

K-Means Clustering

Importing the libraries

In []:

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

Importing the dataset

In []:

```
dataset = pd.read_csv('Mall_Customers.csv')
dataset.head()
```

Out[2]:

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40

In []:

```
dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
#   Column                                Non-Null Count  Dtype  
---  -
0   CustomerID                            200 non-null   int64  
1   Genre                                 200 non-null   object  
2   Age                                   200 non-null   int64  
3   Annual Income (k$)                    200 non-null   int64  
4   Spending Score (1-100)                 200 non-null   int64  
dtypes: int64(4), object(1)
memory usage: 7.9+ KB
```

In []:

```
dataset.shape
```

Out[3]:

```
(200, 5)
```

In []:

```
X = dataset.iloc[:,[3,4]].values
```

In []:

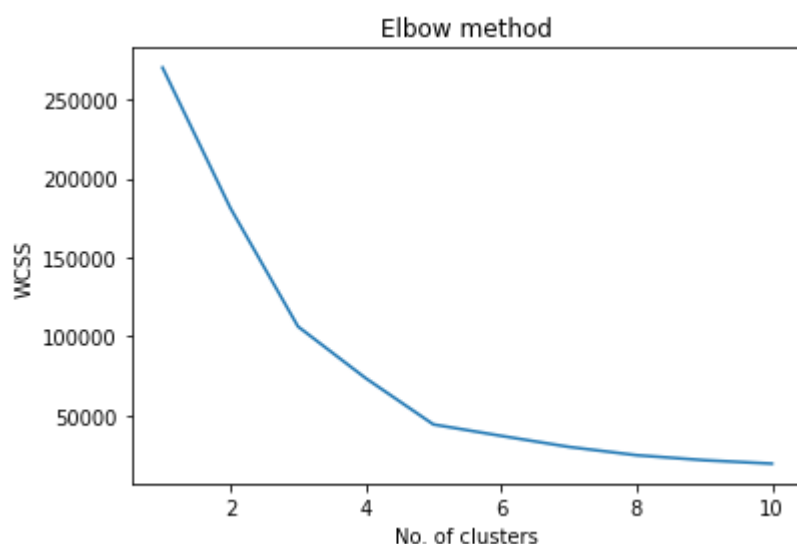
```
print(X[0:5,:])
```

```
[[15 39]
 [15 81]
 [16  6]
 [16 77]
 [17 40]]
```

Using the elbow method to find the optimal number of clusters

In []:

```
from sklearn.cluster import KMeans
wcss = []
for i in range(1,11):
    kmeans = KMeans(n_clusters= i, init='k-means++',random_state= 42)
    kmeans.fit(X)
    wcss.append(kmeans.inertia_)
plt.plot(range(1,11),wcss)
plt.title("Elbow method")
plt.xlabel("No. of clusters")
plt.ylabel("WCSS")
plt.show()
```



Type *Markdown* and LaTeX: α^2

Training the K-Means model on the dataset

In []:

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

Out[8]:

```
KMeans(n_clusters=5, random_state=42)
```

Statistics from the initialization run with the lowest SSE are available as attributes of kmeans after calling .fit()

In []:

```
# The lowest SSE value
kmeans.inertia_
```

Out[24]:

```
44448.45544793371
```

In []:

```
# Final locations of the centroid
kmeans.cluster_centers_
```

Out[25]:

```
array([[25.72727273, 79.36363636],
       [55.2962963 , 49.51851852],
       [86.53846154, 82.12820513],
       [26.30434783, 20.91304348],
       [88.2        , 17.11428571]])
```

In []:

```
# The number of iterations required to converge
kmeans.n_iter_
```

Out[26]:

```
4
```

In []:

```
#Finally, the cluster assignments are stored as a one-dimensional NumPy array in kmeans.  
kmeans.labels_
```

Out[27]:

[illegible]

Creating Output labels for Generating Graph

In []:

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

In []:

```
print(y_kmeans)
```

[illegible]

In []:

```
from sklearn.metrics import accuracy_score
print(accuracy_score(kmeans.labels_, y_kmeans))
```

1.0

Visualising the clusters

In []:

```
plt.scatter(X[y_kmeans == 0,0],X[y_kmeans == 0,1],s=100, c = 'red', label ="Cluster 1")
plt.scatter(X[y_kmeans == 1,0],X[y_kmeans == 1,1],s=100, c = 'blue', label ="Cluster 2")
plt.scatter(X[y_kmeans == 2,0],X[y_kmeans == 2,1],s=100, c = 'lightgreen', label ="Clust
plt.scatter(X[y_kmeans == 3,0],X[y_kmeans == 3,1],s=100, c = 'black', label ="Cluster 4"
plt.scatter(X[y_kmeans == 4,0],X[y_kmeans == 4,1],s=100, c = 'magenta', label ="Cluster
plt.scatter(kmeans.cluster_centers[:,0],kmeans.cluster_centers[:,1],s = 300, c = 'Yell
plt.title("Clusters of Customers",size = 25)
plt.xlabel("Annual Income {k$}")
plt.ylabel("Spending Scores{1 to 100}")
plt.legend()
plt.show()
```

