K-Means clustering (Unsupervised) import required package import numpy as np

```
import pandas as pd
```

import seaborn as sns

Column

Age

CustomerID

Gender

___ 0

1

100

import matplotlib.pyplot as plt

read data from source and describing

df = pd.read csv('mall customers.csv') df.head()

CustomerID Gender Age Annual Income (k\$) Spending Score (1-100) 0 1 Male 19 15

2 1 Male 21 15 2 Female 20 16

3 Female 23 16

4 5 Female 17 31

40 print(df.columns)

39

81

6

77

200.000000

50.200000

25.823522

1.000000

34.750000

50.000000

73.000000

99.000000

 $\label{localization} Index(['CustomerID', 'Gender', 'Age', 'Annual Income (k\$)',$ 'Spending Score (1-100)'],

dtype='object')

print(df.describe())

In [4]: CustomerID Age Annual Income (k\$) Spending Score (1-100) count 200.000000 200.000000 200.000000 60.560000 mean 100.500000 38.850000 26.264721 std 57.879185 13.969007 15.000000 min 1.000000 18.000000

50.750000 28.750000 25% 41.500000 100.500000 50% 36.000000 61.500000 49.000000 78.000000 150.250000 75%

max 137.000000 200.000000 70.000000

df.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 200 entries, 0 to 199 Data columns (total 5 columns):

Annual Income (k\$) 200 non-null int64 Spending Score (1-100) 200 non-null int64 dtypes: int64(4), object(1) memory usage: 7.9+ KB check the relation between variables In [6]: ## state value is excluded as it is categorical corr = df.corr() corr.style.background gradient(cmap = 'Greens')

200 non-null

Non-Null Count Dtype

200 non-null int64

200 non-null object

int64

CustomerID Age Annual Income (k\$) Spending Score (1-100) 0.013835 CustomerID -0.026763 1.000000 0.977548 1.000000 -0.327227 Age -0.026763 -0.012398 Annual Income (k\$) 0.977548 -0.012398 1.000000 0.009903 1.000000 Spending Score (1-100) 0.013835 -0.327227 0.009903

sns.pairplot(df) <seaborn.axisgrid.PairGrid at 0x2ac8c453760> 200 150 CustomerID

50 70 60 Age 140 120 Annual Income (k\$) 100 80 60 40 20 100 Spending Score (1-100) 80 60 100 100 CustomerID Spending Score (1-100) Annual Income (k\$) Age select input and op variable ### since there is no op/predicted variable y is not required

for ele in range(1, 11): model = KMeans(n clusters= ele) model.fit(x) wss.append(model.inertia)

plt.plot(values, wss, label = 'Elbow line')

from sklearn.cluster import KMeans

wss = []

print(wss)

values = np.arange(1, 11)

Out[10]: <matplotlib.legend.Legend at 0x2ac908fe6d0>

2

model = KMeans(n clusters=5)

plt.xlabel('Annual income') plt.ylabel('Spending Score')

Out[14]: Text(0, 0.5, 'Spending Score')

from sklearn.cluster import KMeans

plt.scatter(values, wss)

plt.ylabel('WSS') plt.xlabel('K Values')

plt.legend()

[269981.28000000014, 181363.59595959607, 106348.37306211119, 73679.78903948837, 44448. 45544793369, 37265.86520484345, 30241.34361793659, 24999.3682586117, 22143.22210076743 8, 19676.612585602812]

x = df.drop(['CustomerID', 'Age', 'Gender'], axis=1)

Selecting optimal value of k using elbow method

Elbow line 250000 Elbow point 200000 150000 100000 50000

K Values

plt.scatter([values[4]], [wss[4]], c= 'red', label = 'Elbow point')

KMeans(n_clusters=5)

fit the data model.fit(x)

creating a model

print(model.labels) $[1\ 3\ 1\$

8

10

2 0 2 0 2 0 2 0 2 0 2 0 2 0 2] In [14]: plt.scatter(df['Annual Income (k\$)'][model.labels_ == 0], df['Spending Score (1-100)' plt.scatter(df['Annual Income (k\$)'][model.labels == 1], df['Spending Score (1-100)' plt.scatter(df['Annual Income (k\$)'][model.labels == 2], df['Spending Score (1-100)' plt.scatter(df['Annual Income (k\$)'][model.labels == 3], df['Spending Score (1-100)' plt.scatter(df['Annual Income (k\$)'][model.labels == 4], df['Spending Score (1-100)' plt.scatter(model.cluster centers [:,0], model.cluster centers [:,1], c = 'black', s

100 80 Score 60 Spending 100 120 140 Annual income