

ASSIGNMENT - 4

Logistic regression, Decision tree and random forest classifiers on Employee Attrition dataset

Data Preprocessing.

```
In [1]: 1 #Importing necessary libraries.
        2 import numpy as np
        3 import pandas as pd
        4 import matplotlib.pyplot as plt
        5 import seaborn as sns
```

```
In [2]: 1 #Importing the dataset.
        2 df=pd.read_csv("Employee-Attrition.csv")
```

```
In [3]: 1 df.head()
```

```
Out[3]:
```

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

5 rows × 35 columns



In [4]: 1 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

```

```
In [5]: 1 #Checking for Null Values.  
        2 df.isnull().any()
```

```
Out[5]: 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 [6]: 1 df.isnull().sum()
```

```
Out[6]: 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 [7]: 1 #Data Visualization.  
2 sns.distplot(df["Age"])
```

C:\Users\Admin\AppData\Local\Temp\ipykernel_39480\2400079689.py:2: UserWarning:

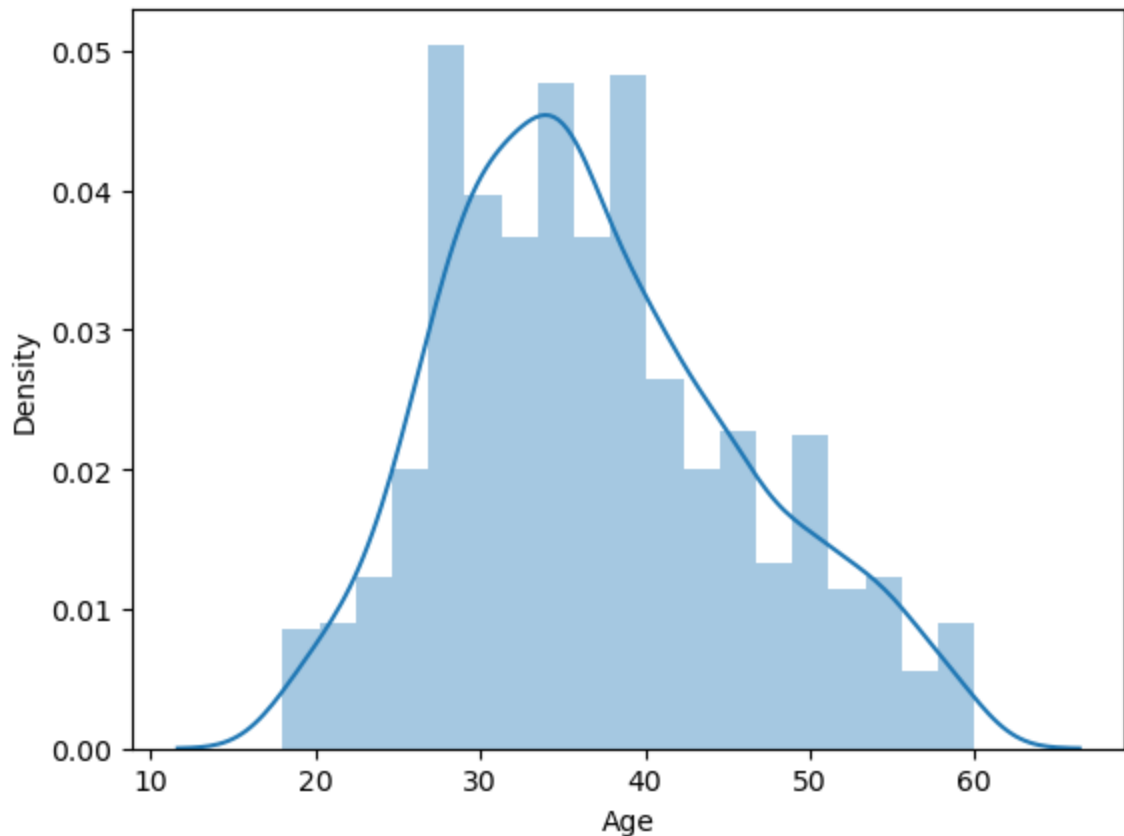
`distplot` is a deprecated function and will be removed in seaborn v0.14.0.

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> (<https://gist.github.com/mwaskom/de44147ed2974457ad6372750bbe5751>)

```
sns.distplot(df["Age"])
```

Out[7]: <Axes: xlabel='Age', ylabel='Density'>



```
In [8]: 1 attrition_count = pd.DataFrame(df['Attrition'].value_counts())  
2 plt.pie(attrition_count['Attrition'], labels = ['No', 'Yes'], explode = (
```

```
Out[8]: ([<matplotlib.patches.Wedge at 0x2634cd0cb80>,  
<matplotlib.patches.Wedge at 0x2634ce105e0>],  
[Text(-1.136781068348268, 0.6306574368426737, 'No'),  
Text(0.961891673217765, -0.5336332157899547, 'Yes')])
```

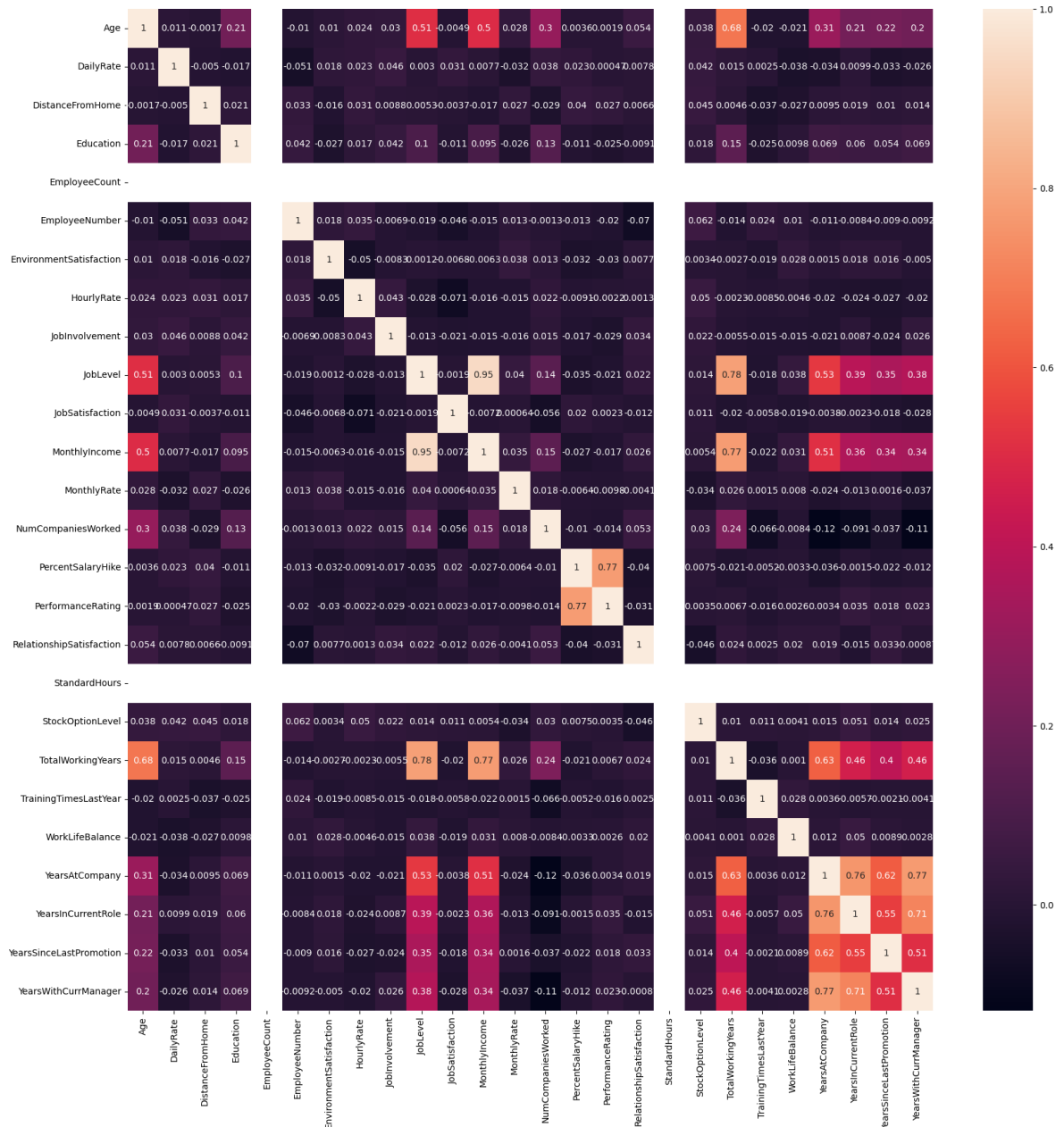


```
In [9]: 1 plt.figure(figsize=[20,20])
        2 sns.heatmap(df.corr(),annot=True)
```

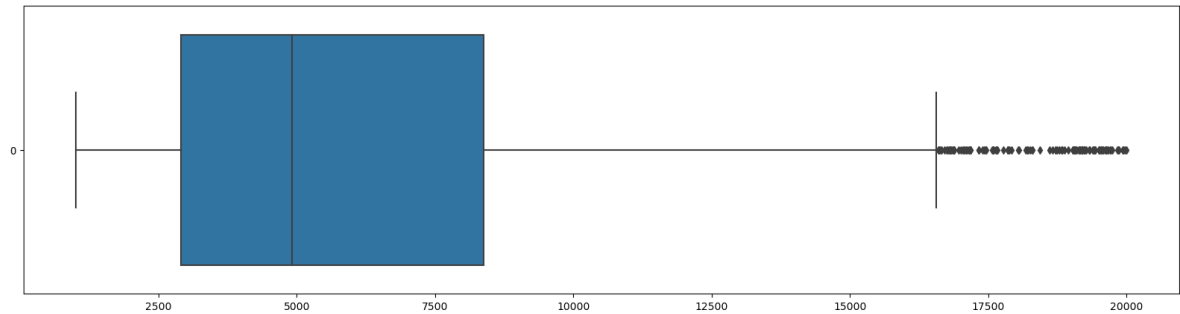
C:\Users\Admin\AppData\Local\Temp\ipykernel_39480\3113117044.py:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

```
sns.heatmap(df.corr(),annot=True)
```

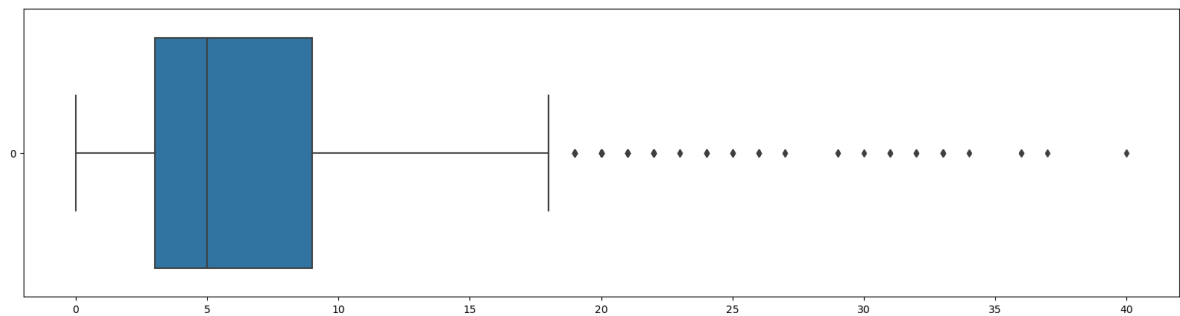
Out[9]: <Axes: >



```
In [10]: 1 #Outlier detection
2 plt.figure(figsize=[20,5])
3 sns.boxplot(df['MonthlyIncome'],orient='h')
4 plt.show()
```



```
In [11]: 1 plt.figure(figsize=[20,5])
2 sns.boxplot(df['YearsAtCompany'],orient='h')
3 plt.show()
```



```
In [12]: 1 # Label Encoding
2 categories = ['BusinessTravel','Department','Education','EducationField',
3             'EnvironmentSatisfaction','JobInvolvement','JobLevel','JobRole',
4             'PerformanceRating','RelationshipSatisfaction','StockOption',
5             'TotalCompensation']
6 categorical = df[categories].astype('object')
7 categorical = pd.get_dummies(df[categories], drop_first = True)
```

```
In [13]: 1 # Splitting Dependent and Independent variables
2 independent = ['Attrition','Over18','EmployeeCount','StandardHours','EmployeeNumber']
3 continuous = df.drop(columns= categories)
4 continuous = continuous.drop(columns= independent)
```

```
In [14]: 1 # X - Features, Y- Target variables
2 X = pd.concat([categorical,continuous],axis=1)
3 Y = df['Attrition'].replace({'Yes': 1, 'No': 0}).values.reshape(-1,1)
```



```
In [15]: 1 # Feature scaling
2 from sklearn.preprocessing import StandardScaler
3
4 scaler = StandardScaler()
5
6 continuous_variables = list(continuous.columns)
7
8 X = X.reset_index()
9 del X['index']
10 X[continuous_variables] = pd.DataFrame(scaler.fit_transform(X[continuous_variables]),
```

```
In [16]: 1 #Splitting Data into Train and Test.
2 from sklearn.model_selection import train_test_split
3 x_train,x_test,y_train,y_test=train_test_split(X,Y,test_size=0.2,random_state=42)
```

```
In [17]: 1 x_train.shape,x_test.shape,y_train.shape,y_test.shape
```

```
Out[17]: ((1176, 44), (294, 44), (1176, 1), (294, 1))
```

Logistic Regression model

```
In [18]: 1 #Importing necessary libraries
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.metrics import accuracy_score,precision_score, recall_score,
```

```
In [19]: 1 #Initializing the model
2 lr = LogisticRegression()
```

```
In [20]: 1 #Training the model
         2 lr.fit(x_train,y_train)
```

C:\Users\Admin\anaconda3\lib\site-packages\sklearn\utils\validation.py:1143: 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().

y = column_or_1d(y, warn=True)
C:\Users\Admin\anaconda3\lib\site-packages\sklearn\linear_model_logistic.py:458: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: 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> (<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 (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)
n_iter_i = _check_optimize_result(

Out[20]: 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 [21]: 1 #Testing the model
         2 y_pred = lr.predict(x_test)
```

```
In [22]: 1 # Evaluation of model
         2 # Accuracy score
         3 print("Accuracy of Logistic regression model:",accuracy_score(y_test,y_pred))
```

Accuracy of Logistic regression model: 0.8843537414965986

```
In [23]: 1 # Precision score
         2 precision_yes = precision_score(y_test, y_pred, pos_label=1)
         3 print("Precision (Yes): " + str(round(precision_yes, 2)))
         4 precision_no = precision_score(y_test, y_pred, pos_label=0)
         5 print("Precision (No): " + str(round(precision_no, 2)))
```

Precision (Yes): 0.76

Precision (No): 0.9

```
In [24]: 1 # Recall score
         2 recall_yes = recall_score(y_test, y_pred, pos_label=1)
         3 print("Recall (Yes): " + str(round(recall_yes, 2)))
         4 recall_no = recall_score(y_test, y_pred, pos_label=0)
         5 print("Recall (No): " + str(round(recall_no, 2)))
```

Recall (Yes): 0.45

Recall (No): 0.97

```
In [25]: 1 # F1 score
2 f1_score_yes = f1_score(y_test, y_pred, pos_label=1)
3 print("F1 Score (Yes): " + str(round(f1_score_yes, 2)))
4 f1_score_no = f1_score(y_test, y_pred, pos_label=0)
5 print("F1 Score (No): " + str(round(f1_score_no, 2)))
```

F1 Score (Yes): 0.56

F1 Score (No): 0.93

```
In [26]: 1 # Confusion matrix
2 print("Confusion matrix:\n\n",confusion_matrix(y_test,y_pred))
```

Confusion matrix:

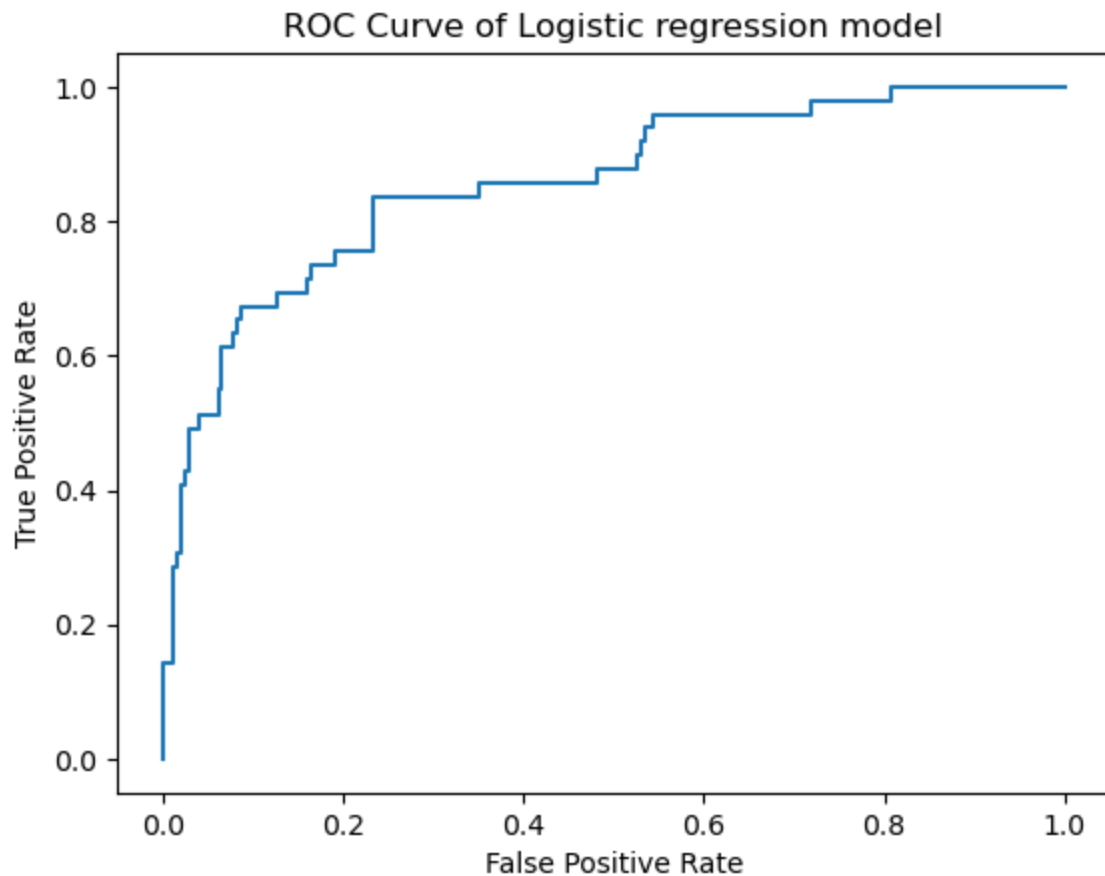
```
[[238  7]
 [ 27 22]]
```

```
In [27]: 1 # Classification Report
2 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	0.90	0.97	0.93	245
1	0.76	0.45	0.56	49
accuracy			0.88	294
macro avg	0.83	0.71	0.75	294
weighted avg	0.87	0.88	0.87	294

```
In [28]: 1 # ROC curve
2 probability = lr.predict_proba(x_test)[:,-1]
3 fpr, tpr, thresholds = roc_curve(y_test, probability)
4 plt.plot(fpr, tpr)
5 plt.xlabel('False Positive Rate')
6 plt.ylabel('True Positive Rate')
7 plt.title('ROC Curve of Logistic regression model')
8 plt.show()
```



Decision Tree Classifier

```
In [29]: 1 # Importing necessary packages
2 from sklearn.tree import DecisionTreeClassifier
```

```
In [30]: 1 # Initializing the model
2 dtc = DecisionTreeClassifier(random_state=30)
```

```
In [31]: 1 # Training the model
        2 dtc.fit(x_train, y_train)
```

Out[31]: DecisionTreeClassifier(random_state=30)

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 [32]: 1 # Testing the model
        2 y_pred1 = dtc.predict(x_test)
```

```
In [33]: 1 # Evaluation metrics
        2 # Accuracy score
        3 accuracy = accuracy_score(y_test, y_pred1)
        4 print("Accuracy of Decision tree model: ",accuracy)
```

Accuracy of Decision tree model: 0.7517006802721088

```
In [34]: 1 # Precision score
        2 precision_yes = precision_score(y_test, y_pred1, pos_label=1)
        3 print("Precision (Yes): " , str(round(precision_yes,2)))
        4 precision_no = precision_score(y_test, y_pred1, pos_label=0)
        5 print("Precision (No): " + str(round(precision_no, 2)))
```

Precision (Yes): 0.27

Precision (No): 0.86

```
In [35]: 1 # Recall score
        2 recall_yes = recall_score(y_test, y_pred1, pos_label=1)
        3 print("Recall (Yes): " + str(round(recall_yes, 2)))
        4 recall_no = recall_score(y_test, y_pred1, pos_label=0)
        5 print("Recall (No): " + str(round(recall_no, 2)))
```

Recall (Yes): 0.29

Recall (No): 0.84

```
In [36]: 1 # F1 score
        2 f1_score_yes = f1_score(y_test, y_pred1, pos_label=1)
        3 print("F1 Score (Yes): " + str(round(f1_score_yes, 2)))
        4 f1_score_no = f1_score(y_test, y_pred1, pos_label=0)
        5 print("F1 Score (No): " + str(round(f1_score_no, 2)))
```

F1 Score (Yes): 0.28

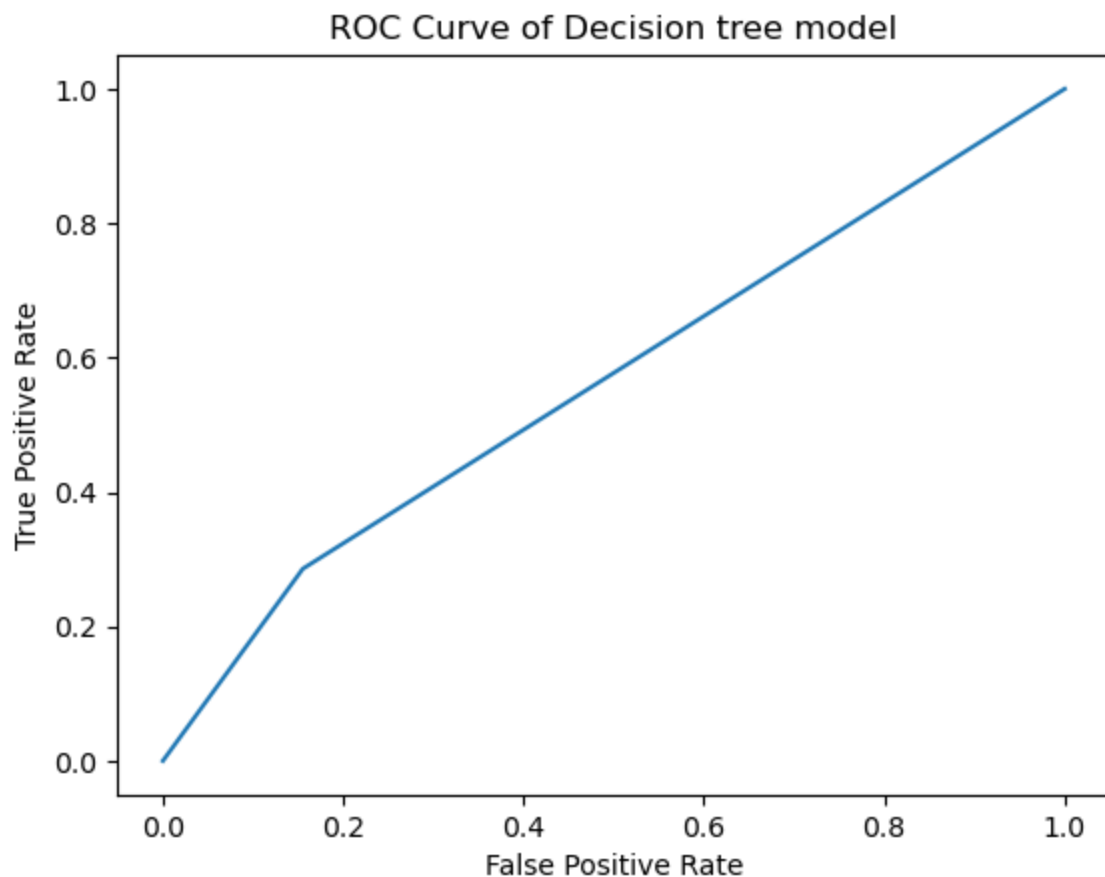
F1 Score (No): 0.85

```
In [37]: 1 # Classification report
2 print("Classification report of Decision tree model:\n\n",classification_
```

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			0.75	294
macro avg	0.56	0.57	0.56	294
weighted avg	0.76	0.75	0.75	294

```
In [38]: 1 # ROC curve
2 probability = dtc.predict_proba(x_test)[: ,1]
3 fpr, tpr, thresholds = roc_curve(y_test, probability)
4 plt.plot(fpr, tpr)
5 plt.xlabel('False Positive Rate')
6 plt.ylabel('True Positive Rate')
7 plt.title('ROC Curve of Decision tree model')
8 plt.show()
```



Random Forest Classifier

```
In [39]: 1 # Importing necessary packages
          2 from sklearn.ensemble import RandomForestClassifier
          3 from sklearn.metrics import accuracy_score
```

```
In [40]: 1 # Initializing the model
          2 rf = RandomForestClassifier(n_estimators=10, criterion='entropy', random_
```

```
In [41]: 1 # Training the model
          2 rf.fit(x_train, y_train)
```

C:\Users\Admin\AppData\Local\Temp\ipykernel_39480\391630832.py:2: 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().
rf.fit(x_train, y_train)

Out[41]: RandomForestClassifier(criterion='entropy', n_estimators=10, random_state=30)
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 [42]: 1 rf.score(x_train, y_train)
```

Out[42]: 0.983843537414966

```
In [43]: 1 # Testing the model
          2 y_pred2 = rf.predict(x_test)
```

```
In [44]: 1 # Evaluation metrics
          2 # Accuracy score
          3 accuracy = accuracy_score(y_test, y_pred2)
          4 print("Accuracy of Random forest model: ",accuracy)
```

Accuracy of Random forest model: 0.8435374149659864

```
In [45]: 1 # Precision score
          2 precision_yes = precision_score(y_test, y_pred2, pos_label=1)
          3 print("Precision (Yes): " , str(round(precision_yes,2)))
          4 precision_no = precision_score(y_test, y_pred2, pos_label=0)
          5 print("Precision (No): " + str(round(precision_no, 2)))
```

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

```
In [46]: 1 # Recall score
2 recall_yes = recall_score(y_test, y_pred2, pos_label=1)
3 print("Recall (Yes): " + str(round(recall_yes, 2)))
4 recall_no = recall_score(y_test, y_pred2, pos_label=0)
5 print("Recall (No): " + str(round(recall_no, 2)))
```

Recall (Yes): 0.1

Recall (No): 0.99

```
In [47]: 1 # F1 score
2 f1_score_yes = f1_score(y_test, y_pred2, pos_label=1)
3 print("F1 Score (Yes): " + str(round(f1_score_yes, 2)))
4 f1_score_no = f1_score(y_test, y_pred2, pos_label=0)
5 print("F1 Score (No): " + str(round(f1_score_no, 2)))
```

F1 Score (Yes): 0.18

F1 Score (No): 0.91

```
In [48]: 1 # Classification Report
2 print("Classification report of Random Forest model:\n\n", classification_
```

Classification report of Random Forest model:

	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	0.79	294


```
In [49]: 1 # ROC curve
2 probability = rf.predict_proba(x_test)[:,-1]
3 fpr, tpr, thresholds = roc_curve(y_test, probability)
4 plt.plot(fpr, tpr)
5 plt.xlabel('False Positive Rate')
6 plt.ylabel('True Positive Rate')
7 plt.title('ROC Curve of Random forest model')
8 plt.show()
```

