

Assignment -4

Python Programming

Customer Segmentation Analysis

In [14]: `# Import required libraries`

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.api as sm
```

In [15]: `# Loading the dataset`

```
df=pd.read_csv('Mall_Customers.csv')
df.head()
```

Out[15]:

	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 [16]: `df.shape`

Out[16]: (200, 5)

In [17]: `df.info()`

```
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 [18]: `df.isnull().any()`

```
Out[18]: CustomerID      False
Gender      False
Age         False
Annual Income (k$)  False
Spending Score (1-100)  False
dtype: bool
```

In [19]: `df.describe()`

Out[19]:

	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)
count	200.000000	200.000000	200.000000	200.000000
mean	100.500000	34.850000	60.560000	50.200000
std	57.879185	13.969907	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	34.000000	61.500000	50.000000
75%	150.250000	49.000000	78.000000	71.000000
max	200.000000	70.000000	137.000000	99.000000

Univariate Analysis

```
In [10]: df['Age'].mean()
```

```
Out[10]: 38.85
```

```
In [11]: df['Age'].median()
```

```
Out[11]: 36.0
```

```
In [12]: df['Age'].std()
```

```
Out[12]: 13.969007331558883
```

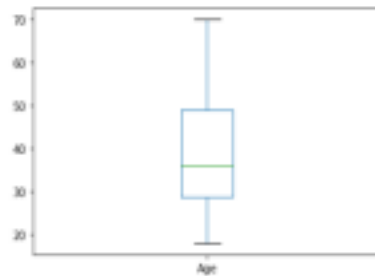
```
In [13]: df['Annual Income (k$)'].value_counts()
```

```
Out[13]: 54    12
       78    12
       48     6
       71     6
       63     6
       ..
       58     2
       59     2
       35     2
       64     2
       37     2
Name: Annual Income (k$), Length: 64, dtype: int64
```

Visualization

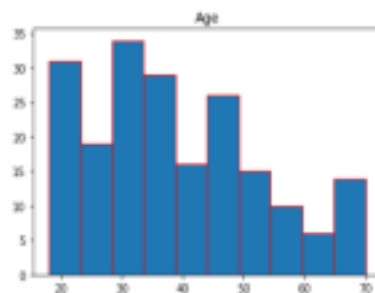
```
In [14]: df.boxplot(column='Age', grid=False)
```

```
Out[14]:
```



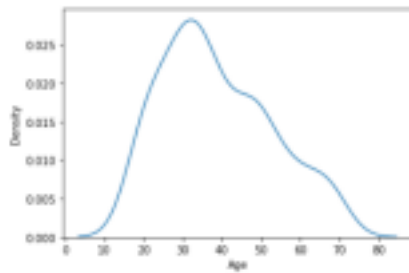
```
In [15]: df.hist(column='Age', grid=False, edgecolor='Red')
```

```
Out[15]: array([], dtype=object)
```



```
In [16]: sns.kdeplot(df['Age'])
```

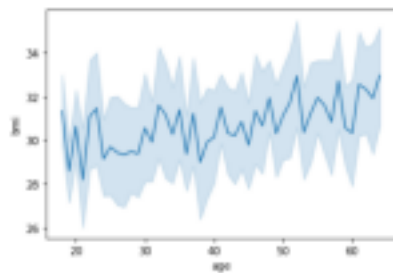
```
Out[16]:
```



```
In [18]: sns.lineplot(df.age, df.bmi)
```

C:\Users\Saanya\Anaconda3\lib\site-packages\seaborn\decorators.py:16: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn(

```
Out[18]:
```

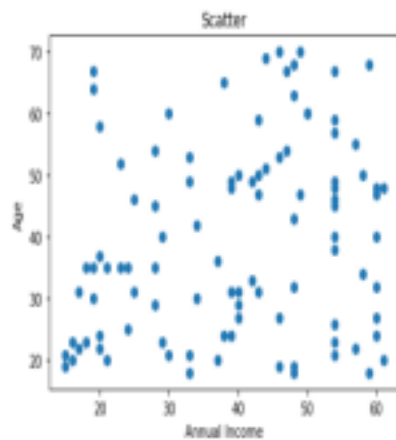


Bi - Variate Analysis

1. Scatterplots

```
In [17]: plt.scatter(x=df['Annual Income (k$)'].head(100), y=df.Age.head(100))  
plt.title('Scatter')  
plt.xlabel('Annual Income')  
plt.ylabel('Age')
```

```
Out[17]: Text(0, 0.5, 'Age')
```



2. Correlation Coefficients

```
In [19]: df.corr()
```

```
Out[19]:
```

	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)
CustomerID	1.000000	-0.026763	0.977548	0.013835
Age	-0.026763	1.000000	-0.012398	-0.327227
Annual Income (k\$)	0.977548	-0.012398	1.000000	0.009903
Spending Score (1-100)	0.013835	-0.327227	0.009903	1.000000

```
In [19]: y = df['Annual Income (k$)']  
x = df['Spending Score (1-100)']  
x = sm.add_constant(x)  
model = sm.OLS(y,x).fit()  
model.summary()
```

```
Out[19]: OLS Regression Results
```

Dep. Variable:	Annual Income (k\$)	R-squared:	0.998				
Model:	OLS	Adj. R-squared:	-0.005				
Method:	Least Squares	F-statistic:	0.01942				
Date:	Sat, 29 Oct 2022	Prob (F-statistic):	0.889				
Time:	10:45:38	Log-Likelihood:	-936.92				
No. Observations:	200	AIC:	1878				
DF Residuals:	198	BIC:	1894				
DF Model:	1						
Covariance Type:	nonrobust						
	coef	std err	t	P> t	[0.025	0.975]	
	const	60.0544	4.078	14.726	0.000	52.012	68.097
Spending Score (1-100)	0.0101	0.072	0.139	0.889	-0.132	0.153	
Omnibus:	3.518	Durbin-Watson:	0.005				
Prob(Omnibus):	0.173	Jarque-Bera (JB):	3.531				
Skew:	0.319	Prob(JB):	0.171				
Kurtosis:	2.875	Cond. No.	124.				

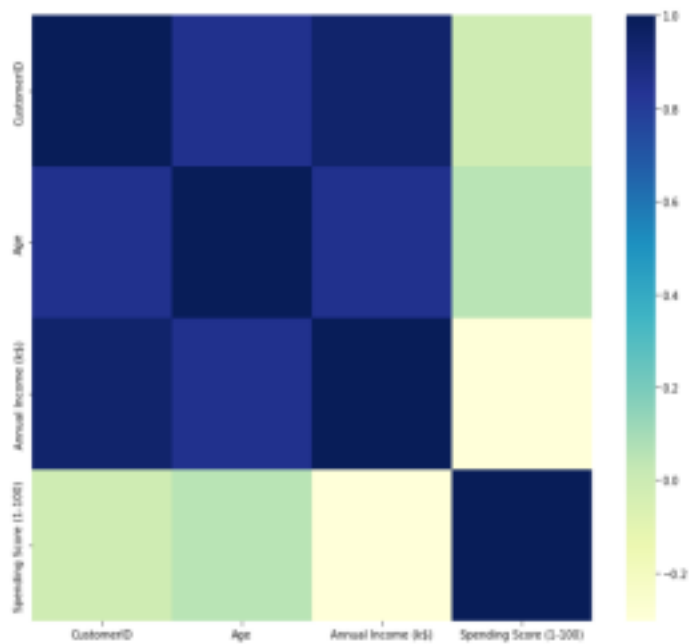
Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Multi - Variate Analysis

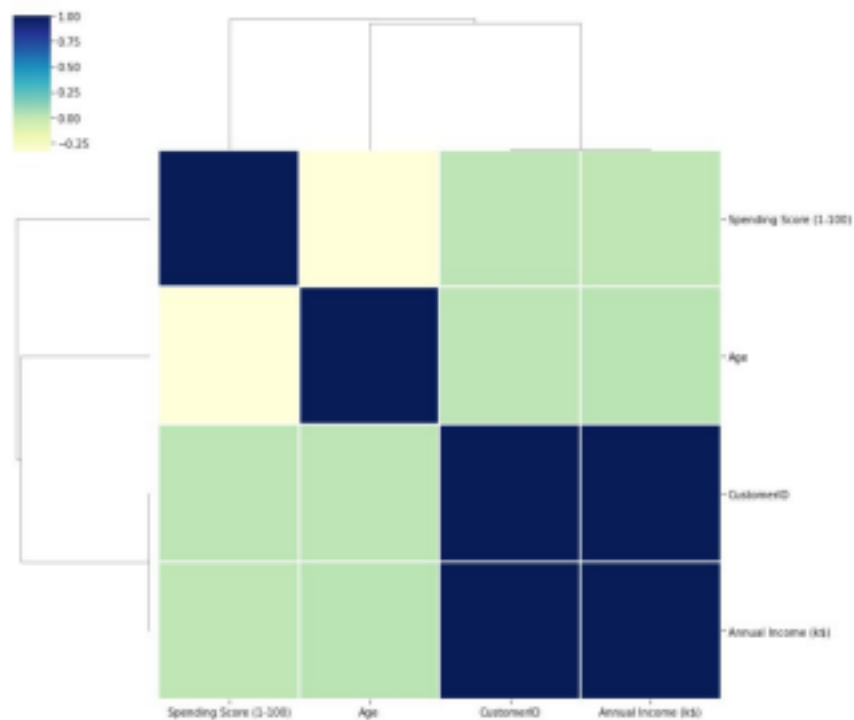
```
In [20]: f = plt.subplots(figsize=(12,10))
sns.heatmap(df.head().corr(), cmap="YlGnBu")
```

Out[20]:



```
In [21]: corrmat = df.corr(method='spearman')
cg = sns.clustermap(corrmat, cmap="YlGnBu", linewidths=0.1)
```

Out[21]:



4. Perform descriptive statistics on the dataset.

```
In [22]: df.shape
```

```
Out[22]: (200, 5)
```

```
In [23]: df.info()
```

```
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 [24]: df.describe()
```

```
Out[24]:
```

	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)
count	200.000000	200.000000	200.000000	200.000000
mean	100.500000	36.850000	60.560000	50.200000
std	57.679105	13.969007	26.264721	25.821522
min	1.000000	18.000000	15.000000	1.000000
25%	50.750000	26.750000	41.500000	34.750000
50%	100.500000	36.000000	61.500000	50.000000
75%	150.250000	46.000000	76.000000	73.000000
max	200.000000	70.000000	137.000000	99.000000

```
In [25]: df.head()
```

```
Out[25]:
```

	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 [27]: df.tail()
```

```
Out[27]:
```

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
195	196	Female	35	120	79
196	197	Female	45	126	28
197	198	Male	32	126	74
198	199	Male	32	137	18
199	200	Male	30	137	83

```
In [28]: df["Annual Income (k$)"].mean()
```

```
Out[28]: 60.56
```

```
In [29]: df["Annual Income (k$)"].median()
```

```
Out[29]: 61.5
```

```
In [30]: df["Annual Income (k$)"].mode()
```

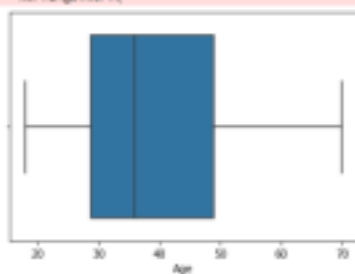
```
Out[30]: 0    54  
        1    76  
        Name: Annual Income (k$), dtype: int64
```

```
In [31]: df["Annual Income (k$)"].var()
```

```
Out[31]: 689.8355778894478
```

```
In [32]: sns.boxplot(df["Age"])
import warnings
warnings.filterwarnings('ignore')
```

C:\Users\sanda\anaconda3\lib\site-packages\seaborn_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be "data", and passing other arguments without an explicit keyword will result in an error or misinterpretation.



5. Handle the Missing values.

```
In [33]: print(df.isnull())
```

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	False	False	False	False	False
1	False	False	False	False	False
2	False	False	False	False	False
3	False	False	False	False	False
4	False	False	False	False	False
...
195	False	False	False	False	False
196	False	False	False	False	False
197	False	False	False	False	False
198	False	False	False	False	False
199	False	False	False	False	False

[200 rows x 5 columns]

```
In [34]: print(df.isnull().sum())
```

CustomerID	0
Gender	0
Age	0
Annual Income (k\$)	0
Spending Score (1-100)	0

dtype: int64

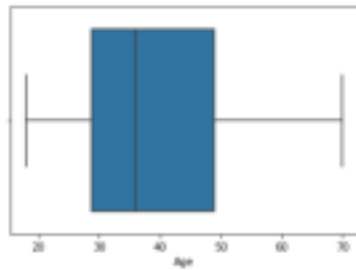
```
In [35]: df.isna().any()
```

```
Out[35]: CustomerID      False
Gender      False
Age         False
Annual Income (k$)      False
Spending Score (1-100)  False
dtype: bool
```

6. Find the outliers and replace the outliers

```
In [36]: x = ses.boxplot(df["Age"])  
x
```

Out[36]:



```
In [37]: ses.boxplot(df["Age"])
```

Out[37]:

