

```
In [15]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: df = pd.read_csv("./Mall_Customers.csv")
```

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

```
Out[3]:
```

| | CustomerID | Gender | 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 [4]: df.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   Gender                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 [7]: df.describe()
```

Out[7]:

| | CustomerID | Age | Annual Income (k\$) | Spending Score (1-100) |
|--------------|------------|------------|---------------------|------------------------|
| count | 200.000000 | 200.000000 | 200.000000 | 200.000000 |
| mean | 100.500000 | 38.850000 | 60.560000 | 50.200000 |
| std | 57.879185 | 13.969007 | 26.264721 | 25.823522 |
| min | 1.000000 | 18.000000 | 15.000000 | 1.000000 |
| 25% | 50.750000 | 28.750000 | 41.500000 | 34.750000 |
| 50% | 100.500000 | 36.000000 | 61.500000 | 50.000000 |
| 75% | 150.250000 | 49.000000 | 78.000000 | 73.000000 |
| max | 200.000000 | 70.000000 | 137.000000 | 99.000000 |

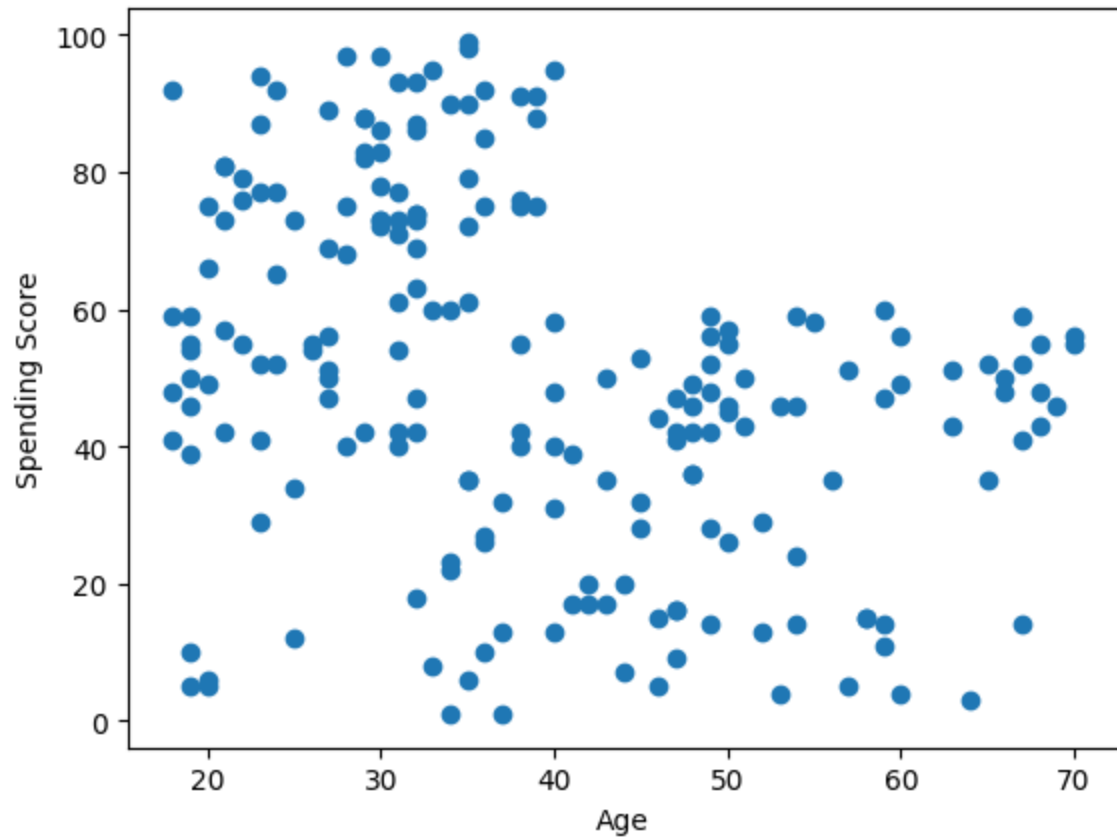
In [8]: `df.shape`

Out[8]: (200, 5)

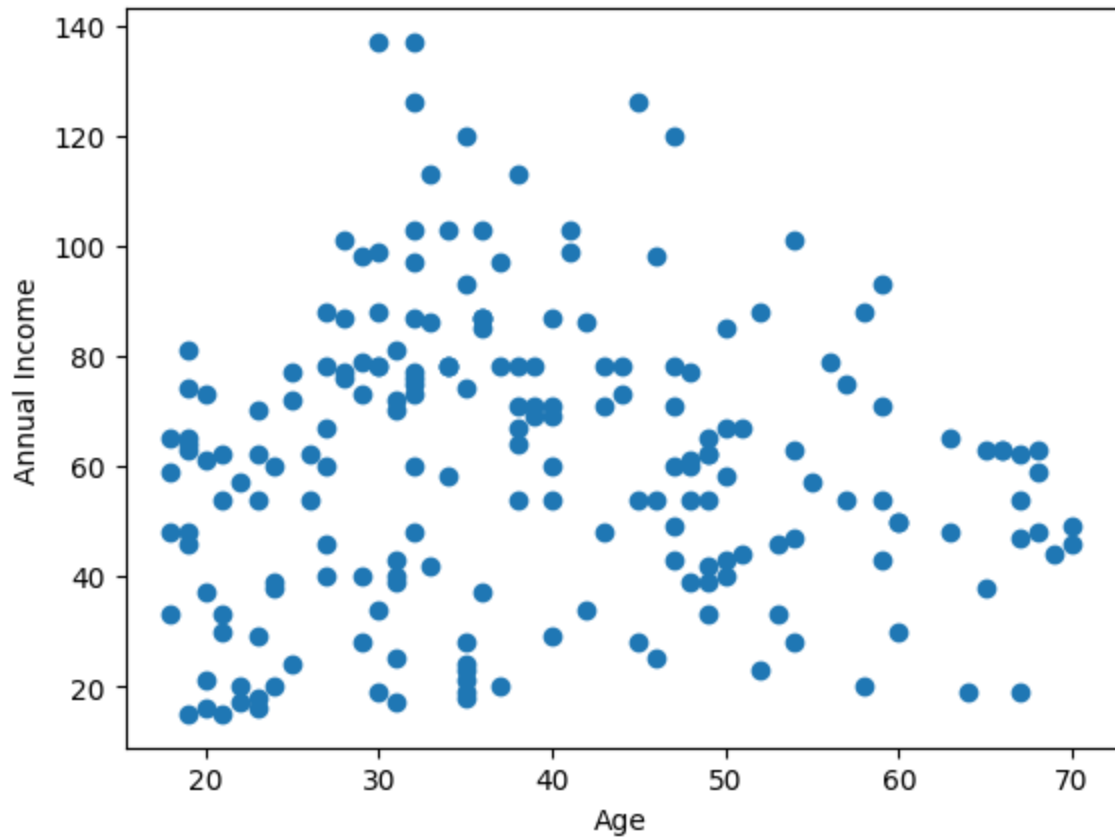
In [9]: `df.columns`

Out[9]: Index(['CustomerID', 'Gender', 'Age', 'Annual Income (k\$)',
'Spending Score (1-100)'],
dtype='object')

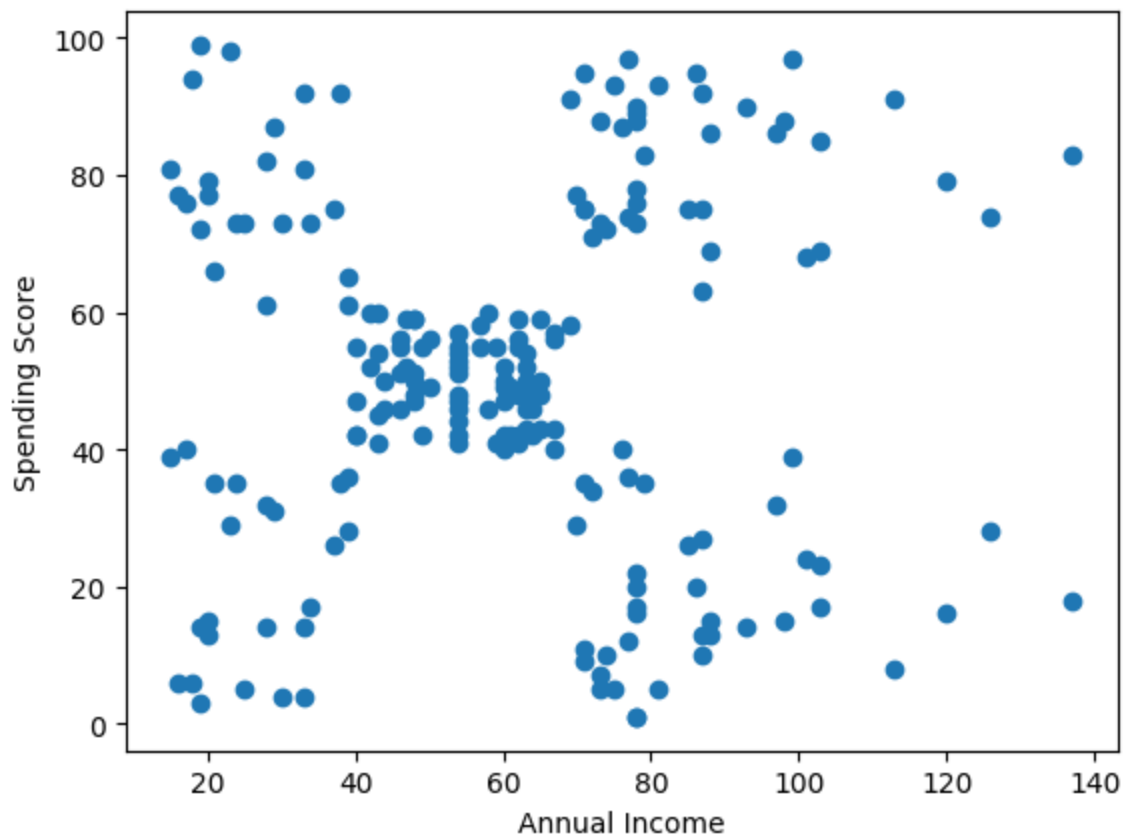
In [10]: `plt.scatter(df['Age'], df['Spending Score (1-100)'])`
`plt.xlabel('Age')`
`plt.ylabel('Spending Score')`
`plt.show()`



