

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import os
import warnings
warnings.filterwarnings('ignore')
%matplotlib inline

os.chdir('C:\Users\utpala mohapatra\Documents\python_folder\python datasets')
```

Load the dataset

```
In [2]: diabetes_df = pd.read_csv('health_care_diabetes.csv')
diabetes_df.head()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

Task 1 - Perform descriptive analysis. Understand the variables and their corresponding values. On the columns below, a value of zero does not make sense and thus indicates missing value.

```
In [3]: diabetes_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Pregnancies            768 non-null    int64
 1   Glucose                768 non-null    int64
 2   BloodPressure          768 non-null    int64
 3   SkinThickness          768 non-null    int64
 4   Insulin                768 non-null    float64
 5   BMI                    768 non-null    float64
 6   DiabetesPedigreeFunction 768 non-null    float64
 7   Age                    768 non-null    int64
 8   Outcome                768 non-null    int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB
```

check for null values

```
In [6]: diabetes_df.isna().sum()

Out[6]: Pregnancies    0
         Glucose        0
         BloodPressure  0
         SkinThickness  0
         Insulin        0
         BMI            0
         DiabetesPedigreeFunction 0
         Age           0
         Outcome       0
         dtype: int64

No null values present in the database.
```

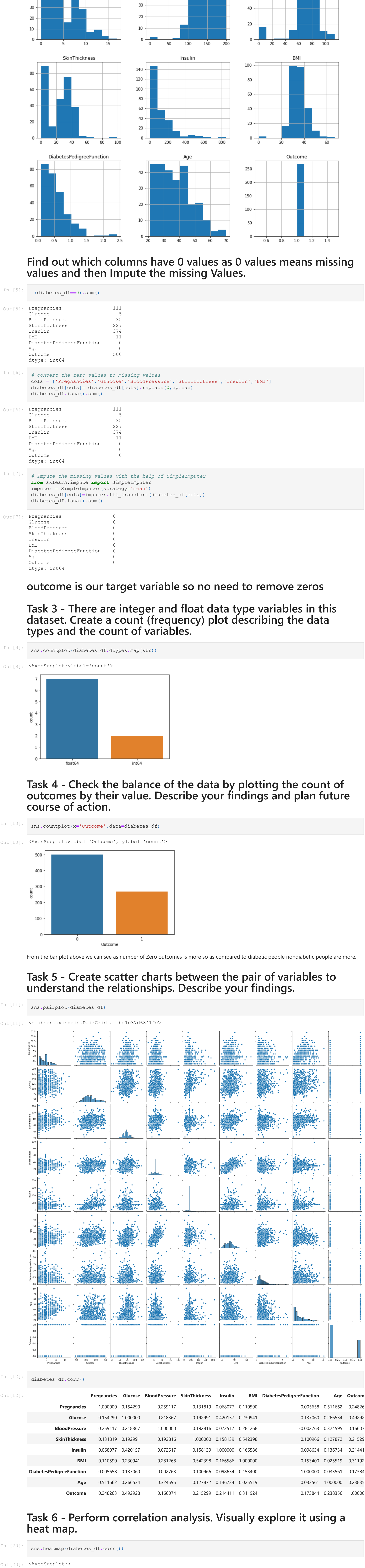
Descriptive Analysis

```
In [7]: diabetes_df.describe()

Out[7]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	0.471876	33.240885	0.348958
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.078000	21.000000	0.000000
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	24.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.372500	29.000000	0.000000
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.626250	41.000000	1.000000
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000	1.000000

Task 2 - Visually explore these variables using histograms. Treat the missing values accordingly.



Find out which columns have 0 values as 0 values means missing values and then Impute the missing Values.

```
In [5]: (diabetes_df==0).sum()

Out[5]: Pregnancies    111
         Glucose        5
         BloodPressure  35
         SkinThickness  227
         Insulin       374
         BMI           11
         DiabetesPedigreeFunction 0
         Age           0
         Outcome       500
         dtype: int64

In [6]: # convert the zero values to missing values
cols = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']
diabetes_df[cols] = diabetes_df[cols].replace(0, np.nan)
diabetes_df.isna().sum()

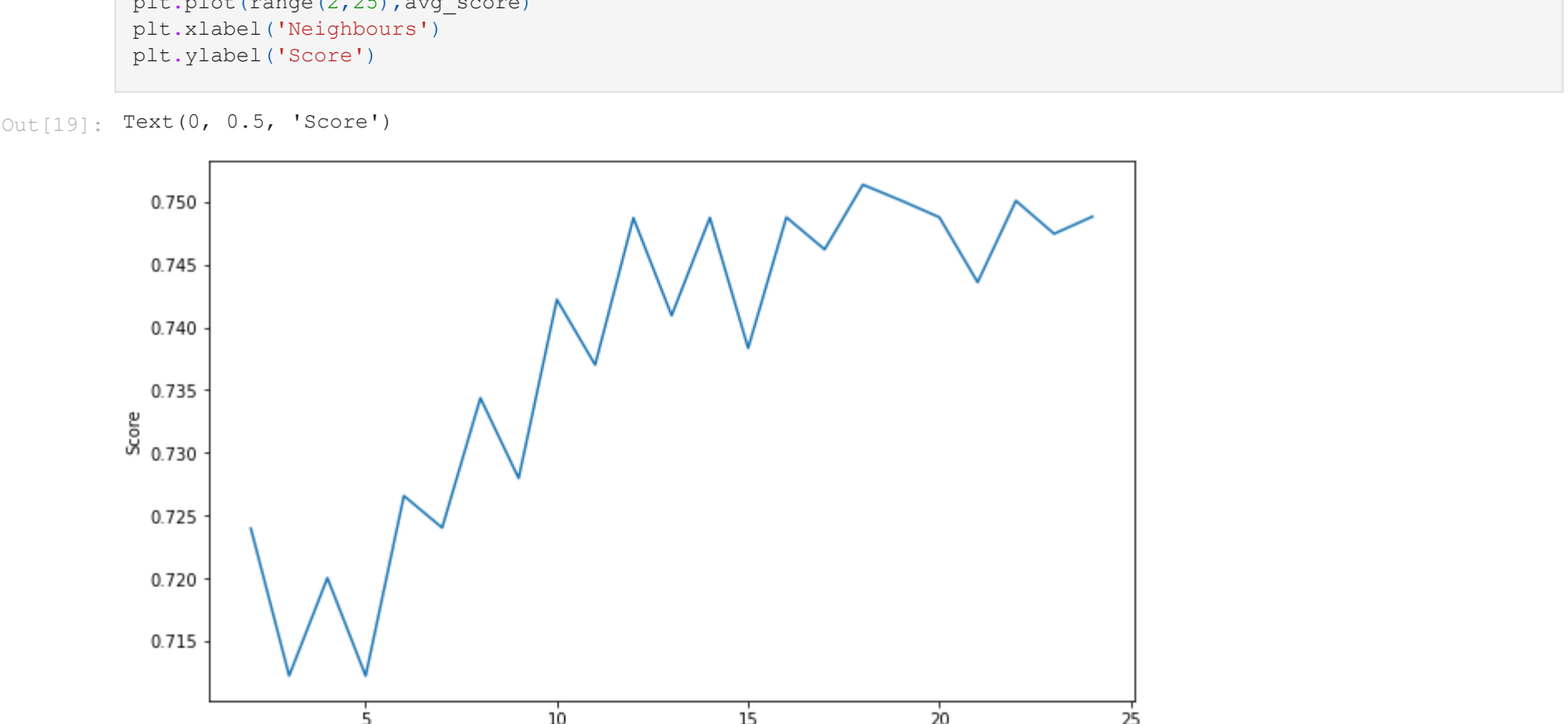
Out[6]: Pregnancies    111
         Glucose        5
         BloodPressure  35
         SkinThickness  227
         Insulin       374
         BMI           11
         DiabetesPedigreeFunction 0
         Age           0
         Outcome       0
         dtype: int64

In [7]: # Impute the missing values with the help of SimpleImputer
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy='mean')
diabetes_df[cols] = imputer.fit_transform(diabetes_df[cols])
diabetes_df.isna().sum()
```

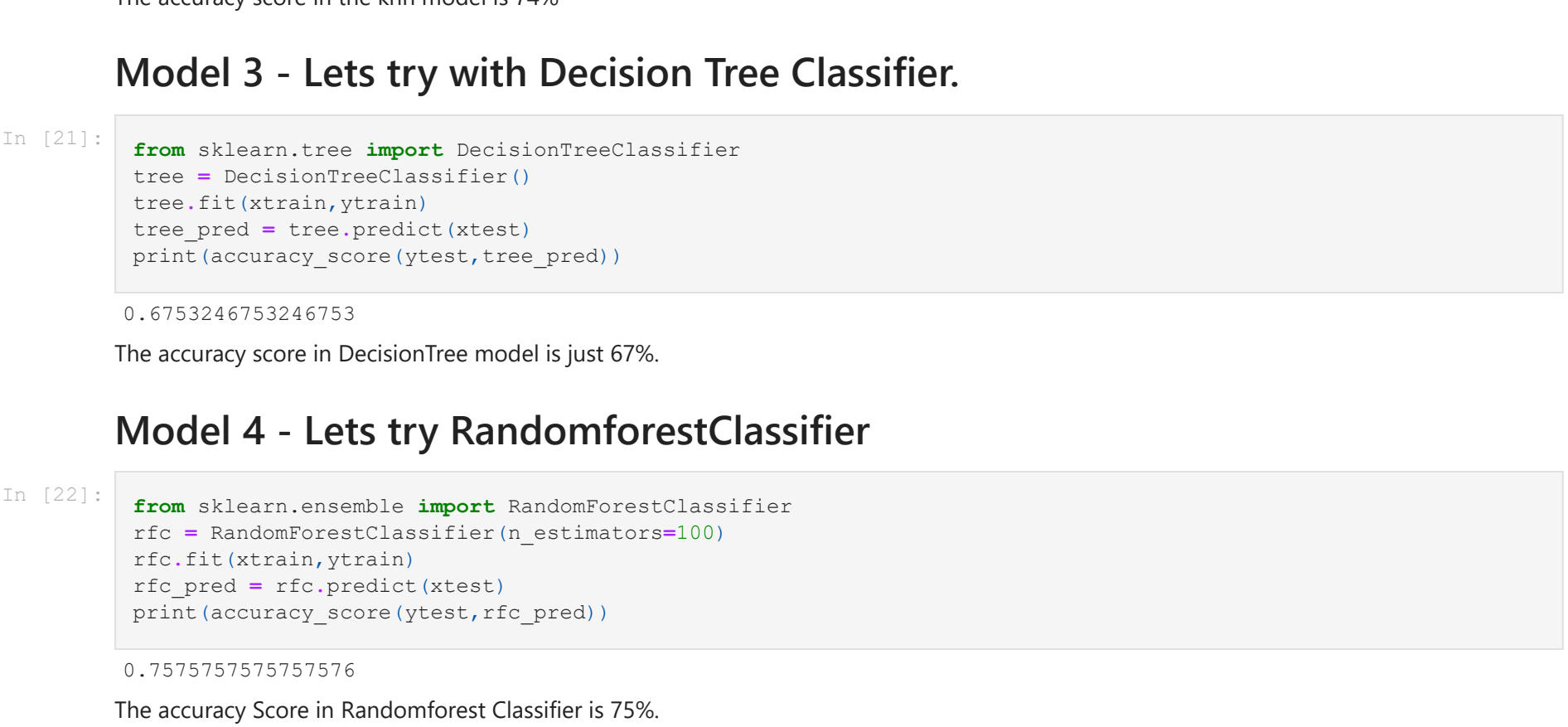
```
Out[7]: Pregnancies    0
         Glucose        0
         BloodPressure  0
         SkinThickness  0
         Insulin        0
         BMI            0
         DiabetesPedigreeFunction 0
         Age           0
         Outcome       0
         dtype: int64
```

outcome is our target variable so no need to remove zeros

Task 3 - There are integer and float data type variables in this dataset. Create a count (frequency) plot describing the data types and the count of variables.

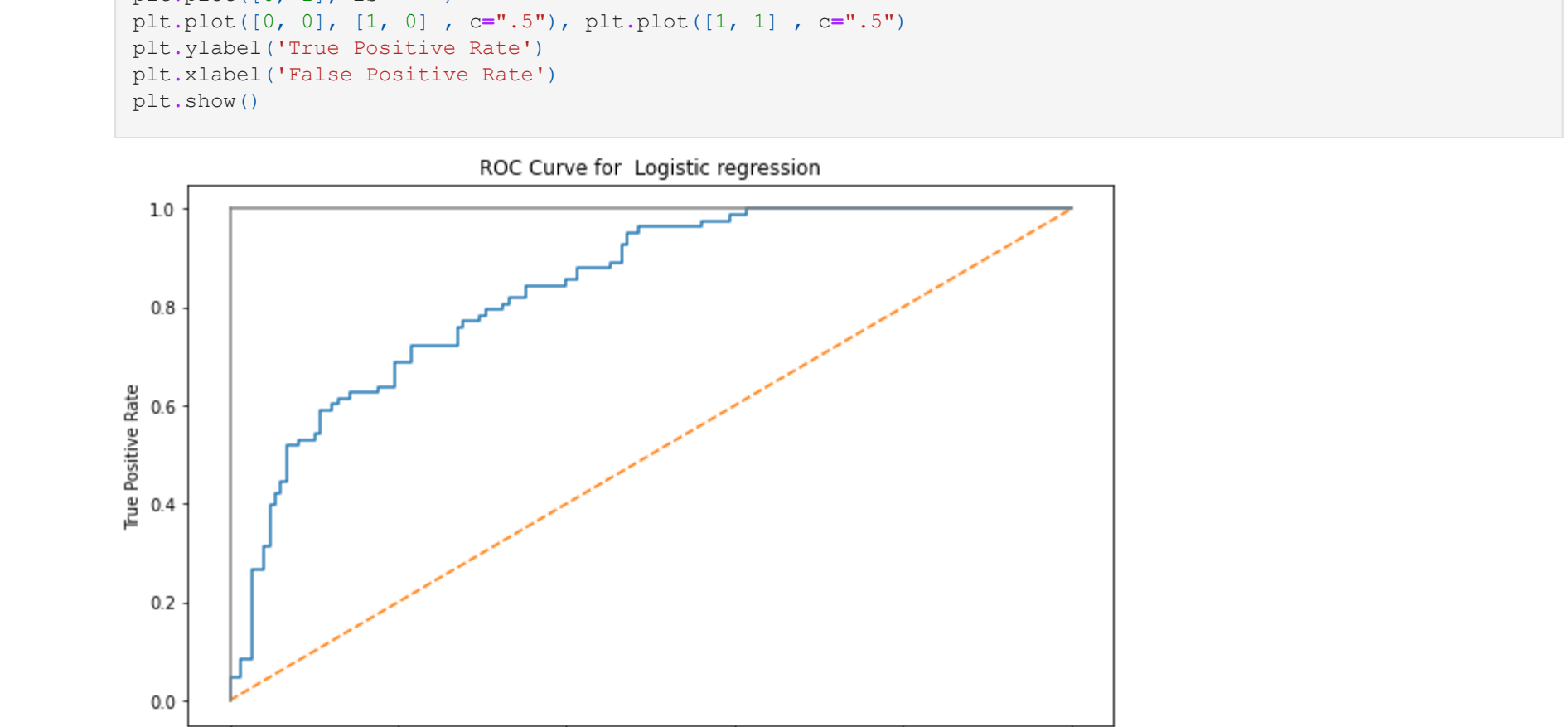
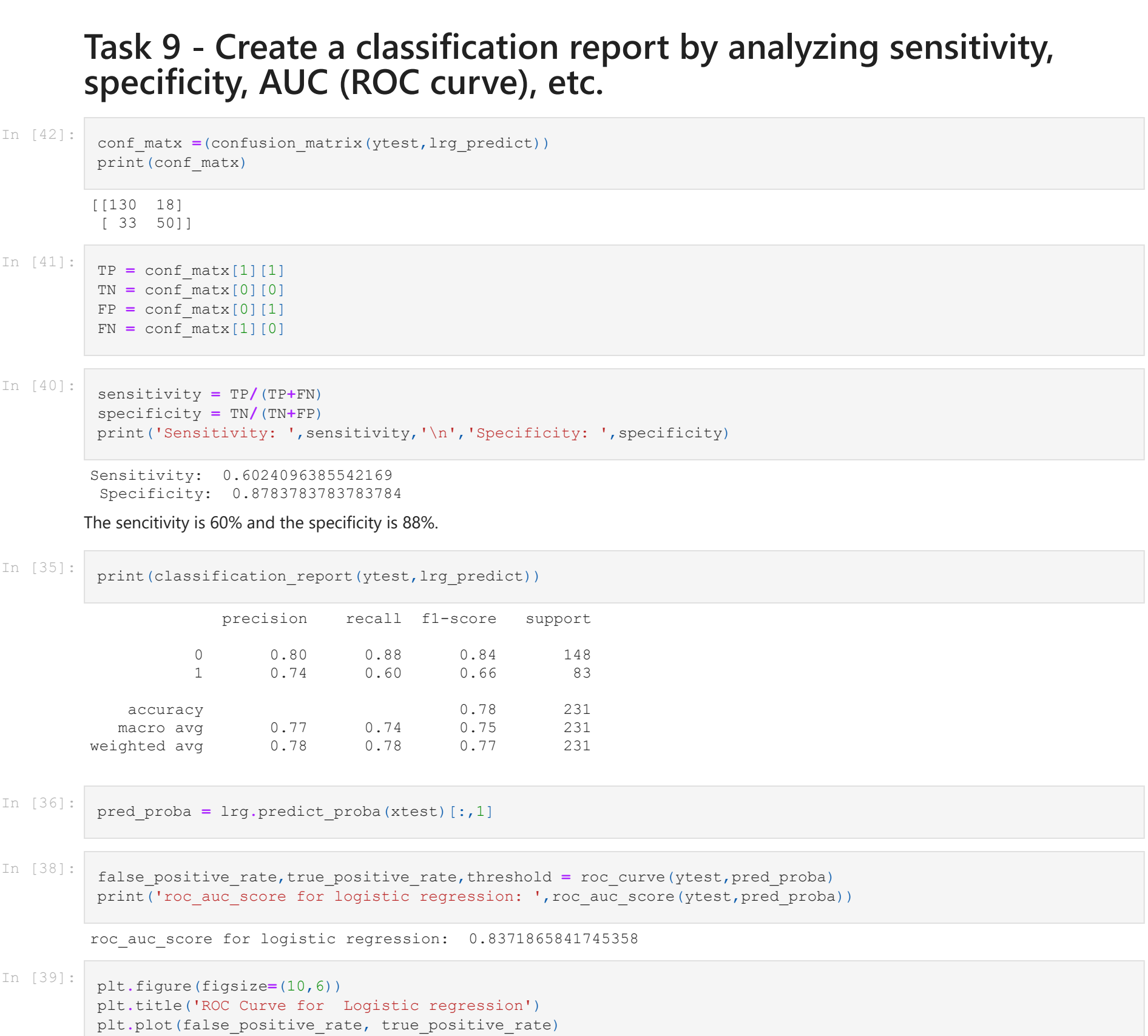


Task 4 - Check the balance of the data by plotting the count of outcomes by their value. Describe your findings and plan future course of action.



From the bar plot above we can see as many Zero outcomes as compared to diabetic people nondiabetic people are more.

Task 5 - Create scatter charts between the pair of variables to understand the relationships. Describe your findings.



Task 6 - Perform correlation analysis. Visually explore it using a heat map.

From the scatter plot and the heatmap we can visualize that outcome of Diabetic data is more correlated with the Glucose level.

Task 7 - Apply an appropriate classification algorithm to build a model.

Model 1 - As the target variable that is Outcome values are binary, so the best possible algorithm should be Logistic Regression.

```
In [13]: # Convert the data into feature and target variables.
X = diabetes_df.drop(['Outcome'], axis=1)
y = diabetes_df['Outcome']

In [14]: # Convert the data to train and test data
from sklearn.model_selection import train_test_split, cross_val_score, KFold
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, roc_auc_score, roc_curve
x_train, x_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=23)
print(x_train.shape, y_train.shape, x_test.shape, y_test.shape)

(537, 8) (537, 1) (231, 8) (231, 1)

In [15]: from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()
lr.fit(x_train, y_train)

Out[15]: LogisticRegression()

In [17]: lrq_predict = lr.predict(x_test)
accuracy_score(y_test, lrq_predict)

Out[17]: 0.7792207792207793

The accuracy score in Logistic Regression is 78%.
```

Task 8 - Compare various models with the results from KNN algorithm.

Model 2 - Lets try KNeighborsClassifier

```
In [18]: from sklearn.neighbors import KNeighborsClassifier

kfold = KFold(n_splits=10)
avg_score = []
for k in range(2, 25):
    knn = KNeighborsClassifier(n_neighbors=k, n_jobs=-1)
    score = cross_val_score(knn, X, y, cv=kfold, scoring='accuracy')
    avg_score.append(score.mean())

In [19]: plt.figure(figsize=(10,6))
plt.plot(range(2,25), avg_score)
plt.xlabel('Neighbours')
plt.ylabel('Score')

Out[19]: Text(0, 0.5, 'Score')
```

We can see in the graph the KNN model gives the best score with the neighbours=17

```
In [20]: knn_17 = KNeighborsClassifier(n_neighbors=17, n_jobs=-1)
knn_17.fit(x_train, y_train)
knn_pred = knn_17.predict(x_test)
print(accuracy_score(y_test, knn_pred))

0.7402597402597403

The accuracy score in the knn model is 74%

Model 3 - Lets try with Decision Tree Classifier.
```

```
In [21]: from sklearn.tree import DecisionTreeClassifier
tree = DecisionTreeClassifier()
tree.fit(x_train, y_train)
tree_pred = tree.predict(x_test)
print(accuracy_score(y_test, tree_pred))

0.6753246753246753

The accuracy score in DecisionTree model is just 67%.
```

Model 4 - Lets try RandomForestClassifier

```
In [22]: from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier(n_estimators=100)
rfc.fit(x_train, y_train)
rfc_pred = rfc.predict(x_test)
print(accuracy_score(y_test, rfc_pred))

0.7575757575757576

The accuracy score in RandomForest Classifier is 75%.
```

Lets try AdaBoost Classifier.

```
In [29]: from sklearn.ensemble import AdaBoostClassifier
ada = AdaBoostClassifier(n_estimators=100)
ada.fit(x_train, y_train)
ada_pred = ada.predict(x_test)
accuracy_score(y_test, ada_pred)

Out[29]: 0.7316017316017316

In AdaBoost model also the Accuracy score is 73%

From all these models we can see RandomForestClassifier is giving better Accuracy than KNN model.
```

Task 9 - Create a classification report by analyzing sensitivity, specificity, AUC (ROC curve), etc.

```
In [42]: conf_matx = (confusion_matrix(y_test, lrq_predict))
print(conf_matx)

[[130  18]
 [ 33  50]]

In [41]: TP = conf_matx[1][1]
TN = conf_matx[0][0]
FP = conf_matx[0][1]
FN = conf_matx[1][0]

In [40]: sensitivity = TP/(TP+FN)
specificity = TN/(TN+FP)
print('Sensitivity: ', sensitivity, '\n', 'Specificity: ', specificity)

Sensitivity: 0.602409695542169
Specificity: 0.8783783783783784

The sensitivity is 60% and the specificity is 88%.
```