In [114]:

```python
import numpy as np
import pandas as pd
import sklearn
import matplotlib.pyplot as plt
import seaborn as sns
```

In [115]:

```python
df = pd.read_csv('mall_customer.csv.zip')
```

In [116]:

```python
df.head()
```

Out[116]:

| | 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 [117]:

```python
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   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 [118]:

```
df.describe()
```

Out[118]:

|       | 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 [ ]:

In [119]:

```
df.dtypes
```

Out[119]:

```
CustomerID              int64
Genre                   object
Age                     int64
Annual Income (k$)      int64
Spending Score (1-100)  int64
dtype: object
```

In [120]:

```
df = df.rename(columns = {'Genre':'Gender'})
```

In [121]:

```
df.dtypes
```

Out[121]:

```
CustomerID              int64
Gender                  object
Age                     int64
Annual Income (k$)      int64
Spending Score (1-100)  int64
dtype: object
```

In [ ]:

In [ ]:

In [122]:

```python
#df.head()
```

In [123]:

```python
df.isnull().sum()
```

Out[123]:

```
CustomerID              0
Gender                  0
Age                     0
Annual Income (k$)      0
Spending Score (1-100)  0
dtype: int64
```
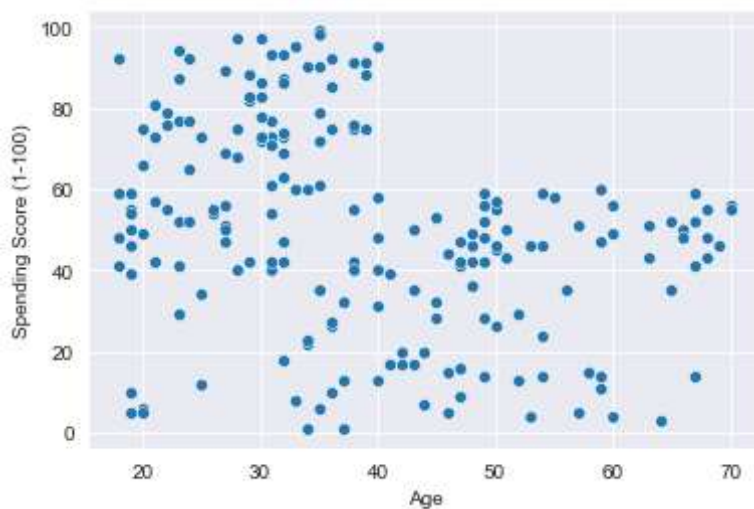
In [124]:

```python
df = df.drop('CustomerID',axis=1)
#dropping it as it's not useful for our analysis
```

In [125]:

```python
sns.scatterplot(x='Age',y='Spending Score (1-100)',data =df)
```
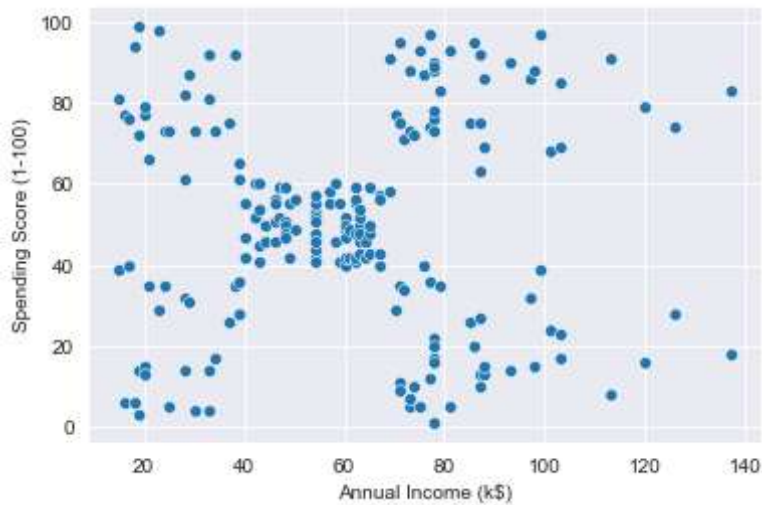
Out[125]:

```
<AxesSubplot:xlabel='Age', ylabel='Spending Score (1-100)'>
```

In [126]:

```python
sns.scatterplot(x='Annual Income (k$)',y='Spending Score (1-100)',data=df)
```
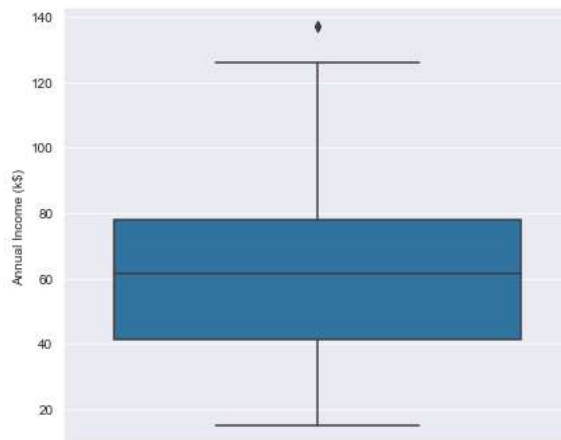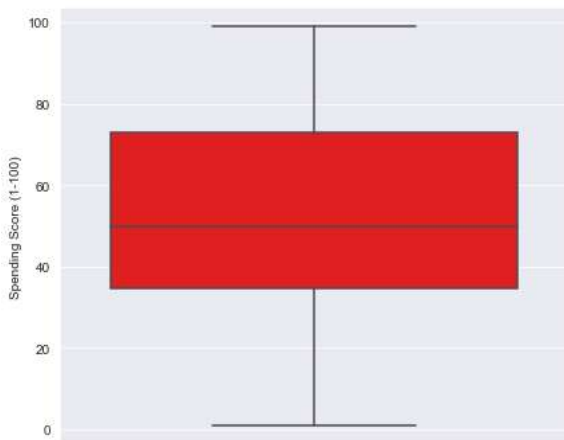
Out[126]:

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



In [127]:

```python
plt.figure(figsize=(15,6))
plt.subplot(1,2,1)
sns.boxplot(y=df["Spending Score (1-100)"], color="red")
plt.subplot(1,2,2)
sns.boxplot(y=df["Annual Income (k$)"])
plt.show()
#spending score is more than income
```
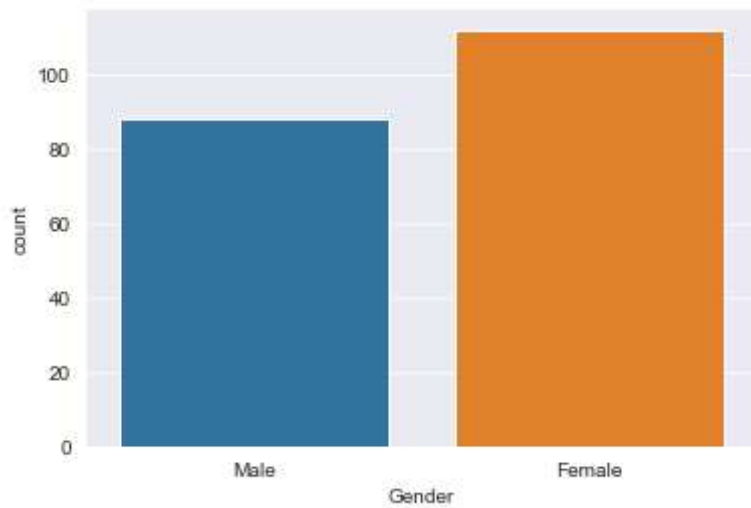


In [ ]:

In [128]:

```python
sns.countplot(x='Gender',data= df)
```

Out[128]:

```
<AxesSubplot:xlabel='Gender', ylabel='count'>
```

In [129]:

```python
age18_25 = df.Age[(df.Age <= 25) & (df.Age >= 18)]
age26_35 = df.Age[(df.Age <= 35) & (df.Age >= 26)]
age36_45 = df.Age[(df.Age <= 45) & (df.Age >= 36)]
age46_55 = df.Age[(df.Age <= 55) & (df.Age >= 46)]
age55above = df.Age[df.Age >= 56]

x = ["18-25","26-35","36-45","46-55","55+"]
y = [len(age18_25.values),len(age26_35.values),len(age36_45.values),len(age46_55.values),le

plt.figure(figsize=(15,6))
sns.barplot(x=x, y=y, palette="rocket")
plt.title("Number of Customer and Ages")
plt.xlabel("Age")
plt.ylabel("Number of Customer")
plt.show()
#checking age grups and spendings
```
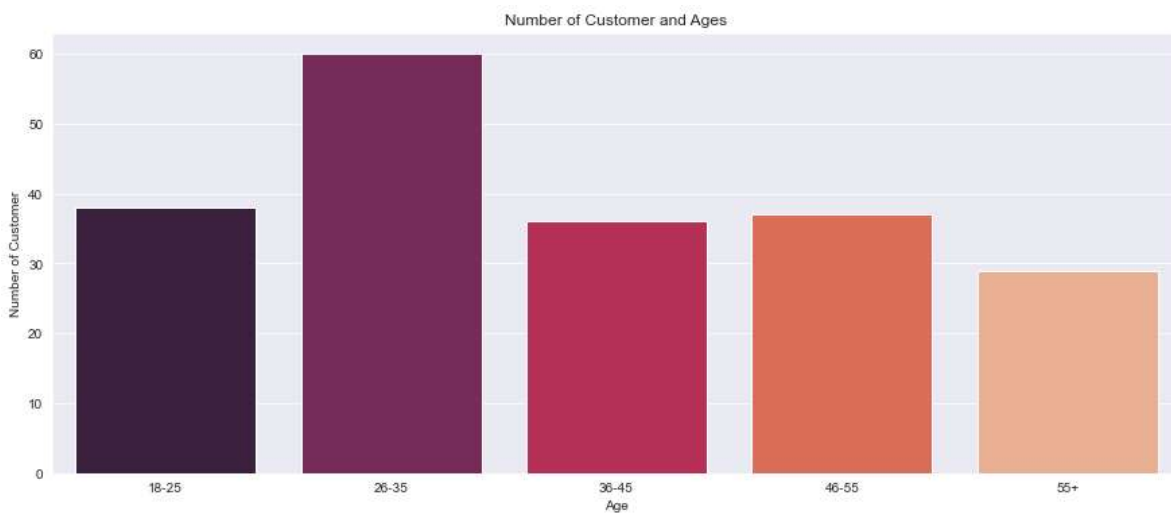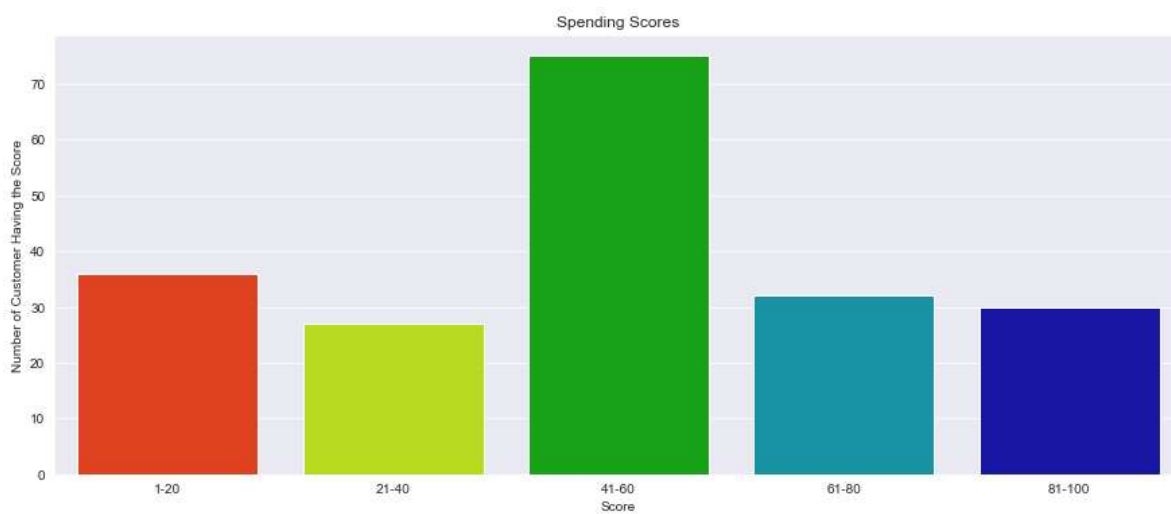
In [130]:

```python
ss1_20 = df["Spending Score (1-100)"][(df["Spending Score (1-100)"] >= 1) & (df["Spending S
ss21_40 = df["Spending Score (1-100)"][(df["Spending Score (1-100)"] >= 21) & (df["Spending
ss41_60 = df["Spending Score (1-100)"][(df["Spending Score (1-100)"] >= 41) & (df["Spending
ss61_80 = df["Spending Score (1-100)"][(df["Spending Score (1-100)"] >= 61) & (df["Spending
ss81_100 = df["Spending Score (1-100)"][(df["Spending Score (1-100)"] >= 81) & (df["Spendin

ssx = ["1-20", "21-40", "41-60", "61-80", "81-100"]
ssy = [len(ss1_20.values), len(ss21_40.values), len(ss41_60.values), len(ss61_80.values), l

plt.figure(figsize=(15,6))
sns.barplot(x=ssx, y=ssy, palette="nipy_spectral_r")
plt.title("Spending Scores")
plt.xlabel("Score")
plt.ylabel("Number of Customer Having the Score")
plt.show()
```
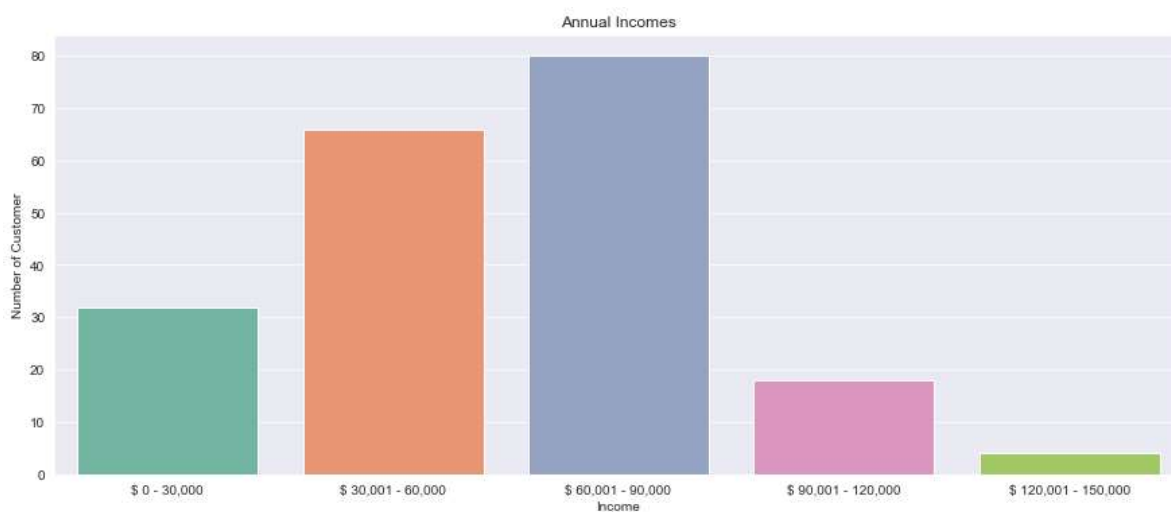
In [131]:

```python
ai0_30 = df["Annual Income (k$)"][(df["Annual Income (k$)"] >= 0) & (df["Annual Income (k$)
ai31_60 = df["Annual Income (k$)"][(df["Annual Income (k$)"] >= 31) & (df["Annual Income (k
ai61_90 = df["Annual Income (k$)"][(df["Annual Income (k$)"] >= 61) & (df["Annual Income (k
ai91_120 = df["Annual Income (k$)"][(df["Annual Income (k$)"] >= 91) & (df["Annual Income (
ai121_150 = df["Annual Income (k$)"][(df["Annual Income (k$)"] >= 121) & (df["Annual Income

aix = ["$ 0 - 30,000", "$ 30,001 - 60,000", "$ 60,001 - 90,000", "$ 90,001 - 120,000", "$ 1
aiy = [len(ai0_30.values), len(ai31_60.values), len(ai61_90.values), len(ai91_120.values),

plt.figure(figsize=(15,6))
sns.barplot(x=aix, y=aiy, palette="Set2")
plt.title("Annual Incomes")
plt.xlabel("Income")
plt.ylabel("Number of Customer")
plt.show()
```



In [132]:

```
ate optimal value of k as it varies with data and here we have 200 datapoints so we use WCSS
```

In [133]:

```
#Next I plotted Within Cluster Sum Of Squares (WCSS) against the the number of clusters (K

#where Yi is centroid for observation Xi. The main goal is to maximize number of clusters a
```
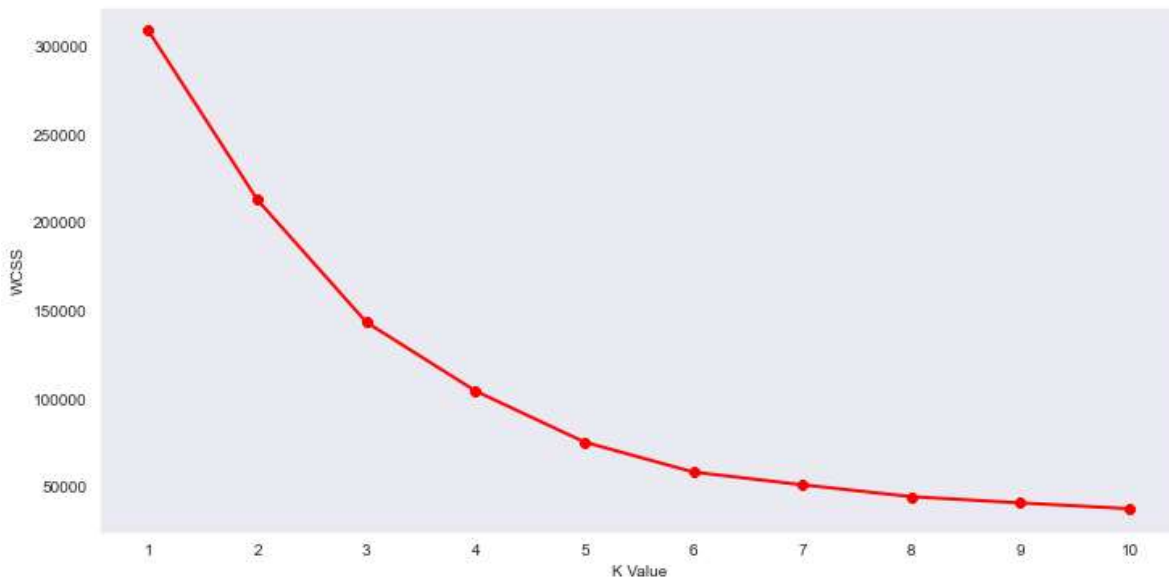
In [134]:

```python
from sklearn.cluster import KMeans
wcss = []
for k in range(1,11):
    kmeans = KMeans(n_clusters=k, init="k-means++")
    kmeans.fit(df.iloc[:,1:])
    wcss.append(kmeans.inertia_)
plt.figure(figsize=(12,6))
plt.grid()
plt.plot(range(1,11),wcss, linewidth=2, color="red", marker ="8")
plt.xlabel("K Value")
plt.xticks(np.arange(1,11,1))
plt.ylabel("WCSS")
plt.show()
```

C:\Users\20ai1\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:881: U
serWarning: KMeans is known to have a memory leak on Windows with MKL, when
there are less chunks than available threads. You can avoid it by setting th
e environment variable OMP_NUM_THREADS=1.
  warnings.warn(



In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

# hint

we can not find the best value for K but we can choose near value by elbow carev The Elbow method is a very popular technique and the idea is to run k-means clustering for a range of clusters k (let's say from 1 to 10) and for each value, we are calculating the sum of squared distances from each point to its assigned center(distortions).

When the distortions are plotted and the plot looks like an arm then the "elbow"(the point of inflection on the curve) is the best value of k.
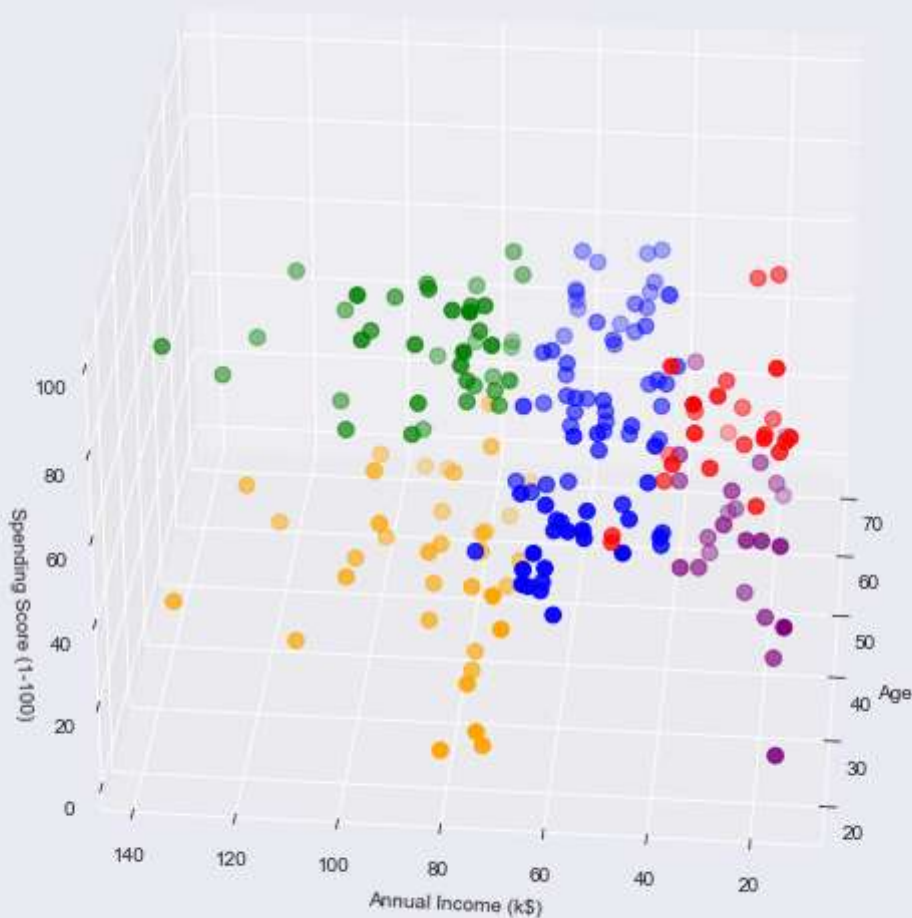
In [ ]:

In [ ]:

In [135]:

```python
km = KMeans(n_clusters=5)
clusters = km.fit_predict(df.iloc[:,1:])
df["label"] = clusters

from mpl_toolkits.mplot3d import Axes3D
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd

fig = plt.figure(figsize=(20,10))
ax = fig.add_subplot(111, projection='3d')
ax.scatter(df.Age[df.label == 0], df["Annual Income (k$)"][df.label == 0], df["Spending Sco
ax.scatter(df.Age[df.label == 1], df["Annual Income (k$)"][df.label == 1], df["Spending Sco
ax.scatter(df.Age[df.label == 2], df["Annual Income (k$)"][df.label == 2], df["Spending Sco
ax.scatter(df.Age[df.label == 3], df["Annual Income (k$)"][df.label == 3], df["Spending Sco
ax.scatter(df.Age[df.label == 4], df["Annual Income (k$)"][df.label == 4], df["Spending Sco
ax.view_init(30, 185)
plt.xlabel("Age")
plt.ylabel("Annual Income (k$)")
ax.set_zlabel('Spending Score (1-100)')
plt.show()
```

In [ ]:

In [ ]: