

**ULLI PAWAN KALYAN**

**21BCE9090**

**AIML ASSIGNMENT-4**

**VIT-AP**

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

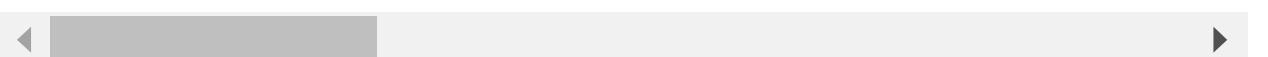
```
In [2]: data=pd.read_csv("Employee-Attrition.csv")
```

```
In [3]: data.head()
```

Out[3]:

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	Education
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sc
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sc
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sc
4	27	No	Travel_Rarely	591	Research & Development	2	1	

5 rows × 35 columns

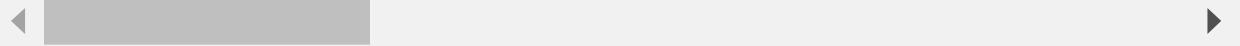


In [4]: `data.tail()`

Out[4]:

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	Ed
1465	36	No	Travel_Frequently	884	Research & Development	23	2	
1466	39	No	Travel_Rarely	613	Research & Development	6	1	
1467	27	No	Travel_Rarely	155	Research & Development	4	3	Lif
1468	49	No	Travel_Frequently	1023	Sales	2	3	
1469	34	No	Travel_Rarely	628	Research & Development	8	3	

5 rows × 35 columns



In [5]: `data.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 35 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Age              1470 non-null    int64  
 1   Attrition        1470 non-null    object  
 2   BusinessTravel   1470 non-null    object  
 3   DailyRate        1470 non-null    int64  
 4   Department       1470 non-null    object  
 5   DistanceFromHome 1470 non-null    int64  
 6   Education        1470 non-null    int64  
 7   EducationField   1470 non-null    object  
 8   EmployeeCount    1470 non-null    int64  
 9   EmployeeNumber   1470 non-null    int64  
 10  EnvironmentSatisfaction 1470 non-null    int64  
 11  Gender            1470 non-null    object  
 12  HourlyRate       1470 non-null    int64  
 13  JobInvolvement   1470 non-null    int64  
 14  JobLevel          1470 non-null    int64  
 15  JobRole           1470 non-null    object  
 16  JobSatisfaction  1470 non-null    int64  
 17  MaritalStatus     1470 non-null    object  
 18  MonthlyIncome     1470 non-null    int64  
 19  MonthlyRate       1470 non-null    int64  
 20  NumCompaniesWorked 1470 non-null    int64  
 21  Over18            1470 non-null    object  
 22  OverTime          1470 non-null    object  
 23  PercentSalaryHike 1470 non-null    int64  
 24  PerformanceRating 1470 non-null    int64  
 25  RelationshipSatisfaction 1470 non-null    int64  
 26  StandardHours     1470 non-null    int64  
 27  StockOptionLevel   1470 non-null    int64  
 28  TotalWorkingYears  1470 non-null    int64  
 29  TrainingTimesLastYear 1470 non-null    int64  
 30  WorkLifeBalance   1470 non-null    int64  
 31  YearsAtCompany    1470 non-null    int64  
 32  YearsInCurrentRole 1470 non-null    int64  
 33  YearsSinceLastPromotion 1470 non-null    int64  
 34  YearsWithCurrManager 1470 non-null    int64  
dtypes: int64(26), object(9)
memory usage: 402.1+ KB
```

In [6]: `data.describe()`

Out[6]:

	Age	DailyRate	DistanceFromHome	Education	EmployeeCount	EmployeeNumber
<b>count</b>	1470.000000	1470.000000	1470.000000	1470.000000	1470.0	1470.000000
<b>mean</b>	36.923810	802.485714	9.192517	2.912925	1.0	1024.86530
<b>std</b>	9.135373	403.509100	8.106864	1.024165	0.0	602.02433
<b>min</b>	18.000000	102.000000	1.000000	1.000000	1.0	1.00000
<b>25%</b>	30.000000	465.000000	2.000000	2.000000	1.0	491.25000
<b>50%</b>	36.000000	802.000000	7.000000	3.000000	1.0	1020.50000
<b>75%</b>	43.000000	1157.000000	14.000000	4.000000	1.0	1555.75000
<b>max</b>	60.000000	1499.000000	29.000000	5.000000	1.0	2068.00000

8 rows × 26 columns



## HANDLING NULL VALUES

```
In [7]: data.isnull().any()
```

```
Out[7]: Age           False
Attrition      False
BusinessTravel  False
DailyRate       False
Department     False
DistanceFromHome False
Education       False
EducationField  False
EmployeeCount   False
EmployeeNumber  False
EnvironmentSatisfaction False
Gender          False
HourlyRate      False
JobInvolvement  False
JobLevel        False
JobRole         False
JobSatisfaction False
MaritalStatus   False
MonthlyIncome   False
MonthlyRate     False
NumCompaniesWorked False
Over18          False
OverTime        False
PercentSalaryHike False
PerformanceRating False
RelationshipSatisfaction False
StandardHours   False
StockOptionLevel False
TotalWorkingYears False
TrainingTimesLastYear False
WorkLifeBalance  False
YearsAtCompany  False
YearsInCurrentRole False
YearsSinceLastPromotion False
YearsWithCurrManager False
dtype: bool
```

```
In [8]: data.isnull().sum()
```

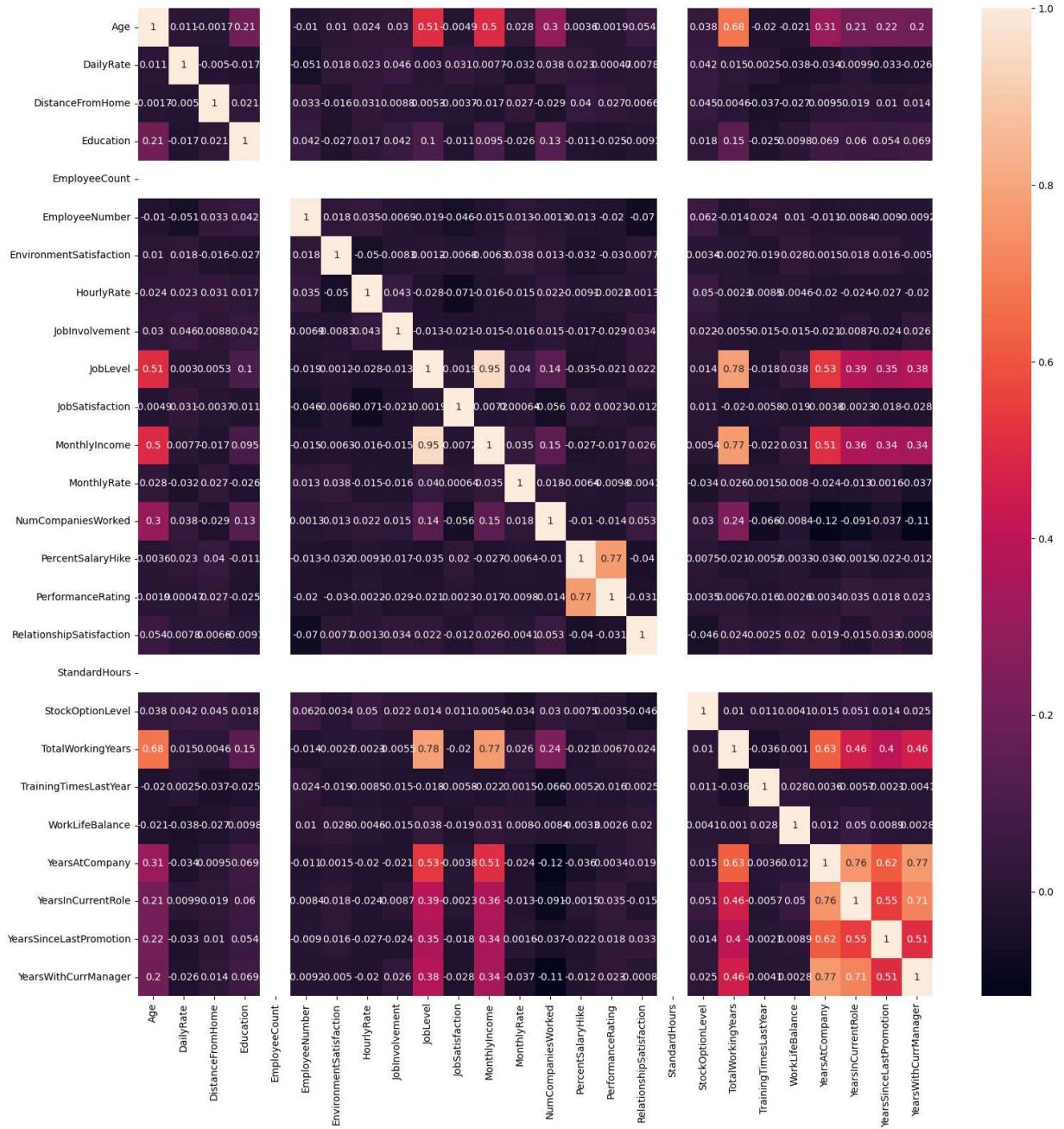
```
Out[8]: Age          0  
Attrition      0  
BusinessTravel  0  
DailyRate       0  
Department      0  
DistanceFromHome 0  
Education        0  
EducationField    0  
EmployeeCount     0  
EmployeeNumber    0  
EnvironmentSatisfaction 0  
Gender          0  
HourlyRate       0  
JobInvolvement    0  
JobLevel         0  
JobRole          0  
JobSatisfaction   0  
MaritalStatus     0  
MonthlyIncome      0  
MonthlyRate        0  
NumCompaniesWorked 0  
Over18           0  
OverTime          0  
PercentSalaryHike 0  
PerformanceRating 0  
RelationshipSatisfaction 0  
StandardHours      0  
StockOptionLevel    0  
TotalWorkingYears   0  
TrainingTimesLastYear 0  
WorkLifeBalance     0  
YearsAtCompany      0  
YearsInCurrentRole  0  
YearsSinceLastPromotion 0  
YearsWithCurrManager 0  
dtype: int64
```

```
In [9]: cor=data.corr()
```

```
C:\Users\srich\AppData\Local\Temp\ipykernel_26344\1426905697.py:1: FutureWarning:  
g: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.  
cor=data.corr()
```

```
In [10]: fig=plt.figure(figsize=(18,18))
sns.heatmap(cor,annot=True)
```

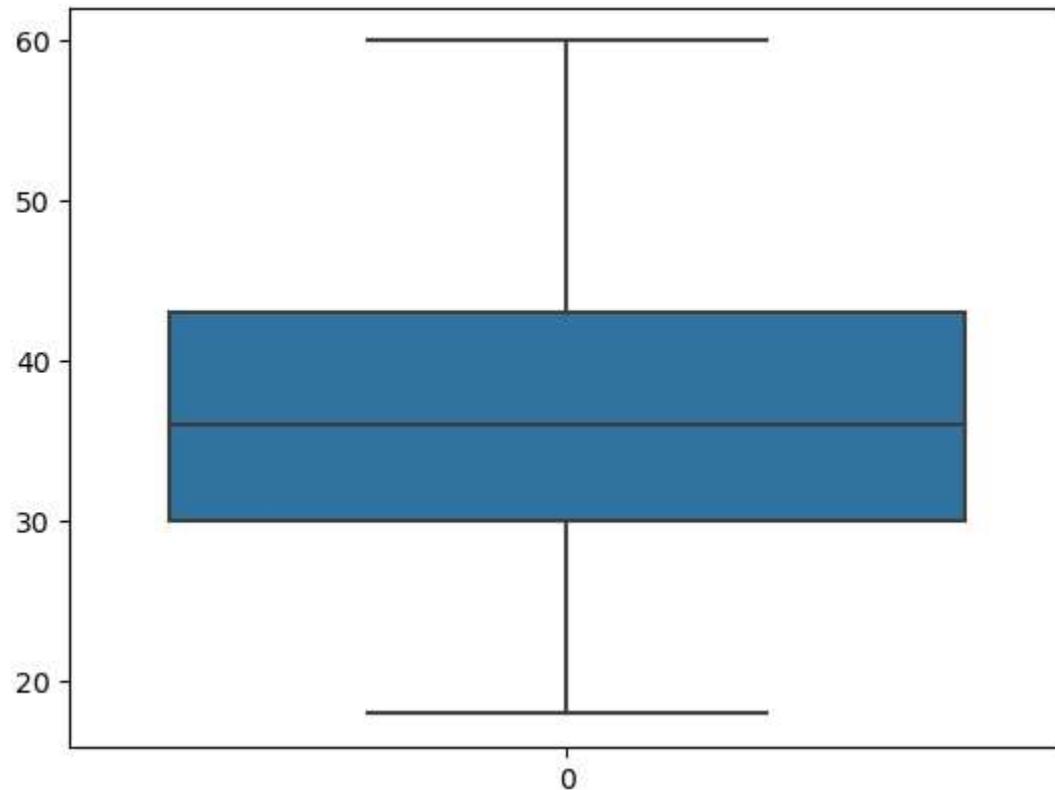
Out[10]: <Axes: >



## OUTLIERS

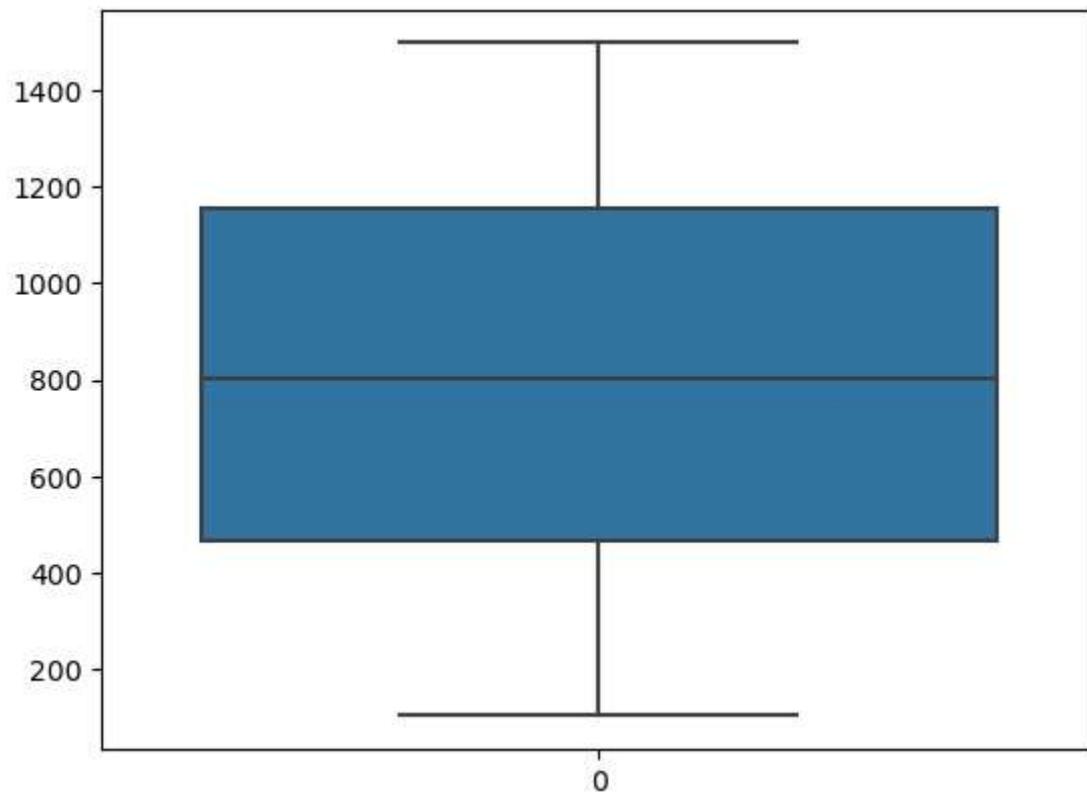
In [11]: `sns.boxplot(data[ "Age" ])`

Out[11]: <Axes: >



In [12]: `sns.boxplot(data["DailyRate"])`

Out[12]: <Axes: >



In [13]: `data.describe()`

Out[13]:

	Age	DailyRate	DistanceFromHome	Education	EmployeeCount	EmployeeNumber
count	1470.000000	1470.000000	1470.000000	1470.000000	1470.0	1470.000000
mean	36.923810	802.485714	9.192517	2.912925	1.0	1024.86530
std	9.135373	403.509100	8.106864	1.024165	0.0	602.02433
min	18.000000	102.000000	1.000000	1.000000	1.0	1.00000
25%	30.000000	465.000000	2.000000	2.000000	1.0	491.25000
50%	36.000000	802.000000	7.000000	3.000000	1.0	1020.50000
75%	43.000000	1157.000000	14.000000	4.000000	1.0	1555.75000
max	60.000000	1499.000000	29.000000	5.000000	1.0	2068.00000

8 rows × 26 columns

In [14]: `data.head()`

Out[14]:

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationalField
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Healthcare
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences
4	27	No	Travel_Rarely	591	Research & Development	2	1	Healthcare

5 rows × 35 columns

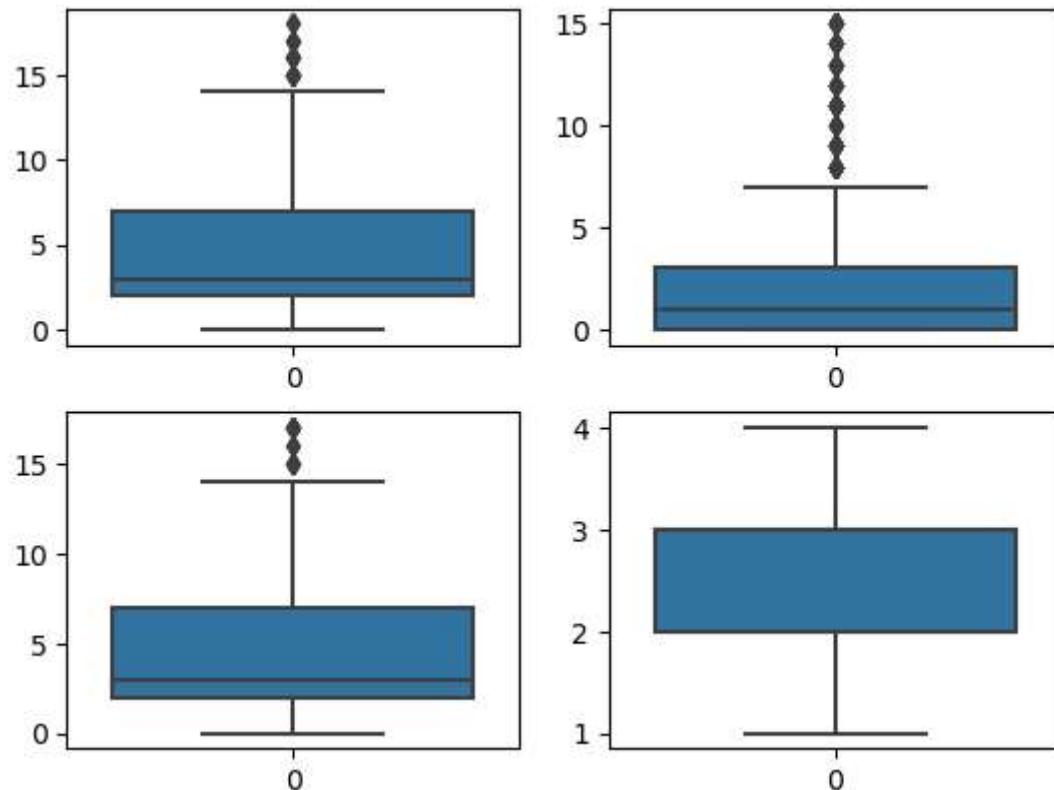


In [ ]:

In [15]:

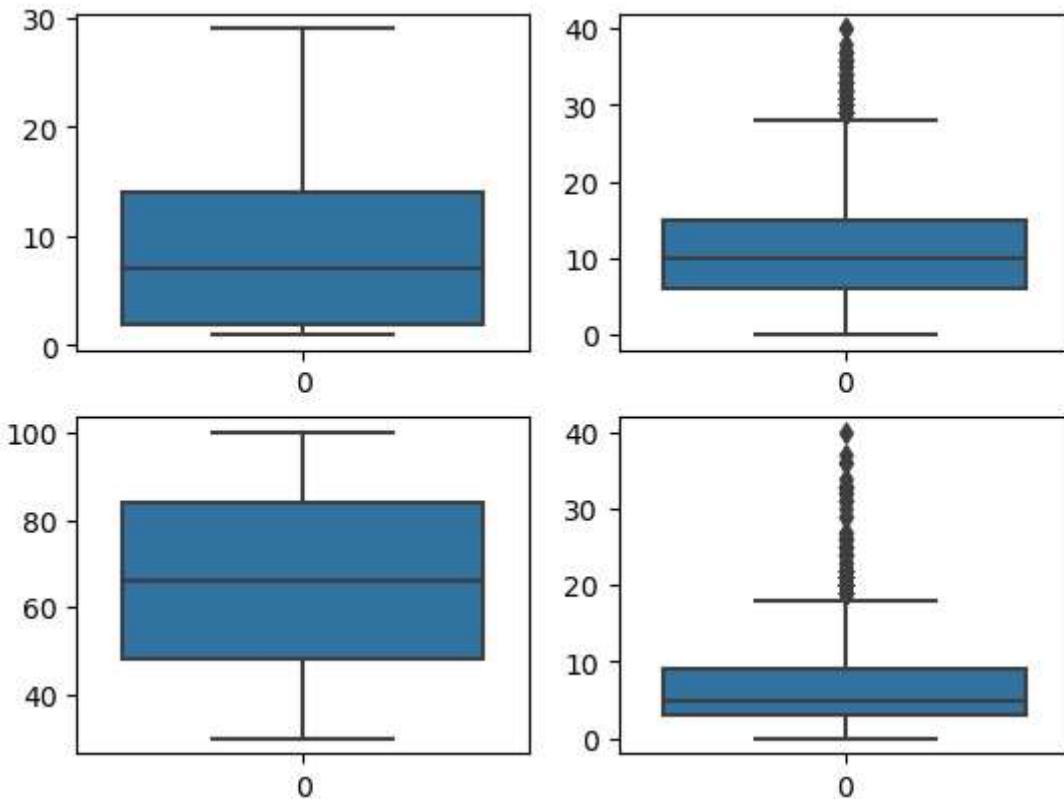
```
fig, axes = plt.subplots(2,2)
sns.boxplot(data=data["YearsInCurrentRole"],ax=axes[0,0])
sns.boxplot(data=data["YearsSinceLastPromotion"],ax=axes[0,1])
sns.boxplot(data=data["YearsWithCurrManager"],ax=axes[1,0])
sns.boxplot(data=data["WorkLifeBalance"],ax=axes[1,1])
```

Out[15]: <Axes: >



```
In [16]: fig, axes = plt.subplots(2,2)
sns.boxplot(data=data[ "DistanceFromHome"],ax=axes[0,0])
sns.boxplot(data=data[ "TotalWorkingYears"],ax=axes[0,1])
sns.boxplot(data=data[ "HourlyRate"],ax=axes[1,0])
sns.boxplot(data=data[ "YearsAtCompany"],ax=axes[1,1])
```

Out[16]: <Axes: >



## HANDLING THE OUTLIERS

```
In [17]: YearsInCurrentRole_q1 = data.YearsInCurrentRole.quantile(0.25)
YearsInCurrentRole_q3 = data.YearsInCurrentRole.quantile(0.75)
IQR_YearsInCurrentRole=YearsInCurrentRole_q3-YearsInCurrentRole_q1
upperlimit_YearsInCurrentRole=YearsInCurrentRole_q3+1.5*IQR_YearsInCurrentRole
lower_limit_YearsInCurrentRole =YearsInCurrentRole_q1-1.5*IQR_YearsInCurrentRole
median_YearsInCurrentRole=data[ "YearsInCurrentRole"].median()
data[ 'YearsInCurrentRole'] = np.where(
    (data[ 'YearsInCurrentRole'] > upperlimit_YearsInCurrentRole),
    median_YearsInCurrentRole,
    data[ 'YearsInCurrentRole']
)
```

```
In [18]: YearsSinceLastPromotion_q1 = data.YearsSinceLastPromotion.quantile(0.25)
YearsSinceLastPromotion_q3 = data.YearsSinceLastPromotion.quantile(0.75)
IQR_YearsSinceLastPromotion=YearsSinceLastPromotion_q3-YearsSinceLastPromotion_q1
upperlimit_YearsSinceLastPromotion=YearsSinceLastPromotion_q3+1.5*IQR_YearsSinceLastPromotion
lower_limit_YearsSinceLastPromotion =YearsSinceLastPromotion_q1-1.5*IQR_YearsSinceLastPromotion
median_YearsSinceLastPromotion=data[ "YearsSinceLastPromotion"].median()
data[ 'YearsSinceLastPromotion' ] = np.where(
    (data[ 'YearsSinceLastPromotion' ] > upperlimit_YearsSinceLastPromotion),
    median_YearsSinceLastPromotion,
    data[ 'YearsSinceLastPromotion' ]
)
```

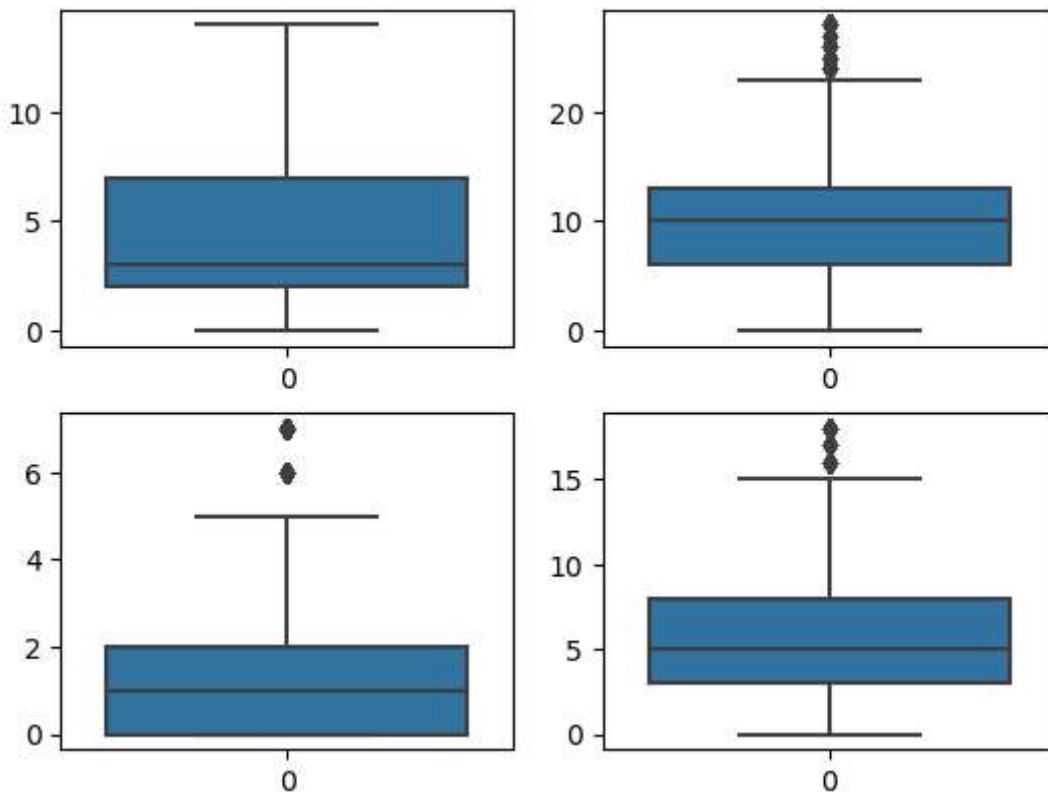
```
In [19]: YearsWithCurrManager_q1 = data.YearsWithCurrManager.quantile(0.25)
YearsWithCurrManager_q3 = data.YearsWithCurrManager.quantile(0.75)
IQR_YearsWithCurrManager=YearsWithCurrManager_q3-YearsWithCurrManager_q1
upperlimit_YearsWithCurrManager=YearsWithCurrManager_q3+1.5*IQR_YearsWithCurrManager
lower_limit_YearsWithCurrManager =YearsWithCurrManager_q1-1.5*IQR_YearsWithCurrManager
median_YearsWithCurrManager=data[ "YearsWithCurrManager"].median()
data[ 'YearsWithCurrManager' ] = np.where(
    (data[ 'YearsWithCurrManager' ] > upperlimit_YearsWithCurrManager),
    median_YearsWithCurrManager,
    data[ 'YearsWithCurrManager' ]
)
```

```
In [20]: TotalWorkingYears_q1 = data.TotalWorkingYears.quantile(0.25)
TotalWorkingYears_q3 = data.TotalWorkingYears.quantile(0.75)
IQR_TotalWorkingYears=TotalWorkingYears_q3-TotalWorkingYears_q1
upperlimit_TotalWorkingYears=TotalWorkingYears_q3+1.5*IQR_TotalWorkingYears
lower_limit_TotalWorkingYears=TotalWorkingYears_q1-1.5*IQR_TotalWorkingYears
median_TotalWorkingYears=data[ "TotalWorkingYears"].median()
data[ 'TotalWorkingYears' ] = np.where(
    (data[ 'TotalWorkingYears' ] > upperlimit_TotalWorkingYears),
    median_TotalWorkingYears,
    data[ 'TotalWorkingYears' ]
)
```

```
In [21]: YearsAtCompany_q1 = data.YearsAtCompany.quantile(0.25)
YearsAtCompany_q3 = data.YearsAtCompany.quantile(0.75)
IQR_YearsAtCompany=YearsAtCompany_q3-YearsAtCompany_q1
upperlimit_YearsAtCompany=YearsAtCompany_q3+1.5*IQR_YearsAtCompany
lower_limit_YearsAtCompany =YearsAtCompany_q1-1.5*IQR_YearsAtCompany
median_YearsAtCompany=data[ "YearsAtCompany"].median()
data[ 'YearsAtCompany' ] = np.where(
    (data[ 'YearsAtCompany' ] > upperlimit_YearsAtCompany),
    median_YearsAtCompany,
    data[ 'YearsAtCompany' ]
)
```

```
In [22]: fig, axes = plt.subplots(2,2)
sns.boxplot(data=data[ "YearsWithCurrManager"],ax=axes[0,0])
sns.boxplot(data=data[ "TotalWorkingYears"],ax=axes[0,1])
sns.boxplot(data=data[ "YearsSinceLastPromotion"],ax=axes[1,0])
sns.boxplot(data=data[ "YearsAtCompany"],ax=axes[1,1])
```

Out[22]: <Axes: >



```
In [23]: data.head()
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	Education
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sc
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sc
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sc
4	27	No	Travel_Rarely	591	Research & Development	2	1	M

5 rows × 35 columns

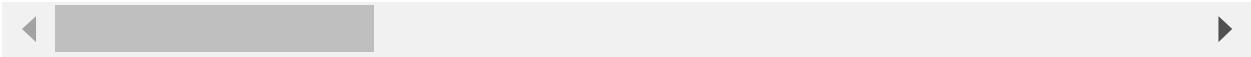
```
In [24]: data.drop("EducationField",axis=1,inplace=True)
```

```
In [25]: data.head(2)
```

```
Out[25]:
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	Employee
0	41	Yes	Travel_Rarely	1102	Sales		1	2
1	49	No	Travel_Frequently	279	Research & Development		8	1

2 rows × 34 columns



```
In [26]: data["BusinessTravel"].unique()
```

```
Out[26]: array(['Travel_Rarely', 'Travel_Frequently', 'Non-Travel'], dtype=object)
```

## SPLITTING THE DATA

```
In [27]: y=data["Attrition"]
```

```
In [28]: y.head()
```

```
Out[28]:
```

0	Yes
1	No
2	Yes
3	No
4	No

Name: Attrition, dtype: object

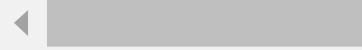
```
In [29]: data.drop("Attrition",axis=1,inplace=True)
```

In [30]: `data.head()`

Out[30]:

	Age	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EmployeeCount	E
0	41	Travel_Rarely	1102	Sales		1	2	1
1	49	Travel_Frequently	279	Research & Development		8	1	1
2	37	Travel_Rarely	1373	Research & Development		2	2	1
3	33	Travel_Frequently	1392	Research & Development		3	4	1
4	27	Travel_Rarely	591	Research & Development		2	1	1

5 rows × 33 columns



## ENCODING

In [31]: `from sklearn.preprocessing import LabelEncoder`

In [32]: `le=LabelEncoder()`

In [33]: `data["BusinessTravel"] = le.fit_transform(data["BusinessTravel"])`

In [34]: `data["Department"] = le.fit_transform(data["Department"])`

In [35]: `data["Gender"] = le.fit_transform(data["Gender"])`

In [36]: `y=le.fit_transform(y)`

In [37]: `y`

Out[37]: `array([1, 0, 1, ..., 0, 0, 0])`

In [38]: `data["JobRole"] = le.fit_transform(data["JobRole"])`

In [39]: `data["Over18"] = le.fit_transform(data["Over18"])`

In [40]: `data["MaritalStatus"] = le.fit_transform(data["MaritalStatus"])`

```
In [41]: data["OverTime"] = le.fit_transform(data["OverTime"])
```

```
In [42]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 33 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Age              1470 non-null    int64  
 1   BusinessTravel   1470 non-null    int32  
 2   DailyRate        1470 non-null    int64  
 3   Department       1470 non-null    int32  
 4   DistanceFromHome 1470 non-null    int64  
 5   Education        1470 non-null    int64  
 6   EmployeeCount    1470 non-null    int64  
 7   EmployeeNumber   1470 non-null    int64  
 8   EnvironmentSatisfaction 1470 non-null    int64  
 9   Gender            1470 non-null    int32  
 10  HourlyRate       1470 non-null    int64  
 11  JobInvolvement   1470 non-null    int64  
 12  JobLevel          1470 non-null    int64  
 13  JobRole           1470 non-null    int32  
 14  JobSatisfaction  1470 non-null    int64  
 15  MaritalStatus    1470 non-null    int32  
 16  MonthlyIncome    1470 non-null    int64  
 17  MonthlyRate      1470 non-null    int64  
 18  NumCompaniesWorked 1470 non-null    int64  
 19  Over18           1470 non-null    int32  
 20  OverTime          1470 non-null    int32  
 21  PercentSalaryHike 1470 non-null    int64  
 22  PerformanceRating 1470 non-null    int64  
 23  RelationshipSatisfaction 1470 non-null    int64  
 24  StandardHours    1470 non-null    int64  
 25  StockOptionLevel 1470 non-null    int64  
 26  TotalWorkingYears 1470 non-null    float64 
 27  TrainingTimesLastYear 1470 non-null    int64  
 28  WorkLifeBalance   1470 non-null    int64  
 29  YearsAtCompany   1470 non-null    float64 
 30  YearsInCurrentRole 1470 non-null    float64 
 31  YearsSinceLastPromotion 1470 non-null    float64 
 32  YearsWithCurrManager 1470 non-null    float64 
dtypes: float64(5), int32(7), int64(21)
memory usage: 338.9 KB
```

## TRAIN TEST SPLIT

```
In [43]: from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(data, y, test_size=0.3, random_state=
```

In [44]: `x_train.shape,x_test.shape,y_train.shape,y_test.shape`

Out[44]: `((1029, 33), (441, 33), (1029,), (441,))`

## FEATURE SCALING

In [45]: `from sklearn.preprocessing import StandardScaler`

In [46]: `sc=StandardScaler()`

In [47]: `x_train=sc.fit_transform(x_train)`

In [48]: `x_test=sc.fit_transform(x_test)`

## BUILDING THE MODEL

### MULTI LINEAR REGRESSION

In [49]: `from sklearn.linear_model import LinearRegression`

In [50]: `lr = LinearRegression()`

In [51]: `lr.fit(x_train,y_train)`

Out[51]: `LinearRegression()`

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**

**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

In [52]: `lr.coef_ #slope(m)`

Out[52]: `array([-3.54982205e-02, 5.83302672e-05, -1.72249253e-02, 3.46305828e-02, 2.45263256e-02, 3.89940919e-03, -8.89214800e+11, -9.43596046e-03, -4.12143211e-02, 1.05731111e-02, -3.10438176e-03, -3.84488534e-02, -1.53491202e-02, -1.57440821e-02, -3.66214700e-02, 3.35650229e-02, -5.85270940e-03, 5.89670852e-03, 3.77895930e-02, -1.09571684e+07, 9.55711035e-02, -2.54155836e-02, 1.99397870e-02, -2.64643443e-02, -2.64138305e+04, -1.79073382e-02, -3.31099987e-02, -1.08374986e-02, -3.09087610e-02, -2.49869759e-02, -1.08495820e-02, 2.10399164e-02, -6.46780261e-03])`

```
In [53]: lr.intercept_ #(c)
```

Out[53]: 0.16229348882410102

```
In [54]: y_pred = lr.predict(x_test)
```

In [55]: y\_pred

```
Out[55]: array([ 1.30796232e-01,  2.17900848e-01,  3.46642078e-01,  5.58365807e-03,
   4.99016900e-01,  1.01768656e-01,  3.45454856e-01,  1.23487843e-01,
  -1.60508211e-01,  4.02774192e-01,  1.43863626e-01,  2.67493211e-01,
  -4.61376545e-02,  5.58563347e-01,  2.81985888e-01,  1.48533683e-02,
  1.78513543e-01,  2.78218063e-01,  9.39501521e-02,  2.16824147e-01,
  2.65686758e-01,  1.40480328e-02,  8.35045241e-02,  9.65278467e-02,
  5.10152592e-01,  2.94947412e-01,  7.87474989e-02,  1.26281046e-01,
  5.05599470e-01,  8.44181518e-02,  -7.94455244e-02,  2.00664655e-02,
  1.07159698e-01,  3.65801296e-01,  1.25083285e-01,  5.15443890e-02,
  1.06691263e-01,  6.10560784e-02,  6.66004307e-02,  4.86894584e-02,
  -1.07548569e-02,  -2.94781175e-02,  5.20898060e-02,  -1.58276250e-02,
  -1.77848612e-02,  4.18379396e-01,  3.66954632e-01,  -2.14407437e-01,
  5.48302265e-01,  4.40485603e-01,  1.97011357e-01,  4.42145584e-01,
  1.46307401e-01,  3.75270817e-01,  4.93018248e-01,  2.95973962e-01,
  -4.67912123e-02,  3.16682815e-01,  -7.70189517e-03,  2.53080649e-01,
  -3.14434764e-02,  2.83109378e-01,  9.08930998e-02,  1.26210137e-01,
  3.60059433e-01,  2.42994338e-02,  3.55930832e-01,  1.96149466e-01,
  1.27766995e-01,  1.19157288e-01,  -2.87607873e-02,  3.18174064e-01,
  1.08182581e-01,  1.25900358e-01,  2.30252390e-01,  9.79248506e-02,
  2.12876232e-02,  2.73223285e-01,  2.53220174e-01,  4.26148277e-02])
```

In [56]: y\_test

## LOGISTIC REGRESSION

```
In [57]: from sklearn.linear_model import LogisticRegression
```

```
In [58]: lg=LogisticRegression()
```

```
In [59]: lg.fit(x_train,y_train)
```

```
Out[59]: LogisticRegression()
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [60]: y_pred_lg=lg.predict(x_test)
```

```
In [61]: y_pred
```

```
1.9723849e-01, 2.06514753e-01, -1.09632049e-01, 3.55228383e-01,  
7.24827035e-02, 2.36363938e-01, -8.16169248e-02, 3.21071947e-01,  
1.62813472e-02, 2.39336106e-01, 3.15931799e-01, 1.22614169e-01,  
9.00083375e-03, 9.34631288e-02, 4.76412767e-01, -5.65172803e-02,  
3.20357008e-01, 1.90009384e-02, 1.84117901e-01, 2.16278454e-01,  
-2.02785555e-01, 2.52564639e-01, -1.50856458e-01, 3.50111406e-01,  
-9.85040851e-02, 2.15620856e-01, 2.45829671e-02, 2.09072817e-01,  
9.49806038e-02, 3.11567457e-01, 4.07284673e-01, -2.19987388e-02,  
7.65184056e-02, 3.07591042e-01, 4.35392606e-02, 1.09065715e-01,  
4.21630293e-01, 1.48897439e-01, 4.42988409e-02, 3.35130176e-02,  
9.96284166e-02, -5.29452820e-02, 3.46797085e-01, 1.03461109e-01,  
-4.11871255e-02, -5.81603769e-02, 3.09315615e-01, 1.01570176e-01,  
2.09448790e-01, 4.10347757e-01, 3.33094014e-01, 1.81228317e-01,  
2.77318425e-01, 2.86403262e-01, 2.62049298e-01, -1.88118573e-02,  
2.35820366e-01, 1.54110630e-01, 6.23994747e-02, 6.70778678e-03,  
1.86183882e-02, 7.74249168e-02, 1.34438409e-01, 1.87764721e-01,  
2.36760139e-01, -1.81894158e-01, 2.98162656e-01, 1.74334561e-01,  
-8.91020361e-02, 3.47082854e-02, 1.36644027e-01, 1.70179690e-01,  
1.70464873e-01, 2.28231706e-01, 2.15402007e-01, 1.04763685e-01,  
-8.16311344e-02])
```

In [62]: y\_test

```
In [63]: score = lg.score(x_test, y_test)
        print(score)
```

0.8820861678004536

## CONFUSION MATRIX

```
In [64]: from sklearn import metrics  
cm = metrics.confusion_matrix(y_test,y_pred_lg)  
print(cm)
```

```
[[ 366    5]
 [ 47   23]]
```

## RIDGE AND LASSO

```
In [65]: from sklearn.linear_model import Ridge  
from sklearn.model_selection import GridSearchCV
```

```
In [66]: rg=Ridge()
```

```
In [67]: parametres={"alpha":[1,2,3,5,10,20,30,40,60,70,80,90]}
ridgecv=GridSearchCV(rg,parametres,scoring="neg_mean_squared_error",cv=5)
ridgecv.fit(x_train,y_train)
```

```
Out[67]: GridSearchCV(cv=5, estimator=Ridge(),
param_grid={'alpha': [1, 2, 3, 5, 10, 20, 30, 40, 60, 70, 80, 90]},
scoring='neg_mean_squared_error')
```

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**

**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

```
In [68]: print(ridgecv.best_params_)

{'alpha': 90}
```

```
In [69]: print(ridgecv.best_score_)

-0.1139062113923418
```

```
In [70]: y_pred_rg=ridgecv.predict(x_test)
```

```
In [71]: y_pred_rg
```

```
Out[71]: array([ 1.34413485e-01,  2.22561818e-01,  3.41692977e-01,  3.88209867e-03,
   4.84617338e-01,  1.16361483e-01,  3.30449743e-01,  1.27358807e-01,
  -1.34442619e-01,  3.77692888e-01,  1.33001445e-01,  2.69898751e-01,
  -2.54707392e-02,  5.25771894e-01,  2.67543514e-01,  2.78725024e-02,
  1.82233111e-01,  2.78896415e-01,  9.12689699e-02,  2.11494641e-01,
  2.70103341e-01,  8.44922044e-03,  8.74746722e-02,  1.05348798e-01,
  4.87749940e-01,  2.83080512e-01,  8.80556209e-02,  1.23817268e-01,
  4.82185624e-01,  9.34824523e-02, -7.16448509e-02,  4.07003104e-02,
  1.08437994e-01,  3.42151399e-01,  1.22270929e-01,  6.85889862e-02,
  1.06690533e-01,  7.08689637e-02,  7.51570276e-02,  6.05829413e-02,
  1.08782897e-02, -6.91368661e-03,  5.83191600e-02, -1.54680056e-02,
 -4.02267475e-03,  4.08010612e-01,  3.43668700e-01, -1.83519405e-01,
  5.29536511e-01,  4.27646098e-01,  1.95234877e-01,  4.25012930e-01,
  1.40754410e-01,  3.52173952e-01,  4.70372694e-01,  2.89240343e-01,
 -3.11642726e-02,  3.04206456e-01,  9.89337674e-03,  2.44569884e-01,
 -1.40249115e-02,  2.75133912e-01,  8.64669565e-02,  1.24214885e-01,
  3.48994545e-01,  3.41026778e-02,  3.40548051e-01,  1.95847356e-01,
  1.30040885e-01,  1.32259137e-01, -2.34680143e-02,  3.04595468e-01,
  1.12452197e-01,  1.30525275e-01,  2.19329505e-01,  9.44722098e-02,
```

In [72]: y\_test

```
In [73]: from sklearn import metrics  
print(metrics.r2_score(y_test,y_pred_rg))  
print(metrics.r2_score(y_train,ridgecv.predict(x_train)))
```

$0.21073458438815917$   
 $0.2061567210285109$

## LASSO

```
In [74]: from sklearn.linear_model import Lasso  
from sklearn.model_selection import GridSearchCV
```

```
In [75]: la=Ridge()
```

```
In [76]: parametres={"alpha":[1,2,3,5,10,20,30,40,60,70,80,90]}\nridgecv=GridSearchCV(la,parametres,scoring="neg_mean_squared_error",cv=5)\nridgecv.fit(x_train,y_train)
```

**In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.**

**On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.**

```
In [77]: print(ridgecv.best_params_)
```

```
{'alpha': 90}
```

```
In [78]: print(ridgecv.best_score_)
```

```
-0.1139062113923418
```

```
In [79]: y_pred_la=ridgecv.predict(x_test)
```

```
In [80]: y_pred_la
```

```
Out[80]: array([ 1.34413485e-01,  2.22561818e-01,  3.41692977e-01,  3.88209867e-03,
   4.84617338e-01,  1.16361483e-01,  3.30449743e-01,  1.27358807e-01,
  -1.34442619e-01,  3.77692888e-01,  1.33001445e-01,  2.69898751e-01,
  -2.54707392e-02,  5.25771894e-01,  2.67543514e-01,  2.78725024e-02,
  1.82233111e-01,  2.78896415e-01,  9.12689699e-02,  2.11494641e-01,
  2.70103341e-01,  8.44922044e-03,  8.74746722e-02,  1.05348798e-01,
  4.87749940e-01,  2.83080512e-01,  8.80556209e-02,  1.23817268e-01,
  4.82185624e-01,  9.34824523e-02,  -7.16448509e-02,  4.07003104e-02,
  1.08437994e-01,  3.42151399e-01,  1.22270929e-01,  6.85889862e-02,
  1.06690533e-01,  7.08689637e-02,  7.51570276e-02,  6.05829413e-02,
  1.08782897e-02,  -6.91368661e-03,  5.83191600e-02,  -1.54680056e-02,
  -4.02267475e-03,  4.08010612e-01,  3.43668700e-01,  -1.83519405e-01,
  5.29536511e-01,  4.27646098e-01,  1.95234877e-01,  4.25012930e-01,
  1.40754410e-01,  3.52173952e-01,  4.70372694e-01,  2.89240343e-01,
  -3.11642726e-02,  3.04206456e-01,  9.89337674e-03,  2.44569884e-01,
  -1.40249115e-02,  2.75133912e-01,  8.64669565e-02,  1.24214885e-01,
  3.48994545e-01,  3.41026778e-02,  3.40548051e-01,  1.95847356e-01,
  1.30040885e-01,  1.32259137e-01,  -2.34680143e-02,  3.04595468e-01,
  1.12452197e-01,  1.30525275e-01,  2.19329505e-01,  9.44722098e-02,
```

```
In [81]: from sklearn import metrics
print(metrics.r2_score(y_test,y_pred_la))
print(metrics.r2_score(y_train,ridgecv.predict(x_train)))
```

```
0.21073458438815917
```

```
0.2061567210285109
```

## DECISION TREE

```
In [82]: from sklearn.tree import DecisionTreeClassifier
dtc=DecisionTreeClassifier()
```

```
In [83]: dtc.fit(x_train,y_train)
```

Out[83]: DecisionTreeClassifier()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [84]: pred=dtc.predict(x_test)
```

In [85]: pred

In [86]: y\_test

```
In [87]: #Accuracy score  
from sklearn.metrics import accuracy_score,confusion_matrix,classification_report
```

```
In [88]: accuracy_score(y_test,pred)
```

Out[88]: 0.7777777777777778

```
In [89]: confusion_matrix(y_test, pred)
```

```
Out[89]: array([[326,  45],  
                 [ 53,  17]], dtype=int64)
```

```
In [90]: pd.crosstab(y_test,pred)
```

```
Out[90]:   col_0    0    1  
           row_0  
_____  
          0    326   45  
          1     53   17
```

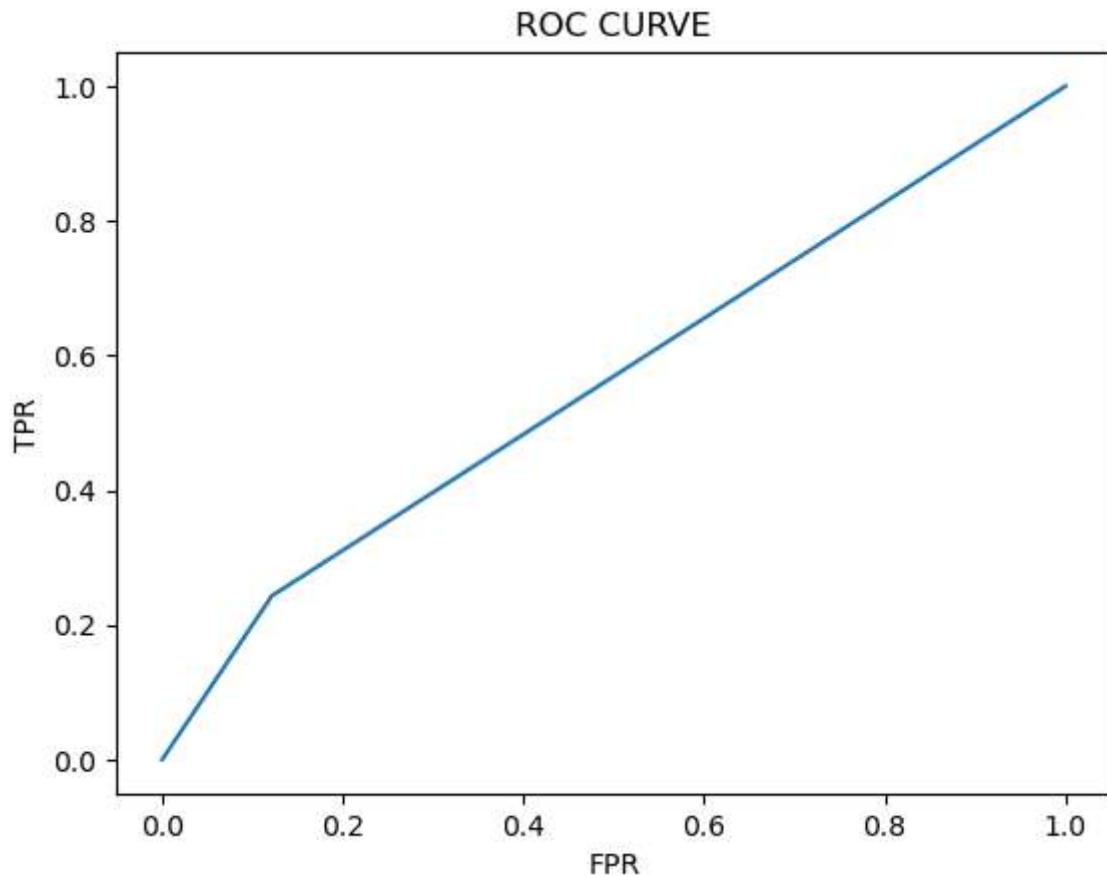
```
In [91]: print(classification_report(y_test,pred))
```

	precision	recall	f1-score	support
0	0.86	0.88	0.87	371
1	0.27	0.24	0.26	70
accuracy			0.78	441
macro avg	0.57	0.56	0.56	441
weighted avg	0.77	0.78	0.77	441

```
In [92]: probability=dtc.predict_proba(x_test)[:,1]
```

```
In [93]: # roc_curve  
fpr,tpr,thresholds = roc_curve(y_test,probability)
```

```
In [94]: plt.plot(fpr,tpr)  
plt.xlabel('FPR')  
plt.ylabel('TPR')  
plt.title('ROC CURVE')  
plt.show()
```



## RANDOM FOREST

```
In [95]: from sklearn.ensemble import RandomForestClassifier  
rfc=RandomForestClassifier()
```

```
In [96]: forest_params = [{'max_depth': list(range(10, 15)), 'max_features': list(range(0,
```

```
In [97]: from sklearn.model_selection import GridSearchCV
```

```
In [98]: rfc_cv= GridSearchCV(rfc,param_grid=forest_params,cv=10,scoring="accuracy")
```

```
In [99]: rfc_cv.fit(x_train,y_train)
```

```
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\model_selection\_validation.py:425: FitFailedWarning:  
50 fits failed out of a total of 700.  
The score on these train-test partitions for these parameters will be set to na  
n.  
If these failures are not expected, you can try to debug them by setting error_  
score='raise'.
```

Below are more details about the failures:

```
-----  
-  
50 fits failed with the following error:  
Traceback (most recent call last):  
  File "C:\ProgramData\anaconda3\Lib\site-packages\sklearn\model_selection\_val  
idation.py", line 732, in _fit_and_score  
    estimator.fit(X_train, y_train, **fit_params)  
  File "C:\ProgramData\anaconda3\Lib\site-packages\sklearn\base.py", line 1144,  
in wrapper  
    estimator._validate_params()  
  File "C:\ProgramData\anaconda3\Lib\site-packages\sklearn\base.py", line 637,  
in _validate_params  
    validate_parameter_constraints()  
  File "C:\ProgramData\anaconda3\Lib\site-packages\sklearn\utils\_param_validat  
ion.py", line 95, in validate_parameter_constraints  
    raise InvalidParameterError(  
sklearn.utils._param_validation.InvalidParameterError: The 'max_features' param  
eter of RandomForestClassifier must be an int in the range [1, inf), a float in  
the range (0.0, 1.0], a str among {'log2', 'sqrt'} or None. Got 0 instead.
```

```
    warnings.warn(some_fits_failed_message, FitFailedWarning)  
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\model_selection\_search.py:9  
76: UserWarning: One or more of the test scores are non-finite: [      nan  0.8  
4256615  0.84644965  0.84935275  0.85420712  0.85615839  
  0.85323625  0.85226537  0.84934323  0.8512945  0.85324576  0.84837236  
  0.85423567  0.85713878      nan  0.84159528  0.84644965  0.85227489  
  0.85226537  0.85422616  0.85228441  0.85227489  0.85226537  0.85130402  
  0.84256615  0.85224634  0.84642109  0.84838188      nan  0.84936227  
  0.8483914   0.847411   0.85322673  0.85519703  0.85324576  0.85421664  
  0.847411   0.85615839  0.84935275  0.84545974  0.84642109  0.8561679  
      nan  0.84645917  0.85032362  0.85228441  0.85033314  0.85422616  
  0.85422616  0.85323625  0.85226537  0.84935275  0.84837236  0.8512945  
  0.8541976   0.847411      nan  0.84255663  0.84644013  0.85226537  
  0.85810965  0.85324576  0.8512945   0.84838188  0.85130402  0.8512945  
  0.85227489  0.8512945   0.84837236  0.84935275]  
  warnings.warn(
```

```
Out[99]: GridSearchCV(cv=10, estimator=RandomForestClassifier(),
                      param_grid=[{'max_depth': [10, 11, 12, 13, 14],
                                   'max_features': [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 1
                                   1,
                                   12, 13]}],
                      scoring='accuracy')
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [100]: pred=rfc_cv.predict(x_test)
```

```
In [101]: print(classification_report(y_test,pred))
```

	precision	recall	f1-score	support
0	0.86	0.99	0.92	371
1	0.73	0.16	0.26	70
accuracy			0.86	441
macro avg	0.80	0.57	0.59	441
weighted avg	0.84	0.86	0.82	441

```
In [102]: rfc_cv.best_params_
```

```
Out[102]: {'max_depth': 14, 'max_features': 4}
```

```
In [103]: rfc_cv.best_score_
```

```
Out[103]: 0.8581096516276412
```

```
In [ ]:
```