

Import Modules

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
import warnings
warnings.filterwarnings('ignore')
```

Load the Dataset

```
In [2]: df = pd.read_csv('/kaggle/input/customer-segmentation-tutorial-in-python/Mall_Customers.csv')
df.head()
```

```
Out[2]:
```

	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 [3]: # statistical info
df.describe()
```

```
Out[3]:
```

	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

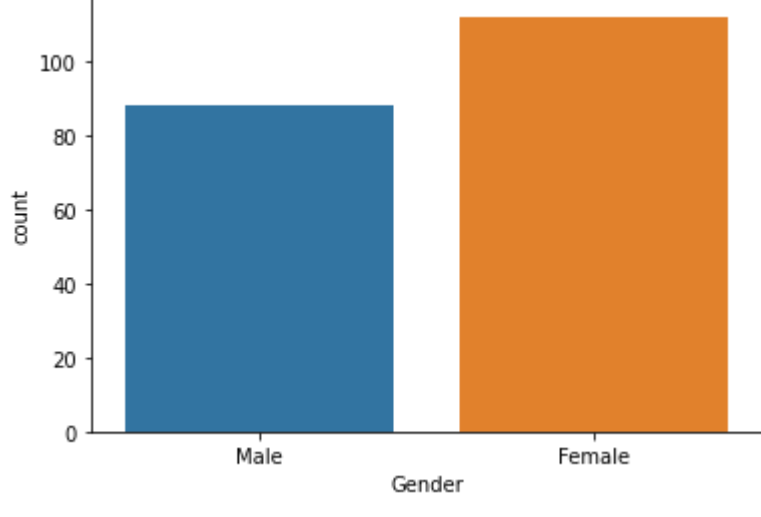
```
In [4]: # datatype info
df.info()
```

```
<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
```

Exploratory Data Analysis

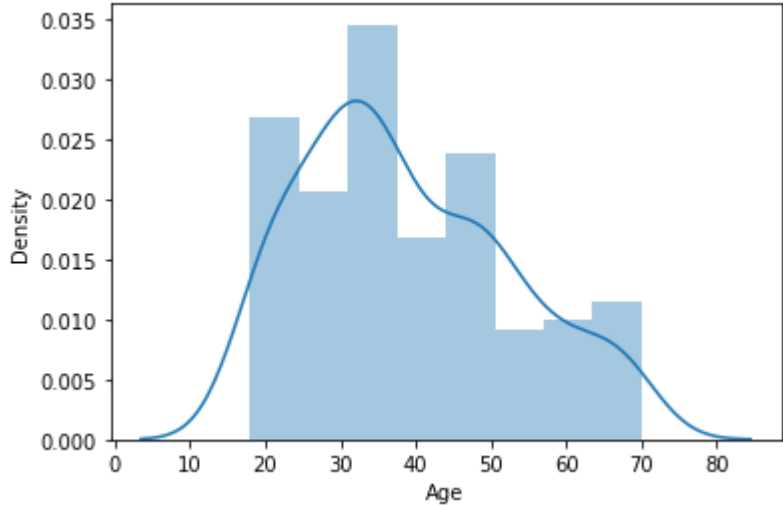
```
In [5]: sns.countplot(df['Gender'])
```

```
Out[5]: <AxesSubplot:xlabel='Gender', ylabel='count'>
```



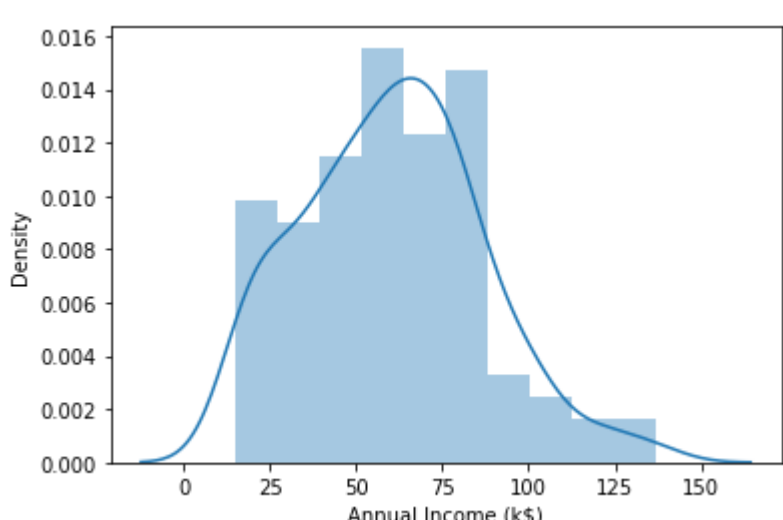
```
In [6]: sns.distplot(df['Age'])
```

```
Out[6]: <AxesSubplot:xlabel='Age', ylabel='Density'>
```



