

Assignment 4

Customer Segmentation Analysis

Problem Statement:-

You own the mall and want to understand the customers who can quickly converge [Target Customers] so that the insight can be given to the marketing team and plan the strategy accordingly.

Clustering the data and performing classification algorithms

1. Perform Below Visualizations.

Univariate Analysis

1. Summary Statistics

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

In [238]:

```
file_data = pd.read_csv('H:\Kavin\Mall_Customers.csv')
file_data
```

In [239]:

Out[239]:

	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
...

	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

200 rows × 5 columns

```
file_data['Age'].mean()
```

In [240]:

```
38.85
```

Out[240]:

```
file_data['Age'].median()
```

In [241]:

```
36.0
```

Out[241]:

```
file_data['Age'].std()
```

In [242]:

```
13.969007331558883
```

Out[242]:

2. Frequency Table

```
file_data['Annual Income (k$)'].value_counts()
```

In [243]:

```
54      12
```

```
78      12
```

```
48       6
```

```
71       6
```

```
63       6
```

```
..
```

```
58       2
```

```
59       2
```

```
16       2
```

```
64       2
```

```
137      2
```

```
Name: Annual Income (k$), Length: 64, dtype: int64
```

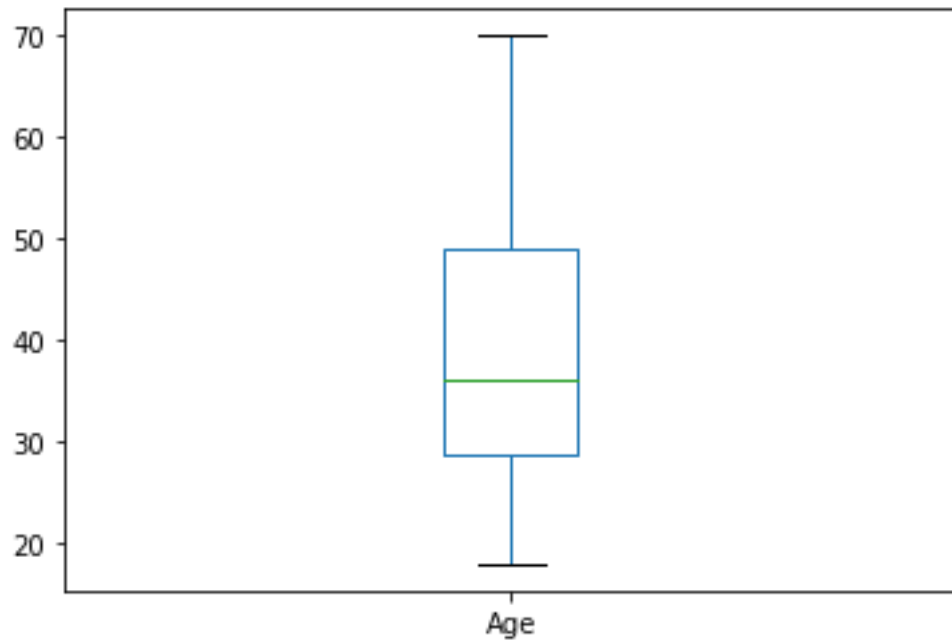
Out[243]:

3. Create Charts

In [244]:

```
file_data.boxplot(column=['Age'], grid=False)
```

Out[244]:

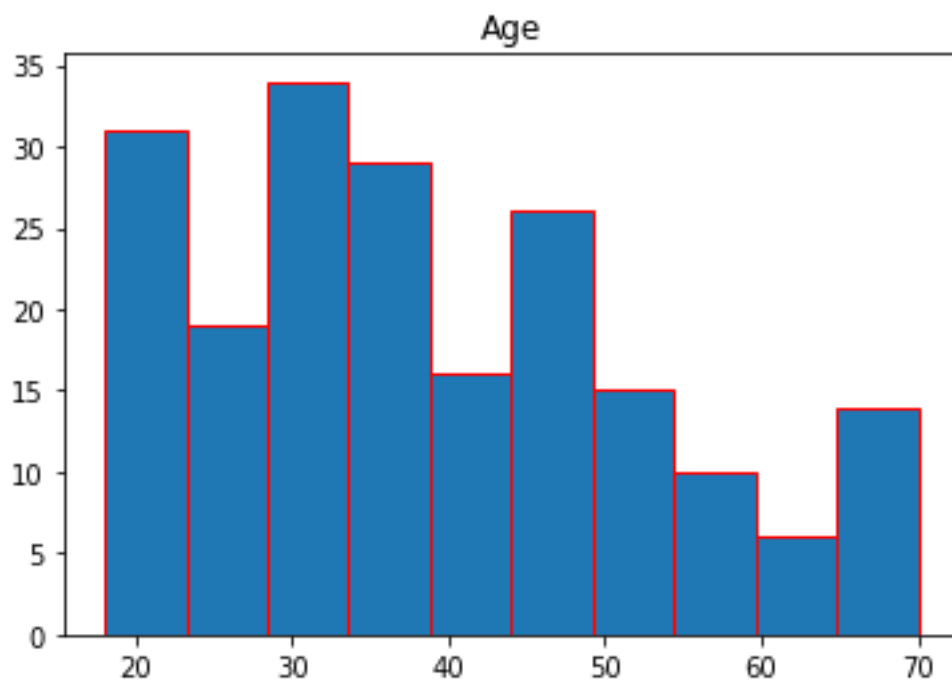


```
file_data.hist(column='Age', grid=False, edgecolor='Red')
```

In [245]:

```
array([[[]], dtype=object)
```

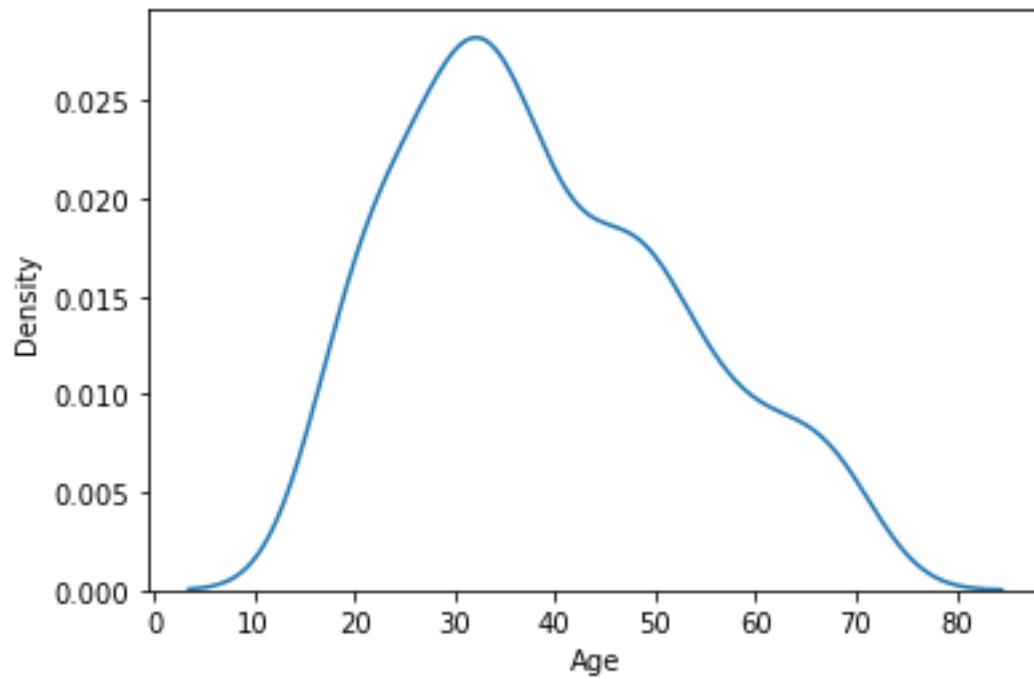
Out[245]:



```
sns.kdeplot(file_data['Age'])
```

In [246]:

Out[246]:



Bi - Variate Analysis

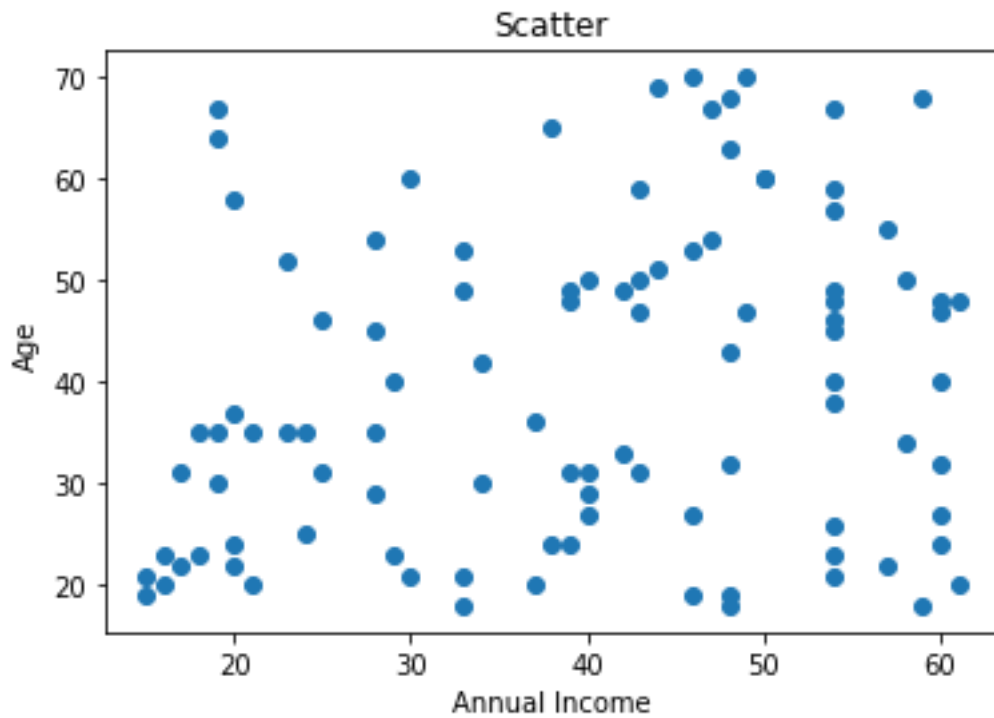
1. Scatterplots

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

In [247]:

```
Text(0, 0.5, 'Age')
```

Out[247]:



2. Correlation Coefficients

```
file_data.corr()
```

In [248]:

Out[248]:

	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

3. Simple Linear Regression

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

In [249]:

Out[249]:

OLS Regression Results

Dep. Variable:	Annual Income (k\$)	R-squared:	0.000
Model:	OLS	Adj. R-squared:	-0.005
Method:	Least Squares	F-statistic:	0.01942
Date:	Thu, 27 Oct 2022	Prob (F-statistic):	0.889
Time:	20:49:22	Log-Likelihood:	-936.92
No. Observations:	200	AIC:	1878.
Df Residuals:	198	BIC:	1884.
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:	3.510	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.

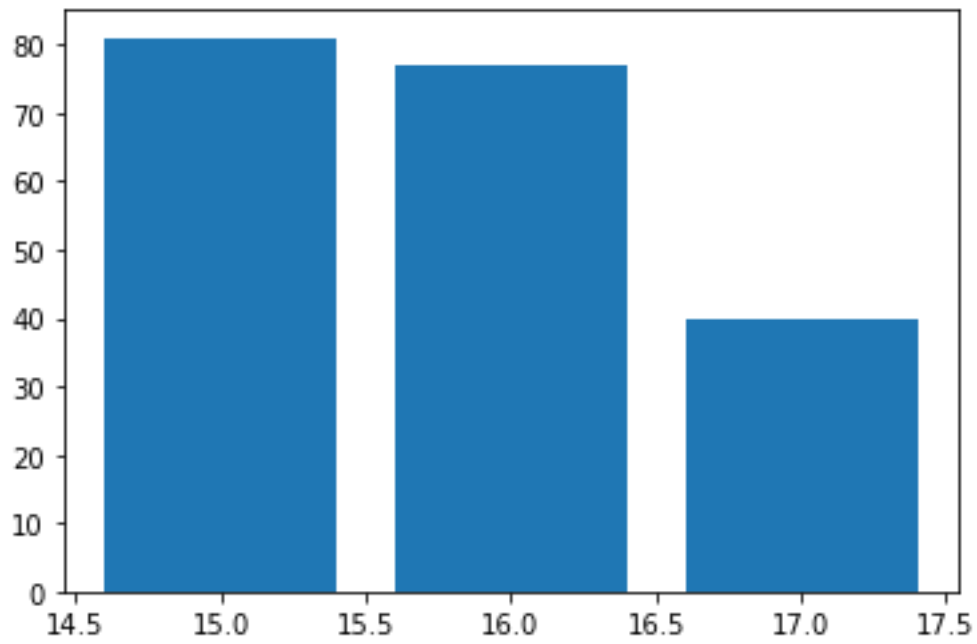
Notes:

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

```
plt.bar(file_data['Annual Income (k$)'].head(), file_data['Spending Score (1-100)'].head(), )
```

In [250]:

Out[250]:

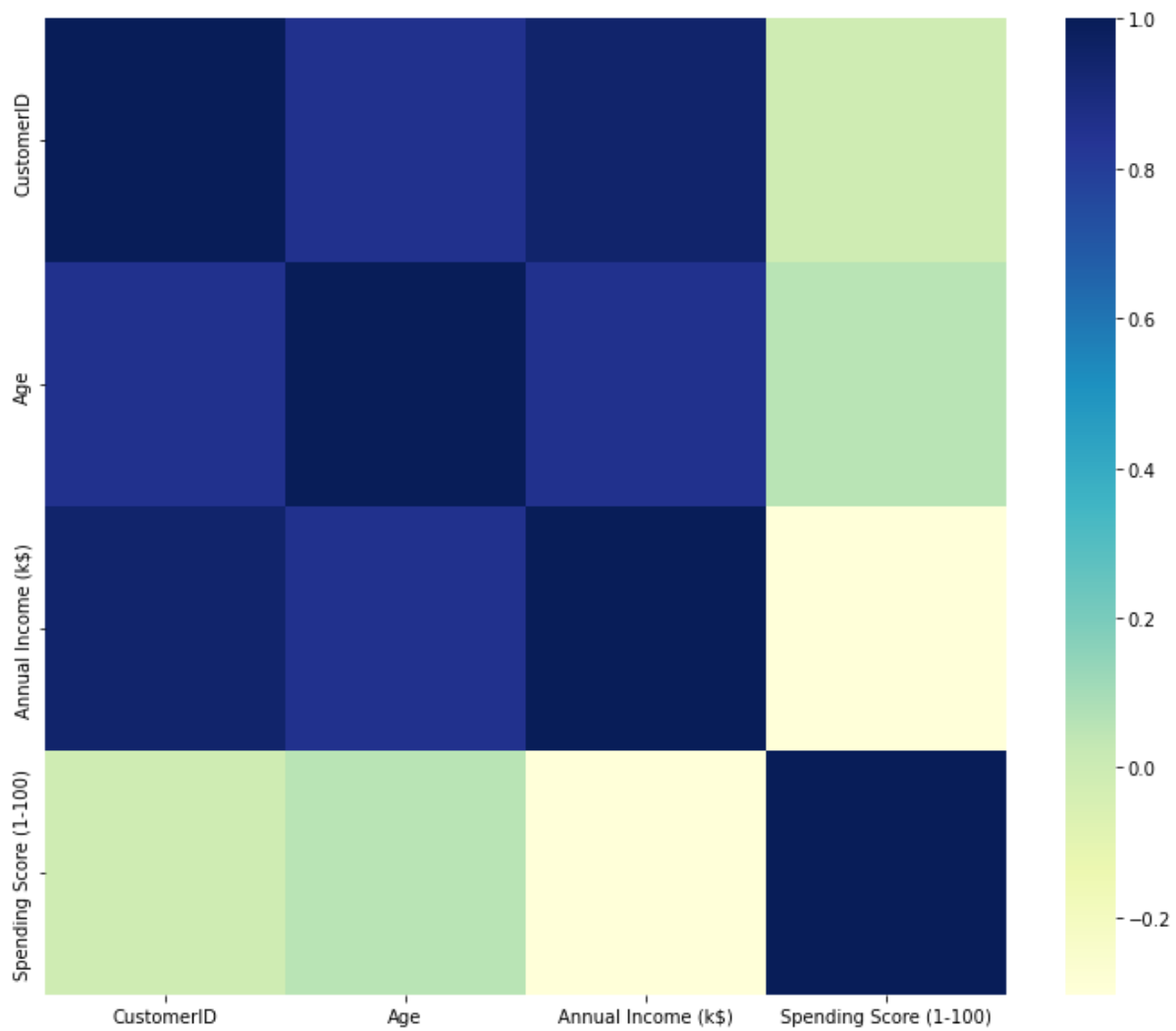


Multi - Variate Analysis

```
f = plt.subplots(figsize=(12,10))  
sns.heatmap(file_data.head().corr(), cmap="YlGnBu")
```

In [251]:

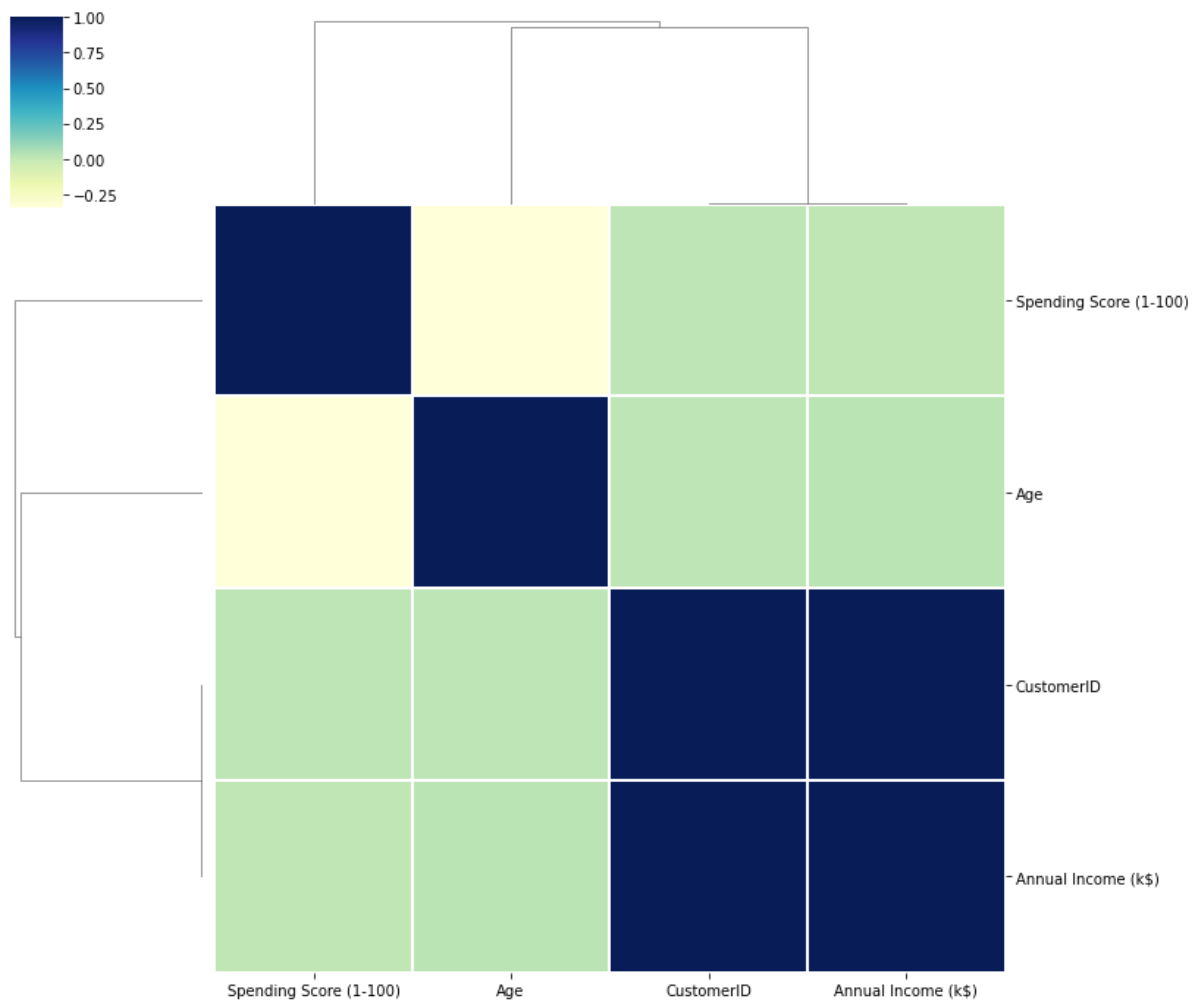
Out[251]:



In [252]:

```
corrmat = file_data.corr(method='spearman')
cg = sns.clustermap(corrmat, cmap="YlGnBu", linewidths=0.1);
plt.setp(cg.ax_heatmap.yaxis.get_majorticklabels(), rotation=0)
cg
```

Out[252]:



4. Perform descriptive statistics on the dataset.

In [253]:
`file_data.shape`

(200, 5)
 Out[253]:

In [254]:

```
file_data.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 [255]:

`file_data.describe()`
 Out[255]:

	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 [256]:

```
file_data.head()
```

Out[256]:

	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 [257]:

```
file_data.tail()
```

Out[257]:

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
195	196	Female	35	120	79

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
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 [258]:
file_data["Annual Income (k$)"].mean()

Out[258]:
60.56

In [259]:
file_data["Annual Income (k$)"].median()

Out[259]:
61.5

In [260]:
file_data["Annual Income (k$)"].mode()

Out[260]:
0    54
1    78
Name: Annual Income (k$), dtype: int64

In [261]:
file_data["Annual Income (k$)"].var()

Out[261]:
689.8355778894478

In [262]:
file_data["Annual Income (k$)"].std()

Out[262]:
26.264721165271254

In [263]:
file_data["Annual Income (k$)"].skew()

Out[263]:
0.3218425498619055

In [264]:
file_data["Annual Income (k$)"].kurt()

Out[264]:
-0.09848708652696203

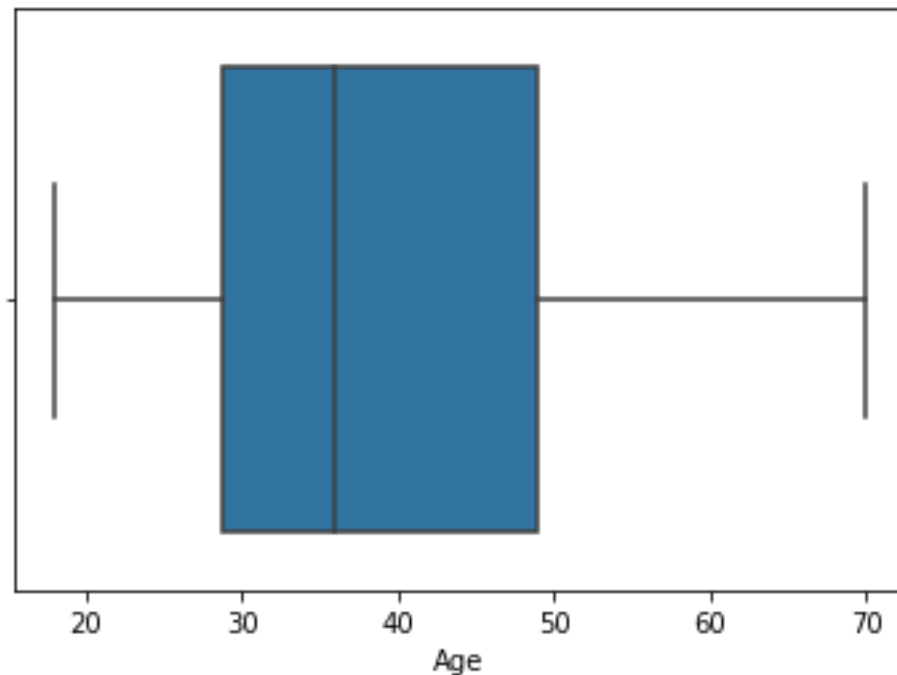
In [265]:
quantile = file_data['Age'].quantile(q=[0.75, 0.25])
quantile

Out[265]:
0.75    49.00
0.25    28.75
Name: Age, dtype: float64

```

In [266]:

```
sns.boxplot(file_data["Age"])
import warnings
warnings.filterwarnings('ignore')
```



5. Handle the Missing values.

In [267]:

```
print(file_data.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 [268]:

```
print(file_data.isnull().sum())
```

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

In [269]:

```
file_data.isna().any()
```

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

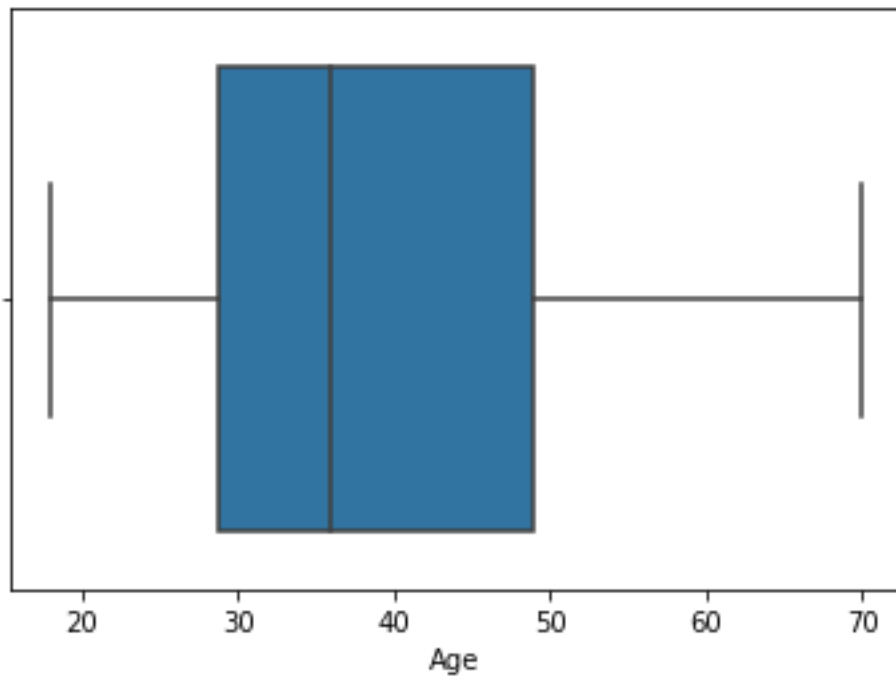
Out[269]:

6. Find the outliers and replace the outliers

```
x = sns.boxplot(x=file_data["Age"])
x
```

In [270]:

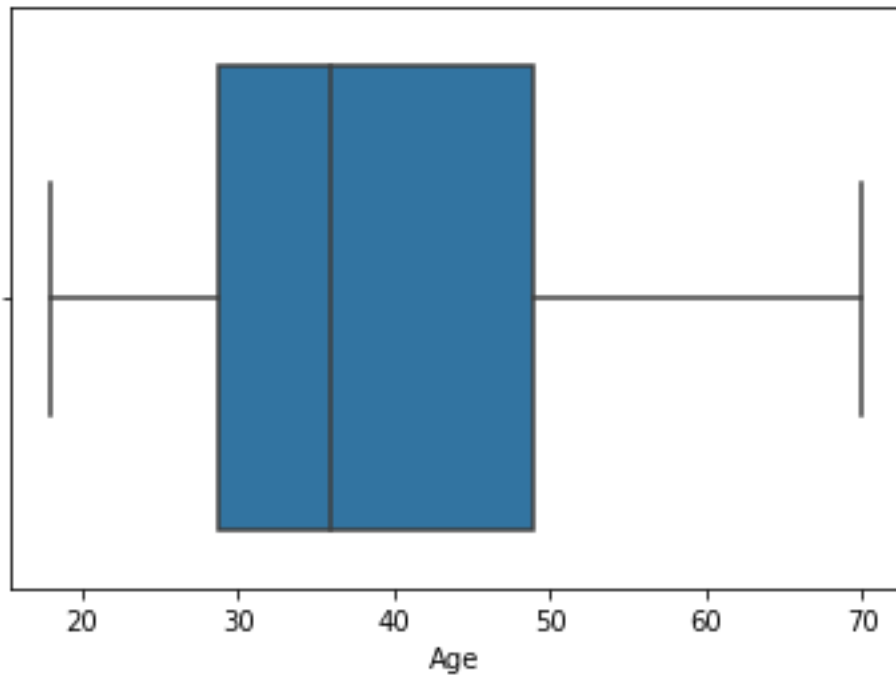
Out[270]:



```
sns.boxplot(file_data['Age'])
```

In [271]:

Out[271]:



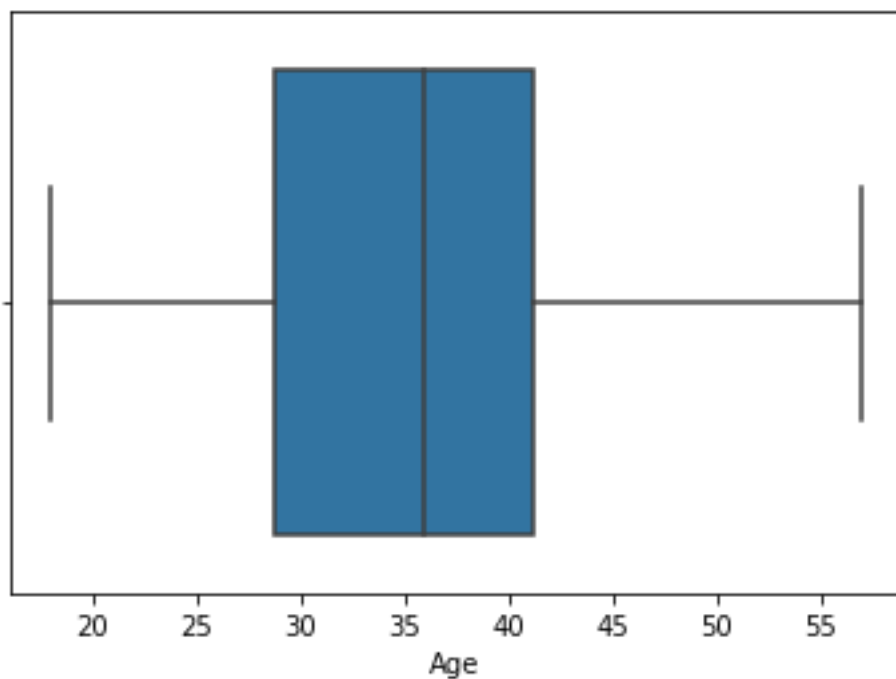
In [272]:

```
file_data['Age']=np.where(file_data['Age']>57,39, file_data['Age'])
```

In [273]:

```
sns.boxplot(file_data['Age'])
```

Out[273]:



7. Check for Categorical columns and perform encoding.

In [274]:

```
pd.Categorical(file_data["Annual Income (k$)"])
```

```
[15, 15, 16, 16, 17, ..., 120, 126, 126, 137, 137]
Length: 200
Categories (64, int64): [15, 16, 17, 18, ..., 113, 120, 126, 137]
```

Out[274]:

```
# One Hot Encoding
```

In [275]:

```
pd.get_dummies(file_data["Annual Income (k$)"]).head(10)
```

Out[275]:

	15	16	17	18	19	20	21	23	24	25	..	93	97	98	99	101	103	113	120	126	137
0	1	0	0	0	0	0	0	0	0	0	..	0	0	0	0	0	0	0	0	0	0
1	1	0	0	0	0	0	0	0	0	0	..	0	0	0	0	0	0	0	0	0	0
2	0	1	0	0	0	0	0	0	0	0	..	0	0	0	0	0	0	0	0	0	0
3	0	1	0	0	0	0	0	0	0	0	..	0	0	0	0	0	0	0	0	0	0
4	0	0	1	0	0	0	0	0	0	0	..	0	0	0	0	0	0	0	0	0	0
5	0	0	1	0	0	0	0	0	0	0	..	0	0	0	0	0	0	0	0	0	0
6	0	0	0	1	0	0	0	0	0	0	..	0	0	0	0	0	0	0	0	0	0
7	0	0	0	1	0	0	0	0	0	0	..	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	1	0	0	0	0	0	..	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	1	0	0	0	0	0	..	0	0	0	0	0	0	0	0	0	0

10 rows × 64 columns

```
pd.get_dummies(file_data).head(10)
```

In [276]:

Out[276]:

	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)	Gender_Female	Gender_Male
0	1	19	15	39	0	1
1	2	21	15	81	0	1
2	3	20	16	6	1	0
3	4	23	16	77	1	0
4	5	31	17	40	1	0
5	6	22	17	76	1	0
6	7	35	18	6	1	0
7	8	23	18	94	1	0
8	9	39	19	3	0	1
9	10	30	19	72	1	0

8. Scaling the data

```
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
```

In [277]:

```
label = LabelEncoder()
label = label.fit_transform(file_data['Gender'])
file_data["Gender"] = label
file_data['Gender'].value_counts()
```

In [278]:

```
X = file_data.drop("Age",axis=1)
Y = file_data['Age']
```

In [279]:

```
object1 = StandardScaler()
scale = object1.fit_transform(X)
scale
```

Out[279]:

```
array([[ -1.7234121,   1.12815215, -1.73899919, -0.43480148],
       [-1.70609137,   1.12815215, -1.73899919,   1.19570407],
       [-1.68877065, -0.88640526, -1.70082976, -1.71591298],
```


[-1.67144992, -0.88640526, -1.70082976, 1.04041783],
[-1.6541292, -0.88640526, -1.66266033, -0.39597992],
[-1.63680847, -0.88640526, -1.66266033, 1.00159627],
[-1.61948775, -0.88640526, -1.62449091, -1.71591298],
[-1.60216702, -0.88640526, -1.62449091, 1.70038436],
[-1.5848463, 1.12815215, -1.58632148, -1.83237767],
[-1.56752558, -0.88640526, -1.58632148, 0.84631002],
[-1.55020485, 1.12815215, -1.58632148, -1.4053405],
[-1.53288413, -0.88640526, -1.58632148, 1.89449216],
[-1.5155634, -0.88640526, -1.54815205, -1.36651894],
[-1.49824268, -0.88640526, -1.54815205, 1.04041783],
[-1.48092195, 1.12815215, -1.54815205, -1.44416206],
[-1.46360123, 1.12815215, -1.54815205, 1.11806095],
[-1.4462805, -0.88640526, -1.50998262, -0.59008772],
[-1.42895978, 1.12815215, -1.50998262, 0.61338066],
[-1.41163905, 1.12815215, -1.43364376, -0.82301709],
[-1.39431833, -0.88640526, -1.43364376, 1.8556706],
[-1.3769976, 1.12815215, -1.39547433, -0.59008772],
[-1.35967688, 1.12815215, -1.39547433, 0.88513158],
[-1.34235616, -0.88640526, -1.3573049, -1.75473454],
[-1.32503543, 1.12815215, -1.3573049, 0.88513158],
[-1.30771471, -0.88640526, -1.24279661, -1.4053405],
[-1.29039398, 1.12815215, -1.24279661, 1.23452563],
[-1.27307326, -0.88640526, -1.24279661, -0.7065524],
[-1.25575253, 1.12815215, -1.24279661, 0.41927286],
[-1.23843181, -0.88640526, -1.20462718, -0.74537397],
[-1.22111108, -0.88640526, -1.20462718, 1.42863343],
[-1.20379036, 1.12815215, -1.16645776, -1.7935561],
[-1.18646963, -0.88640526, -1.16645776, 0.88513158],
[-1.16914891, 1.12815215, -1.05194947, -1.7935561],
[-1.15182818, 1.12815215, -1.05194947, 1.62274124],
[-1.13450746, -0.88640526, -1.05194947, -1.4053405],
[-1.11718674, -0.88640526, -1.05194947, 1.19570407],
[-1.09986601, -0.88640526, -1.01378004, -1.28887582],
[-1.08254529, -0.88640526, -1.01378004, 0.88513158],
[-1.06522456, -0.88640526, -0.89927175, -0.93948177],
[-1.04790384, -0.88640526, -0.89927175, 0.96277471],
[-1.03058311, -0.88640526, -0.86110232, -0.59008772],
[-1.01326239, 1.12815215, -0.86110232, 1.62274124],
[-0.99594166, 1.12815215, -0.82293289, -0.55126616],
[-0.97862094, -0.88640526, -0.82293289, 0.41927286],
[-0.96130021, -0.88640526, -0.82293289, -0.86183865],
[-0.94397949, -0.88640526, -0.82293289, 0.5745591],
[-0.92665877, -0.88640526, -0.78476346, 0.18634349],
[-0.90933804, -0.88640526, -0.78476346, -0.12422899],
[-0.89201732, -0.88640526, -0.78476346, -0.3183368],
[-0.87469659, -0.88640526, -0.78476346, -0.3183368],
[-0.85737587, -0.88640526, -0.70842461, 0.06987881],
[-0.84005514, 1.12815215, -0.70842461, 0.38045129],
[-0.82273442, -0.88640526, -0.67025518, 0.14752193],
[-0.80541369, 1.12815215, -0.67025518, 0.38045129],
[-0.78809297, -0.88640526, -0.67025518, -0.20187212],
[-0.77077224, 1.12815215, -0.67025518, -0.35715836],
[-0.75345152, -0.88640526, -0.63208575, -0.00776431],
[-0.73613079, 1.12815215, -0.63208575, -0.16305055],
[-0.71881007, -0.88640526, -0.55574689, 0.03105725],
[-0.70148935, 1.12815215, -0.55574689, -0.16305055],

[-0.68416862, 1.12815215, -0.55574689, 0.22516505],
[-0.6668479 , 1.12815215, -0.55574689, 0.18634349],
[-0.64952717, -0.88640526, -0.51757746, 0.06987881],
[-0.63220645, -0.88640526, -0.51757746, 0.34162973],
[-0.61488572, 1.12815215, -0.47940803, 0.03105725],
[-0.597565 , 1.12815215, -0.47940803, 0.34162973],
[-0.58024427, -0.88640526, -0.47940803, -0.00776431],
[-0.56292355, -0.88640526, -0.47940803, -0.08540743],
[-0.54560282, 1.12815215, -0.47940803, 0.34162973],
[-0.5282821 , -0.88640526, -0.47940803, -0.12422899],
[-0.51096138, 1.12815215, -0.4412386 , 0.18634349],
[-0.49364065, -0.88640526, -0.4412386 , -0.3183368],
[-0.47631993, -0.88640526, -0.40306917, -0.04658587],
[-0.4589992 , -0.88640526, -0.40306917, 0.22516505],
[-0.44167848, 1.12815215, -0.25039146, -0.12422899],
[-0.42435775, 1.12815215, -0.25039146, 0.14752193],
[-0.40703703, -0.88640526, -0.25039146, 0.10870037],
[-0.3897163 , 1.12815215, -0.25039146, -0.08540743],
[-0.37239558, -0.88640526, -0.25039146, 0.06987881],
[-0.35507485, -0.88640526, -0.25039146, -0.3183368],
[-0.33775413, 1.12815215, -0.25039146, 0.03105725],
[-0.3204334 , 1.12815215, -0.25039146, 0.18634349],
[-0.30311268, 1.12815215, -0.25039146, -0.35715836],
[-0.28579196, -0.88640526, -0.25039146, -0.24069368],
[-0.26847123, -0.88640526, -0.25039146, 0.26398661],
[-0.25115051, 1.12815215, -0.25039146, -0.16305055],
[-0.23382978, -0.88640526, -0.13588317, 0.30280817],
[-0.21650906, -0.88640526, -0.13588317, 0.18634349],
[-0.19918833, -0.88640526, -0.09771374, 0.38045129],
[-0.18186761, -0.88640526, -0.09771374, -0.16305055],
[-0.16454688, -0.88640526, -0.05954431, 0.18634349],
[-0.14722616, 1.12815215, -0.05954431, -0.35715836],
[-0.12990543, 1.12815215, -0.02137488, -0.04658587],
[-0.11258471, -0.88640526, -0.02137488, -0.39597992],
[-0.09526399, -0.88640526, -0.02137488, -0.3183368],
[-0.07794326, 1.12815215, -0.02137488, 0.06987881],
[-0.06062254, -0.88640526, -0.02137488, -0.12422899],
[-0.04330181, -0.88640526, -0.02137488, -0.00776431],
[-0.02598109, 1.12815215, 0.01679455, -0.3183368],
[-0.00866036, 1.12815215, 0.01679455, -0.04658587],
[0.00866036, -0.88640526, 0.05496398, -0.35715836],
[0.02598109, -0.88640526, 0.05496398, -0.08540743],
[0.04330181, 1.12815215, 0.05496398, 0.34162973],
[0.06062254, 1.12815215, 0.05496398, 0.18634349],
[0.07794326, 1.12815215, 0.05496398, 0.22516505],
[0.09526399, -0.88640526, 0.05496398, -0.3183368],
[0.11258471, -0.88640526, 0.09313341, -0.00776431],
[0.12990543, 1.12815215, 0.09313341, -0.16305055],
[0.14722616, 1.12815215, 0.09313341, -0.27951524],
[0.16454688, 1.12815215, 0.09313341, -0.08540743],
[0.18186761, 1.12815215, 0.09313341, 0.06987881],
[0.19918833, -0.88640526, 0.09313341, 0.14752193],
[0.21650906, -0.88640526, 0.13130284, -0.3183368],
[0.23382978, 1.12815215, 0.13130284, -0.16305055],
[0.25115051, -0.88640526, 0.16947227, -0.08540743],
[0.26847123, -0.88640526, 0.16947227, -0.00776431],
[0.28579196, -0.88640526, 0.16947227, -0.27951524],

[0.30311268, -0.88640526, 0.16947227, 0.34162973],
[0.3204334 , -0.88640526, 0.24581112, -0.27951524],
[0.33775413, -0.88640526, 0.24581112, 0.26398661],
[0.35507485, 1.12815215, 0.24581112, 0.22516505],
[0.37239558, -0.88640526, 0.24581112, -0.39597992],
[0.3897163 , -0.88640526, 0.32214998, 0.30280817],
[0.40703703, 1.12815215, 0.32214998, 1.58391968],
[0.42435775, -0.88640526, 0.36031941, -0.82301709],
[0.44167848, -0.88640526, 0.36031941, 1.04041783],
[0.4589992 , 1.12815215, 0.39848884, -0.59008772],
[0.47631993, 1.12815215, 0.39848884, 1.73920592],
[0.49364065, 1.12815215, 0.39848884, -1.52180518],
[0.51096138, 1.12815215, 0.39848884, 0.96277471],
[0.5282821 , 1.12815215, 0.39848884, -1.5994483],
[0.54560282, 1.12815215, 0.39848884, 0.96277471],
[0.56292355, -0.88640526, 0.43665827, -0.62890928],
[0.58024427, -0.88640526, 0.43665827, 0.80748846],
[0.597565 , 1.12815215, 0.4748277 , -1.75473454],
[0.61488572, -0.88640526, 0.4748277 , 1.46745499],
[0.63220645, -0.88640526, 0.4748277 , -1.67709142],
[0.64952717, 1.12815215, 0.4748277 , 0.88513158],
[0.6668479 , 1.12815215, 0.51299713, -1.56062674],
[0.68416862, -0.88640526, 0.51299713, 0.84631002],
[0.70148935, -0.88640526, 0.55116656, -1.75473454],
[0.71881007, 1.12815215, 0.55116656, 1.6615628],
[0.73613079, -0.88640526, 0.58933599, -0.39597992],
[0.75345152, -0.88640526, 0.58933599, 1.42863343],
[0.77077224, 1.12815215, 0.62750542, -1.48298362],
[0.78809297, 1.12815215, 0.62750542, 1.81684904],
[0.80541369, 1.12815215, 0.62750542, -0.55126616],
[0.82273442, -0.88640526, 0.62750542, 0.92395314],
[0.84005514, -0.88640526, 0.66567484, -1.09476801],
[0.85737587, 1.12815215, 0.66567484, 1.54509812],
[0.87469659, 1.12815215, 0.66567484, -1.28887582],
[0.89201732, 1.12815215, 0.66567484, 1.46745499],
[0.90933804, -0.88640526, 0.66567484, -1.17241113],
[0.92665877, -0.88640526, 0.66567484, 1.00159627],
[0.94397949, -0.88640526, 0.66567484, -1.32769738],
[0.96130021, -0.88640526, 0.66567484, 1.50627656],
[0.97862094, 1.12815215, 0.66567484, -1.91002079],
[0.99594166, -0.88640526, 0.66567484, 1.07923939],
[1.01326239, 1.12815215, 0.66567484, -1.91002079],
[1.03058311, -0.88640526, 0.66567484, 0.88513158],
[1.04790384, -0.88640526, 0.70384427, -0.59008772],
[1.06522456, -0.88640526, 0.70384427, 1.27334719],
[1.08254529, 1.12815215, 0.78018313, -1.75473454],
[1.09986601, -0.88640526, 0.78018313, 1.6615628],
[1.11718674, 1.12815215, 0.93286085, -0.93948177],
[1.13450746, -0.88640526, 0.93286085, 0.96277471],
[1.15182818, 1.12815215, 0.97103028, -1.17241113],
[1.16914891, -0.88640526, 0.97103028, 1.73920592],
[1.18646963, -0.88640526, 1.00919971, -0.90066021],
[1.20379036, 1.12815215, 1.00919971, 0.49691598],
[1.22111108, 1.12815215, 1.00919971, -1.44416206],
[1.23843181, 1.12815215, 1.00919971, 0.96277471],
[1.25575253, 1.12815215, 1.00919971, -1.56062674],
[1.27307326, 1.12815215, 1.00919971, 1.62274124],

```
[ 1.29039398, -0.88640526, 1.04736914, -1.44416206],
[ 1.30771471, -0.88640526, 1.04736914, 1.38981187],
[ 1.32503543, 1.12815215, 1.04736914, -1.36651894],
[ 1.34235616, 1.12815215, 1.04736914, 0.72984534],
[ 1.35967688, 1.12815215, 1.23821628, -1.4053405 ],
[ 1.3769976 , 1.12815215, 1.23821628, 1.54509812],
[ 1.39431833, -0.88640526, 1.390894 , -0.7065524 ],
[ 1.41163905, -0.88640526, 1.390894 , 1.38981187],
[ 1.42895978, 1.12815215, 1.42906343, -1.36651894],
[ 1.4462805 , -0.88640526, 1.42906343, 1.46745499],
[ 1.46360123, -0.88640526, 1.46723286, -0.43480148],
[ 1.48092195, 1.12815215, 1.46723286, 1.81684904],
[ 1.49824268, -0.88640526, 1.54357172, -1.01712489],
[ 1.5155634 , 1.12815215, 1.54357172, 0.69102378],
[ 1.53288413, -0.88640526, 1.61991057, -1.28887582],
[ 1.55020485, -0.88640526, 1.61991057, 1.35099031],
[ 1.56752558, -0.88640526, 1.61991057, -1.05594645],
[ 1.5848463 , -0.88640526, 1.61991057, 0.72984534],
[ 1.60216702, 1.12815215, 2.00160487, -1.63826986],
[ 1.61948775, -0.88640526, 2.00160487, 1.58391968],
[ 1.63680847, -0.88640526, 2.26879087, -1.32769738],
[ 1.6541292 , -0.88640526, 2.26879087, 1.11806095],
[ 1.67144992, -0.88640526, 2.49780745, -0.86183865],
[ 1.68877065, 1.12815215, 2.49780745, 0.92395314],
[ 1.70609137, 1.12815215, 2.91767117, -1.25005425],
[ 1.7234121 , 1.12815215, 2.91767117, 1.27334719]])
```

In [280]:

```
X_scaled = pd.DataFrame(scale, columns = X.columns)
X_scaled
```

Out[280]:

	CustomerID	Gender	Annual Income (k\$)	Spending Score (1-100)
0	-1.723412	1.128152	-1.738999	-0.434801
1	-1.706091	1.128152	-1.738999	1.195704
2	-1.688771	-0.886405	-1.700830	-1.715913
3	-1.671450	-0.886405	-1.700830	1.040418
4	-1.654129	-0.886405	-1.662660	-0.395980
...
195	1.654129	-0.886405	2.268791	1.118061
196	1.671450	-0.886405	2.497807	-0.861839

	CustomerID	Gender	Annual Income (k\$)	Spending Score (1-100)
197	1.688771	1.128152	2.497807	0.923953
198	1.706091	1.128152	2.917671	-1.250054
199	1.723412	1.128152	2.917671	1.273347

200 rows × 4 columns

In [281]:

```
X_train, X_test, Y_train, Y_test = train_test_split(X_scaled, Y,
test_size=0.20, random_state=0)
```

Input In [281]

```
X_train, X_test, Y_train, Y_test = train_test_split(X_scaled, Y, test_s
ize=0.20, random_state=0)
```

^

SyntaxError: unexpected EOF while parsing

9. Perform any of the clustering algorithms

In []:

```
from sklearn.cluster import KMeans
```

In []:

```
x = file_data.iloc[:, [3, 4]].values
```

In []:

```
list= []
```

```
for i in range(1, 11):
    kmeans = KMeans(n_clusters=i, init='k-means++', random_state= 42)
    kmeans.fit(x)
    list.append(kmeans.inertia_)
plt.plot(range(1, 11), list)
plt.title('The Elbow Method Graph')
plt.xlabel('Number of clusters(k)')
plt.ylabel('wcss_list')
plt.show()
```

In []:

```
kmeans = KMeans(n_clusters=5, init='k-means++', random_state= 42)
y_predict= kmeans.fit_predict(x)
```

In []:

```
plt.scatter(x[y_predict == 0, 0], x[y_predict == 0, 1], s = 100, c =
'blue', label = 'Cluster 1') #for first cluster
plt.scatter(x[y_predict == 1, 0], x[y_predict == 1, 1], s = 100, c =
'green', label = 'Cluster 2') #for second cluster
plt.scatter(x[y_predict== 2, 0], x[y_predict == 2, 1], s = 100, c = 'red',
label = 'Cluster 3') #for third cluster
plt.scatter(x[y_predict == 3, 0], x[y_predict == 3, 1], s = 100, c =
'cyan', label = 'Cluster 4') #for fourth cluster
```

```
plt.scatter(x[y_predict == 4, 0], x[y_predict == 4, 1], s = 100, c =
'magenta', label = 'Cluster 5') #for fifth cluster
plt.scatter(kmeans.cluster_centers[:, 0], kmeans.cluster_centers[:, 1], s
= 300, c = 'yellow', label = 'Centroid')
plt.title('Clusters of customers')
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.legend()
plt.show()
```

10. Add the cluster data with the primary dataset

```
In [ ]:
file_data['Cluster']=kmeans.labels_
file_data.head()

In [ ]:
file_data.tail()
```

11. Split the data into dependent and independent variables.

```
In [ ]:
X=file_data.drop('Cluster',axis=1)
Y=file_data['Cluster']
y=file_data['Cluster']
Y

In [ ]:
from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test=train_test_split(X,Y,test_size=0.2,random_state=42)

In [ ]:
X_train.shape

In [ ]:
y_train.shape
```

12. Split the data into training and testing

```
In [ ]:
X_train

In [ ]:
X_test

In [ ]:
y_train

In [ ]:
y_test
```

13. Build the Model

```
from sklearn.linear_model import LogisticRegression
model=LogisticRegression()
model.fit(X_train,y_train)
```

In []:

14. Train the Model

```
model.score(X_train,y_train)
```

In []:

15. Test the Model

```
model.score(X_test,y_test)
```

In []:

16. Measure the performance using Evaluation Metrics.

```
from sklearn.metrics import confusion_matrix,classification_report
```

In []:

```
y_pred=model.predict(X_test)
confusion_matrix(y_test,y_pred)
```

In []:

```
print(classification_report(y_test,y_pred))
```

In []: