 0.869548 | 1.8839 |
| 6 | -1.04088 | -0.736484 | 0.799761 | -1.09272 |
| 7 | -0.437875 | 0.81503 | 0.450825 | -0.585834 |
| 8 | 1.76541 | 1.99079 | 1.00912 | 1.15053 |
| 9 | 2.22927 | 0.487757 | -0.386621 | 0.492653 |
| 10 | -0.57703 | -1.51224 | 1.21848 | -0.1113 |
| 11 | -1.20323 | -0.615272 | -0.805344 | -0.758392 |
| 12 | 0.605789 | 0.948363 | 1.21848 | 0.298525 |
| 13 | -0.136372 | -0.700121 | -0.0376851 | -0.0250209 |
| 14 | -1.296 | -1.39103 | -0.595982 | -1.07115 |
| 15 | -0.414682 | -0.675878 | 0.0321021 | -0.348567 |
| 16 | 0.443441 | -0.748605 | -0.944918 | -0.53191 |
| 17 | 1.76541 | 0.948363 | 0.0321021 | 0.104398 |
| 18 | -1.31919 | -1.06376 | -1.01471 | -1.44862 |
| 19 | 0.814521 | 1.56654 | 0.101889 | 0.70835 |
| 20 | -0.785763 | -0.263757 | 1.35806 | -0.53191 |

Format  Resize   ☐ Background color  ☐ Column min/max    Save and Close   Close

---

| Index | Type | Size | Value |
|---|---|---|---|
| 0 | float64 | 1 | 200.0 |
| 1 | float64 | 1 | 104.96163315756871 |
| 2 | float64 | 1 | 80.08569526137276 |
| 3 | float64 | 1 | 57.55425863091105 |
| 4 | float64 | 1 | 50.05119672966492 |
| 5 | float64 | 1 | 44.20084360686918 |
| 6 | float64 | 1 | 40.70266300421998 |
| 7 | float64 | 1 | 34.85432218312674 |
| 8 | float64 | 1 | 30.616017585613793 |
| 9 | float64 | 1 | 26.856423116032715 |

Save and Close   Close

The Elbow Method

# Fitting K-Means to the dataset

kmeans = KMeans(n_clusters = 4, init = 'k-means++', random_state = 42)

y_kmeans = kmeans.fit_predict(crime_rates_standardized)


y_kmeans

print(y_kmeans)
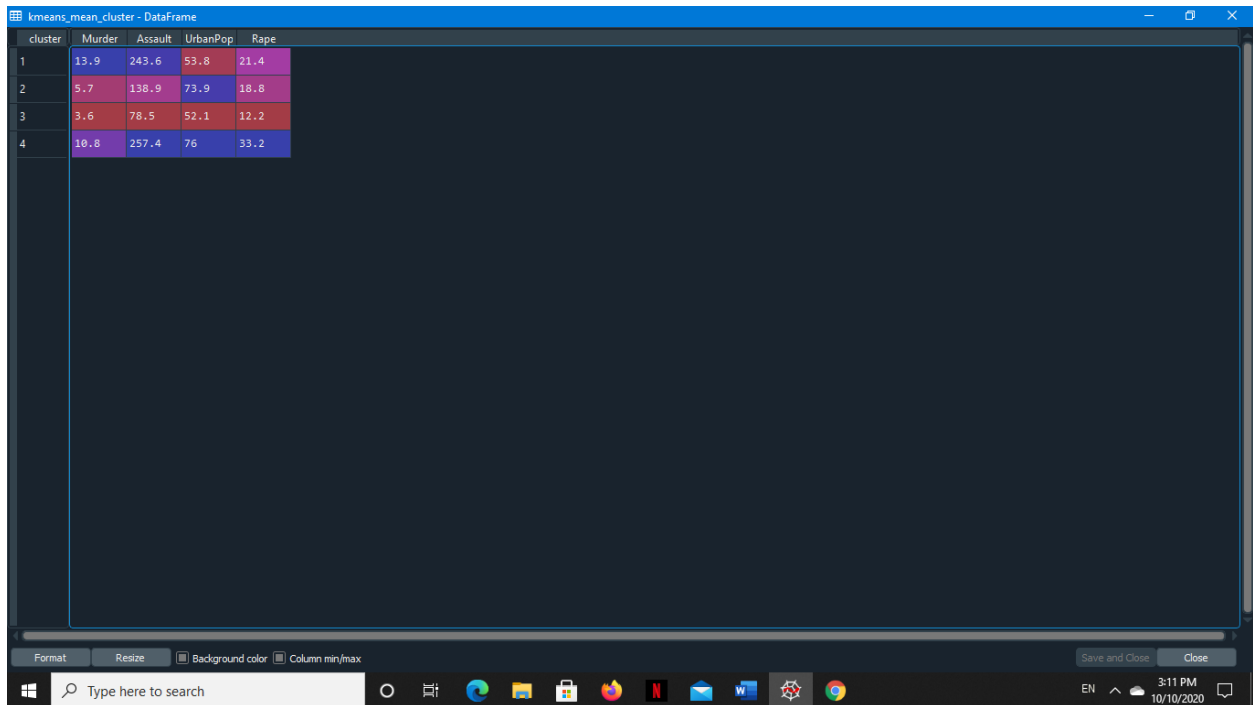

#beginning of  the cluster numbering with 1 instead of 0

y_kmeans1=y_kmeans+1

# New list called cluster
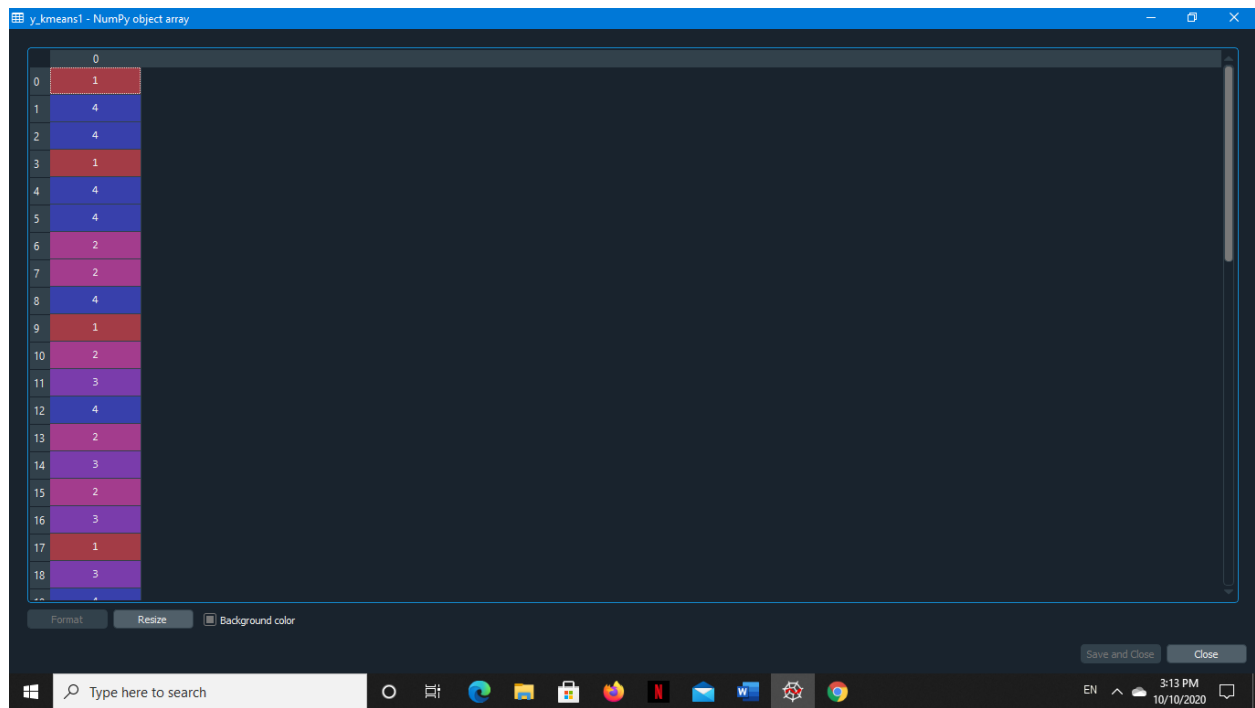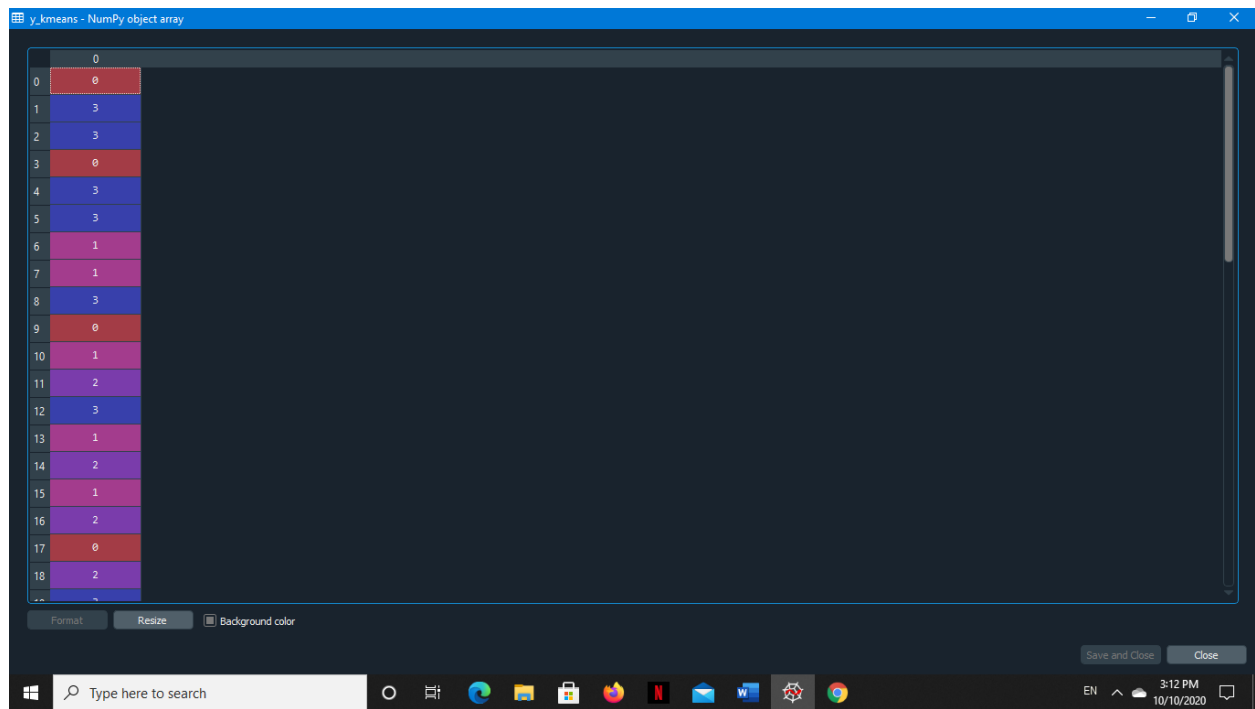
cluster = list(y_kmeans1)

# Adding cluster to our data set

crime_rates['cluster'] = cluster

#Mean of clusters 1 to 4

kmeans_mean_cluster = pd.DataFrame(round(crime_rates.groupby('cluster').mean(),1))

kmeans_mean_cluster

print(kmeans_mean_cluster)

| kmeans_mean_cluster - DataFrame | | | | | — □ × |
|---|---|---|---|---|---|
| cluster | Murder | Assault | UrbanPop | Rape | |
| 1 | 13.9 | 243.6 | 53.8 | 21.4 | |
| 2 | 5.7 | 138.9 | 73.9 | 18.8 | |
| 3 | 3.6 | 78.5 | 52.1 | 12.2 | |
| 4 | 10.8 | 257.4 | 76 | 33.2 | |

Format    Resize    ☐ Background color  ☐ Column min/max

Save and Close    Close

import seaborn as sns

```python
plt.figure(figsize=(12,6))

sns.scatterplot(x=crime_rates['Murder'], y = crime_rates['Assault'],hue=y_kmeans1)


crime_rates[crime_rates['cluster']==1]

print(crime_rates[crime_rates['cluster']==1])
```