

# ASSIGNMENT-04

Assignment Date	2 November 2022
Student Name	Sahithya V
Student Roll Number	113219071034
Maximum Marks	2 Marks

## Problem Statement: Customer Segmentation Analysis

```
[ ]
```

▼ Download the dataset

```
import pandas as pd
import numpy as np
import seaborn as sns
from matplotlib import pyplot as plt
from sklearn.preprocessing import scale
import warnings
warnings.filterwarnings('ignore')
```

## load the dataset into the tool

```
[ ] data=pd.read_csv("drive/MyDrive/CONTENT/Mall_Customers.csv")
data.head()
```

	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

```
[ ] data.shape
```

```
(200, 5)
```

```
[ ] data.size
```

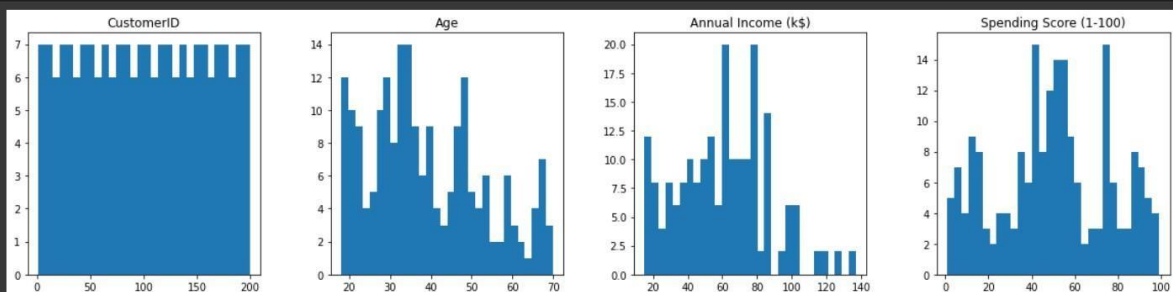
```
1000
```

```
data.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   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
```

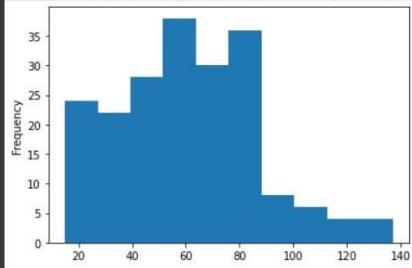
## Perform Below Visualizations

```
[ ] data.hist(figsize=(20,10), grid=False, layout=(2,4),bins=30)
plt.show()
```



```
[ ] data["Annual Income (k$)"].plot(kind='hist')
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fa78eca9290>

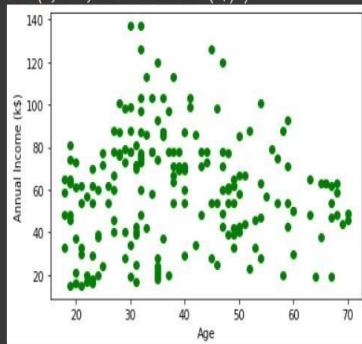


+ Code + Text

Connect Editing

```
[ ] from matplotlib import pyplot as plt
plt.scatter(data['Age'],data['Annual Income (k$)'],color='green')
plt.xlabel("Age")
plt.ylabel("Annual Income (k$)")
```

Text(0, 0.5, 'Annual Income (k\$)')

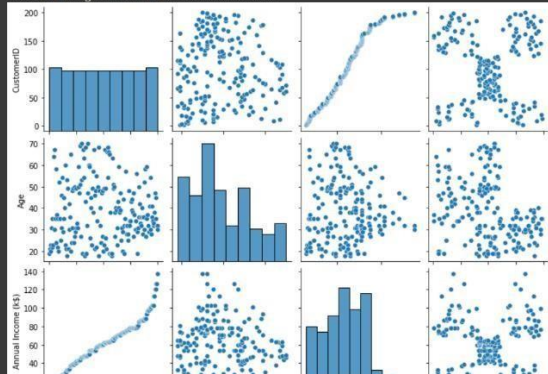


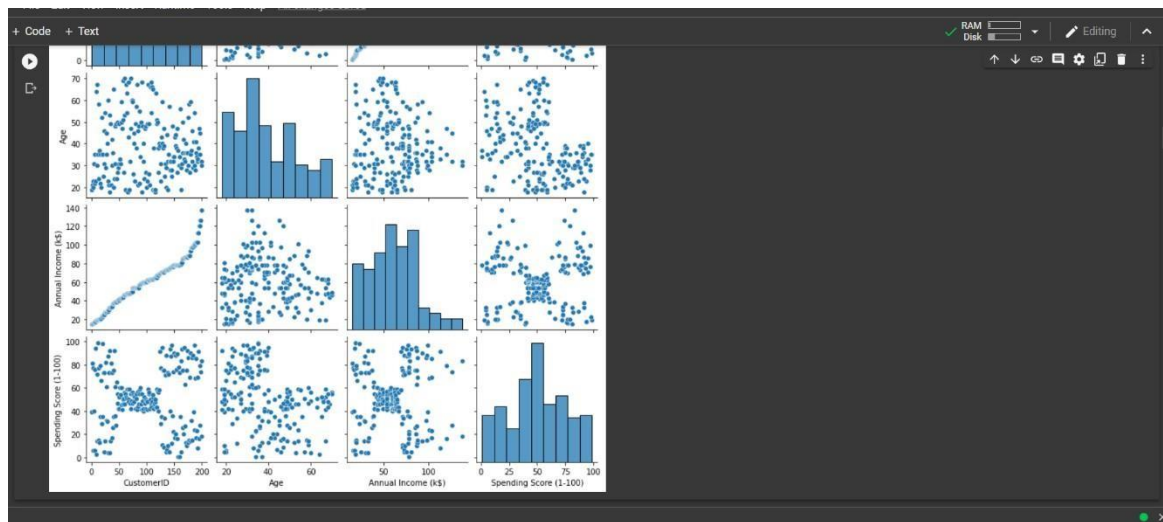
+ Code + Text

Connecting Editing

sns.pairplot(data)

<seaborn.axisgrid.PairGrid at 0x7fa78eb2c1d0>





## Perform descriptive statistics on the dataset

```
[ ] data.describe()
```

	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

+ Code + Text

✓ RAM  Disk  Editing ^

## ▼ Check for Missing values and deal with them

```
[ ] data.isna().sum()
```

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

## ▼ Find the outliers and replace them outliers.

```
[ ] data.skew()
```

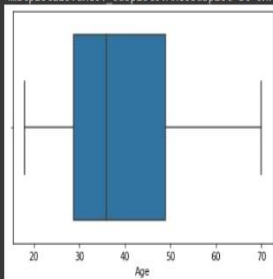
```
CustomerID      0.000000
Age              0.485569
Annual Income (k$)  0.321843
Spending Score (1-100) -0.047220
dtype: float64
```

+ Code + Text

✓ RAM  Disk  Editing ^

```
[ ] sns.boxplot(x=data['Age'],data=data)
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fa78f3fd150>

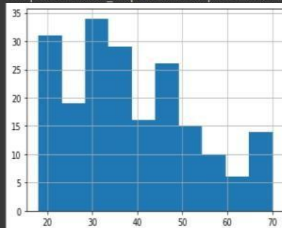


+ Code + Text

✓ RAM  Disk  Editing ^

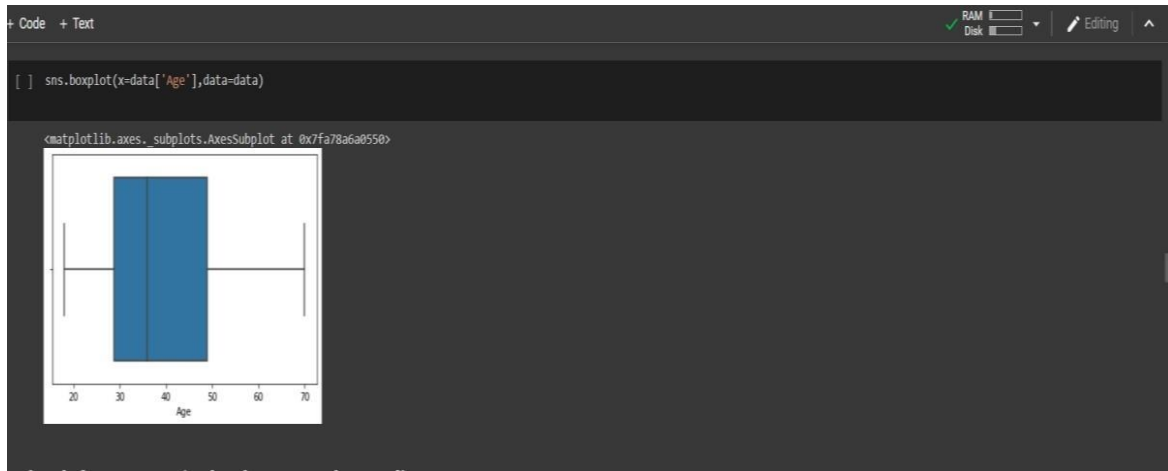
```
[ ] data['Age'].hist()
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7fa78f493b50>






```
[ ] print('skewness value of Age:',data['Age'].skew())
```

skewness value of Age: 0.4855685096681657



+ Code + Text

RAM  Disk  Editing  ^

## Check for Categorical columns and encoding




```
[ ] data.info
```

<bound method DataFrame.info of

			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	
..	...	...	...	...	...	...	...
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 x 5 columns]>

+ Code + Text

RAM  Disk  Editing  ^

```
[ ] from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
data['Gender']=le.fit_transform(data['Gender'])
data.head()
```

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	1	19	15	39
1	2	1	21	15	81
2	3	0	20	16	6
3	4	0	23	16	77
4	5	0	31	17	40

```
[ ] data["Gender"].unique()
```

array([1, 0])







+ Code + Text

RAM 16GB Disk 512GB

Editing

### Build the model

```
[ ] from sklearn.linear_model import LinearRegression
    LR = LinearRegression()
```

### Train the Model

```
[ ] LR.fit(x_train,y_train)

    LinearRegression()
```

```
+ Code + Text
[ ] pred_LR.predict(x_test)
pred

array([[ 41.79651469,  35.44897396,  32.32182941,  62.15230947,
        97.15499 , 102.74527464,  57.52984542, 18.50596884,
        28.90898195,  90.05616474,  90.63951146,  25.17877999,
        21.47607213,  56.15450717,  65.58284431,  68.81365584,
        85.74449988,  31.45756756,  76.51559556,  42.98039276,
        38.70178627,  23.89238204,  36.61738406,  57.67164216,
        29.74845621,  86.65460588,  33.53032334,  29.31235764,
        100.75984295,  28.3364555 ,  37.02836966,  88.57006476,
        101.81449573,  93.23392219,  94.16104415,  58.75918464,
        93.31570423,  49.53263005,  46.78164703,  91.618992 ,
        64.85923756,  63.89021447,  98.96847593,  22.93975353,
        41.82689378,  24.95860094,  65.82297944,  33.18229176,
        96.71878777,  70.4308092 ,  59.76768524,  70.1173078 ,
        69.1581952 ,  40.54244593,  30.19338393,  94.32293272,
        95.33656664,  64.12923371, 102.85955135,  76.19945402]])

[ ] pred.astype(int)

array([[ 41,  35,  32,  62,  97, 102,  57,  18,  28,  90,  90,  25,  21,
        56,  65,  68,  85,  31,  76,  42,  38,  23,  36,  57,  29,  86,
        33,  29, 100,  28,  27,  88, 101,  93,  94,  58,  93,  49,  46,
        91,  64,  63,  98,  22,  41,  24,  65,  33,  96,  70,  59,  70,
        69,  40,  30,  94,  95,  64, 102,  76]])
```

```
+ Code + Text
[ ] y_test

58      46
40      38
24      33
102     62
184     99
198     137
95      60
4       17
29      29
168     87
171     87
18      23
11      19
89      58
110     63
118     67
159     78
35      33
136     73
59      46
51      42
16      21
44      39
94      60
31      30
162     81
38      37
28      29
193     113
```

```
+ Code + Text
[ ]
38      37
28      29
193     113
27      28
47      40
165     85
194     120
177     88
176     88
97      60
174     88
73      50
69      48
172     87
108     63
107     63
189     103
14      20
56      44
19      23
114     65
39      37
185     99
124     70
98      61
123     69
119     67
53      43
33      33
179     93
181     97
106     63
```

## Measure the performance using Evaluation Metrics.

```
from sklearn.metrics import r2_score
score=r2_score(pred,y_test)
score
```

0.9234274149757858