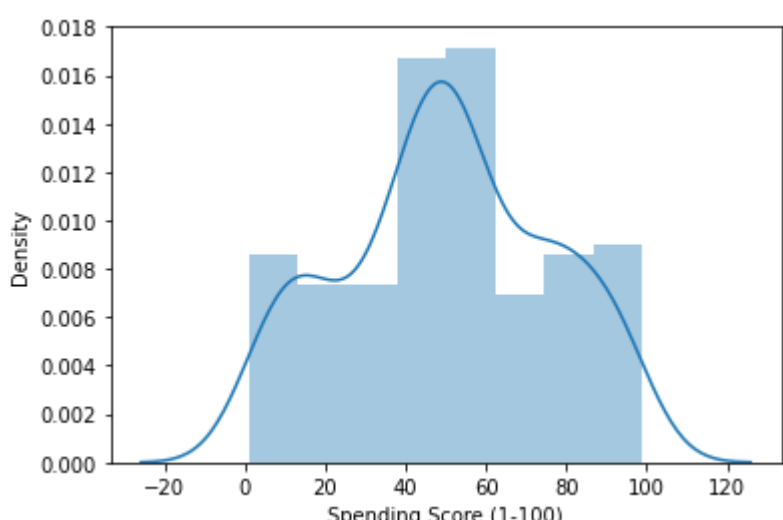
```
In [7]: sns.distplot(df['Annual Income (k$)'])
```

```
Out[7]: <AxesSubplot:xlabel='Annual Income (k$)', ylabel='Density'>
```



```
In [8]: sns.distplot(df['Spending Score (1-100)'])
```

```
Out[8]: <AxesSubplot:xlabel='Spending Score (1-100)', ylabel='Density'>
```



Correlation Matrix

```
In [9]: corr = df.corr()
sns.heatmap(corr, annot=True, cmap='coolwarm')
```

```
Out[9]: <AxesSubplot:>
```



Clustering

```
In [10]: df.head()
```

```
Out[10]:
```

	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

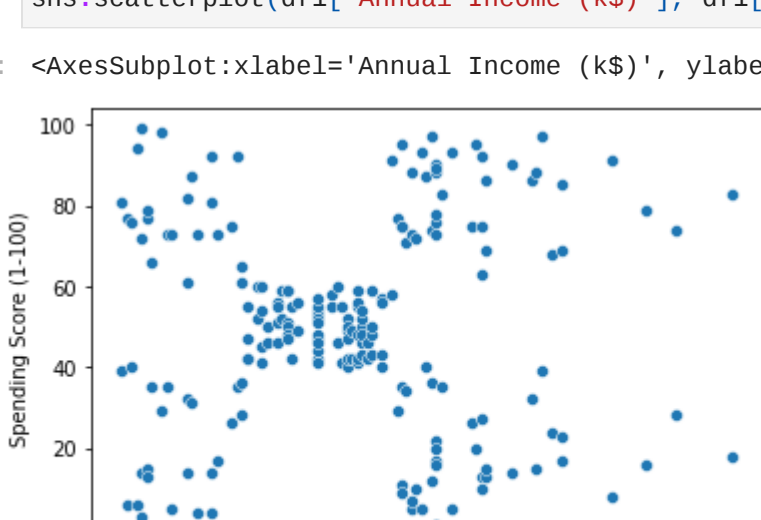
```
In [11]: # cluster on 2 features
df1 = df[['Annual Income (k$)', 'Spending Score (1-100)']]
df1.head()
```

```
Out[11]:
```

	Annual Income (k\$)	Spending Score (1-100)
0	15	39
1	15	81
2	16	6
3	16	77
4	17	40

```
In [12]: # scatter plot
sns.scatterplot(df1['Annual Income (k$)'], df1['Spending Score (1-100)'])
```

```
Out[12]: <AxesSubplot:xlabel='Annual Income (k$)', ylabel='Spending Score (1-100)'>
```



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

```
errors = []
for i in range(1, 11):
```

```
    kmeans = KMeans(n_clusters=i)
    kmeans.fit(df1)
    errors.append(kmeans.inertia_)
```

```
In [14]: # plot the results for elbow method
plt.figure(figsize=(15,6))
plt.plot(range(1,11), errors)
```

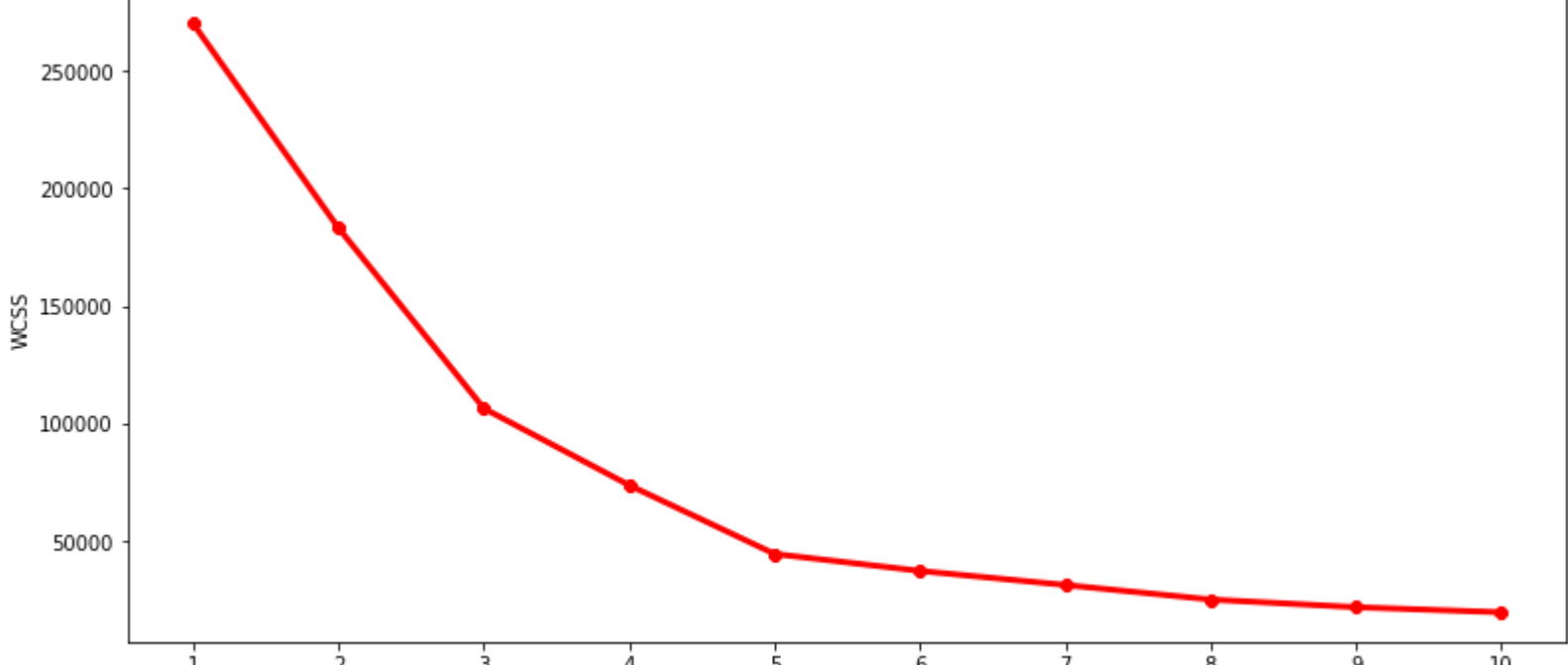
```
plt.plot(range(1,11), errors, linewidth=3, color='red', markers='8')
```

```
plt.xlabel('No. of clusters')
```

```
plt.ylabel('WCSS')
```

```
plt.xticks(np.arange(1,11,1))
```

```
plt.show()
```



```
In [15]: km = KMeans(n_clusters=5)
```

```
km.fit(df1)
```

```
y = km.predict(df1)
```

```
df1['Label'] = y
```

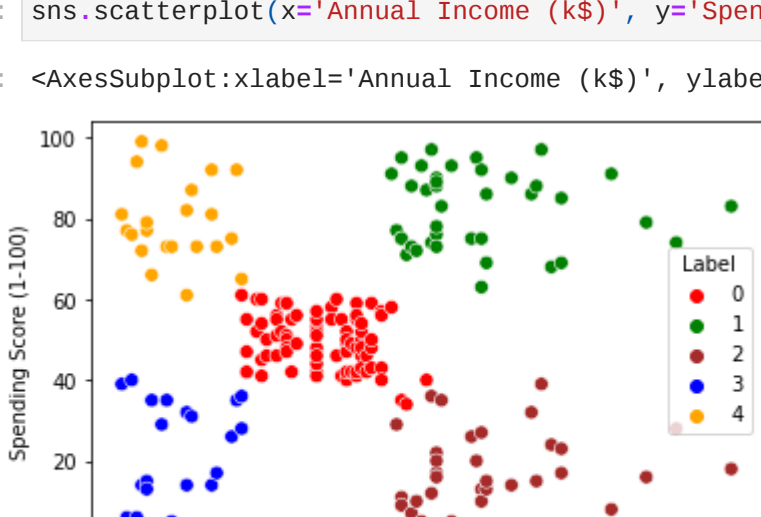
```
df1.head()
```

```
Out[15]:
```

	Annual Income (k\$)	Spending Score (1-100)	Label
0	15	39	3
1	15	81	4
2	16	6	3
3	16	77	4
4	17	40	3

```
In [16]: sns.scatterplot(x='Annual Income (k$)', y='Spending Score (1-100)', data=df1, hue='Label', s=50, palette=['red', 'green', 'brown', 'blue', 'orange'])
```

```
Out[16]: <AxesSubplot:xlabel='Annual Income (k$)', ylabel='Spending Score (1-100)'>
```



```
In [17]: # cluster on 3 features
```

```
df2 = df[['Annual Income (k$)', 'Spending Score (1-100)', 'Age']]
df2.head()
```

```
Out[17]:
```

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

```
In [18]: errors = []
```

```
for i in range(1, 11):
```

```
    kmeans = KMeans(n_clusters=i)
```

```
    kmeans.fit(df2)
    errors.append(kmeans.inertia_)
```

```
In [19]: # plot the results for elbow method
```

```
plt.figure(figsize=(15,6))
```

```
plt.plot(range(1,11), errors)
```

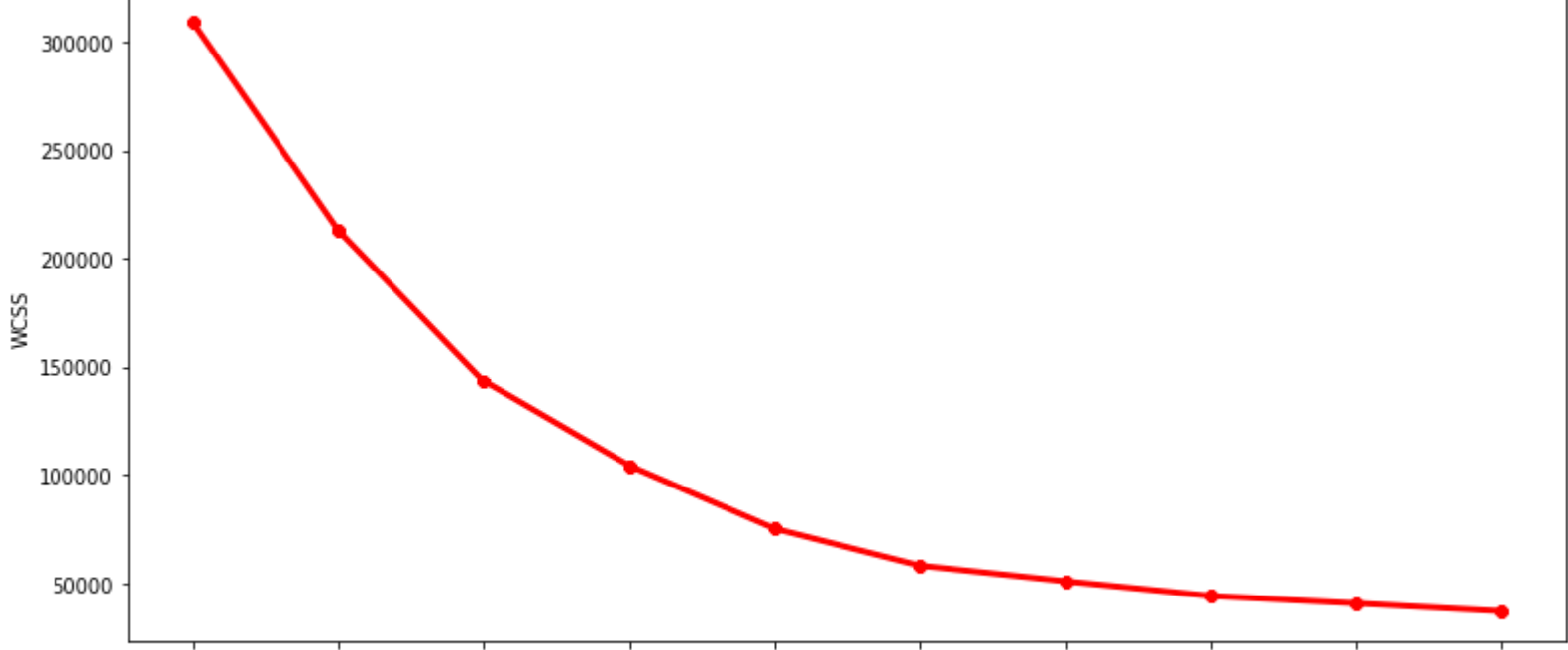
```
plt.plot(range(1,11), errors, linewidth=3, color='red', markers='8')
```

```
plt.xlabel('No. of clusters')
```

```
plt.ylabel('WCSS')
```

```
plt.xticks(np.arange(1,11,1))
```

```
plt.show()
```



```
In [20]: km = KMeans(n_clusters=5)
```

```
km.fit(df2)
```

```
y = km.predict(df2)
```

```
df2['Label'] = y
```

```
df2.head()
```

```
Out[20]:
```

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

```
In [21]: # 3d scatter plot
```

```
fig = plt.figure(figsize=(20,15))
```

```
ax = fig.add_subplot(111, projection='3d')
```

```
ax.scatter(df2['Age'][df2['Label']==0], df2['Annual Income (k$)'][df2['Label']==0], df2['Spending Score (1-100)'][df2['Label']==0], c='red', s=50)
```

```
ax.scatter(df2['Age'][df2['Label']==1], df2['Annual Income (k$)'][df2['Label']==1], df2['Spending Score (1-100)'][df2['Label']==1], c='green', s=50)
```

```
ax.scatter(df2['Age'][df2['Label']==2], df2['Annual Income (k$)'][df2['Label']==2], df2['Spending Score (1-100)'][df2['Label']==2], c='blue', s=50)
```

```
ax.scatter(df2['Age'][df2['Label']==3], df2['Annual Income (k$)'][df2['Label']==3], df2['Spending Score (1-100)'][df2['Label']==3], c='brown', s=50)
```

```
ax.scatter(df2['Age'][df2['Label']==4], df2['Annual Income (k$)'][df2['Label']==4], df2['Spending Score (1-100)'][df2['Label']==4], c='orange', s=50)
```

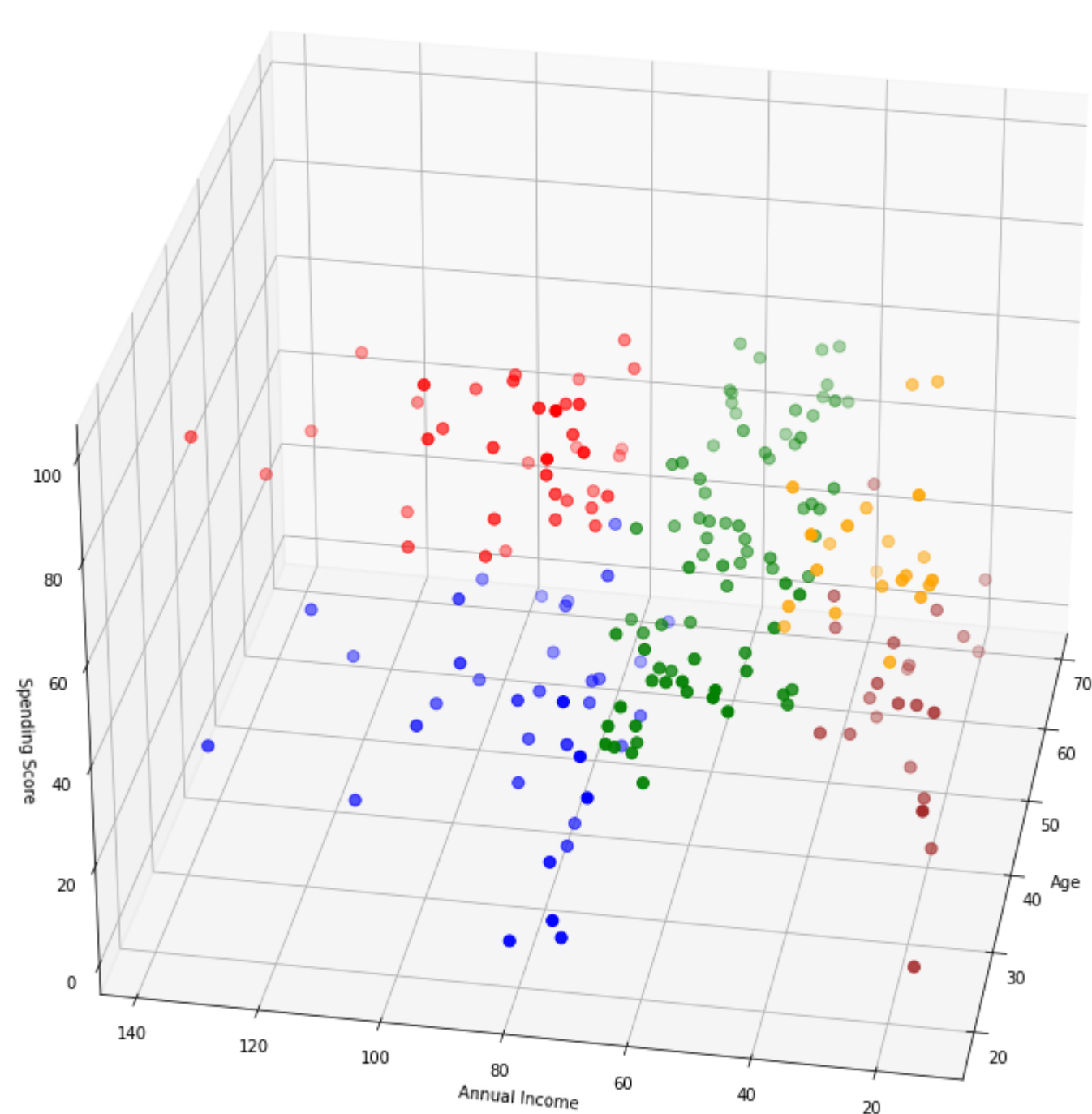
```
ax.view_init(30, 190)
```

```
ax.set_xlabel('Age')
```

```
ax.set_ylabel('Annual Income')
```

```
ax.set_zlabel('Spending Score')
```

```
plt.show()
```



In []: