# Capstone Project for Health Care
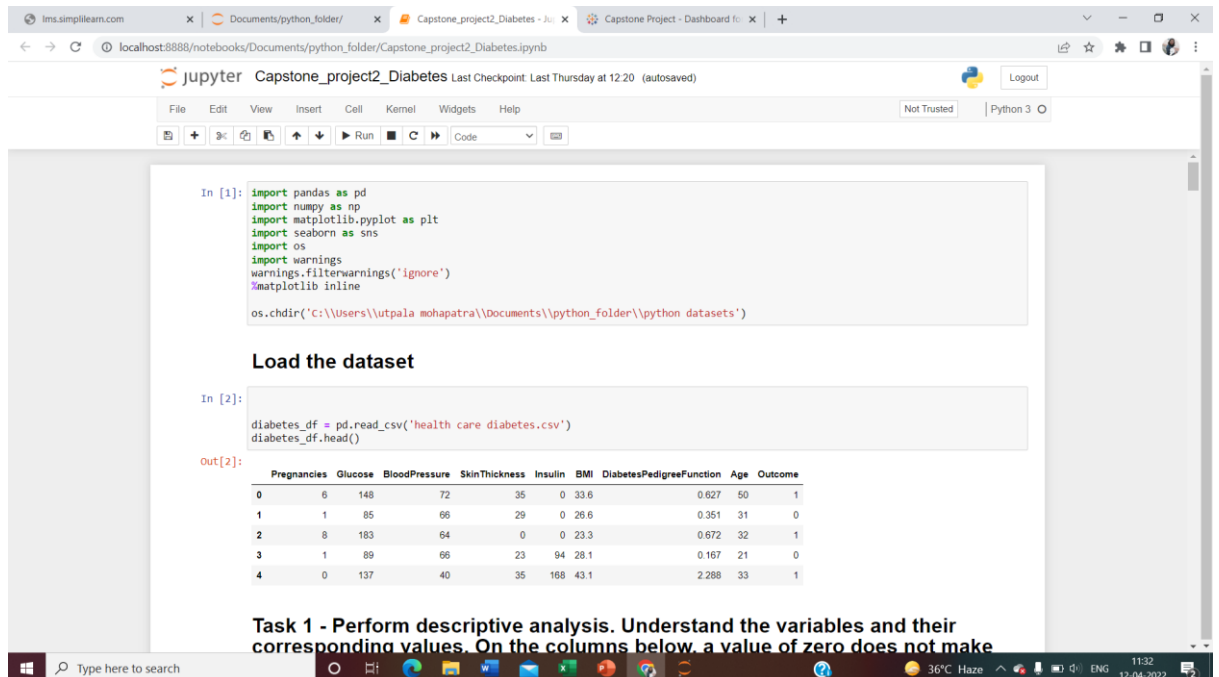
**Description-**NIDDK (National Institute of Diabetes and Digestive and Kidney Diseases) research creates knowledge about and treatments for the most chronic, costly, and consequential diseases.

- The dataset used in this project is originally from NIDDK. The objective is to predict whether or not a patient has diabetes, based on certain diagnostic measurements included in the dataset.
- Build a model to accurately predict whether the patients in the dataset have diabetes or not.

## Data Exploration-

**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:

## Load Dataset

# Descriptive Analysis of the data.



```
Out[6]: Pregnancies                 0
        Glucose                     0
        BloodPressure               0
        SkinThickness               0
        Insulin                     0
        BMI                         0
        DiabetesPedigreeFunction    0
        Age                         0
        Outcome                     0
        dtype: int64
```
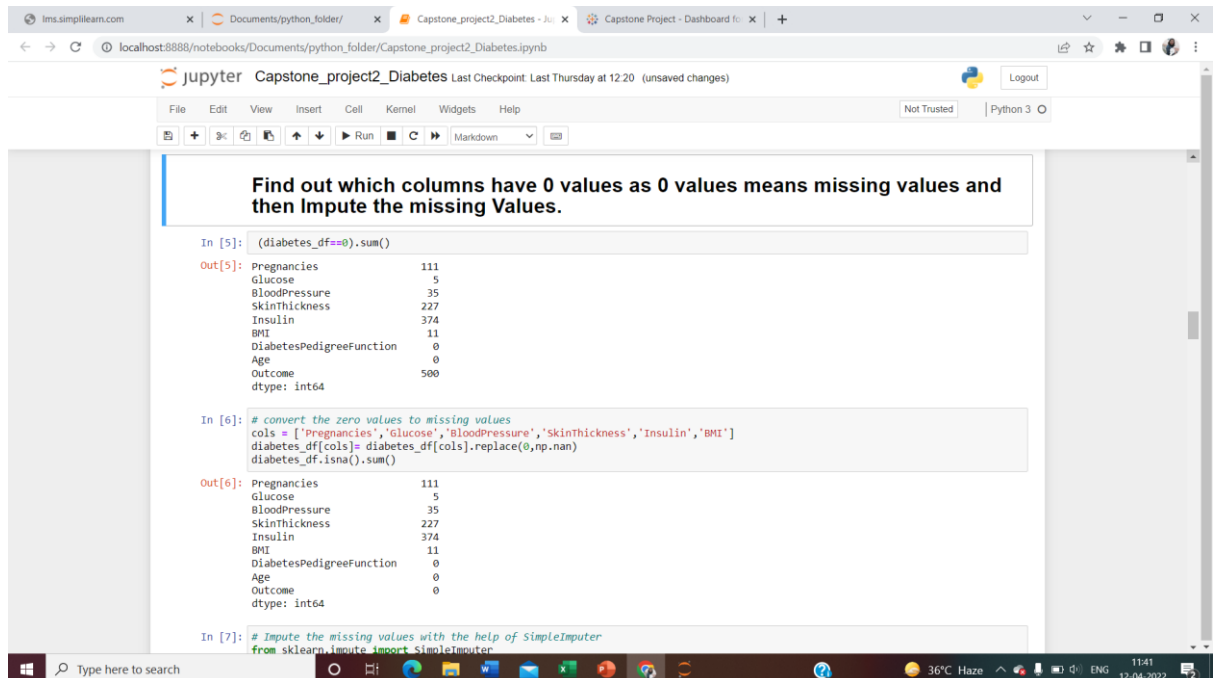
No null values present in the database.

## 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

# Check for missing values.



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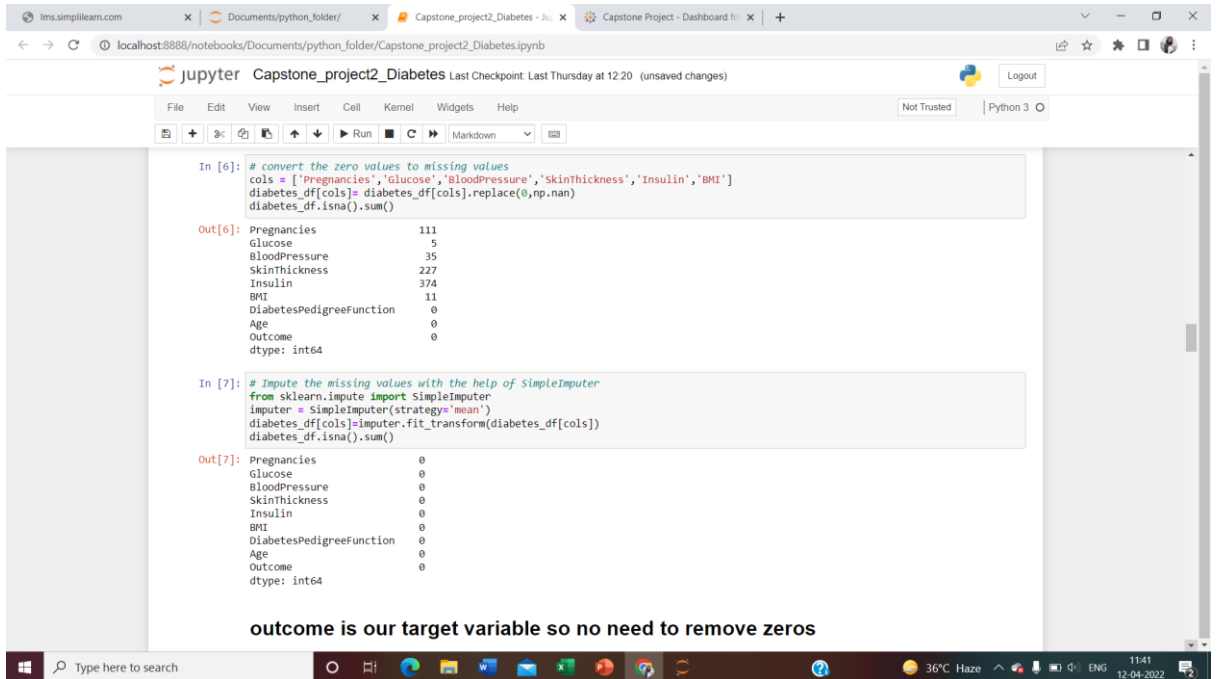
In [6]:
```
# convert the zero values to missing values
cols = ['Pregnancies','Glucose','BloodPressure','SkinThickness','Insulin','BMI']
diabetes_df[cols]= diabetes_df[cols].replace(0,np.nan)
diabetes_df.isna().sum()
```

```
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
```

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

# Impute missing values



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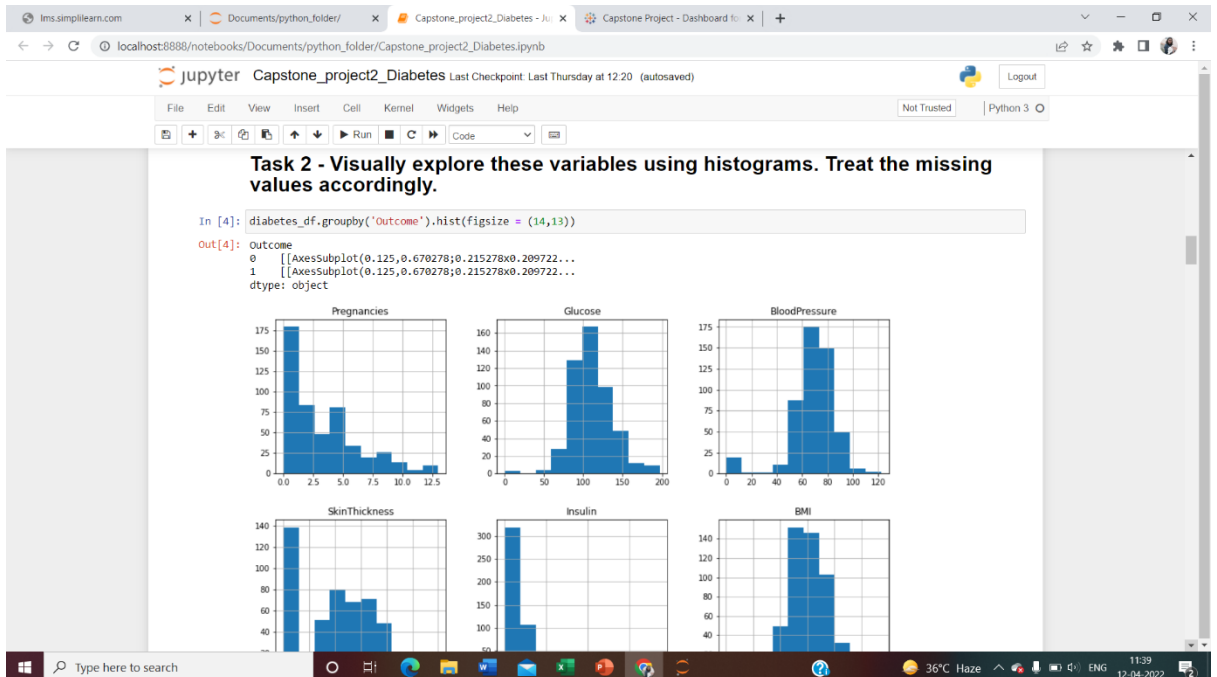
```
In [7]: # Impute the missing values with the help of SimpleImputer
        from sklearn.impute import SimpleImputer
        imputer = SimpleImputer(strategy='mean')
        diabetes_df[cols]=imputer.fit_transform(diabetes_df[cols])
        diabetes_df.isna().sum()

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

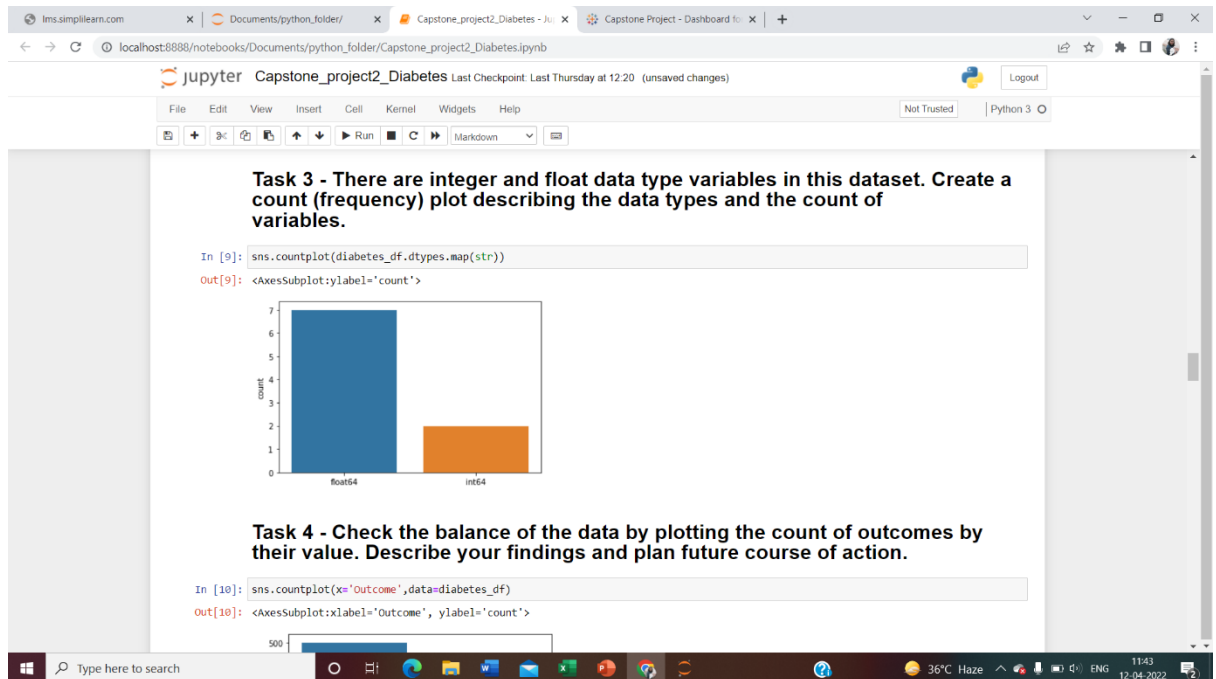**outcome is our target variable so no need to remove zeros**

# Visual analysis by Histograms



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

```
In [4]: diabetes_df.groupby('Outcome').hist(figsize = (14,13))

Out[4]: Outcome
        0    [[AxesSubplot(0.125,0.670278;0.215278x0.209722...
        1    [[AxesSubplot(0.125,0.670278;0.215278x0.209722...
        dtype: object
```
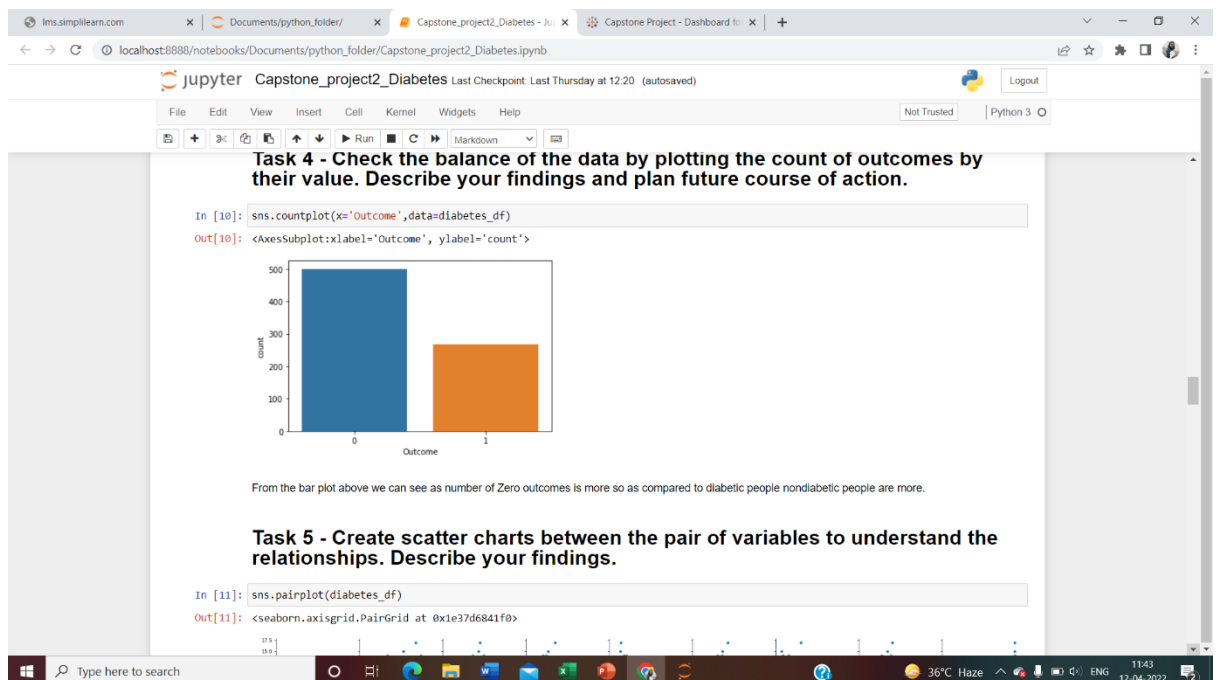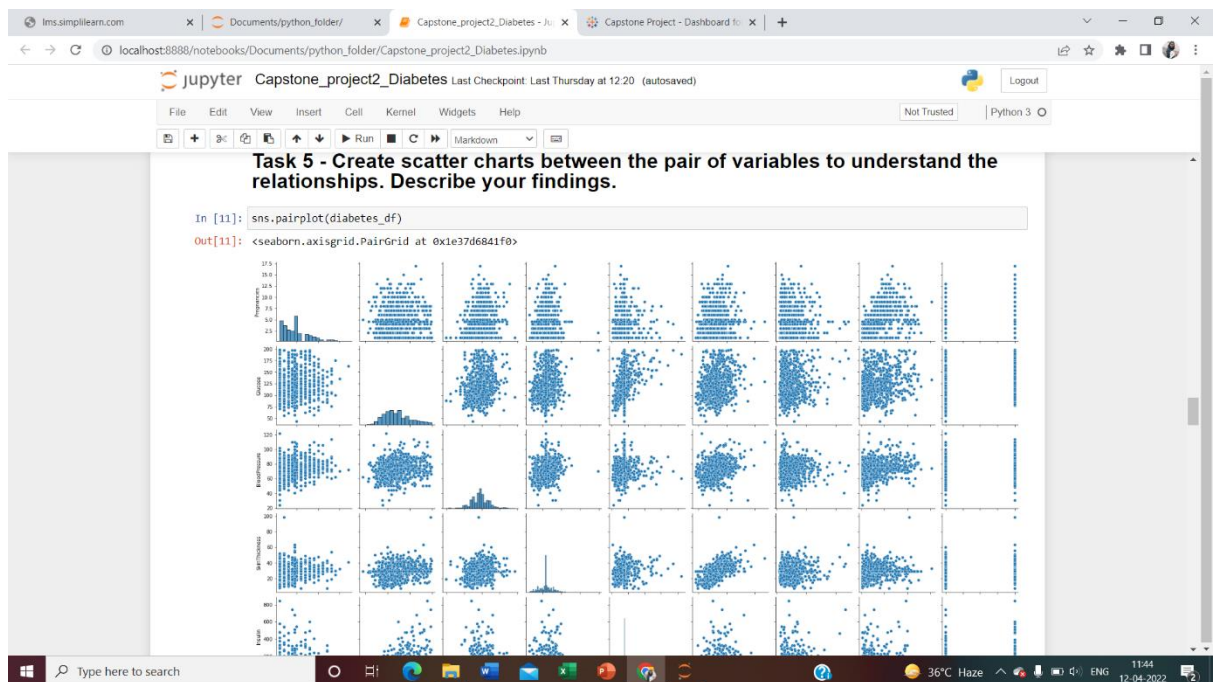
## Task 3 – There are integer and float data type variables in this dataset. Create a count (frequency) plot describing the data types and the count of variables.
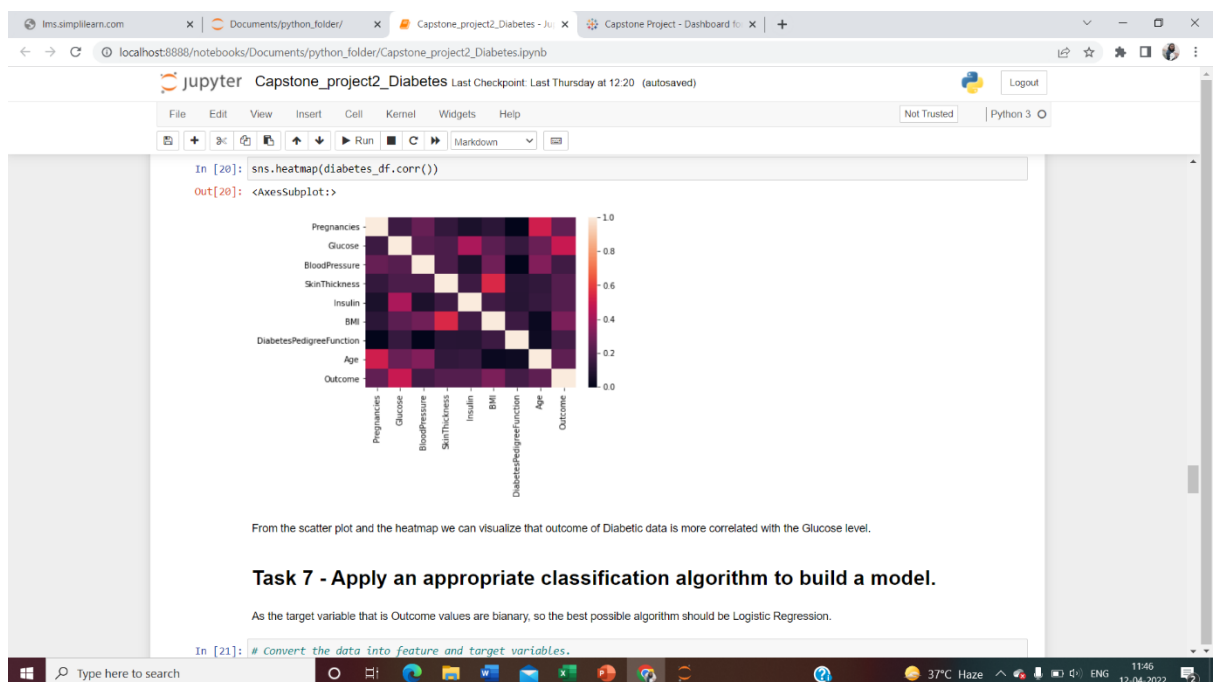


## Task 4 – Check the balance of the data by plotting the count of outcomes by their value. Describe your findings and plan future course of action.

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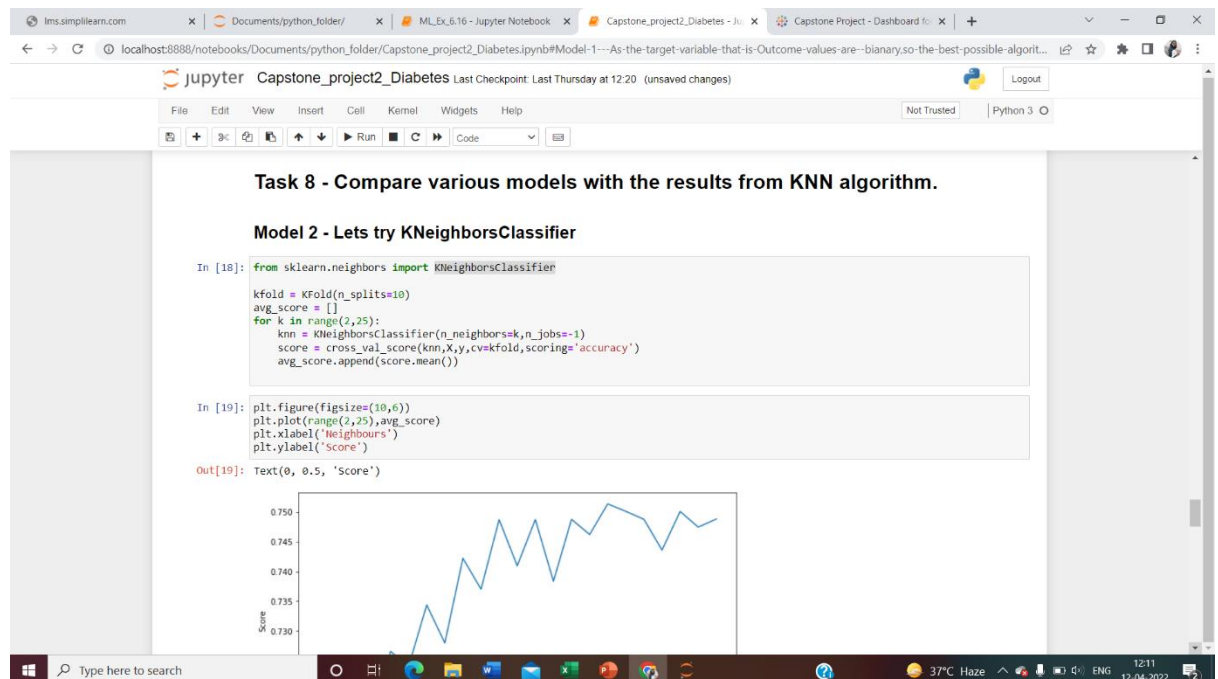


**Task 6 –** Perform correlation analysis. Visually explore it using a heat map.

**Task 7 –** Devise strategies for model building. It is important to decide the right validation framework. Express your thought process.
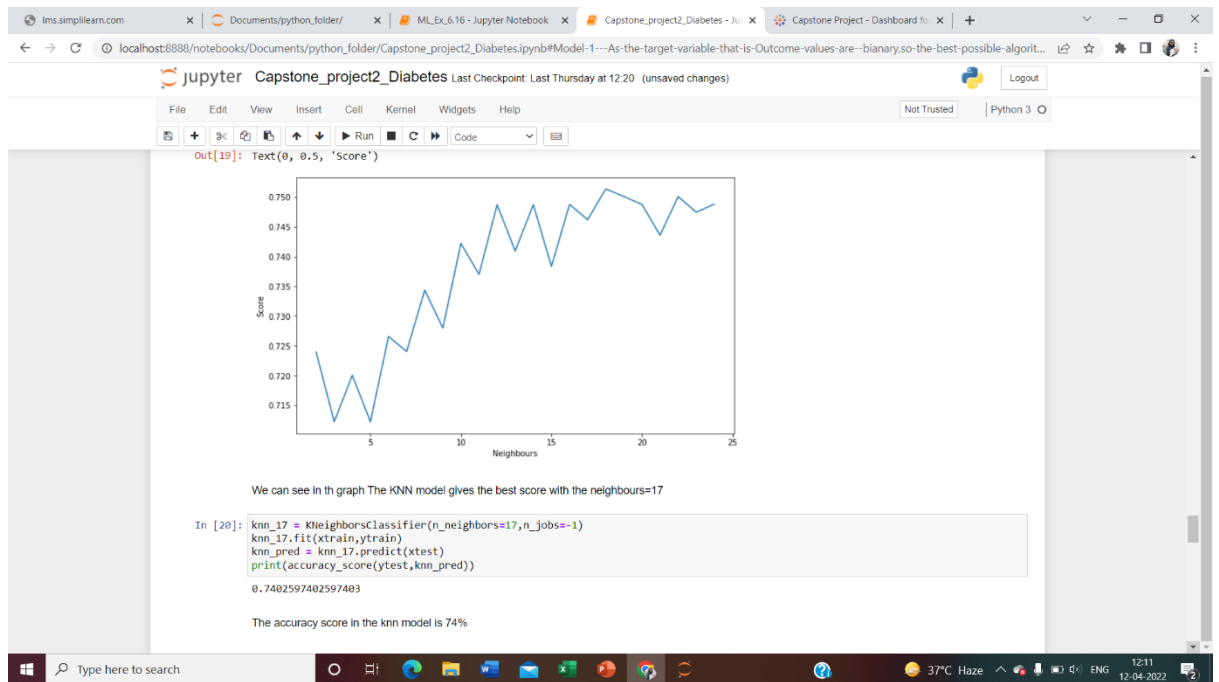
ANS - As the target variable is binary so Logistic regression should be the suitable model for this data to predict best results.



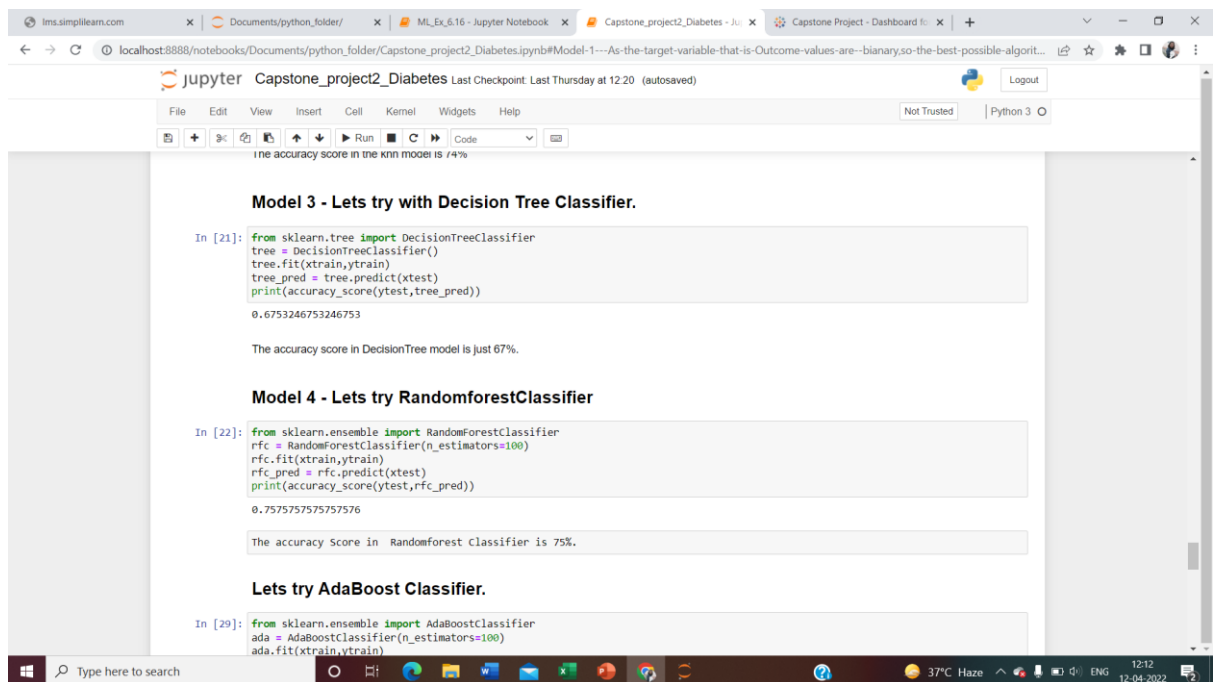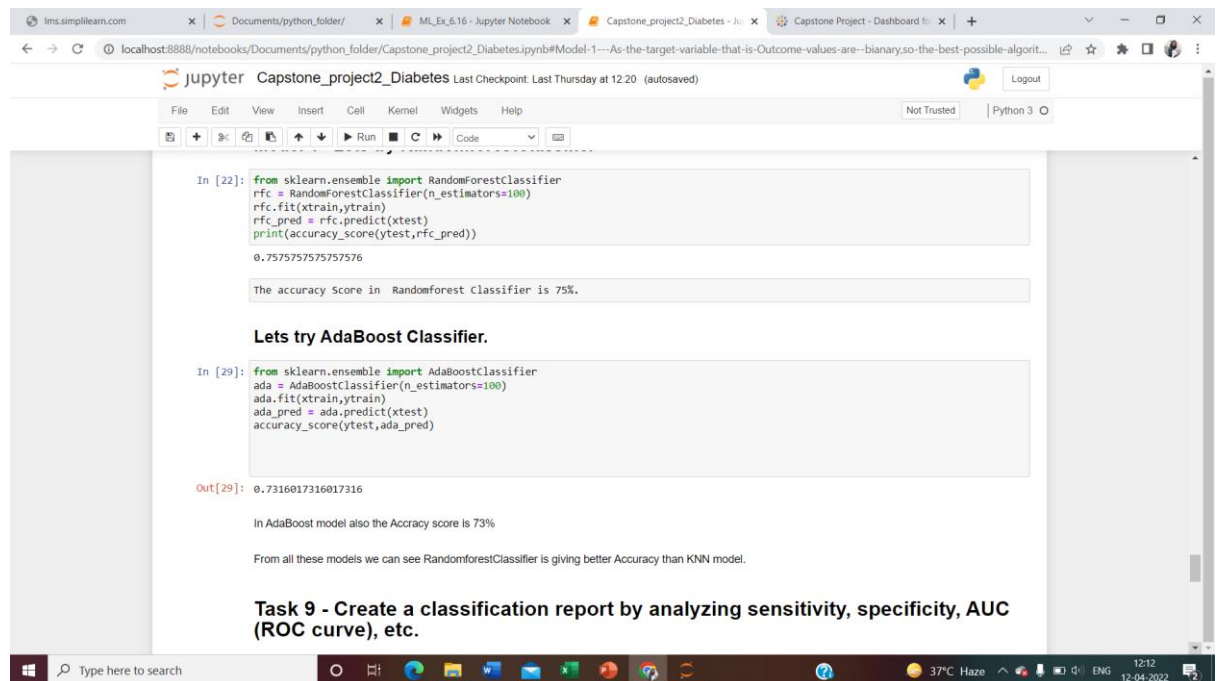**Task 7 –** Compare various models with the results from KNN algorithm.

# KNN model

Jupyter **Capstone_project2_Diabetes** Last Checkpoint: Last Thursday at 12:20 (unsaved changes)

File   Edit   View   Insert   Cell   Kernel   Widgets   Help

Not Trusted   |   Python 3 ○

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



We can see in th 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(xtrain,ytrain)
         knn_pred = knn_17.predict(xtest)
         print(accuracy_score(ytest,knn_pred))

         0.7402597402597403
```

The accuracy score in the knn model is 74%

---

# DecisionTree and RandomForest Classifier

Jupyter **Capstone_project2_Diabetes** Last Checkpoint: Last Thursday at 12:20 (autosaved)

File   Edit   View   Insert   Cell   Kernel   Widgets   Help

Not Trusted   |   Python 3 ○

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(xtrain,ytrain)
         tree_pred = tree.predict(xtest)
         print(accuracy_score(ytest,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(xtrain,ytrain)
         rfc_pred = rfc.predict(xtest)
         print(accuracy_score(ytest,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(xtrain,ytrain)
```

# Adaboost Classifier



## Task 8 – Create a classification report by analyzing sensitivity, specificity, AUC (ROC curve), etc.

# Sensitivity, Specificity and Classification Report

# ROC curve



**Task 5 -** Create a dashboard in tableau by choosing appropriate chart types and metrics useful for the business. The dashboard must entail the following:

- Pie chart to describe the diabetic or non-diabetic population

- Scatter charts between relevant variables to analyze the relationships

- Histogram or frequency charts to analyze the distribution of the data

- Heatmap of correlation analysis among the relevant variables

.