Chinmay Rajan Mhatre

22102B2001

CMPN B

https://github.com/chinmay0910/ML-Lab---VIT/blob/main/EXP%208/ML_LAB_8.ipynb

## ∨ Importing required libraries

```
 1 import pandas as pd
 2 import numpy as np
 3 import seaborn as sb
 4 import matplotlib.pyplot as pt
 5
 6 from sklearn.cluster import DBSCAN
 7 from sklearn.preprocessing import StandardScaler
 8 from sklearn import metrics
 9
10 import warnings
11 warnings.simplefilter("ignore")
```

## ∨ i) Exploratory Data Analysis

```
1 # Loading dataset or dataframe
2 segment=pd.read_csv("/content/Mall_Customers.csv")
```

```
1 # Looking for shape of dataframe
2 segment.shape
```

    (200, 5)

```
1 # Viewing columns
2 segment.columns
```

    Index(['CustomerID', 'Gender', 'Age', 'Annual Income (k$)',
           'Spending Score (1-100)'],
          dtype='object')

```
1 # Head of the dataframe
2 segment.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 |

Next steps:  [ Generate code with `segment` ]   [ ⬤ View recommended plots ]   [ New interactive sheet ]

```
1 # Tail of the dataframe
2
3 segment.tail()
```

|     | 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 |

```
1 # Datatypes involved..
2 segment.dtypes.value_counts()
```

|         | count |
|---------|-------|
| **int64** | 4     |
| **object** | 1     |

**dtype:** int64

```
1 # Information about Dataframe
2 segment.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
```

```
1 # Removing white spaces and remaining columns
2 segment.columns=segment.columns.str.replace(" ","")
3 segment.columns
```

```
Index(['CustomerID', 'Gender', 'Age', 'AnnualIncome(k$)',
       'SpendingScore(1-100)'],
      dtype='object')
```

```
1 # renaming columns
2 segment.columns=segment.rename(columns={'AnnualIncome(k$)':'AnnualIncome',
3                              'SpendingScore(1-100)':'SpendingScore',"Genre":"Gender"}).columns
4 segment.columns
```

```
Index(['CustomerID', 'Gender', 'Age', 'AnnualIncome', 'SpendingScore'], dtype='object')
```

```
1 # Viewing int columns
2 int_col=segment.select_dtypes(include="int64").columns.tolist()
3 int_col
```

```
['CustomerID', 'Age', 'AnnualIncome', 'SpendingScore']
```

```
1 # Viewing categorical columns
2 cat_col=segment.select_dtypes(include="O").columns.tolist()
3 cat_col
```

```
['Gender']
```

```
1 # Drop the id column
2 copy_segment=segment.copy()
3 segment.drop("CustomerID",axis=1,inplace=True)
```

```
1 # Summary statistics
2 segment.describe()
```

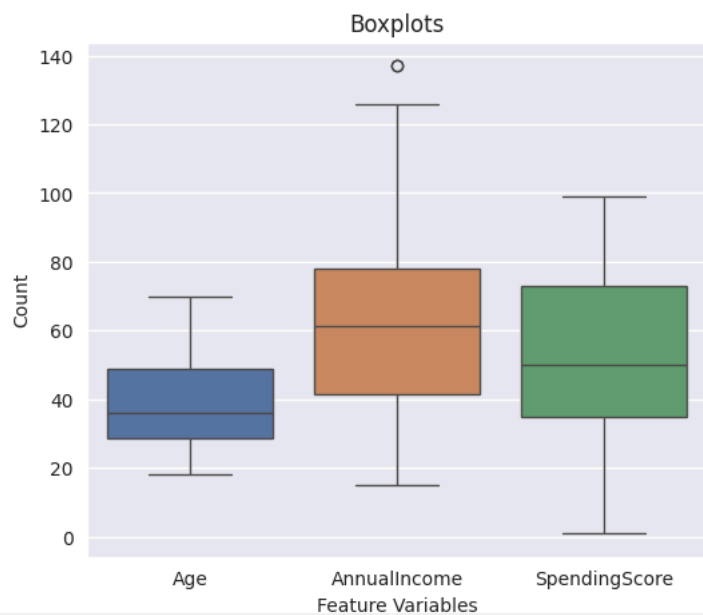|       | Age | AnnualIncome | SpendingScore |
|-------|-----|--------------|---------------|
| **count** | 200.000000 | 200.000000 | 200.000000 |
| **mean** | 38.850000 | 60.560000 | 50.200000 |
| **std** | 13.969007 | 26.264721 | 25.823522 |
| **min** | 18.000000 | 15.000000 | 1.000000 |
| **25%** | 28.750000 | 41.500000 | 34.750000 |
| **50%** | 36.000000 | 61.500000 | 50.000000 |
| **75%** | 49.000000 | 78.000000 | 73.000000 |
| **max** | 70.000000 | 137.000000 | 99.000000 |

```
1 Start coding or generate with AI.
```

```
1 # looking for outliers through boxplot
2 sb.set({"figure.figsize":(6,5)})
3 sb.boxplot(segment)
4 pt.title("Boxplots")
5 pt.xlabel("Feature Variables")
6 pt.ylabel("Count")
7
```

Text(0, 0.5, 'Count')



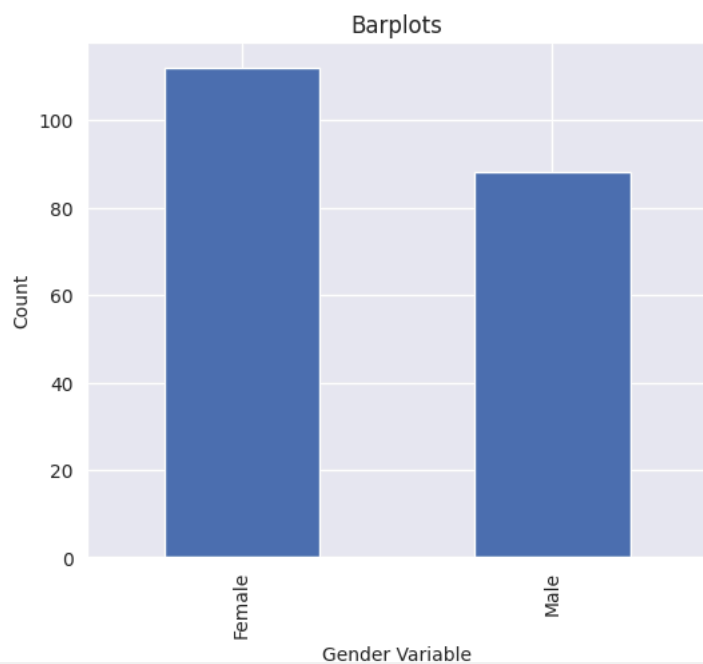since,we have one outlier in annual income,we neglect it as it may be the person with higher annual income.

```
1 # Plotting gender for count
2 segment.Gender.value_counts().plot(kind="bar")
3 pt.title("Barplots")
4 pt.xlabel("Gender Variable")
5 pt.ylabel("Count")
6
```
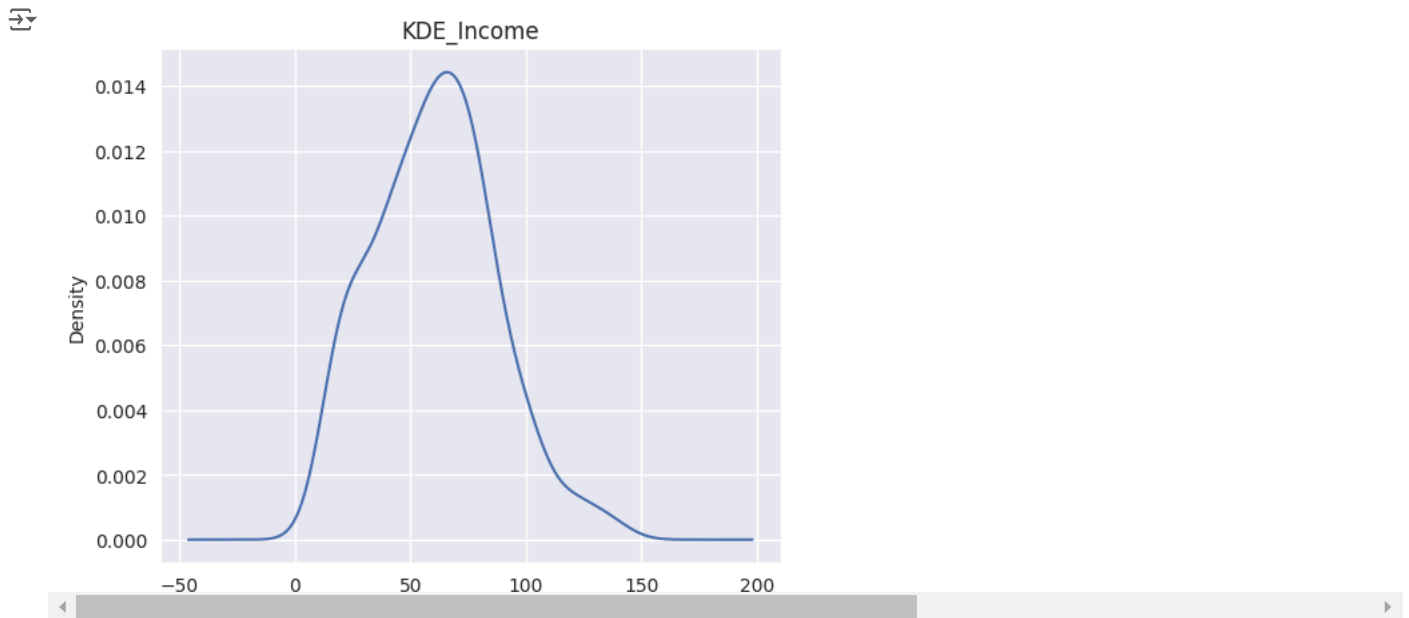
Text(0, 0.5, 'Count')



```
1 # KDE plot for Annual income
2 segment.AnnualIncome.plot(kind="kde")
3 pt.title('KDE_Income')
4 pt.show()
```

KDE_Income

```
1 # Scatter plot for the AnnualIncome & SpendingScore
2 sb.set({"figure.figsize":(10,5)})
3 sb.scatterplot(data=segment,x='AnnualIncome',y='SpendingScore',hue="Gender")
4 pt.title('Income vs Score')
5 pt.show()
```



Income vs Score

~: We can see that most of the customer data points lies at annual income(40-70) and spending score (40-60).

```
1 # Scatter plot for the Age & SpendingScore
2 sb.scatterplot(data=segment,x='Age',y='SpendingScore',hue="Gender")
3 pt.title('Score vs Age')
4 pt.show()
```

Score vs Age

~: We interpret that the age (40-60) of having spending score around (20-60),the age (20-40) of having higher spending score around (40-100) and the age (60-70) of having balanced spending score around (40-60) respectively .

## ii) Data Preprocessing

```
1 #converting gender to binary
2 segment.Gender=np.where(segment["Gender"]=="Male",1,0)
3 segment.Gender
```

|     | Gender |
| --- | --- |
| 0   | 1 |
| 1   | 1 |
| 2   | 0 |
| 3   | 0 |
| 4   | 0 |
| ... | ... |
| 195 | 0 |
| 196 | 0 |
| 197 | 1 |
| 198 | 1 |
| 199 | 1 |

200 rows × 1 columns

```
1 segment.head()
```

|     | Gender | Age | AnnualIncome | SpendingScore |
| --- | --- | --- | --- | --- |
| 0   | 1 | 19 | 15 | 39 |
| 1   | 1 | 21 | 15 | 81 |
| 2   | 0 | 20 | 16 | 6 |
| 3   | 0 | 23 | 16 | 77 |
| 4   | 0 | 31 | 17 | 40 |

Next steps:   Generate code with `segment`      ⬤ View recommended plots      New interactive sheet

```
1 # Scaling desired columns for modelling
2 scaler=StandardScaler()
3 scaled_val=scaler.fit_transform(segment[["AnnualIncome","SpendingScore"]])
```

```
1 scaled_val[:5]
```

```
array([[-1.73899919, -0.43480148],
       [-1.73899919,  1.19570407],
       [-1.70082976, -1.71591298],
       [-1.70082976,  1.04041783],
       [-1.66266033, -0.39597992]])
```

```
1 # Creating a new dataframe with out Gender variable
2 features=pd.DataFrame(scaled_val,columns=segment.columns[2:4].tolist())
3 features
```

|     | AnnualIncome | SpendingScore |
| --- | --- | --- |
| 0   | -1.738999 | -0.434801 |
| 1   | -1.738999 | 1.195704 |
| 2   | -1.700830 | -1.715913 |
| 3   | -1.700830 | 1.040418 |
| 4   | -1.662660 | -0.395980 |
| ... | ... | ... |
| 195 | 2.268791 | 1.118061 |
| 196 | 2.497807 | -0.861839 |
| 197 | 2.497807 | 0.923953 |
| 198 | 2.917671 | -1.250054 |
| 199 | 2.917671 | 1.273347 |

200 rows × 2 columns

Next steps: | Generate code with `features` | View recommended plots | New interactive sheet |

## iii) Model Building & Evaluation

~: DBSCAN groups observations into clusters of high density ,which does not make use of k-clusters.

```
1 #DBSCAN Algo...
2 # eps is maximum distance b/w datapoints , min_samples for core datapoints.
3 dbscan=DBSCAN(eps=0.5,min_samples=10,metric="euclidean")
4 dbscan.fit(features)
5 ypre=dbscan.fit_predict(features)
```

```
1 dbscan
```

```
▼        DBSCAN        ⓘ ?
DBSCAN(min_samples=10)
```

```
1 ypre
```

```
array([-1,  0,  1,  0, -1,  0,  1, -1,  1,  0,  1, -1,  1,  0,  1,  0, -1,
        0, -1, -1, -1,  0,  1,  0,  1,  0,  0,  0,  0,  0,  1,  0,  1,  0,
        1,  0,  1,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,
        0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,
        0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,
        0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,
        0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,
        0,  0,  0,  0,  2,  0,  2,  0,  2,  3,  2,  3,  2,  0,  2,  3,  2,
        3,  2,  3,  2,  3,  2,  0,  2,  3,  2,  0,  2,  3,  2,  3,  2,  3,
        2,  3,  2,  3,  2,  3,  2,  0,  2,  3,  2,  3,  2,  3,  2,  3,  2,
        3,  2,  3,  2,  3,  2,  3,  2,  3,  2, -1,  2,  3,  2, -1,  2,  3,
       -1,  3, -1,  3, -1, -1, -1, -1, -1, -1, -1, -1, -1])
```
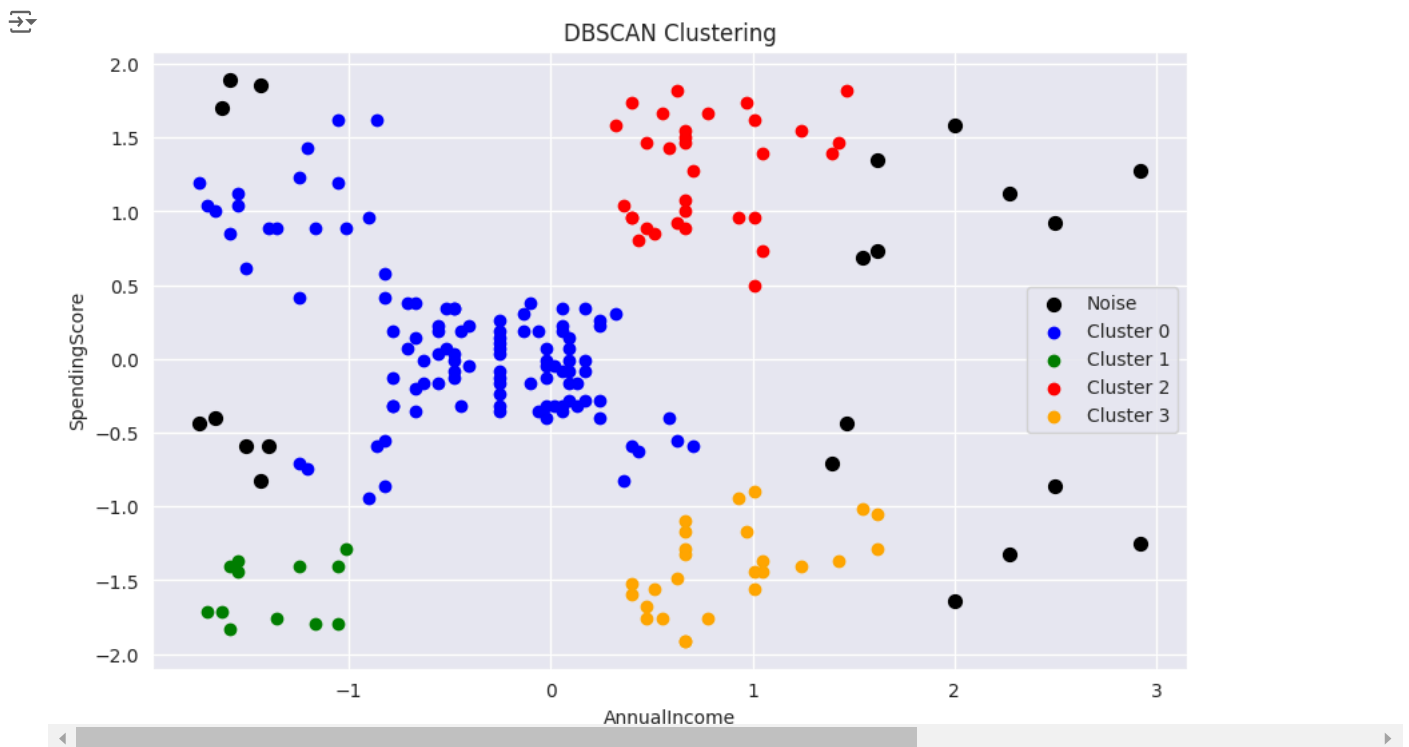
~:Outliers are denoted with " -1 "

```
1 dbscan.labels_
```

```
array([-1,  0,  1,  0, -1,  0,  1, -1,  1,  0,  1, -1,  1,  0,  1,  0, -1,
        0, -1, -1, -1,  0,  1,  0,  1,  0,  0,  0,  0,  0,  1,  0,  1,  0,
        1,  0,  1,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,
        0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,
        0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,
        0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,
        0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,  0,
        0,  0,  0,  0,  2,  0,  2,  0,  2,  3,  2,  3,  2,  0,  2,  3,  2,
        3,  2,  3,  2,  3,  2,  0,  2,  3,  2,  0,  2,  3,  2,  3,  2,  3,
        2,  3,  2,  3,  2,  3,  2,  0,  2,  3,  2,  3,  2,  3,  2,  3,  2,
        3,  2,  3,  2,  3,  2,  3,  2,  3,  2, -1,  2,  3,  2, -1,  2,  3,
       -1,  3, -1,  3, -1, -1, -1, -1, -1, -1, -1, -1, -1])
```

```python
1 # Plotting for clusters visualization
2 clr=["blue","green","red","orange","purple"]
3
4 pt.figure(figsize=(10,6))
5 for i in np.unique(ypre):
6     if i==-1:
7         pt.scatter(features[ypre== i]['AnnualIncome'],features[ypre== i]['SpendingScore'],s=50,c="black",label="Noise")
8     else:
9         pt.scatter(features[ypre== i]['AnnualIncome'],features[ypre== i]['SpendingScore'],label=f"Cluster {i}",c=clr[i])
10
11 pt.xlabel('AnnualIncome')
12 pt.ylabel('SpendingScore')
13 pt.title('DBSCAN Clustering')
14 pt.legend()
15 pt.show()
```



~: It's to be interpreted that black color datapoints represents outliers/noise in the clusters.

```python
1 #Evaluation Metrics..
2 ss=metrics.silhouette_score(features,dbscan.labels_)
3 print("Silhouette_Score Coefficient : {:.2f}".format(ss))
```

```
Silhouette_Score Coefficient : 0.41
```

~:Silhouette_Score ranges from -1 to 1,which near to one is best and near to -1 is worst. Since we got coffecient as *0.41* in which datapoints are very Moderately compact with the clusters.

Double-click (or enter) to edit

Double-click (or enter) to edit

```
1 Start coding or generate with AI.
```

END