

SOURCE CODE

```
#Importing the necessary libraries

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

from mpl_toolkits.mplot3d import Axes3D

%matplotlib inline


data=pd.read_csv("Mall_Customers.csv")


data.head()


data.corr()


#Distribution of Annual Income
plt.figure(figsize=(10, 6))
sns.set(style = 'whitegrid')
sns.distplot(data['Annual Income (k$)'])
plt.title('Distribution of Annual Income (k$)', fontsize = 20)
plt.xlabel('Range of Annual Income (k$)')
plt.ylabel('Count')

Text(0, 0.5, 'Count')


#Distribution of age
plt.figure(figsize=(10, 6))
sns.set(style = 'whitegrid')
sns.distplot(data['Age'])
plt.title('Distribution of Age', fontsize = 20)
```

```
plt.xlabel('Range of Age')
plt.ylabel('Count')
```

```
#Distribution of spending score
plt.figure(figsize=(10, 6))
sns.set(style = 'whitegrid')
sns.distplot(data['Spending Score (1-100)'])
plt.title('Distribution of Spending Score (1-100)', fontsize = 20)
plt.xlabel('Range of Spending Score (1-100)')
plt.ylabel('Count')
```

```
genders = data.Gender.value_counts()
sns.set_style("darkgrid")
plt.figure(figsize=(10,4))
sns.barplot(x=genders.index, y=genders.values)
plt.show()
```

```
df1=data[["CustomerID","Gender","Age","Annual Income (k$)","Spending Score (1-100)"]]
X=df1[["Annual Income (k$)","Spending Score (1-100)"]]
```

```
X.head()
```

```
#Scatterplot of the input data
plt.figure(figsize=(10,6))
sns.scatterplot
```

```
(x = 'Annual Income (k$)', y = 'Spending Score (1-100)', data = X ,s = 60 )
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.title('Spending Score (1-100) vs Annual Income (k$)')
plt.show()
```

```
#Importing KMeans from sklearn
from sklearn.cluster import KMeans
```

```
wcss=[]
for i in range(1,11):
    km=KMeans(n_clusters=i)
    km.fit(X)
    wcss.append(km.inertia_)
```

```
#The elbow curve
plt.figure(figsize=(12,6))
plt.plot(range(1,11),wcss)
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()
```

```
#Taking 5 clusters
km1=KMeans(n_clusters=5)
#Fitting the input data
km1.fit(X)
#predicting the labels of the input data
y=km1.predict(X)
#adding the labels to a column named label
df1["label"] = y
#The new dataframe with the clustering done
df1.head()
```

```
#Scatterplot of the clusters
```

```
plt.figure(figsize=(10,6))
sns.scatterplot(x = 'Annual Income (k$)',y = 'Spending Score (1-100)',hue="label",
               palette=['green','orange','brown','dodgerblue','red'], legend='full',data = df1 ,s = 60 )
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.title('Spending Score (1-100) vs Annual Income (k$)')
plt.show()
```

 NameError Traceback (most recent call last)

<ipython-input-20-6473c0047571> in <module>

```
1 #We choose the k for which WSS starts to diminish
2 km2 = KMeans(n_clusters=5)
----> 3 y2 = km.fit_predict(X2)
4 df2["label"] = y2
5 #The data with labels
```

NameError: name 'X2' is not defined

#Taking the features

```
X2=df1[["Age","Annual Income (k$)","Spending Score (1-100)"]]
```

#Now we calculate the Within Cluster Sum of Squared Errors (WSS) for different values of k.

```
wcss = []
```

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

```
    kmeans = KMeans(n_clusters=k, init="k-means++")
```

```
    kmeans.fit(X2)
```

```
    wcss.append(kmeans.inertia_)
```

```
plt.figure(figsize=(12,6))
```

```
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()
```

```
#We choose the k for which WSS starts to diminish
```

```
km2 = KMeans(n_clusters=5)
```

```
y2 = km.fit_predict(X2)
```

```
df1["label"] = y2
```

```
#The data with labels
```

```
df1.head()
```

```
#3D Plot as we did the clustering on the basis of 3 input features
```

```
fig = plt.figure(figsize=(20,10))
```

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

```
ax.scatter(df1.Age[df1.label == 0], df1["Annual Income (k$)"][df1.label == 0], df1["Spending Score (1-100)"][df1.label == 0], c='purple', s=60)
```

```
ax.scatter(df1.Age[df1.label == 1], df1["Annual Income (k$)"][df1.label == 1], df1["Spending Score (1-100)"][df1.label == 1], c='red', s=60)
```

```
ax.scatter(df1.Age[df1.label == 2], df1["Annual Income (k$)"][df1.label == 2], df1["Spending Score (1-100)"][df1.label == 2], c='blue', s=60)
```

```
ax.scatter(df1.Age[df1.label == 3], df1["Annual Income (k$)"][df1.label == 3], df1["Spending Score (1-100)"][df1.label == 3], c='green', s=60)
```

```
ax.scatter(df1.Age[df1.label == 4], df1["Annual Income (k$)"][df1.label == 4], df1["Spending Score (1-100)"][df1.label == 4], c='yellow', s=60)
```

```
ax.view_init(35, 185)
```

```
plt.xlabel("Age")
```

```
plt.ylabel("Annual Income (k$)")
```

```
ax.set_zlabel('Spending Score (1-100)')
```

```
plt.show()
```

```
cust1=df1[df1["label"]==1]
```

```
print('Number of customer in 1st group=', len(cust1))
```

```
print('They are -', cust1["CustomerID"].values)
```

```
print("-----")
```

```

cust2=df1[df1["label"]==2]
print('Number of customer in 2nd group=', len(cust2))
print('They are -', cust2["CustomerID"].values)
print("-----")
cust3=df1[df1["label"]==0]
print('Number of customer in 3rd group=', len(cust3))
print('They are -', cust3["CustomerID"].values)
print("-----")
cust4=df1[df1["label"]==3]
print('Number of customer in 4th group=', len(cust4))
print('They are -', cust4["CustomerID"].values)
print("-----")
cust5=df1[df1["label"]==4]
print('Number of customer in 5th group=', len(cust5))
print('They are -', cust5["CustomerID"].values)
print("-----")

```

Number of customer in 1st group= 27

They are - [77 78 80 84 86 90 93 94 97 99 102 105 107 108 109 110 111 113
117 118 119 120 122 123 127 147 161]

Number of customer in 2nd group= 22

They are - [2 4 6 8 10 12 14 16 18 20 22 24 26 28 30 32 34 36 38 40 42 46]

Number of customer in 3rd group= 30

They are - [44 48 52 53 59 62 66 69 70 76 79 82 85 88 89 92 95 96
98 100 101 104 106 112 114 115 116 121 133 143]

Number of customer in 4th group= 11

They are - [180 182 184 186 188 190 192 194 196 198 200]

Number of customer in 5th group= 11

They are - [179 181 183 185 187 189 191 193 195 197 199]
