

```
[1]: # Logistic regression, Decision Tree and random forest classifiers on Employee Attrition dataset
```

```
[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("W:\Python\HR-Employee-Attrition.csv")
```

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

Out[4]:	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EmployeeNumber	...	RelationshipSatisfaction	StandardHours	StockOptionLevel	TotalWorkingYears	Tri
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	1	1	...	1	80	0	1	8
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	1	2	...	4	80	0	1	10
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	1	4	...	2	80	0	0	7
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	1	5	...	3	80	0	0	8
4	27	No	Travel_Rarely	581	Research & Development	2	1	Medical	1	7	...	4	80	1	0	6

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   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   object
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   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  YearsCurrentRole    1470 non-null   int64
33  YearsSinceLastPromotion 1470 non-null   int64
34  YearsWithCurrManager 1470 non-null   int64
dtypes: int64(26), object(9)
memory usage: 451.1+ 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
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
YearsCurrentRole  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
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
YearsCurrentRole  0
YearsSinceLastPromotion 0
YearsWithCurrManager 0
dtype: int64
```

```
In [11]: #Data visualization.
sns.distplot(df['Age'])

C:\Users\Jokes\AppData\Local\Temp\ipykernel_7684\248007698.py:2: UserWarning:
'distplot' is a deprecated function and will be removed in seaborn v0.14.0.

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

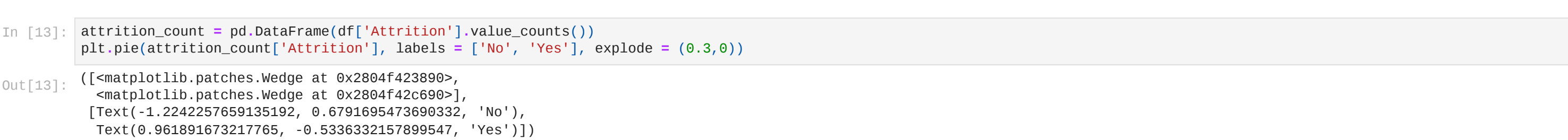
For a guide to updating your code to use the new functions, please see
https://gist.github.com/mwaskon/8e4417e02974578d9372750bbe5761
```

```
Out[11]:
sns.distplot(df['Age'])
<Axes: xlabel='Age', ylabel='Density'>
```



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

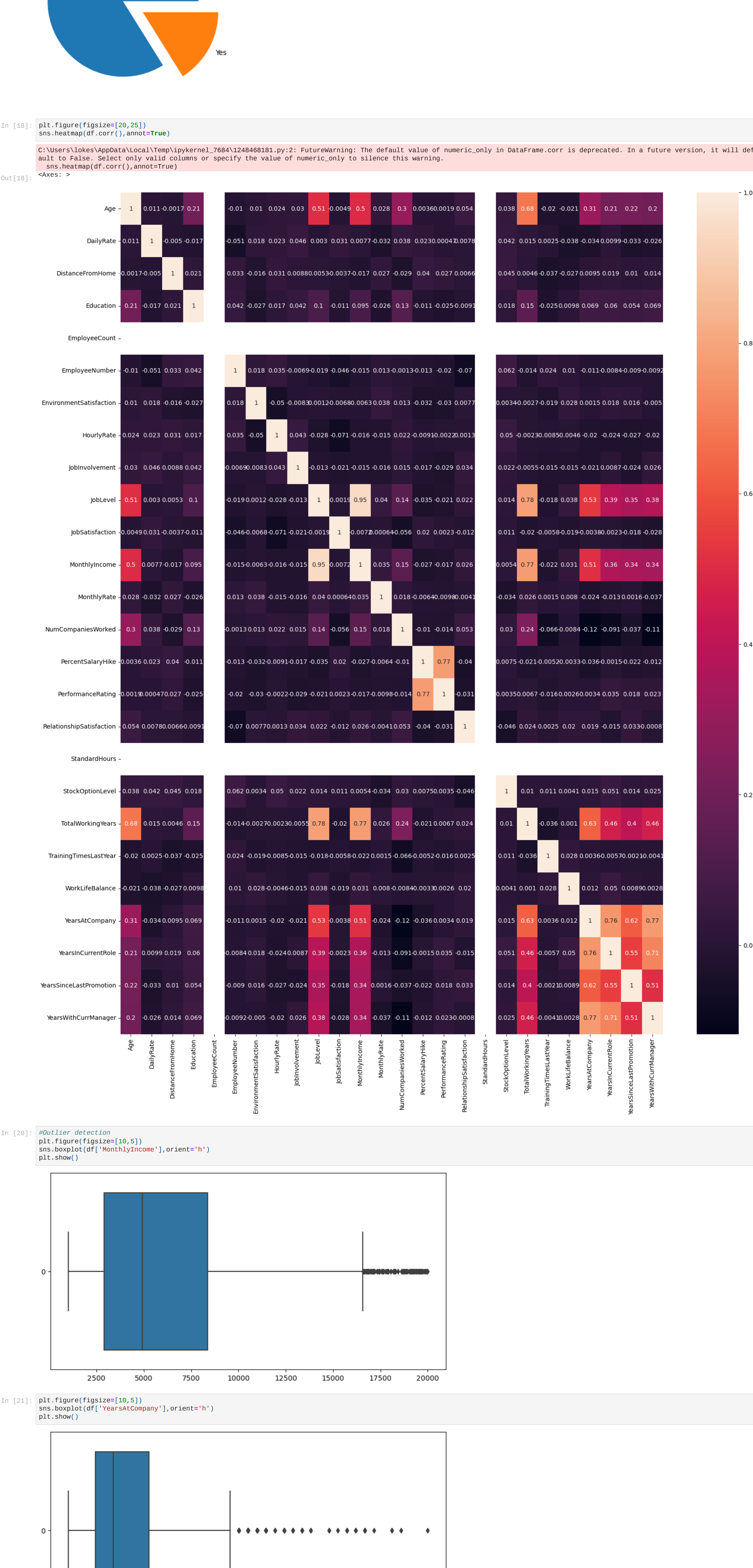
```
Out[13]:
((matplotlib.patches.Wedge at 0x2884f423899d,
matplotlib.patches.Wedge at 0x2884f42c099d),
[Text(-1.2242257693913192, 8.679169647269932, 'No'),
Text(0.8683691231752, -0.93831157695947, 'Yes')])
```



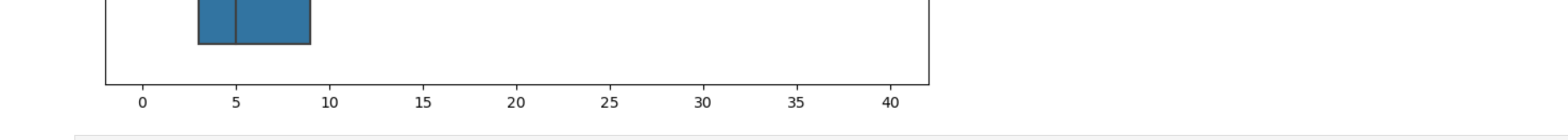
```
In [18]: plt.figure(figsize=(20,25))
sns.heatmap(df.corr(),annot=True)
```

```
C:\Users\Jokes\AppData\Local\Temp\ipykernel_7684\3248408181.py:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will def
ault to None, which means to exclude non-numeric columns. Specify the value of numeric_only to silence this warning.
sns.heatmap(df.corr(),annot=True)
```

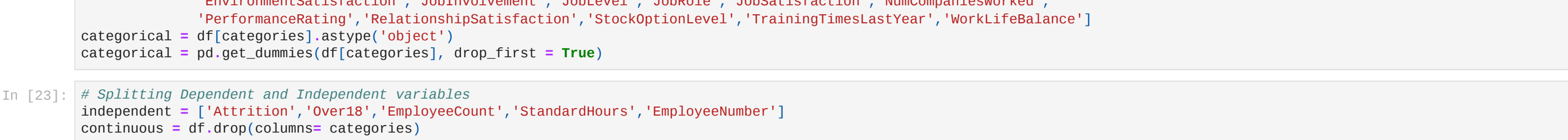
```
Out[18]:
<Axes: >
```



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



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



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

```
In [23]: # Splitting dependent and independent variables
dependent = 'Attrition'
independent = df.drop(columns=dependent)
continuous = df.drop(columns=dependent)
```

```
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_train,X_test = train_test_split(X, Y, test_size=0.2, random_state=0)
```

```
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]: X_train.shape,X_test.shape,y_train.shape,y_test.shape
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\Jokes\venv\workspace\jupyter\lib\site-packages\sklearn\utils\validation.py:1184: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please
change the shape of y to (n_samples,) for example using.ravel().
C:\Users\Jokes\venv\workspace\jupyter\lib\site-packages\sklearn\linear_model\_logistic.py:468: ConvergenceWarning: lbfgs failed to converge (status=1):
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_ = 100
LogisticRegression
```

```
Out[30]:
LogisticRegression()
```

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

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

```
Accuracy of Logistic regression model: 0.8843537414965986
```

```
In [34]: # 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)))
```

```
Precision (Yes): 0.76
Precision (No): 0.9
```

```
In [35]: # 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)))
```

```
Recall (Yes): 0.45
Recall (No): 0.97
```

```
In [36]: # 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)))
```

```
F1 Score (Yes): 0.58
F1 Score (No): 0.93
```

```
In [37]: # Confusion matrix
print('Confusion matrix:\n\n',confusion_matrix(y_test,y_pred))
```

```
[ [28  2]
  [ 2 21] ]
```

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

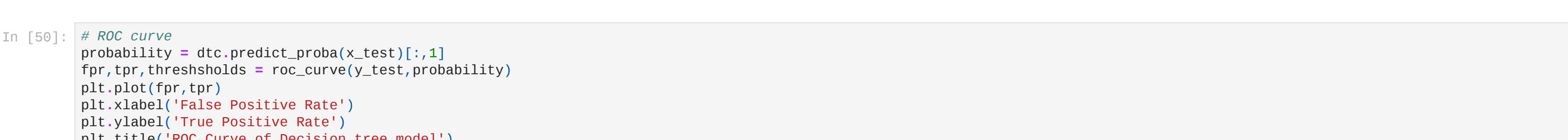
```
Classification report of Logistic Regression model:

```

```
precision    recall  f1-score   support
0           0.86      0.97      0.93      245
1           0.76      0.45      0.56       49
accuracy          0.82      0.71      0.88      294
macro avg          0.83      0.75      0.84      294
weighted avg          0.87      0.88      0.87      294
```

```
In [40]: # ROC curve
probability = lr.predict_proba(X_test)[:,1]
for tpr, thresholds = roc_curve(y_test,probability)
plt.plot(tpr, thresholds)
```

```
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
from sklearn.metrics import accuracy_score
```

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

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

```
Out[43]:
DecisionTreeClassifier
DecisionTreeClassifier(random_state=30)
```

```
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 Decision tree model: 0.7517868682721888
```

```
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): 0.67
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))
```

```
Classification report of Decision Tree model:

```

```
precision    recall  f1-score   support
0           0.86      0.84      0.85      245
1           0.27      0.29      0.28       49
accuracy          0.56      0.57      0.56      294
macro avg          0.56      0.57      0.56      294
weighted avg          0.76      0.75      0.75      294
```

```
In [61]: # ROC curve
probability = rf.predict_proba(X_test)[:,1]
for tpr, thresholds = roc_curve(y_test,probability)
plt.plot(tpr, thresholds)
```

```
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve of Random forest model')
plt.show()
```

