

▼ Assignment 5

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

```
df=pd.read_csv('/content/Mall_Customers.csv')
df
```

	CustomerID	Genre	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 × 5 columns

```
print(df.head())
print(df.tail())
```

	CustomerID	Genre	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

```
print(df.shape)
```

(200, 5)

```
print(df.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   Genre                 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
None
```

```
df.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

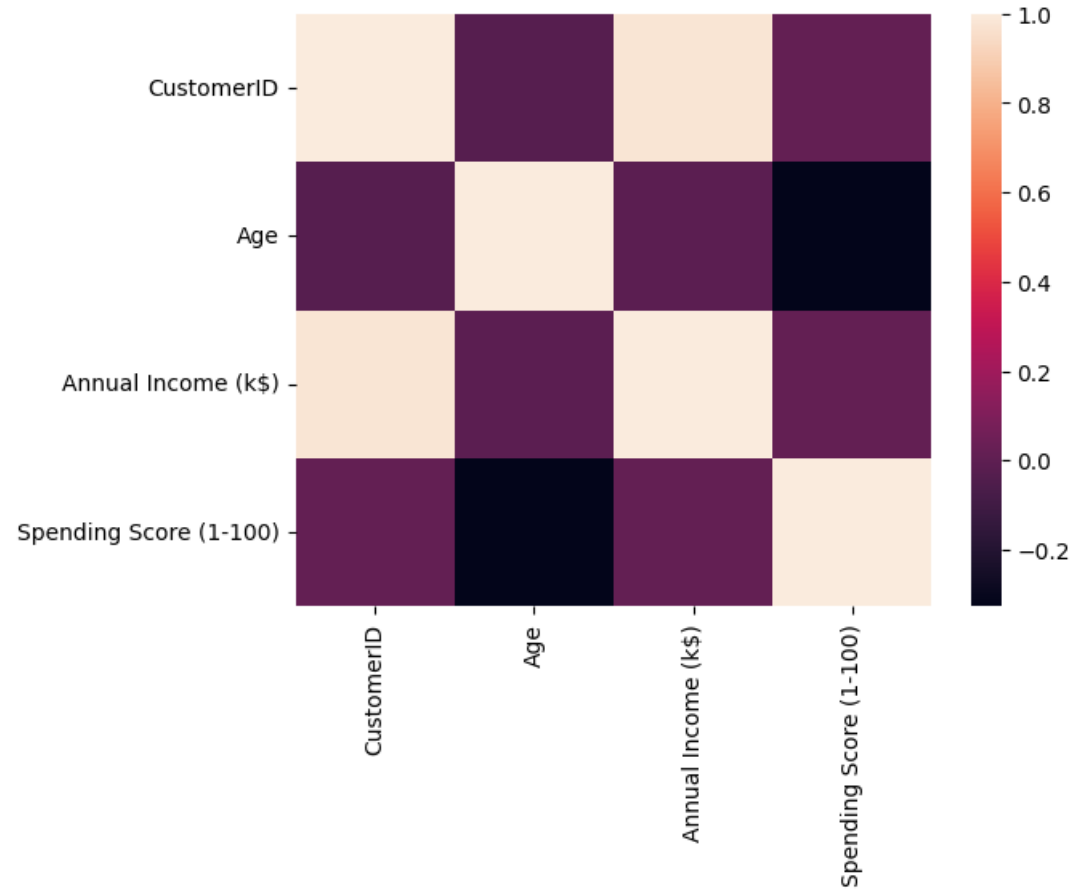
```
df.corr()
```

```
<ipython-input-9-2f6f6606aa2c>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid column
df.corr()
```

	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

```
sns.heatmap(df.corr())
```

```
<ipython-input-10-aa4f4450a243>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid column
sns.heatmap(df.corr())
<Axes: >
```



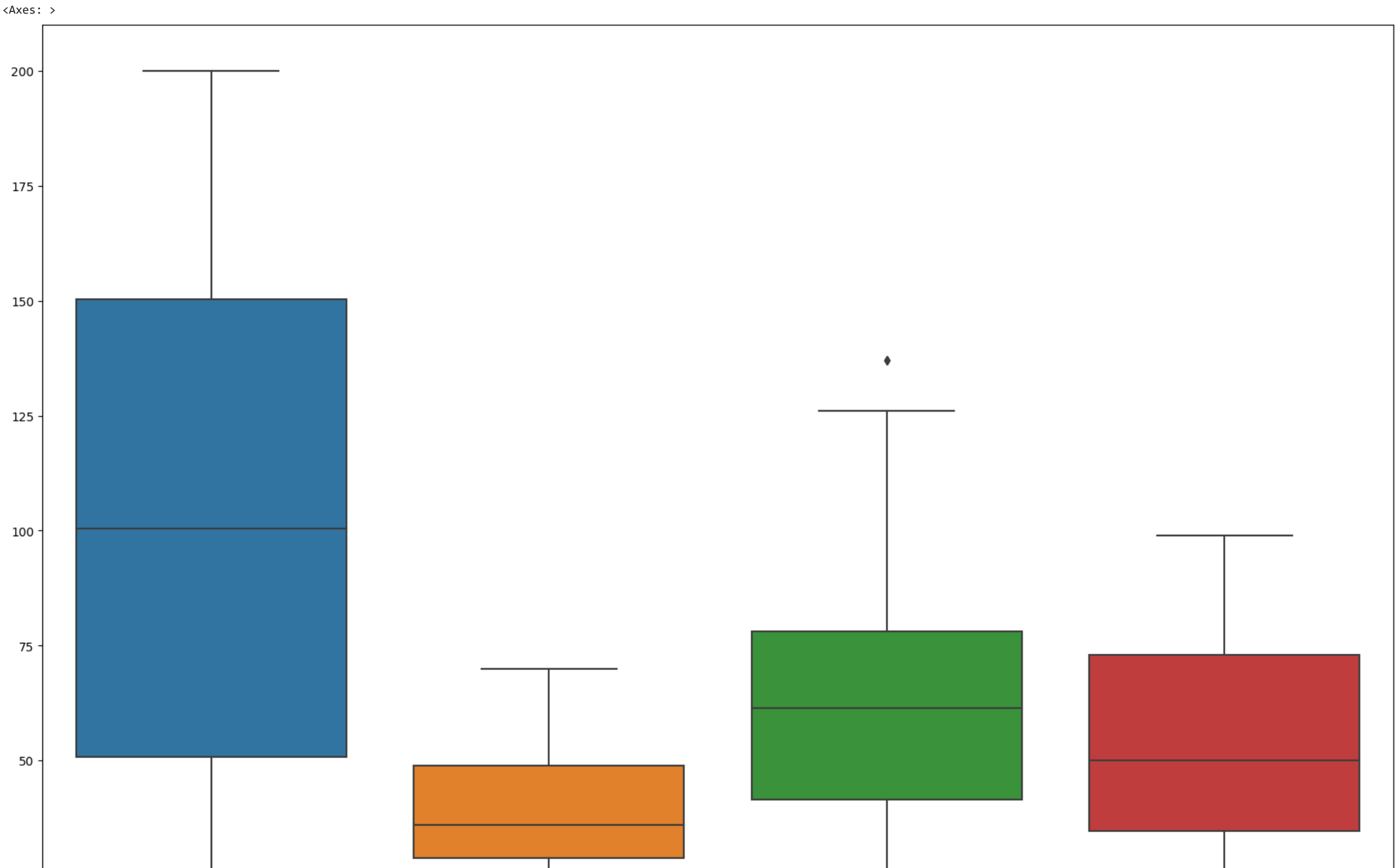
```
df.isnull().any()
```

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

```
df.isnull().sum()
```

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

```
plt.subplots(figsize=(20,15))
sns.boxplot(df)
```



```
x=df.iloc[:, :3]
y=df.iloc[:, 3:4]
```

```
x.head()
```

	CustomerID	Genre	Age
0	1	Male	19
1	2	Male	21
2	3	Female	20
3	4	Female	23
4	5	Female	31

```
y.head()
```

	Annual Income (k\$)
0	15
1	15
2	16
3	16
4	17

```
print(x.shape)
print(y.shape)
```

```
(200, 3)
(200, 1)
```

```
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
x['Genre']=le.fit_transform(x['Genre'])
```

```
x[['Genre']]
```



```
[0.11891892 0. 0.51923077]
[0.58918919 1. 0.01923077]
[0.88648649 0. 0.34615385]
[0.31891892 0. 0.69230769]
[0.02162162 1. 0.88461538]
[0.38378378 1. 0.15384615]
[0.61621622 0. 0.63461538]
[0.75135135 0. 0.26923077]
[0.36216216 0. 0.55769231]
[0.64864865 0. 0.09615385]
[0.97297297 0. 0.44230769]
[0.5027027 0. 0.17307692]
[0.78378378 1. 0.30769231]
[0.10810811 0. 0.69230769]
[0.14054054 1. 0.80769231]
[0.84334334 0. 0.33376033]
```