

```
In [1]: # Customer Segmentation and Analysis
# Steps to solve the problem :
# Importing Libraries.
# Exploration of data.
# Data Visualization.
# Clustering using K-Means.
# Selection of Clusters.
# Plotting the Cluster Boundary and Clusters.
# 3D Plot of Clusters.

In [2]: # Importing Libraries.
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import matplotlib.pyplot as plt
import seaborn as sns
import plotly as py
import plotly.graph_objs as go
from sklearn.cluster import KMeans
import warnings
import os
warnings.filterwarnings("ignore")
py.offline.init_notebook_mode(connected = True)
#%print(os.listdir("./input"))

In [3]: # Data Exploration
df = pd.read_csv(r'Mall_Customers.csv')
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.shape
Out[4]: (2000, 5)

In [5]: df.describe()

Out[5]:
```

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

```
In [6]: df.dtypes
```

```
Out[6]: CustomerID      int64
Gender          object
Age            int64
Annual Income (k$)    int64
Spending Score (1-100) int64
dtype: object
```

```
In [7]: df.isnull().sum()
Out[7]: CustomerID      0
Gender          0
Age            0
Annual Income (k$) 0
Spending Score (1-100) 0
dtype: int64
```

```
In [8]: # Data Visualization
plt.style.use('fivethirtyeight')
```

```
In [9]: # Histograms
plt.figure(1 , figsize = (15 , 6))
n = 0
for x in ['Age' , 'Annual Income (k$)' , 'Spending Score (1-100)']:
    n += 1
    plt.subplot(1 , 3 , n)
    plt.subplots_adjust(hspace = 0.5 , wspace = 0.5)
    sns.distplot(df[x] , bins = 20)
    plt.title('Distplot of {}'.format(x))
plt.show()
```

```
In [10]: # Count Plot of Gender
plt.figure(1 , figsize = (15 , 5))
sns.countplot(y = 'Gender' , data = df)
plt.show()
```

```
In [11]: # Plotting the Relation between Age , Annual Income and Spending Score
plt.figure(1 , figsize = (15 , 7))
n = 0
for x in ['Age' , 'Annual Income (k$)' , 'Spending Score (1-100)']:
    n += 1
    plt.subplot(3 , 3 , n)
    plt.subplots_adjust(hspace = 0.5 , wspace = 0.5)
    sns.regplot(x = x , y = y , data = df)
    plt.ylabel(y.split()[0]+''+y.split()[-1] if len(y.split()) > 1 else y)
plt.show()
```

```
In [12]: plt.figure(1 , figsize = (15 , 6))
for gender in ['Male' , 'Female']:
    plt.scatter(x = 'Age' , y = 'Annual Income (k$)' , data = df[df['Gender'] == gender] ,
                s = 200 , alpha = 0.5 , label = gender)
    plt.xlabel('Age vs Annual Income w.r.t Gender')
    plt.title('Age vs Annual Income w.r.t Gender')
    plt.legend()
    plt.show()
```

```
In [13]: plt.figure(1 , figsize = (15 , 6))
for gender in ['Male' , 'Female']:
    data = df[df['Gender'] == gender]
    plt.plot(x = 'Annual Income (k$)' , y = 'Spending Score (1-100)' ,
             data = data , s = 200 , alpha = 0.5 , label = gender)
    plt.xlabel('Annual Income (k$)') , plt.ylabel('Spending Score (1-100)')
    plt.title('Annual Income vs Spending Score w.r.t Gender')
    plt.legend()
    plt.show()
```

```
In [14]: # Distribution of values in Age , Annual Income and Spending Score according to Gender
plt.figure(1 , figsize = (15 , 7))
n = 0
for cols in ['Age' , 'Annual Income (k$)' , 'Spending Score (1-100)']:
    n += 1
    plt.subplot(3 , 3 , n)
    plt.subplots_adjust(hspace = 0.5 , wspace = 0.5)
    sns.violinplot(x = cols , y = 'Gender' , data = df , palette = 'vlag')
    sns.swarmplot(x = cols , y = 'Gender' , data = df)
    plt.ylabel('Gender' if n == 1 else '')
    plt.title('Boxplots & Swarmplots' if n == 2 else '')
plt.show()
```

```
In [15]: # Clustering using K-means
# 1. Segmentation using Age and Spending Score
'''Age and spending Score'''
X1 = df[['Age' , 'Spending Score (1-100)']].iloc[:, :].values
inertia = []
for n in range(1 , 11):
    algorithm = (KMeans(n_clusters = n, init='k-means++', n_init = 10 , max_iter=300,
                        tol=0.0001, random_state= 111 , algorithm='elkan'))
    inertia.append(algorithm.inertia_)
```

```
In [16]: # Selecting N Clusters based on Inertia (Squared Distance between Centroids and data points, should be less)
plt.plot(np.arange(1 , 11) , inertia , 'o')
plt.xlabel('Number of Clusters') , plt.ylabel('Inertia')
plt.show()
```

```
In [27]: # Import the necessities libraries
```

```
import plotly.offline as pyo
```

```
import plotly.graph_objs as go
```

```
# Set notebook mode to work in offline
```

```
pyo.init_notebook_mode()
```

```
df['label3'] = labels3
```

```
trace1 = go.Scatter3d(
```

```
x=df['Age'],
```

```
y=df['Spending Score (1-100)'],
```

```
z=df['Annual Income (k$)'],
```

```
mode='markers',
```

```
marker=dict(
```

```
color=df['label3'],
```

```
size=20,
```

```
line=dict(
```

```
color=df['label3'],
```

```
width=12
```

```
),
```

```
opacity=0.8
```

```
)
```

```
)
```

```
data = [trace1]
```

```
layout = go.Layout(
```

```
margin=dict(
```

```
l=0,
```

```
r=0,
```

```
b=0,
```

```
t=0
```

```
),
```

```
title= 'Clusters',
```

```
scene = dict(
```

```
xaxis = dict(title = 'Age'),
```

```
yaxis = dict(title = 'Spending Score'),
```

```
zaxis = dict(title = 'Annual Income')
```

```
)
```

```
)
```

```
fig = go.Figure(data=data, layout=layout)
```

```
pyo.offline.iplot(fig)
```

```
-----
```

```
In [ ]:
```