

day59-k-means-clustering

December 7, 2023

Day59 K Means Clustering By: Loga Aswin

```
[1]: # import libraries
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import KMeans
```

```
[2]: # load datasets
df = pd.read_csv('/content/Mall_Customers.csv')
```

Exploratory Data Analysis(EDA):

```
[3]: df.head()
```

```
[3]:
```

	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

```
[19]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 200 entries, 0 to 199
```

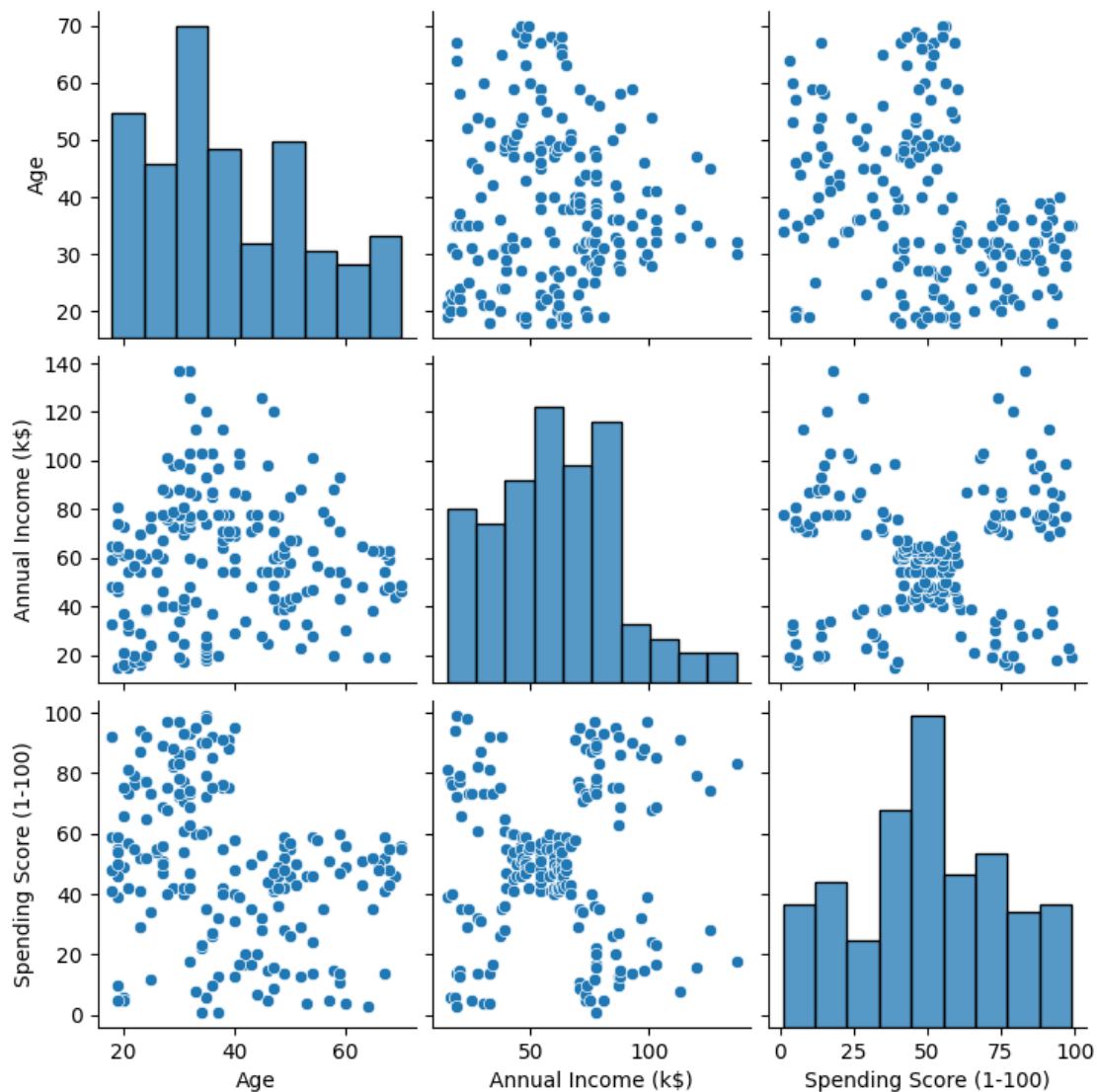
```
Data columns (total 6 columns):
```

#	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
5	Cluster	200 non-null	int32

```
dtypes: int32(1), int64(4), object(1)
```

```
memory usage: 8.7+ KB
```

```
[4]: sns.pairplot(df[['Age', 'Annual Income (k$)', 'Spending Score (1-100)']])
plt.show()
```



```
[5]: # Feature Selection
X = df[['Annual Income (k$)', 'Spending Score (1-100)']]
```

Choosing Number of Clusters (Elbow Method)

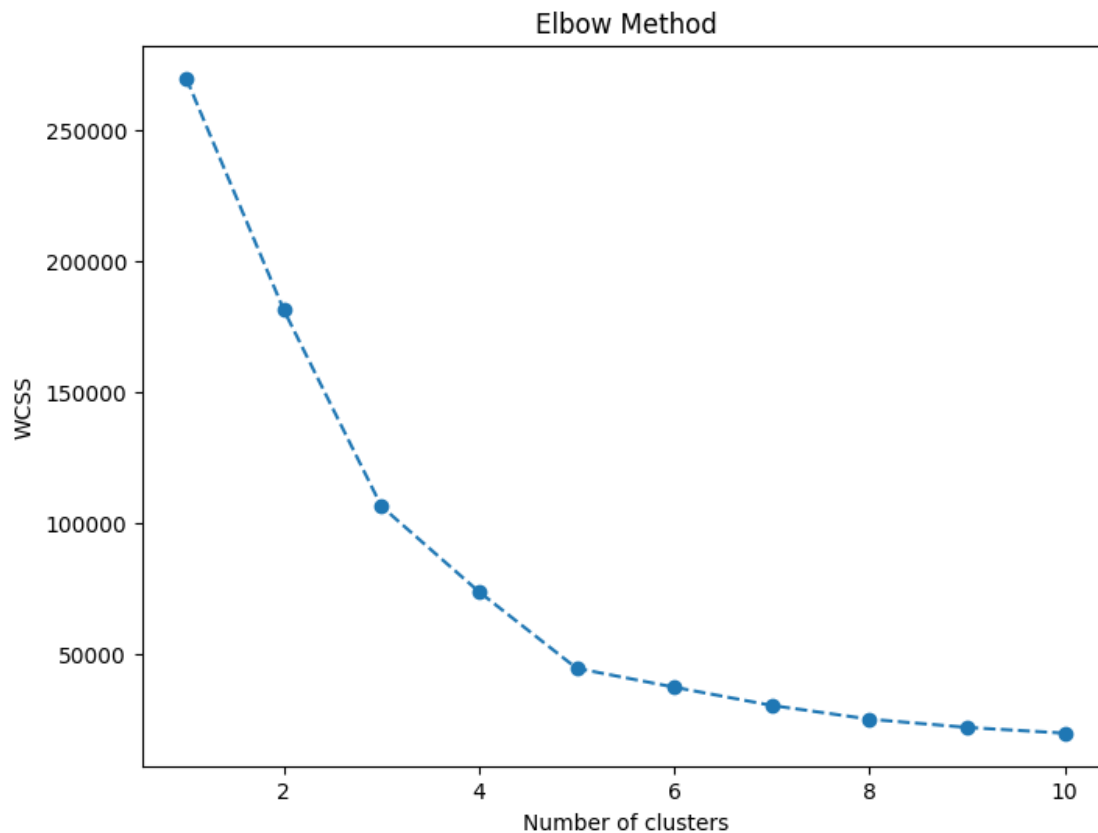
Using the Elbow Method, we'll plot the Within-Cluster Sum of Squares (WCSS) against different numbers of clusters.

After calculating WCSS for different numbers of clusters, we'll plot the number of clusters against their respective WCSS values.

```
[6]: 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_)
```

[illegible]

```
[7]: plt.figure(figsize=(8, 6))
plt.plot(range(1, 11), wcss, marker='o', linestyle='--')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.title('Elbow Method')
plt.show()
```



Applying KMeans Clustering:

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

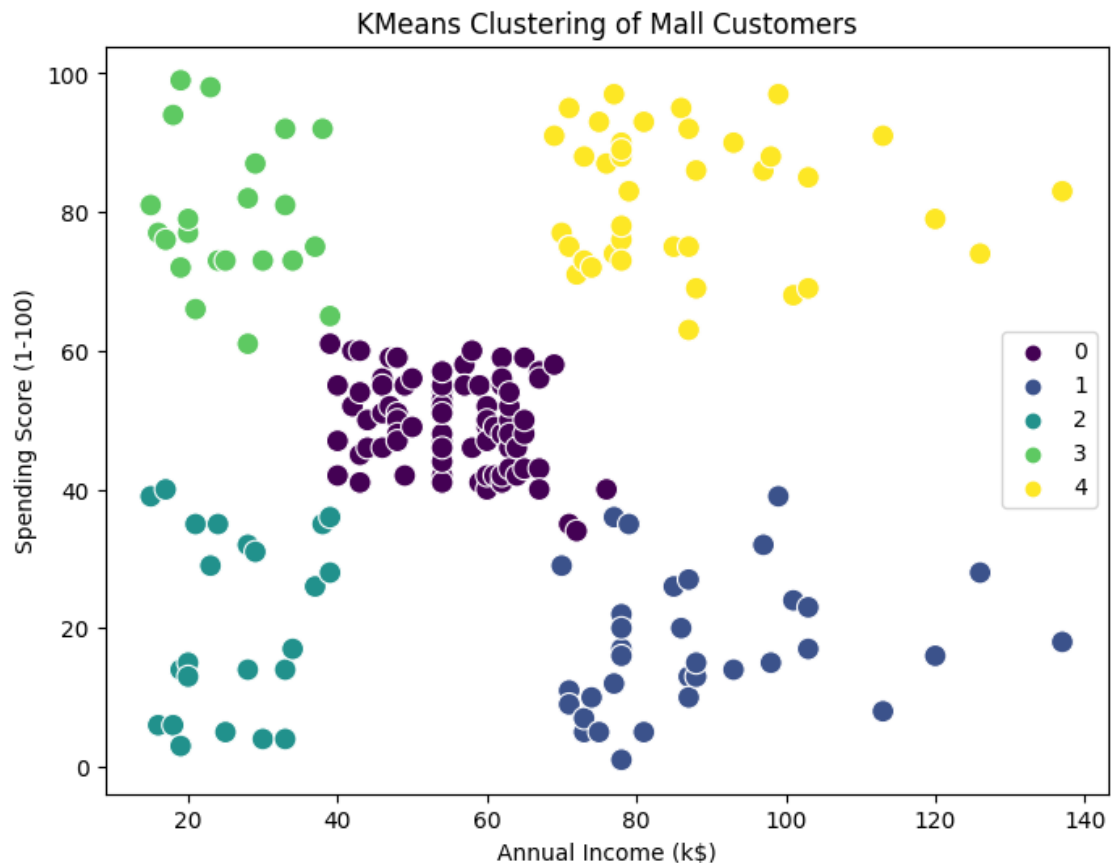
```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870:
FutureWarning: The default value of `n_init` will change from 10 to 'auto' in
1.4. Set the value of `n_init` explicitly to suppress the warning
warnings.warn(
```

```
[8]: KMeans(n_clusters=5, random_state=42)
```

```
[9]: df['Cluster'] = kmeans.labels_
```

Visualization of Clusters:

```
[16]: plt.figure(figsize=(8, 6))
sns.scatterplot(x='Annual Income (k$)', y='Spending Score (1-100)', data=df,
               hue='Cluster', palette='viridis', s=100)
plt.title('KMeans Clustering of Mall Customers')
plt.legend()
plt.show()
```



Model Evaluation Metrics:

```
[15]: from sklearn.metrics import silhouette_score

y = kmeans.labels_
silhouette = silhouette_score(X, y)
print(f"Silhouette Score: {silhouette}")
```

Silhouette Score: 0.553931997444648