

NAME: B PAVAN KUMAR

REG NO: 21BCE8241

1.Download the Employee Attrition Dataset

<https://www.kaggle.com/datasets/patelprashant/employee-attrition>

2.Perfrom Data Preprocessing

3.Model Building using Logistic Regression and Decision Tree and Random Forest

4.Calculate Performance metrics

```
#Import the Libraries.
```

```
import numpy as np
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

```
#Importing the dataset.
```

```
df=pd.read_csv("Employee-Attrition.csv")
```

```
df.head()
```

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

5 rows × 35 columns

```
df.shape
```

```
(1470, 35)
```

```
df.Age.value_counts()
```

```

35    78
34    77
36    69
31    69
29    68
32    61
30    60
33    58
38    58
40    57
37    50
27    48
28    48
42    46
39    42
45    41
41    40
26    39
44    33
46    33
43    32
50    30
25    26
24    26
49    24
47    24
55    22
51    19
53    19
48    19
54    18
52    18

```

```
22 16
56 14
23 14
58 14
21 13
20 11
59 10
19 9
18 8
60 5
57 4
Name: Age, dtype: int64
```

```
df.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
```

```
df.describe()

      Age  DailyRate  DistanceFromHome  Education  EmployeeCount  EmployeeNumber  EnvironmentSatisfaction  HourlyRate  J
count 1470.000000  1470.000000      1470.000000  1470.000000          1470.0      1470.000000      1470.000000  1470.000000
mean   36.923810   802.485714        9.192517    2.912925            1.0      1024.865306        2.721769    65.891156
std     9.135373   403.509100        8.106864    1.024165            0.0        602.024335        1.093082    20.329428
min    18.000000   102.000000        1.000000    1.000000            1.0         1.000000        1.000000    30.000000
25%    30.000000   465.000000        2.000000    2.000000            1.0        491.250000        2.000000    48.000000
50%    36.000000   802.000000        7.000000    3.000000            1.0       1020.500000        3.000000    66.000000
75%    43.000000  1157.000000       14.000000    4.000000            1.0       1555.750000        4.000000    83.750000
max    60.000000  1499.000000       29.000000    5.000000            1.0       2068.000000        4.000000   100.000000

8 rows × 26 columns
```

```
#Checking for Null Values.
df.isnull().any()

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

```

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

```

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

```

```

#Data Visualization.
sns.distplot(df["Age"])

```

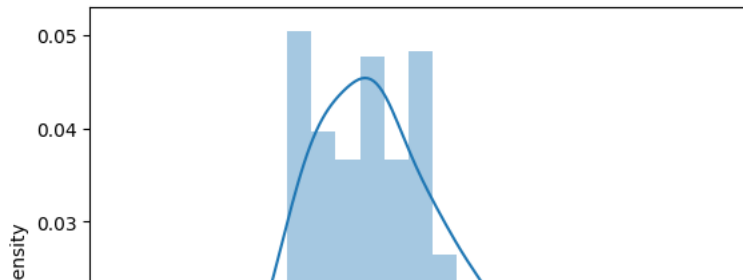
```
<ipython-input-84-25fc8198007f>:2: UserWarning:
```

```
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.
```

Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

For a guide to updating your code to use the new functions, please see <https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>

```
sns.distplot(df["Age"])
<Axes: xlabel='Age', ylabel='Density'>
```



```
df.corr()
```

```
<ipython-input-85-2f6f6606aa2c>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future ver
df.corr()
```

| | Age | DailyRate | DistanceFromHome | Education | EmployeeCount | EmployeeNumber | EnvironmentSatisfaction | Hc |
|---------------------------------|-----------|-----------|------------------|-----------|---------------|----------------|-------------------------|----|
| Age | 1.000000 | 0.010661 | -0.001686 | 0.208034 | NaN | -0.010145 | 0.010146 | |
| DailyRate | 0.010661 | 1.000000 | -0.004985 | -0.016806 | NaN | -0.050990 | 0.018355 | |
| DistanceFromHome | -0.001686 | -0.004985 | 1.000000 | 0.021042 | NaN | 0.032916 | -0.016075 | |
| Education | 0.208034 | -0.016806 | 0.021042 | 1.000000 | NaN | 0.042070 | -0.027128 | |
| EmployeeCount | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| EmployeeNumber | -0.010145 | -0.050990 | 0.032916 | 0.042070 | NaN | 1.000000 | 0.017621 | |
| EnvironmentSatisfaction | 0.010146 | 0.018355 | -0.016075 | -0.027128 | NaN | 0.017621 | 1.000000 | |
| HourlyRate | 0.024287 | 0.023381 | 0.031131 | 0.016775 | NaN | 0.035179 | -0.049857 | |
| JobInvolvement | 0.029820 | 0.046135 | 0.008783 | 0.042438 | NaN | -0.006888 | -0.008278 | |
| JobLevel | 0.509604 | 0.002966 | 0.005303 | 0.101589 | NaN | -0.018519 | 0.001212 | |
| JobSatisfaction | -0.004892 | 0.030571 | -0.003669 | -0.011296 | NaN | -0.046247 | -0.006784 | |
| MonthlyIncome | 0.497855 | 0.007707 | -0.017014 | 0.094961 | NaN | -0.014829 | -0.006259 | |
| MonthlyRate | 0.028051 | -0.032182 | 0.027473 | -0.026084 | NaN | 0.012648 | 0.037600 | |
| NumCompaniesWorked | 0.299635 | 0.038153 | -0.029251 | 0.126317 | NaN | -0.001251 | 0.012594 | |
| PercentSalaryHike | 0.003634 | 0.022704 | 0.040235 | -0.011111 | NaN | -0.012944 | -0.031701 | |
| PerformanceRating | 0.001904 | 0.000473 | 0.027110 | -0.024539 | NaN | -0.020359 | -0.029548 | |
| RelationshipSatisfaction | 0.053535 | 0.007846 | 0.006557 | -0.009118 | NaN | -0.069861 | 0.007665 | |
| StandardHours | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| StockOptionLevel | 0.037510 | 0.042143 | 0.044872 | 0.018422 | NaN | 0.062227 | 0.003432 | |
| TotalWorkingYears | 0.680381 | 0.014515 | 0.004628 | 0.148280 | NaN | -0.014365 | -0.002693 | |
| TrainingTimesLastYear | -0.019621 | 0.002453 | -0.036942 | -0.025100 | NaN | 0.023603 | -0.019359 | |
| WorkLifeBalance | -0.021490 | -0.037848 | -0.026556 | 0.009819 | NaN | 0.010309 | 0.027627 | |
| YearsAtCompany | 0.311309 | -0.034055 | 0.009508 | 0.069114 | NaN | -0.011240 | 0.001458 | |
| YearsInCurrentRole | 0.212901 | 0.009932 | 0.018845 | 0.060236 | NaN | -0.008416 | 0.018007 | |
| YearsSinceLastPromotion | 0.216513 | -0.033229 | 0.010029 | 0.054254 | NaN | -0.009019 | 0.016194 | |
| YearsWithCurrManager | 0.202089 | -0.026363 | 0.014406 | 0.069065 | NaN | -0.009197 | -0.004999 | |

26 rows × 26 columns

```
df.head()
```

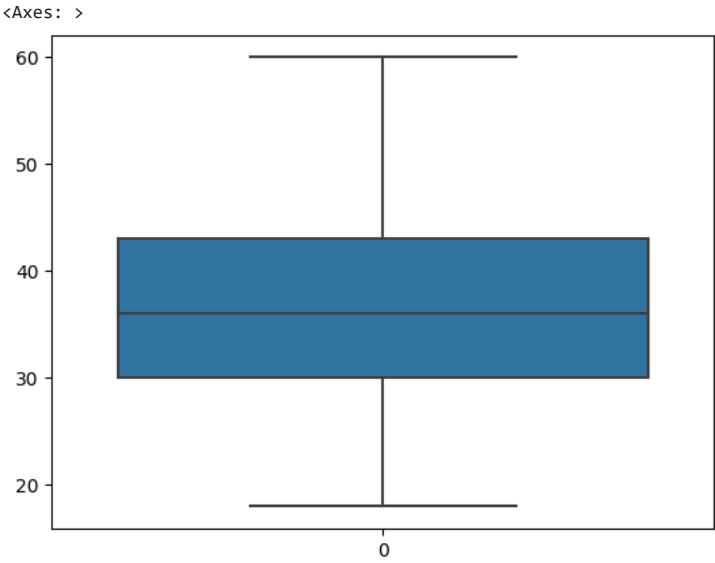
| | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | EmployeeCount | EmployeeNumber |
|---|-----|-----------|-------------------|-----------|------------------------|------------------|-----------|----------------|---------------|----------------|
| 0 | 41 | Yes | Travel_Rarely | 1102 | Sales | 1 | 2 | Life Sciences | 1 | 1 |
| 1 | 49 | No | Travel_Frequently | 279 | Research & Development | 8 | 1 | Life Sciences | 1 | 2 |
| 2 | 37 | Yes | Travel_Rarely | 1373 | Research & Development | 2 | 2 | Other | 1 | 3 |
| 3 | 33 | No | Travel_Frequently | 1392 | Research & Development | 3 | 4 | Life Sciences | 1 | 4 |
| 4 | 32 | No | Travel_Frequently | 1392 | Research & Development | 3 | 4 | Life Sciences | 1 | 5 |

```
plt.subplots(figsize = (25,25))
sns.heatmap(df.corr(),annot=True)
```

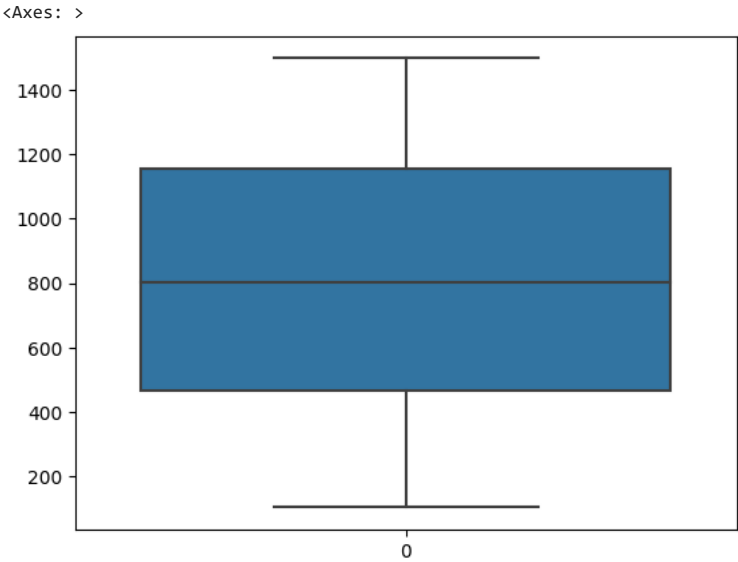
```
<ipython-input-87-9329d5e70af4>:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future ver
sns.heatmap(df.corr(),annot=True)
<Axes: >
```



```
sns.boxplot(df. Age)
```



```
sns.boxplot(df.DailyRate )
```



```
df.head()
```

| | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | EmployeeCount | EmployeeNumber |
|---|-----|-----------|-------------------|-----------|------------------------|------------------|-----------|----------------|---------------|----------------|
| 0 | 41 | Yes | Travel_Rarely | 1102 | Sales | | 1 | 2 | Life Sciences | 1 |
| 1 | 49 | No | Travel_Frequently | 279 | Research & Development | | 8 | 1 | Life Sciences | 1 |
| 2 | 37 | Yes | Travel_Rarely | 1373 | Research & Development | | 2 | 2 | Other | 1 |

```
x=df.iloc[:,1:4]
x.head()
```

| | Attrition | BusinessTravel | DailyRate |
|---|-----------|-------------------|-----------|
| 0 | Yes | Travel_Rarely | 1102 |
| 1 | No | Travel_Frequently | 279 |
| 2 | Yes | Travel_Rarely | 1373 |
| 3 | No | Travel_Frequently | 1392 |
| 4 | No | Travel_Rarely | 591 |

```
y=df.Attrition
y.head()

0    Yes
1     No
2     Yes
3     No
4     No
Name: Attrition, dtype: object
```

```
#label encoding
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
y=le.fit_transform(y)
```

```
#label encoding
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
y_test=le.fit_transform(y_test)
```

```
y

array([1, 0, 1, ..., 0, 0, 0])
```

```
y_test

array([[0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0,
        0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0,
        1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 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, 1, 0, 0, 0, 0, 0, 0, 0, 1,
        1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0,
        1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0,
        0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0,
        0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0,
        1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
        0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0,
        0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0,
        0, 1, 0, 0, 0, 1, 0, 0])
```

```
#label encoding
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
x.Attrition=le.fit_transform(x.Attrition)
x.head()
```

```

Attrition    BusinessTravel    DailyRate
#label encoding
from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
x.BusinessTravel    =le.fit_transform(x.BusinessTravel )
x.head()

```

| | Attrition | BusinessTravel | DailyRate |
|---|-----------|----------------|-----------|
| 0 | 1 | 2 | 1102 |
| 1 | 0 | 1 | 279 |
| 2 | 1 | 2 | 1373 |
| 3 | 0 | 1 | 1392 |
| 4 | 0 | 2 | 591 |

```

#feature scaling
from sklearn.preprocessing import MinMaxScaler
ms=MinMaxScaler()
x_scaled=pd.DataFrame(ms.fit_transform(x),columns=x.columns)

```

x_scaled

| | Attrition | BusinessTravel | DailyRate |
|------|-----------|----------------|-----------|
| 0 | 1.0 | 1.0 | 0.715820 |
| 1 | 0.0 | 0.5 | 0.126700 |
| 2 | 1.0 | 1.0 | 0.909807 |
| 3 | 0.0 | 0.5 | 0.923407 |
| 4 | 0.0 | 1.0 | 0.350036 |
| ... | ... | ... | ... |
| 1465 | 0.0 | 0.5 | 0.559771 |
| 1466 | 0.0 | 1.0 | 0.365784 |
| 1467 | 0.0 | 1.0 | 0.037938 |
| 1468 | 0.0 | 0.5 | 0.659270 |
| 1469 | 0.0 | 1.0 | 0.376521 |

1470 rows × 3 columns

```

#Splitting Data into Train and Test.
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x_scaled,y,test_size=0.2,random_state=0)

```

x_train.shape,x_test.shape,y_train.shape,y_test.shape

((1176, 3), (294, 3), (1176,), (294,))

x_train.head()

| | Attrition | BusinessTravel | DailyRate |
|------|-----------|----------------|-----------|
| 1374 | 0.0 | 1.0 | 0.360057 |
| 1092 | 0.0 | 1.0 | 0.607015 |
| 768 | 0.0 | 1.0 | 0.141732 |
| 569 | 0.0 | 0.0 | 0.953472 |
| 911 | 1.0 | 0.5 | 0.355762 |

```

from sklearn.linear_model import LogisticRegression
model=LogisticRegression()

```

```

model.fit(x_train,y_train)
pred=model.predict(x_test)
pred

```



```
array([0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0,
       0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 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, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1,
       1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1,
       0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0,
       0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0,
       1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0,
       0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0,
       0, 1, 0, 0, 0, 1, 0, 0, 0])
```

```
y_test
array([0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
       0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0,
       1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 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, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1,
       1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1,
       0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0,
       0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0,
       1, 0, 0, 0, 1, 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, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0,
       0, 1, 0, 0, 0, 1, 0, 0, 0])
```

df

| | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | EmployeeCount | EmployeeNur |
|------|-----|-----------|-------------------|-----------|------------------------|------------------|-----------|----------------|---------------|-------------|
| 0 | 41 | Yes | Travel_Rarely | 1102 | Sales | 1 | 2 | Life Sciences | 1 | |
| 1 | 49 | No | Travel_Frequently | 279 | Research & Development | 8 | 1 | Life Sciences | 1 | |
| 2 | 37 | Yes | Travel_Rarely | 1373 | Research & Development | 2 | 2 | Other | 1 | |
| 3 | 33 | No | Travel_Frequently | 1392 | Research & Development | 3 | 4 | Life Sciences | 1 | |
| 4 | 27 | No | Travel_Rarely | 591 | Research & Development | 2 | 1 | Medical | 1 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 1465 | 36 | No | Travel_Frequently | 884 | Research & Development | 23 | 2 | Medical | 1 | : |
| 1466 | 39 | No | Travel_Rarely | 613 | Research & Development | 6 | 1 | Medical | 1 | : |
| 1467 | 27 | No | Travel_Rarely | 155 | Research & Development | 4 | 3 | Life Sciences | 1 | : |
| 1468 | 49 | No | Travel_Frequently | 1023 | Sales | 2 | 3 | Medical | 1 | : |
| 1469 | 34 | No | Travel_Rarely | 628 | Research & Development | 8 | 3 | Medical | 1 | : |

1470 rows × 35 columns

▼ Evaluation of classification model

```
#Accuracy score
from sklearn.metrics import accuracy_score,confusion_matrix,classification_report,roc_auc_score,roc_curve

accuracy_score(y_test,pred)

1.0

confusion_matrix(y_test,pred)

array([[245,  0],
       [ 0, 49]])
```

```
pd.crosstab(y_test,pred)
```

```
col_0    0    1
row_0
0      245    0
1         0   49
```

▼ Roc-AUC curve

```
probability=model.predict_proba(x_test)[:,1]
probability
```

```
array([0.00812244, 0.00828669, 0.95811557, 0.00736818, 0.95915491,
0.00868515, 0.95769994, 0.00747535, 0.00864836, 0.95948831,
0.00746992, 0.9598929 , 0.00780439, 0.95840564, 0.00780955,
0.00754884, 0.00819051, 0.95859568, 0.00855622, 0.00773983,
0.00821274, 0.00802763, 0.00836879, 0.00802746, 0.0074146 ,
0.0082729 , 0.00813215, 0.00757403, 0.00836558, 0.0080338 ,
0.00847012, 0.00758941, 0.00776549, 0.9586776 , 0.00739941,
0.00786175, 0.00813215, 0.00829093, 0.95832836, 0.00824523,
0.00869423, 0.007611 , 0.00869312, 0.00878484, 0.00793765,
0.96025833, 0.00826973, 0.00803054, 0.95868102, 0.95757947,
0.0075865 , 0.95463824, 0.00865168, 0.00737196, 0.95822857,
0.00851361, 0.00811741, 0.95863069, 0.00767957, 0.00844072,
0.0086454 , 0.00748238, 0.00846579, 0.00733505, 0.00836005,
0.00790619, 0.95612966, 0.9572522 , 0.00791531, 0.95601035,
0.00766876, 0.0081645 , 0.00857157, 0.00818544, 0.00738205,
0.00796719, 0.00874445, 0.00813944, 0.95790261, 0.00842686,
0.00872992, 0.00804101, 0.00847989, 0.00855075, 0.00823679,
0.00784465, 0.00816345, 0.00759541, 0.00780755, 0.00776466,
0.00733489, 0.00755271, 0.00855184, 0.00857157, 0.00805819,
0.00843441, 0.00832604, 0.00879947, 0.00741825, 0.0077944 ,
0.95458789, 0.95561224, 0.00849492, 0.00767072, 0.00804925,
0.0075498 , 0.00793054, 0.0074594 , 0.00817461, 0.00736803,
0.00738331, 0.00795088, 0.00826761, 0.00815388, 0.00795495,
0.0086574 , 0.00792309, 0.00795292, 0.00805441, 0.00771804,
0.00753627, 0.00811238, 0.00844504, 0.96115113, 0.00879384,
0.00842434, 0.00770898, 0.00773685, 0.00784666, 0.00861778,
0.00786578, 0.95829225, 0.95454309, 0.00787284, 0.00745558,
0.00791835, 0.96048519, 0.00783143, 0.00767547, 0.00813111,
0.0080699 , 0.0076776 , 0.00843009, 0.00799787, 0.00742126,
0.95698192, 0.008544 , 0.00833244, 0.95903925, 0.00750062,
0.00777229, 0.00822431, 0.00825157, 0.95411538, 0.00773091,
0.00855385, 0.00774876, 0.00852761, 0.00861999, 0.00739736,
0.95904348, 0.00753144, 0.00806439, 0.95982325, 0.00838131,
0.00793257, 0.00787687, 0.00772101, 0.00803672, 0.0079285 ,
0.00825685, 0.95951339, 0.00814153, 0.00756045, 0.95461587,
0.00788899, 0.00776864, 0.00829925, 0.00776168, 0.0086233 ,
0.0076032 , 0.00779839, 0.00738978, 0.95567792, 0.00752758,
0.00875341, 0.00766974, 0.00827502, 0.95890716, 0.00857377,
0.008401 , 0.95498365, 0.00836558, 0.00871429, 0.00757112,
0.00863987, 0.00857926, 0.0074165 , 0.95391145, 0.00763705,
0.0074947 , 0.00849184, 0.95977841, 0.00814361, 0.00820013,
0.00756692, 0.00838596, 0.00799582, 0.00867422, 0.00838596,
0.00792749, 0.00745367, 0.0081645 , 0.95357441, 0.95804902,
0.00736535, 0.00774876, 0.00792224, 0.00778857, 0.00743363,
0.00799479, 0.00869201, 0.00830971, 0.008613 , 0.00876014,
0.00776366, 0.00831539, 0.00787082, 0.95823976, 0.00770405,
0.00813406, 0.00850598, 0.00794766, 0.00825263, 0.00753337,
0.00834953, 0.00839778, 0.0073748 , 0.0075498 , 0.00880397,
0.9593064 , 0.0082782 , 0.00763053, 0.00821169, 0.00832604,
0.00837201, 0.00765567, 0.00753321, 0.00799684, 0.00761783,
0.00846579, 0.00825456, 0.95988628, 0.00827184, 0.00741841,
0.00762857, 0.00760515, 0.95708899, 0.00771409, 0.95641043,
0.00818649, 0.00765976, 0.00806869, 0.00825914, 0.00782959,
0.00804719, 0.00825597, 0.95475459, 0.9540758 , 0.00775869,
0.0086998 , 0.00820223, 0.00777528, 0.00837523, 0.00835702,
0.00828775, 0.00794766, 0.00740416, 0.00850054, 0.00785672,
0.95430063, 0.00766777, 0.00796413, 0.00814987, 0.00841823,
0.00811862, 0.00757468, 0.96004177, 0.00753337, 0.00830581,
```

```
y
```

```
array([1, 0, 1, ..., 0, 0, 0])
```

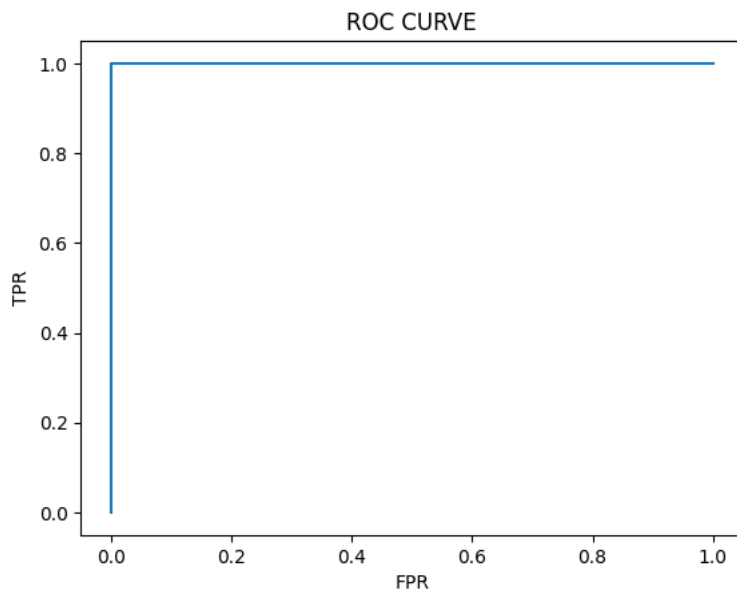
```
y_test
```

```
array([0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0,
1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 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, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1,
1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1,
0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0,
0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0,
1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0,
0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0,
0, 1, 0, 0, 0, 1, 0, 0])
```

```
# roc_curve
fpr, tpr, thresholds = roc_curve(y_test, probability)
```

```
plt.plot(fpr, tpr)
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.title('ROC CURVE')
plt.show()
```



DecisionTreeClassifier

```
from sklearn.tree import DecisionTreeClassifier
dtc=DecisionTreeClassifier()
```

```
dtc.fit(x_train,y_train)
```

```
DecisionTreeClassifier
```

```
pred=dtc.predict(x_test)
```

```
pred
```

```
array([0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0,
1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 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, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1,
1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1,
0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0,
0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0,
1, 0, 0, 0, 1, 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, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0,
0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
0, 1, 0, 0, 0, 0, 1, 0, 0])
```

```
y_test
```

```
array([0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
0, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
0, 1, 0, 0, 0, 0, 1, 0, 0])
```

df

1470 rows × 35 columns

```

/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but MinMaxScaler was
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but DecisionTreeClas
warnings.warn(
array([1])

```

- ▼ Evaluation of classification model

y_test

12/18

```
1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0,
0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0,
0, 1, 0, 0, 0, 1, 0, 0])
```

```
accuracy_score(y_test,pred)
```

```
1.0
```

```
confusion_matrix(y_test,pred)
```

```
array([[245,  0],
       [ 0, 49]])
```

```
pd.crosstab(y_test,pred)
```

```
col_0    0    1
row_0
0      245    0
1         0   49
```

```
print(classification_report(y_test,pred))
```

```
              precision    recall  f1-score   support

0               1.00         1.00         1.00         245
1               1.00         1.00         1.00          49

 accuracy               1.00         1.00         1.00         294
 macro avg              1.00         1.00         1.00         294
 weighted avg           1.00         1.00         1.00         294
```

▼ Roc-AUC curve

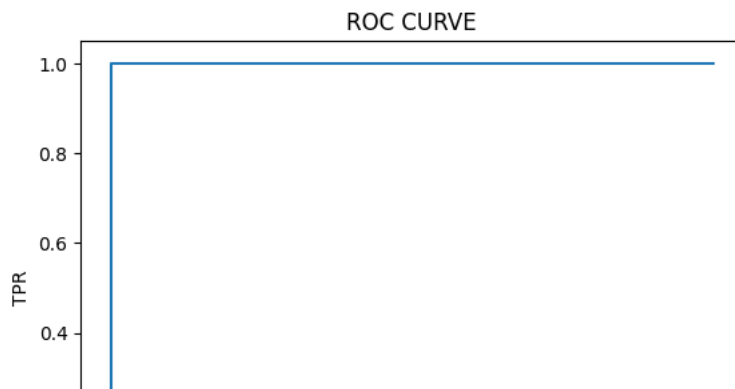
```
probability=dtc.predict_proba(x_test)[: ,1]
```

```
probability
```

```
array([0., 0., 1., 0., 1., 0., 1., 0., 0., 1., 0., 1., 0., 1., 0., 0., 0.,
1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1.,
0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 1., 0., 0., 1., 1., 0.,
1., 0., 0., 1., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 1., 1.,
0., 1., 0., 0., 0., 0., 0., 0., 0., 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., 1., 0., 0., 0., 0., 0., 0., 0., 1., 1., 0., 0., 0.,
1., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 1., 0., 0., 0., 0.,
1., 0., 0., 0., 0., 0., 1., 0., 0., 1., 0., 0., 0., 0., 0.,
0., 1., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0.,
0., 1., 0., 0., 1., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 1., 0.,
0., 0., 0., 0., 0., 0., 0., 0., 1., 1., 0., 0., 0., 0., 0.,
0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0.,
0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0., 1., 0., 0.,
0., 0., 1., 0., 1., 0., 0., 0., 0., 0., 0., 1., 1., 0., 0., 0.,
0., 0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 1., 0.,
0., 0., 1., 0., 0.]
```

```
fpr, tpr, thresholds = roc_curve(y_test, probability)
```

```
plt.plot(fpr, tpr)
plt.xlabel('FPR')
plt.ylabel('TPR')
plt.title('ROC CURVE')
plt.show()
```



```
from sklearn import tree
plt.figure(figsize=(25,15))
tree.plot_tree(dtc,filled=True)
```

```
[Text(0.5, 0.75, 'x[0] <= 0.5\ngini = 0.269\nsamples = 1176\nvalue = [988, 188]'),
Text(0.25, 0.25, 'gini = 0.0\nsamples = 988\nvalue = [988, 0]'),
Text(0.75, 0.25, 'gini = 0.0\nsamples = 188\nvalue = [0, 188]')]
```

$x[0] \leq 0.5$
gini = 0.269
samples = 1176
value = [988, 188]

gini = 0.0
samples = 988
value = [988, 0]

```
from sklearn.model_selection import GridSearchCV
parameter={
    'criterion':['gini','entropy'],
    'splitter':['best','random'],
    'max_depth':[1,2,3,4,5],
    'max_features':['auto', 'sqrt', 'log2']
}
```

```
grid_search=GridSearchCV(estimator=dtc,param_grid=parameter,cv=5,scoring="accuracy")
```

```
grid_search.fit(x_train,y_train)
```

[illegible]


```

/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
grid_search.best_params_

{'criterion': 'entropy',
 'max_depth': 5,
 'max_features': 'log2',
 'splitter': 'best'}
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
dtc_cv=DecisionTreeClassifier(criterion='entropy',
                             max_depth=3,
                             max_features='sqrt',
                             splitter='best')
dtc_cv.fit(x_train,y_train)

```

```

DecisionTreeClassifier
DecisionTreeClassifier(criterion='entropy', max_depth=3, max_features='sqrt')

```

```

/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
pred=dtc_cv.predict(x_test)
warnings.warn(
print(classification_report(y_test,pred))

```

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 1.00 | 1.00 | 1.00 | 245 |
| 1 | 1.00 | 1.00 | 1.00 | 49 |
| accuracy | | | 1.00 | 294 |
| macro avg | 1.00 | 1.00 | 1.00 | 294 |
| weighted avg | 1.00 | 1.00 | 1.00 | 294 |

```

/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(

```

RandomForestClassifier

```

warnings.warn(
from sklearn.ensemble import RandomForestClassifier
rfc=RandomForestClassifier()
warnings.warn(
forest_params = [{'max_depth': list(range(10, 15)), 'max_features': list(range(0,14))}]
/usr/local/lib/python3.10/dist-packages/sklearn/tree/_classes.py:269: FutureWarning: `max_features='auto'` has been deprecated in 1.
  warnings.warn(
rfc_cv= GridSearchCV(rfc,param_grid=forest_params,cv=10,scoring="accuracy")
warnings.warn(
rfc_cv.fit(x_train,y_train)

```

```
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_validation.py:378: FitFailedWarning:
50 fits failed out of a total of 700.
The score on these train-test partitions for these parameters will be set to nan.
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 "/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_validation.py", line 686, in _fit_and_score
    estimator.fit(X_train, y_train, **fit_params)
  File "/usr/local/lib/python3.10/dist-packages/sklearn/ensemble/_forest.py", line 340, in fit
    self._validate_params()
  File "/usr/local/lib/python3.10/dist-packages/sklearn/base.py", line 600, in _validate_params
    validate_parameter_constraints(
  File "/usr/local/lib/python3.10/dist-packages/sklearn/utils/_param_validation.py", line 97, in validate_parameter_constraints
    raise InvalidParameterError(
sklearn.utils._param_validation.InvalidParameterError: The 'max_features' parameter of RandomForestClassifier must be an int in the
warnings.warn(some_fits_failed_message, FitFailedWarning)
/usr/local/lib/python3.10/dist-packages/sklearn/model_selection/_search.py:952: UserWarning: One or more of the test scores are non-
1. 1. 1. 1. 1. 1. 1. 1. 1. 1. nan 1. 1. 1. 1. 1. 1.
1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
| 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
| 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
pred=rfc_cv.predict(x_test)

| 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
| 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.
print(classification_report(y_test,pred))

          precision    recall  f1-score   support

     0         1.00      1.00      1.00        245
     1         1.00      1.00      1.00         49

 accuracy          1.00
 macro avg         1.00      1.00      1.00        294
weighted avg         1.00      1.00      1.00        294

rfc_cv.best_params_
{'max_depth': 10, 'max_features': 1}
```