

In [1]: #Logistic regression, Decisive tree and random forest classifiers on Employee Attrition dataset

```
In [2]: #Importing necessary libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [3]: #Importing the dataset
df=pd.read_csv("Wk_Fn-UseC-HR-Employee-Attrition.csv")

In [4]: df.head()
```

	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	...	RelationshipSatisfaction	StandardHours	
0	41	Yes	Travel_Rarely	1102	Sales		1	2	Life Sciences	1	1	...	1	80
1	49	No	Travel_Frequently	279	Research & Development		8	1	Life Sciences	1	2	...	4	80
2	37	Yes	Travel_Rarely	1392	Research & Development		2	2	Other	1	4	...		80
3	33	No	Travel_Frequently	1392	Research & Development		3	4	Life Sciences	1	5	...	3	80
4	27	No	Travel_Rarely	591	Research & Development		2	1	Medical	1	7	...	4	80

5 rows x 35 columns

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

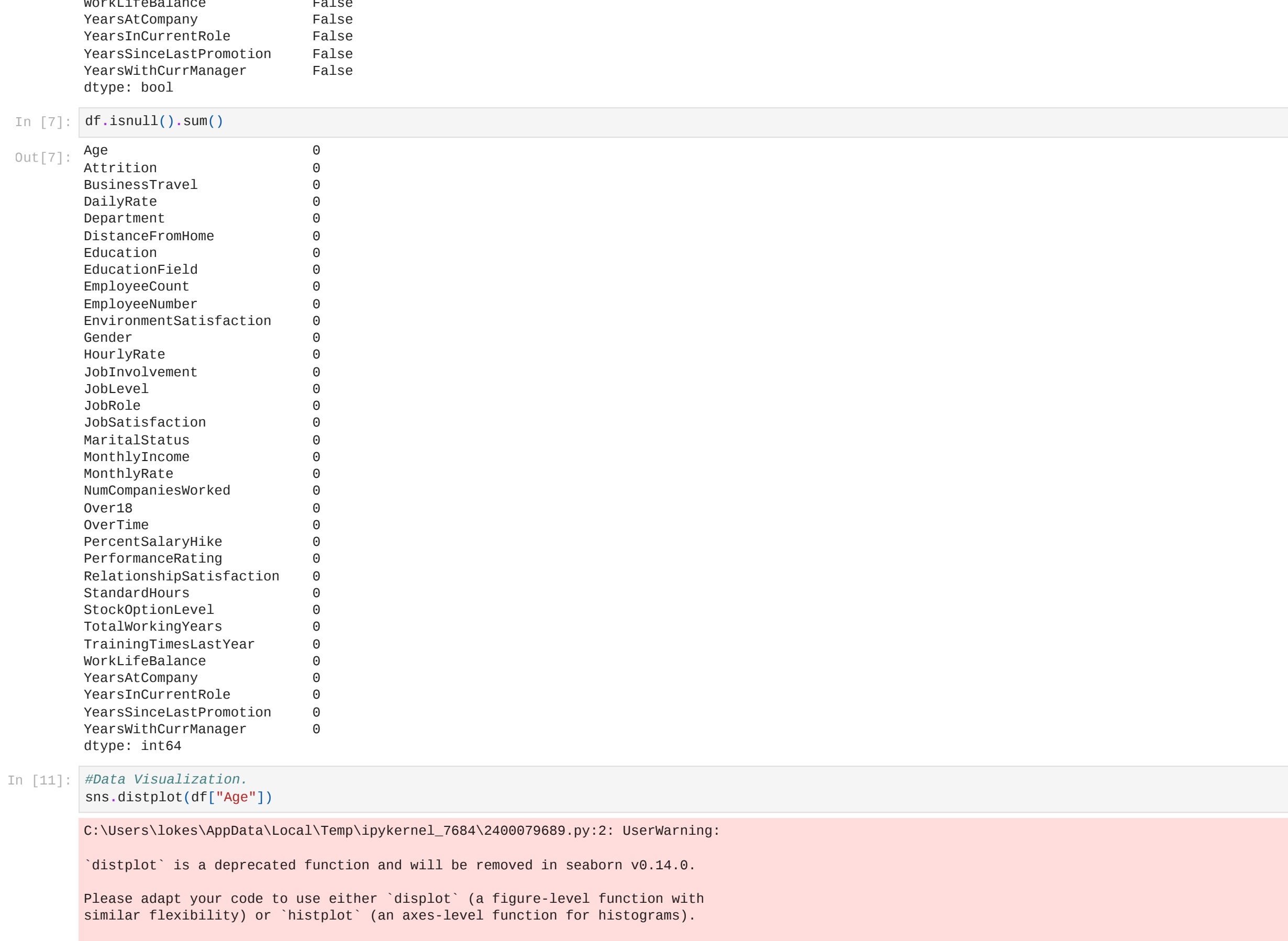
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 35 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Age                   1470 non-null   int64
 1   Attrition              1470 non-null   object
 2   BusinessTravel         1470 non-null   object
 3   DailyRate              1470 non-null   int64
 4   Department             1470 non-null   object
 5   DistanceFromHome       1470 non-null   int64
 6   Education              1470 non-null   object
 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   object
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   object
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: 482.1+ KB
```

```
In [6]: #Checking for Null Values.
df.isnull().any()
```

```
Out[6]:
Age                   False
Attrition              False
BusinessTravel         False
DailyRate              False
Department             False
DistanceFromHome       False
Education              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
```

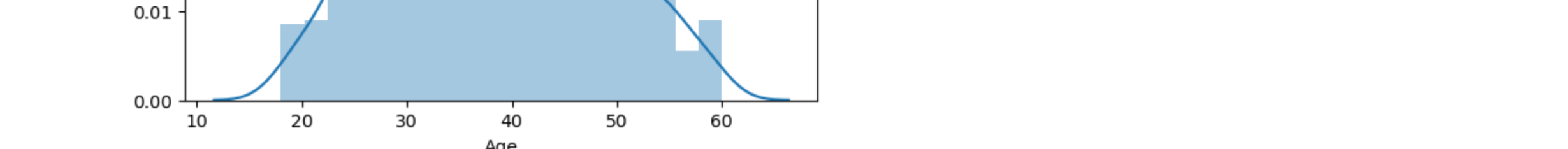
```
In [7]: df.isnull().sum()

Out[7]:
Age                   0
Attrition              0
BusinessTravel         0
DailyRate              0
Department             0
DistanceFromHome       0
Education              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
```



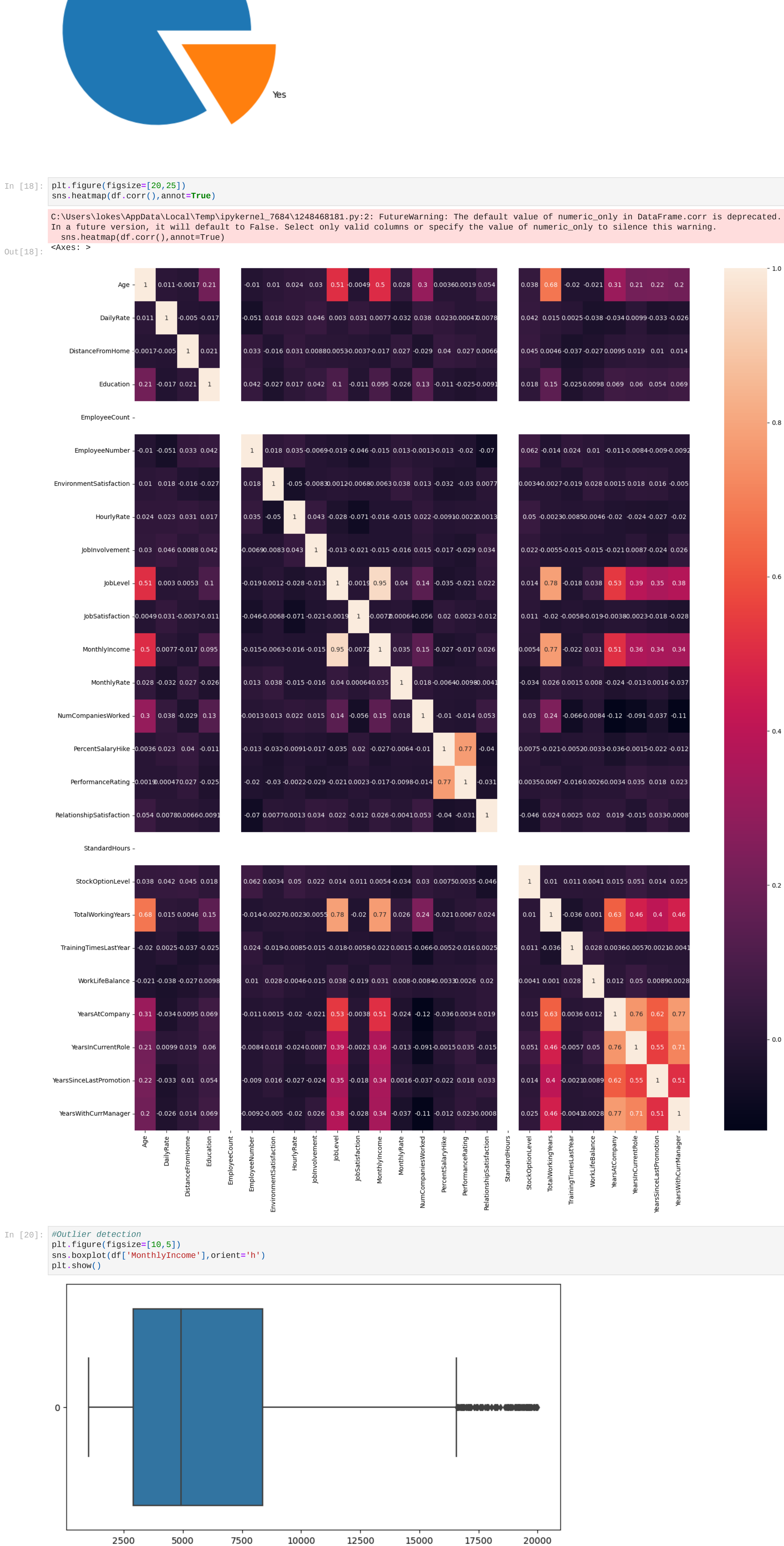
```
In [13]: attrition_count = pd.DataFrame(df['Attrition'].value_counts())
plt.pie(attrition_count['Attrition'], labels = ['No', 'Yes'], explode = (0.3,0))
```

```
Out[13]:
([<matplotlib.patches.Wedge at 0x2844f428689>,
<matplotlib.patches.Wedge at 0x2844f428689>],
[Text(-1.2242257859135192, 0.6791695473696332, 'No'),
Text(0.961891673217765, -0.5336332157899547, 'Yes')])
```



```
In [18]: plt.figure(figsize=[20,25])
sns.heatmap(df.corr(),annot=True)
```

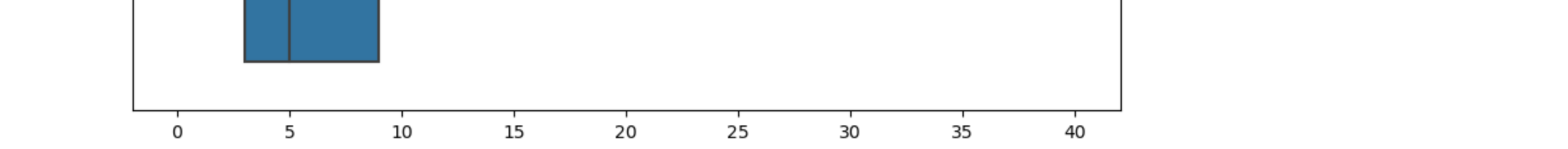
```
Out[18]:
C:\Users\lokes\AppData\Local\Temp\ipykernel_7684\1248468181.py:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
sns.heatmap(df.corr(),annot=True)
```



```
In [20]: #Outlier detection
plt.figure(figsize=[10,5])
sns.boxplot(df['MonthlyIncome'],orient='h')
```



```
In [21]: plt.figure(figsize=[10,5])
sns.boxplot(df['YearsAtCompany'],orient='h')
```



```
In [22]: # Label Encoding
categories = ['BusinessTravel','Department','Education','EducationField','Gender','MaritalStatus','OverTime',
'EnvironmentSatisfaction','JobInvolvement','JobLevel','JobRole','JobSatisfaction','NumCompaniesWorked',
'PerformanceRating','RelationshipSatisfaction','StockOptionLevel','TrainingTimesLastYear','WorkLifeBalance']
categories = df[categories].astype('object')
categorical = pd.get_dummies(df[categories], drop_first = True)
```

```
In [23]: # Splitting Dependent and Independent variables
independent = ['Attrition','Over18','EmployeeCount','StandardHours','EmployeeNumber']
continuous = df.drop(columns=categorical)
continuous = continuous.drop(columns=independent)
```

```
In [24]: # X = features, Y= Target variables
X = pd.concat([categorical,continuous],axis=1)
Y = df['Attrition'].replace({'Yes': 1, 'No': 0}).values.reshape(-1,1)
```

```
In [25]: # Feature scaling
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
continuous_variables = list(continuous.columns)
X = X.reset_index()
del X['index']
X[continuous_variables] = pd.DataFrame(scaler.fit_transform(X[continuous_variables]), columns = continuous_variables)
```

```
In [26]: #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,Y,test_size=0.2,random_state=0)
```

```
In [27]: # (1176, 44), (294, 44), (1176, 1), (294, 1))
Out[27]:
((1176, 44), (294, 44), (1176, 1), (294, 1))
```

Logistic Regression Model

```
In [28]: #Importing necessary libraries
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score,precision_score, recall_score, f1_score,confusion_matrix,classification_report,roc_auc_score,roc_curve
```

```
In [29]: #Initializing the model
lr = LogisticRegression()
```

```
In [30]: #Training the model
lr.fit(x_train,y_train)
```

```
C:\Users\lokes\vscode\workspace\jupyter\Lib\site-packages\sklearn\utils\validation.py:1084: DataConversionWarning: A column-vector y was passed when a 2d array was expected. Please change the shape of y to (n_samples, ), for example using ravel().
y = column_or_1d(y, warn=True)
C:\Users\lokes\vscode\workspace\jupyter\Lib\site-packages\sklearn\linear_model\_logistic.py:468: ConvergenceWarning: lbfgs failed to converge (at step=1): TOTAL NO. of ITERATIONS REACHED LIMIT.
Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_1 = 1, check optimize_result()
n.LogisticRegression
LogisticRegression()
```

```
In [31]: #Testing the model
y_pred = lr.predict(x_test)
```

```
In [32]: # Accuracy score
print("Accuracy of Logistic regression model: ",accuracy_score(y_test,y_pred))
```

```
In [33]: # Precision score
precision_yes = precision_score(y_test, y_pred, pos_label=1)
print("Precision (Yes): " + str(round(precision_yes, 2)))
precision_no = precision_score(y_test, y_pred, pos_label=0)
print("Precision (No): " + str(round(precision_no, 2)))
```

```
In [34]: # Recall score
recall_yes = recall_score(y_test, y_pred, pos_label=1)
print("Recall (Yes): " + str(round(recall_yes, 2)))
recall_no = recall_score(y_test, y_pred, pos_label=0)
print("Recall (No): " + str(round(recall_no, 2)))
```

```
In [35]: # F1 score
f1_score_yes = f1_score(y_test, y_pred, pos_label=1)
print("F1 Score (Yes): " + str(round(f1_score_yes, 2)))
f1_score_no = f1_score(y_test, y_pred, pos_label=0)
print("F1 Score (No): " + str(round(f1_score_no, 2)))
```

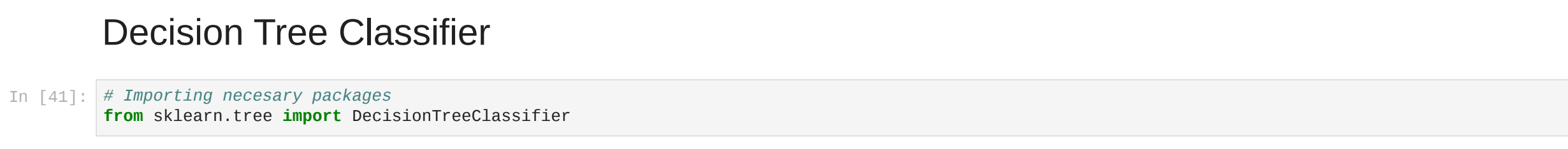
```
In [36]: # Confusion matrix
print("Confusion matrix:\n\n",confusion_matrix(y_test,y_pred))
```

```
In [37]: Confusion matrix:
[[230 7]
 [27 22]]
```

```
In [38]: # Classification Report
print("Classification report of Logistic Regression model:\n\n",classification_report(y_test,y_pred))
```

	precision	recall	f1-score	support
0	0.90	0.97	0.93	245
1	0.76	0.45	0.56	49
accuracy			0.89	294
macro avg	0.83	0.71	0.75	294
weighted avg	0.87	0.88	0.87	294

```
In [40]: # ROC curve
probability = lr.predict_proba(x_test)[:,:1]
fpr,tpr,thresholds = roc_curve(y_test,probability)
plt.plot(fpr,tpr)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of Logistic regression model')
plt.show()
```



Decision Tree Classifier

```
In [41]: # Importing necessary packages
from sklearn.tree import DecisionTreeClassifier
```

```
In [42]: # Initializing the model
dtc = DecisionTreeClassifier(random_state=30)
```

```
In [43]: # Training the model
etc.fit(x_train, y_train)
```

```
Out[43]:
DecisionTreeClassifier
```

```
In [44]: # Testing the model
y_pred2 = dtc.predict(x_test)
```

```
In [45]: # Accuracy metrics
# Accuracy score
accuracy = accuracy_score(y_test, y_pred2)
print("Accuracy of Decision tree model: ",accuracy)
```

```
Accuracy of Random Forest model: 0.7517068682721088
```

```
In [46]: # Precision score
precision_yes = precision_score(y_test, y_pred2, pos_label=1)
print("Precision (Yes): " + str(round(precision_yes, 2)))
precision_no = precision_score(y_test, y_pred2, pos_label=0)
print("Precision (No): " + str(round(precision_no, 2)))
```

```
Precision (Yes): 0.27
Precision (No): 0.86
```

```
In [47]: # Recall score
recall_yes = recall_score(y_test, y_pred2, pos_label=1)
print("Recall (Yes): " + str(round(recall_yes, 2)))
recall_no = recall_score(y_test, y_pred2, pos_label=0)
print("Recall (No): " + str(round(recall_no, 2)))
```

```
Recall (Yes): 0.29
Recall (No): 0.84
```

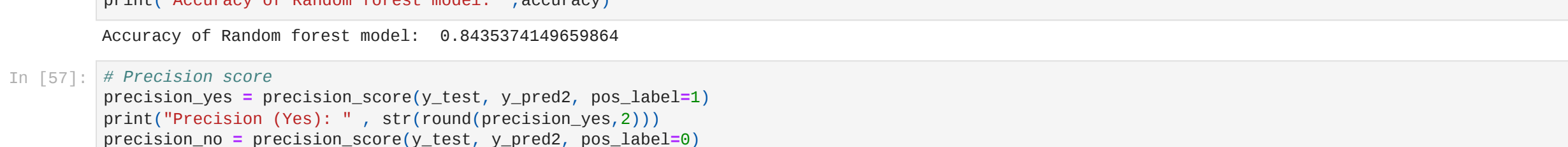
```
In [48]: # F1 score
f1_score_yes = f1_score(y_test, y_pred2, pos_label=1)
print("F1 Score (Yes): " + str(round(f1_score_yes, 2)))
f1_score_no = f1_score(y_test, y_pred2, pos_label=0)
print("F1 Score (No): " + str(round(f1_score_no, 2)))
```

```
F1 Score (Yes): 0.28
F1 Score (No): 0.85
```

```
In [49]: # Classification report
print("Classification report of Decision tree model:\n\n",classification_report(y_test,y_pred2))
```

	precision	recall	f1-score	support
0	0.86	0.84	0.85	245
1	0.71	0.10	0.18	49
accuracy			0.84	294
macro avg	0.78	0.55	0.55	294
weighted avg	0.76	0.57	0.75	294

```
In [50]: # ROC curve
probability = rf.predict_proba(x_test)[:,:1]
fpr,tpr,thresholds = roc_curve(y_test,probability)
plt.plot(fpr,tpr)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of Decision tree model')
plt.show()
```



Random Forest Classifier

```
In [51]: # Importing necessary packages
from sklearn.ensemble import RandomForestClassifier
```

```
In [52]: # Initializing the model
rf = RandomForestClassifier(n_estimators=10, criterion='entropy', random_state=30)
```

```
In [53]: # Training the model
rf.fit(x_train, y_train)
```

```
Out[53]:
RandomForestClassifier
```

```
In [54]: rf.score(x_train, y_train)
```

```
Out[54]:
0.98384337414966
```

```
In [55]: # Testing the model
y_pred2 = rf.predict(x_test)
```

```
In [56]: # Accuracy metrics
# Accuracy score
accuracy = accuracy_score(y_test, y_pred2)
print("Accuracy of Random Forest model: ",accuracy)
```

```
Accuracy of Random Forest model: 0.8435374149659864
```

```
In [57]: # Precision score
precision_yes = precision_score(y_test, y_pred2, pos_label=1)
print("Precision (Yes): " + str(round(precision_yes, 2)))
precision_no = precision_score(y_test, y_pred2, pos_label=0)
print("Precision (No): " + str(round(precision_no, 2)))
```

```
Precision (Yes): 0.71
Precision (No): 0.85
```

```
In [58]: # Recall score
recall_yes = recall_score(y_test, y_pred2, pos_label=1)
print("Recall (Yes): " + str(round(recall_yes, 2)))
recall_no = recall_score(y_test, y_pred2, pos_label=0)
print("Recall (No): " + str(round(recall_no, 2)))
```

```
Recall (Yes): 0.81
Recall (No): 0.89
```

```
In [59]: # F1 score
f1_score_yes = f1_score(y_test, y_pred2, pos_label=1)
print("F1 Score (Yes): " + str(round(f1_score_yes, 2)))
f1_score_no = f1_score(y_test, y_pred2, pos_label=0)
print("F1 Score (No): " + str(round(f1_score_no, 2)))
```

```
F1 Score (Yes): 0.81
F1 Score (No): 0.91
```

```
In [60]: # Classification report
print("Classification report of Random Forest model:\n\n",classification_report(y_test,y_pred2))
```

	precision	recall	f1-score	support
0	0.85	0.99	0.91	245
1	0.71	0.10	0.18	49
accuracy			0.84	294
macro avg	0.78	0.55	0.55	294
weighted avg	0.82	0.84	0.84	294

```
In [61]: # ROC curve
probability = rf.predict_proba(x_test)[:,:1]
fpr,tpr,thresholds = roc_curve(y_test,probability)
plt.plot(fpr,tpr)
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of Random forest model')
plt.show()
```