```
In [11]: plt.scatter(df['Age'], df['Annual Income (k$)'])  
plt.xlabel('Age')  
plt.ylabel('Annual Income')  
plt.show()
```



```
In [12]: plt.scatter(df['Annual Income (k$)'], df['Spending Score (1-100)'])  
plt.xlabel('Annual Income')  
plt.ylabel('Spending Score')  
plt.show()
```



```
In [13]: from sklearn.cluster import KMeans
```

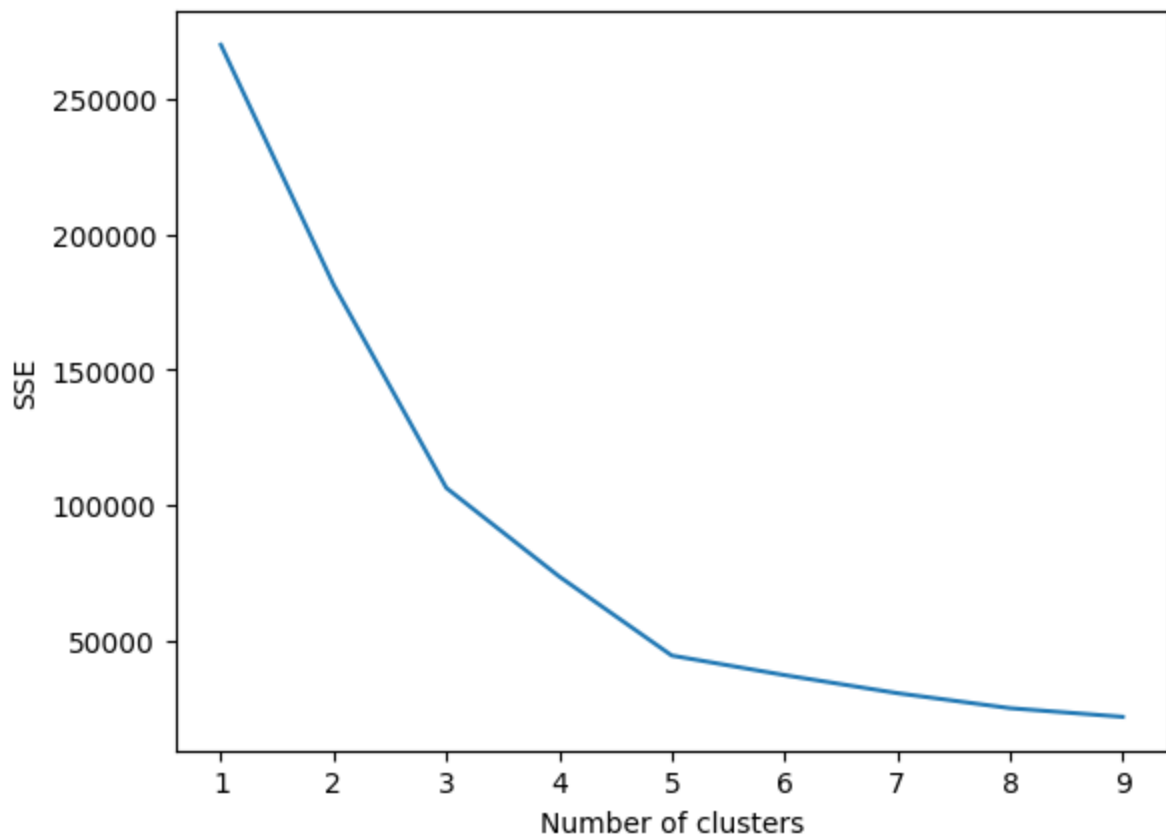
```
In [16]: sse = []
for i in range(1, 10):
    km = KMeans(n_clusters = i)
    km.fit(df[['Annual Income (k$)', 'Spending Score (1-100)']])
    sse.append(km.inertia_)
```

```
In [17]: sse
```

```
Out[17]: [269981.28,
181363.59595959593,
106348.37306211118,
73679.78903948836,
44448.45544793371,
37233.81451071001,
30566.45113025186,
25018.781613414074,
21818.11458845218]
```

```
In [18]: plt.xlabel('Number of clusters')
plt.ylabel("SSE")
plt.plot(range(1,10), sse)
```

```
Out[18]: [<matplotlib.lines.Line2D at 0x20679b02f10>]
```



```
In [19]: km = KMeans(n_clusters = 5)
predicted = km.fit_predict(df[['Annual Income (k$)', 'Spending Score (1-100)']])
predicted
```

```
Out[19]: array([2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4,
                2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 0,
                2, 4, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
                0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 3, 1, 3, 0, 3, 1, 3, 1, 3,
                0, 3, 1, 3, 1, 3, 1, 3, 1, 3, 0, 3, 1, 3, 1, 3, 1, 3, 1, 3,
                1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3,
                1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3, 1, 3,
                1, 3])
```

```
In [20]: df['Cluster'] = predicted
df
```

Out[20]:

| | CustomerID | Gender | Age | Annual Income (k\$) | Spending Score (1-100) | Cluster |
|------------|------------|--------|-----|---------------------|------------------------|---------|
| 0 | 1 | Male | 19 | 15 | 39 | 2 |
| 1 | 2 | Male | 21 | 15 | 81 | 4 |
| 2 | 3 | Female | 20 | 16 | 6 | 2 |
| 3 | 4 | Female | 23 | 16 | 77 | 4 |
| 4 | 5 | Female | 31 | 17 | 40 | 2 |
| ... | ... | ... | ... | ... | ... | ... |
| 195 | 196 | Female | 35 | 120 | 79 | 3 |
| 196 | 197 | Female | 45 | 126 | 28 | 1 |
| 197 | 198 | Male | 32 | 126 | 74 | 3 |
| 198 | 199 | Male | 32 | 137 | 18 | 1 |
| 199 | 200 | Male | 30 | 137 | 83 | 3 |

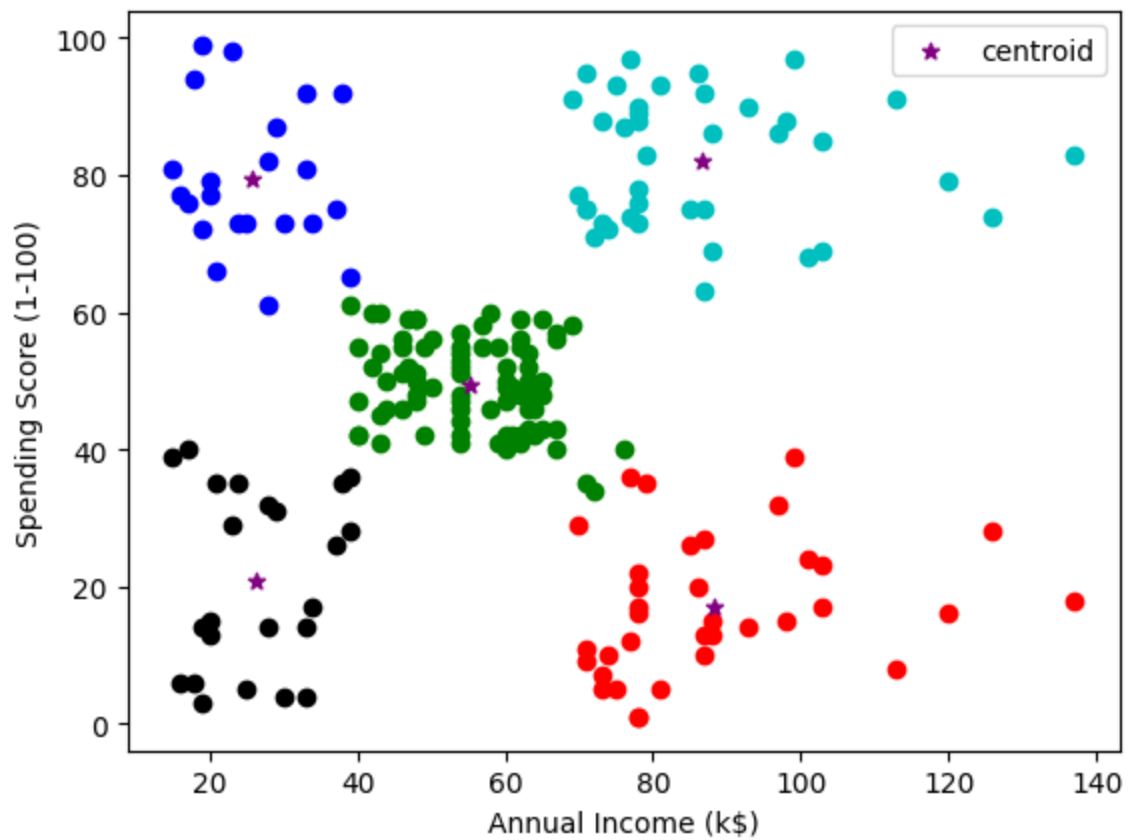
200 rows × 6 columns

```

In [21]: df1 = df[df.Cluster==0]
df2 = df[df.Cluster==1]
df3 = df[df.Cluster==2]
df4 = df[df.Cluster==3]
df5 = df[df.Cluster==4]
plt.scatter(df1['Annual Income (k$)'],df1['Spending Score (1-100)'],color='green')
plt.scatter(df2['Annual Income (k$)'],df2['Spending Score (1-100)'],color='red')
plt.scatter(df3['Annual Income (k$)'],df3['Spending Score (1-100)'],color='black')
plt.scatter(df4['Annual Income (k$)'],df4['Spending Score (1-100)'],color='c')
plt.scatter(df5['Annual Income (k$)'],df5['Spending Score (1-100)'],color='blue')
plt.scatter(km.cluster_centers_[ :,0],km.cluster_centers_[ :,1],color='purple',marker
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.legend()

```

Out[21]: <matplotlib.legend.Legend at 0x2067c44e790>



In []: