

# DIABETES PROJECT DATA ANALYSIS

Submitted to Ms.Shalini Kumari

Submitted By:

Namratha K

Batch Number: 4868

Bangalore

## TABLE OF CONTENTS

- 1) Abstract
- 2) Introduction
- 3) Description
- 4) Handling missing values
- 5) Data visualization
- 6) Analysis of parameter
- 7) Bilateral Analysis
- 8) Box plot
- 9) Conclusion

**Abstract:-** Diabetes is an illness caused because of high glucose level in a human body. Diabetes should not be ignored if it is untreated then Diabetes may cause some major issues in a person like: heart related problems, kidney problem, blood pressure, eye damage and it can also affects other organs of human body. Diabetes can be controlled if it is predicted earlier. To achieve this goal this project work we will do early prediction of Diabetes in a human body or a patient for a higher accuracy through applying Techniques. Provide better result for prediction by constructing models from datasets collected from patients

## Introduction

Data analysis is all about getting and overall understanding of data. It is mainly done to find it's properties , patterns and visualizations

In this blog we are using python as our programming language for the analysis purpose. Python has a wide variety of libraries like pandas , seaborn , numpy , matplotlib which can be used for this purpose. We are using the Pima Indians Dataset where it shows the various diagnostic factors influencing the diabetes.

## Brief Introduction to the used libraries

As discussed above , we are going to use the following libraries to perform different operations on the data

- Numpy
  - NumPy is a library for the Python programming language, adding support for large, multi-dimensional arrays and matrices, along with a large collection of high-level mathematical functions to operate on these arrays.
  - <http://www.numpy.org/>
- Pandas
  - pandas is a python package providing fast, flexible, and expressive data structures designed to make working with “relational” or “labeled” data both easy and intuitive. It aims to be the fundamental high-level building block for doing practical, real world data analysis in Python. Additionally, it has the broader goal of becoming the most powerful and flexible open source data analysis / manipulation tool available in any language. It is already well on its way toward this goal.
  - <https://pandas.pydata.org/>
- Seaborn
  - Seaborn is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.
  - <https://seaborn.pydata.org/>
- Matplotlib
  - Matplotlib is a Python 2D plotting library which produces publication quality figures in a variety of hardcopy formats and interactive environments across platforms.
  - <https://matplotlib.org/>

# Understanding the data :-

This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases. The objective of the dataset is to diagnostically predict whether or not a patient has diabetes, based on certain diagnostic measurements included in the dataset. Several constraints were placed on the selection of these instances from a larger database. In particular, all patients here are females at least 21 years old of Pima Indian heritage.

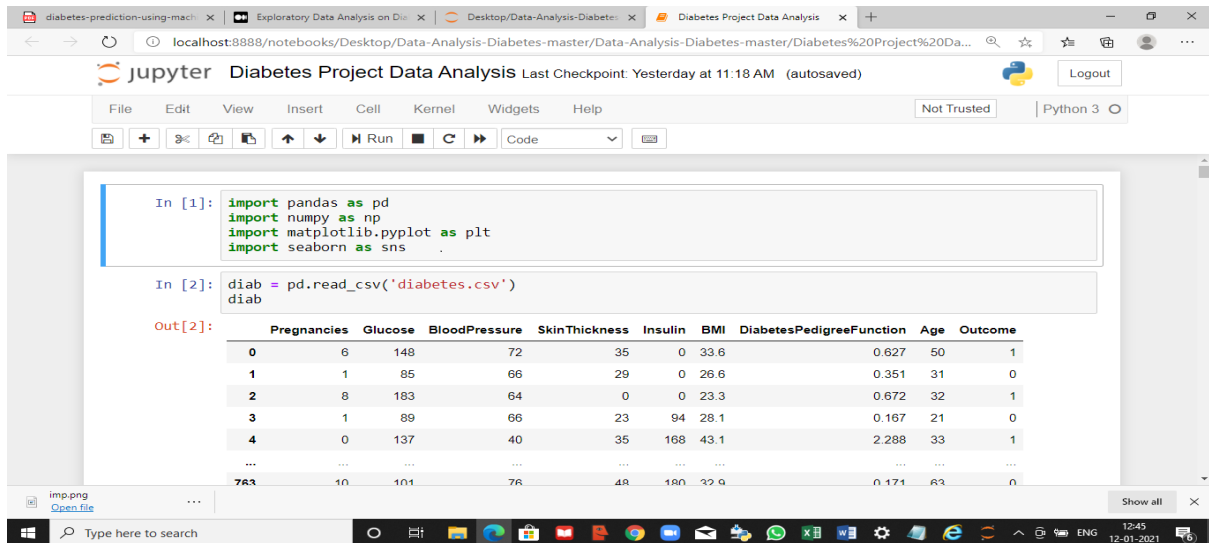
## Content

The datasets consists of several medical predictor variables and one target variable, Outcome. Columns are following :-

1. Pregnancies :- Number of times pregnant
2. Glucose:- Plasma glucose concentration a 2 hours in an oral glucose tolerance test
3. Blood Pressure:- Diastolic blood pressure (mm Hg)
4. Skin Thickness:- Triceps skin fold thickness (mm)
5. Insulin:- 2-Hour serum insulin (mu U/ml)
6. BMI:- Body mass index (weight in kg/(height in m)<sup>2</sup>)
7. Diabetes Pedigree Function:- Diabetes pedigree function
8. Age:-Age in years
9. Outcome:- Class variable (0 or 1) 268 of 768 are 1, the others

are 0

## 1) Importing the necessary libraries and Loading the dataset



The screenshot shows a Jupyter Notebook titled "Diabetes Project Data Analysis" running on a local server at localhost:8888. The notebook has two input cells. The first cell, labeled "In [1]:", contains the following code to import necessary libraries:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

The second cell, labeled "In [2]:", contains the code to load the dataset:

```
diab = pd.read_csv('diabetes.csv')
diab
```

The output of the second cell, labeled "Out[2]:", displays the first few rows of the dataset as a table:

	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
...	...	...	...	...	...	...	...	...	...
763	10	101	76	48	190	32.9	0.171	63	0

The notebook interface includes a menu bar with options like File, Edit, View, Insert, Cell, Kernel, Widgets, and Help. The status bar at the bottom shows the current kernel is Python 3 and the notebook is not trusted. The system tray at the very bottom shows the Windows taskbar with various application icons and the system clock indicating 12:45 on 12-01-2021.

## 2) Collecting Basic Information about the Data

The screenshot shows a Jupyter Notebook interface with the title 'Diabetes Project Data Analysis'. The notebook is running on a local host (localhost:8888). The code cell shows the following code:

```
In [5]: diab50 = diab.head(50)
        diab50
```

The output of the code is a table with 10 columns: Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DiabetesPedigreeFunction, Age, and Outcome. The table displays the first 10 rows of data.

	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
5	5	116	74	0	0	25.6	0.201	30	0
6	3	78	50	32	88	31.0	0.248	26	1
7	10	115	0	0	0	35.3	0.134	29	0
8	2	197	70	45	543	30.5	0.158	53	1
9	0	125	66	0	0	0.0	0.222	54	1

It shows that there are eight independent variables (Pregnancies, Glucose, Blood Pressure, SkinThickness, Insulin, BMI, DiabetesPedigreeFunction, Age) and one dependent variable (Outcome).

## 3) Descriptive statistics using describe method as shown below

diabetes-prediction-us | Exploratory Data Analy | Desktop/Data-Analysis | Diabetes Project Data | Exploratory Data Analy | Exploratory Data Analy | + -

localhost:8888/notebooks/Desktop/Data-Analysis-Diabetes-master/Data-Analysis-Diabetes-master/Diabetes%20Project%20Da... Logout

jupyter Diabetes Project Data Analysis Last Checkpoint: Yesterday at 11:18 AM (autosaved)

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3

Run Code

This tells us that there are no missing values. clear!

## 2. Description

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

Out[7]:

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

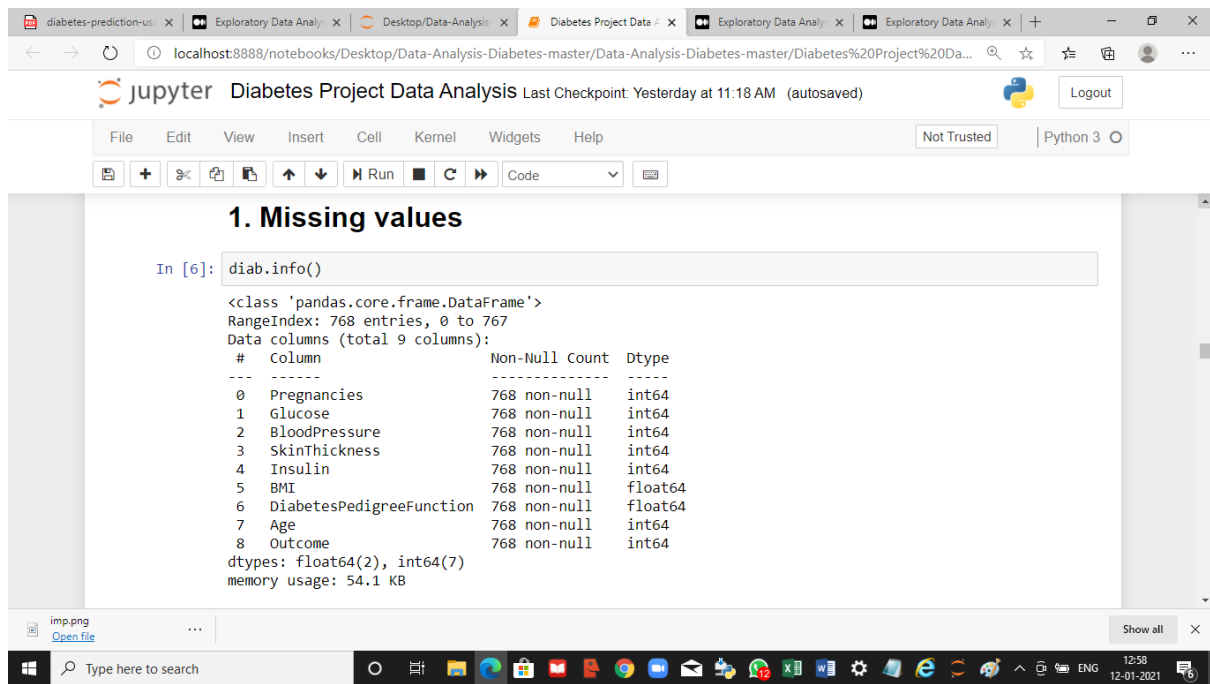
imp.png Open file Show all

Type here to search 12:58 12-01-2021

## 4) Handling the missing values

We need to check the presence of the missing values and need to replace them with mean, median and mode. Sometimes we have null values in the form of 0 , so we need to convert them to NaN and then replace them accordingly. The missing values can be removed also but it should be less than 5 percent of the whole dataset



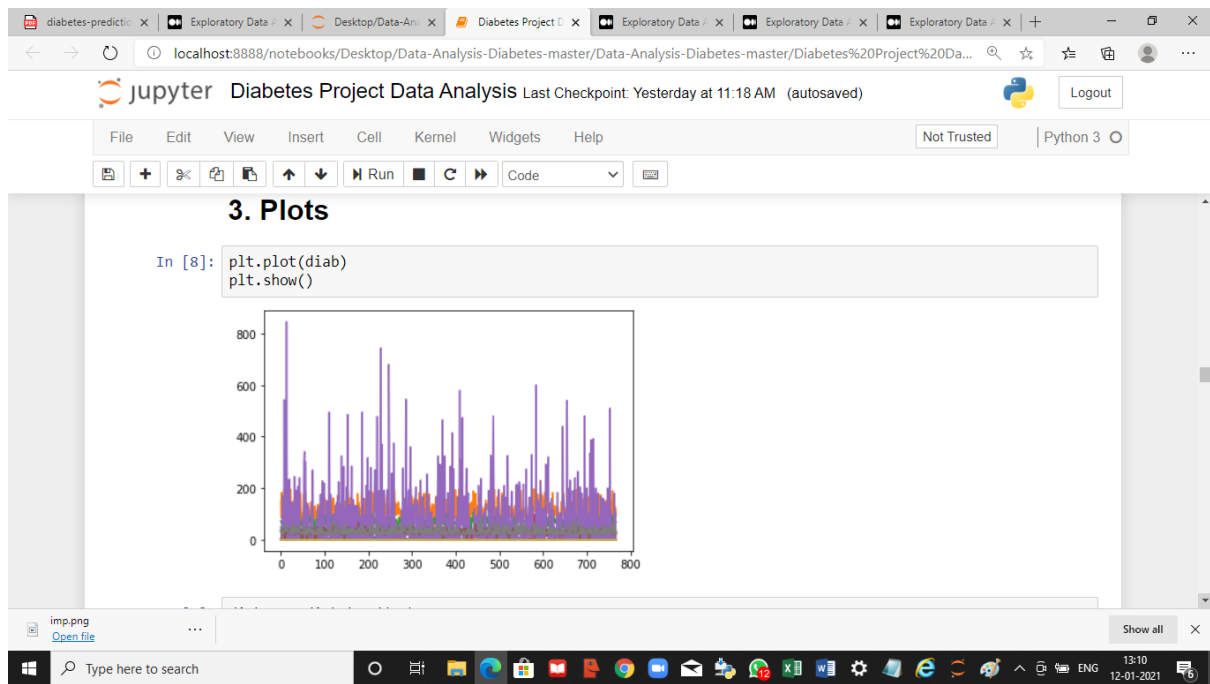


```
In [6]: diab.info()

<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   int64   
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
```

## 5)Data Visualization

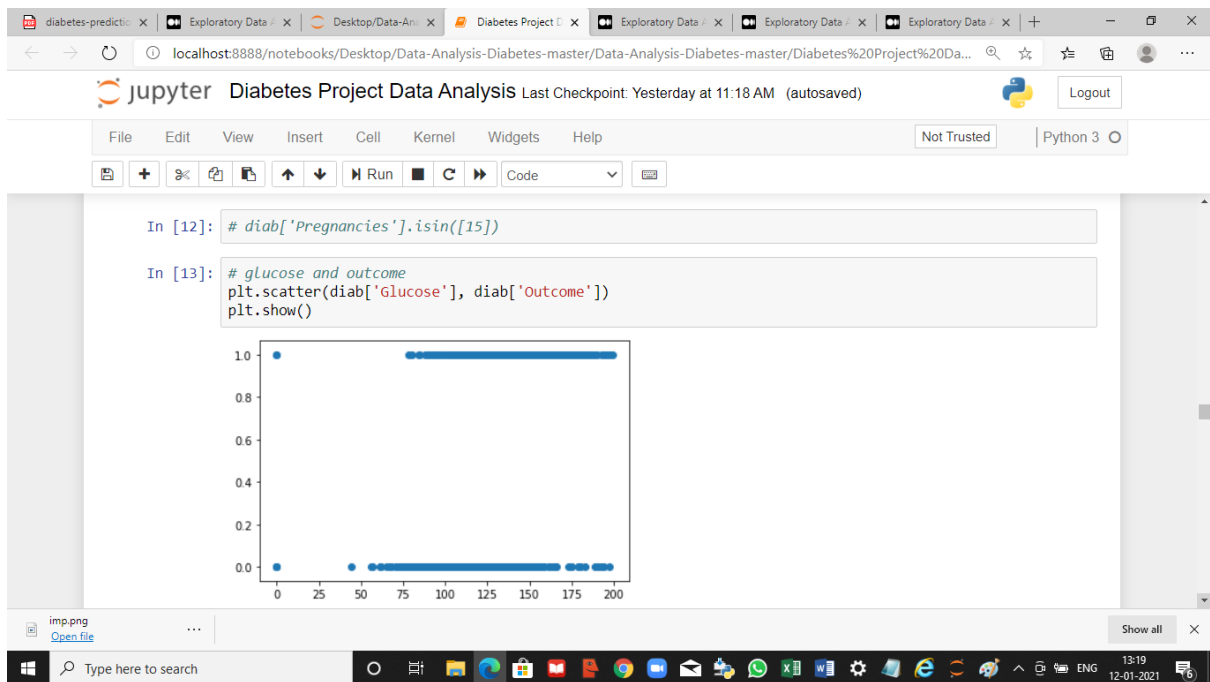
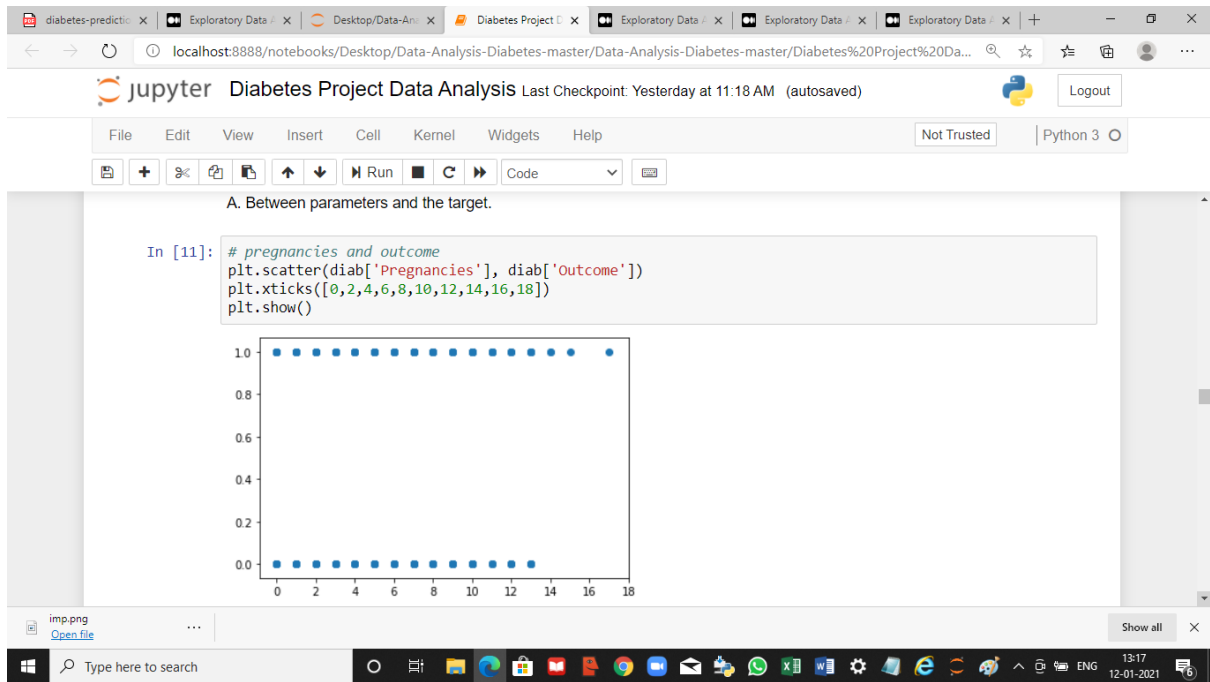
Visualizing data in different type of graphs will provide us with greater insights into our data. We will explore different options on visualizing our data and find out any patterns between within it

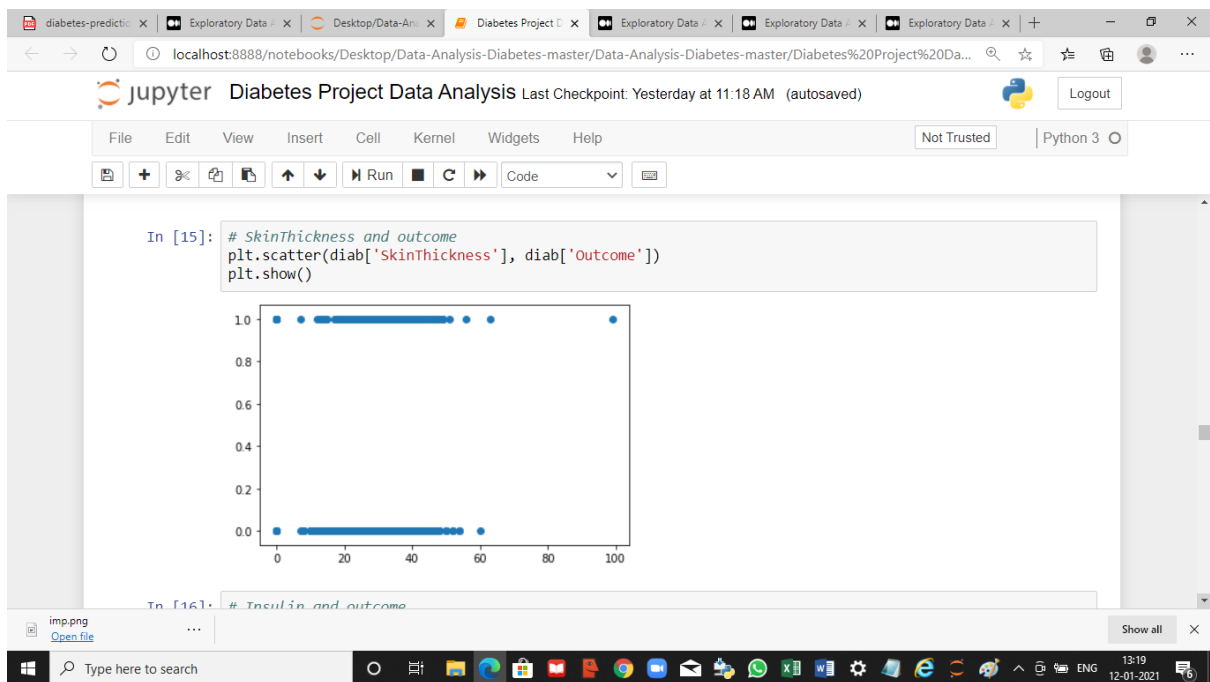
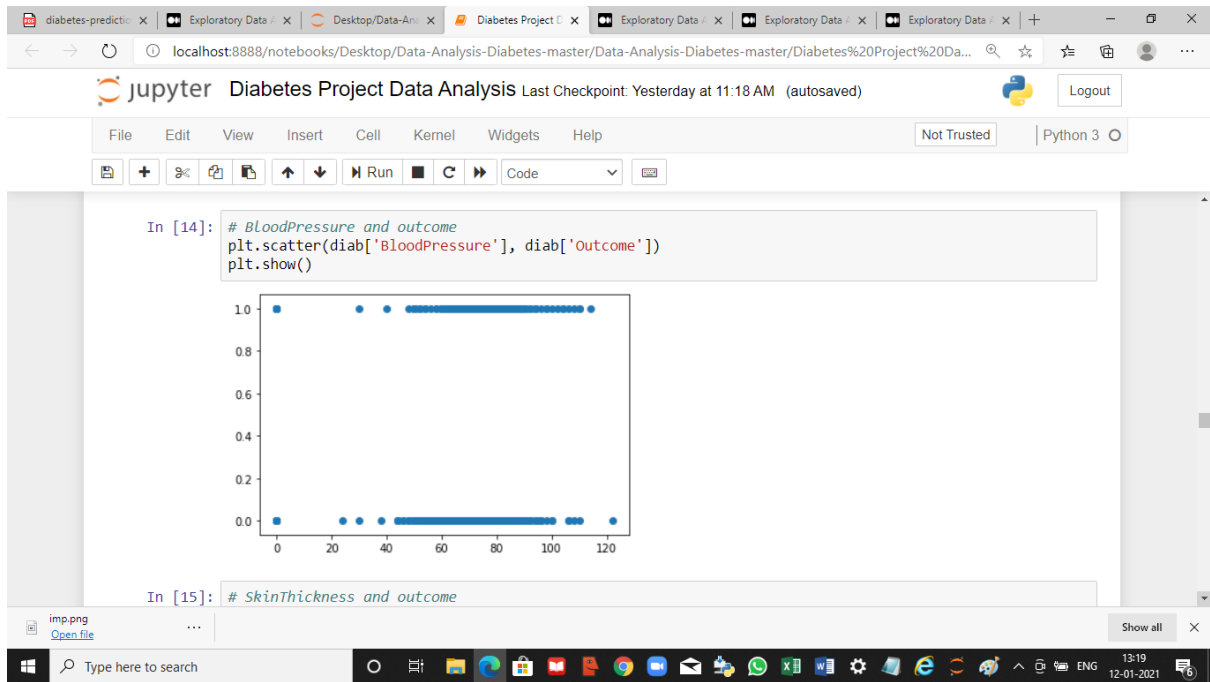


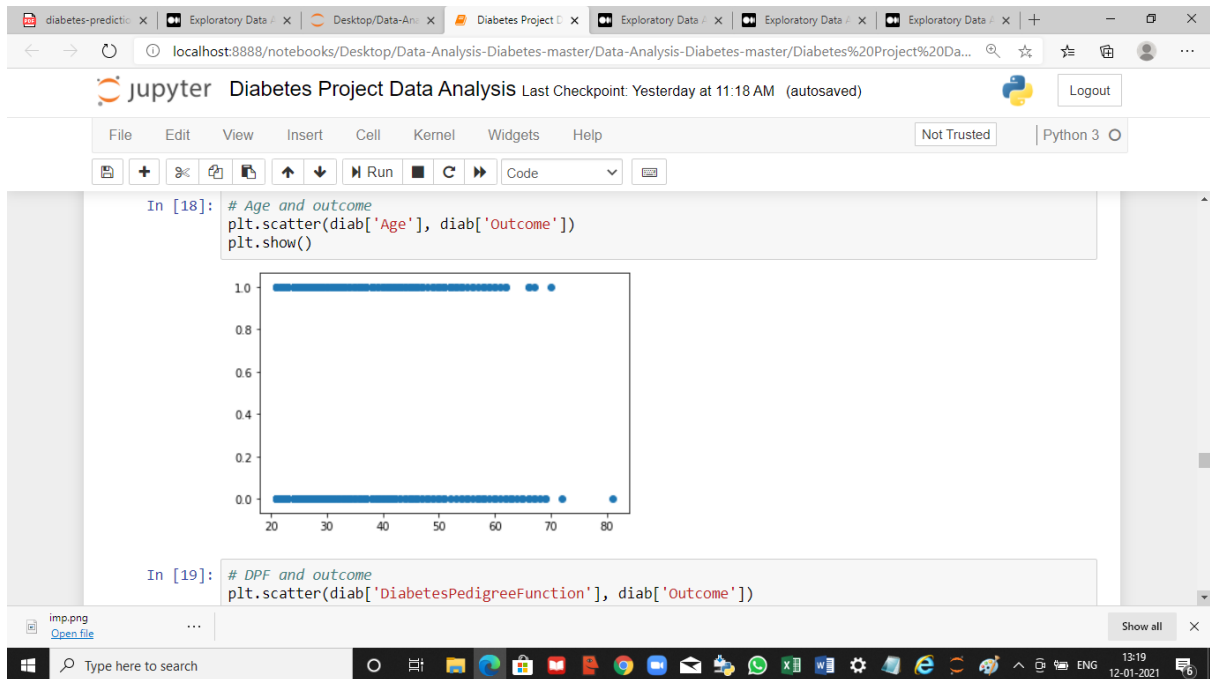
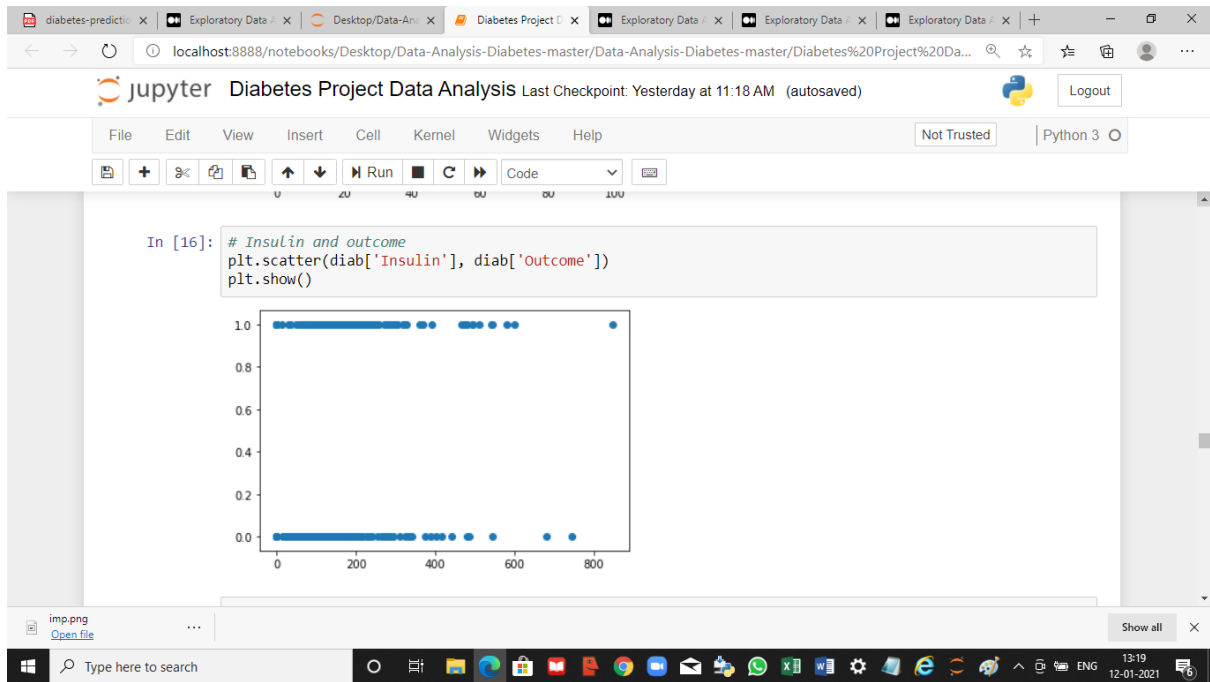
## 6) Analysis of each parameter w.r.t outcome parameter

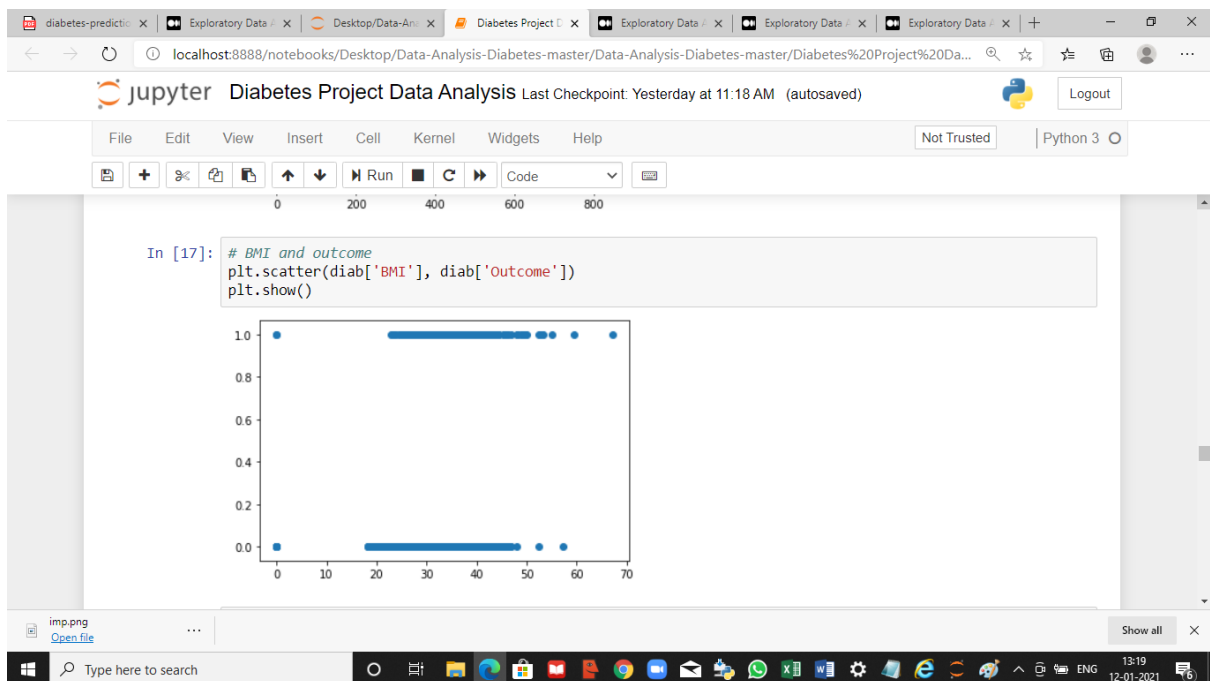
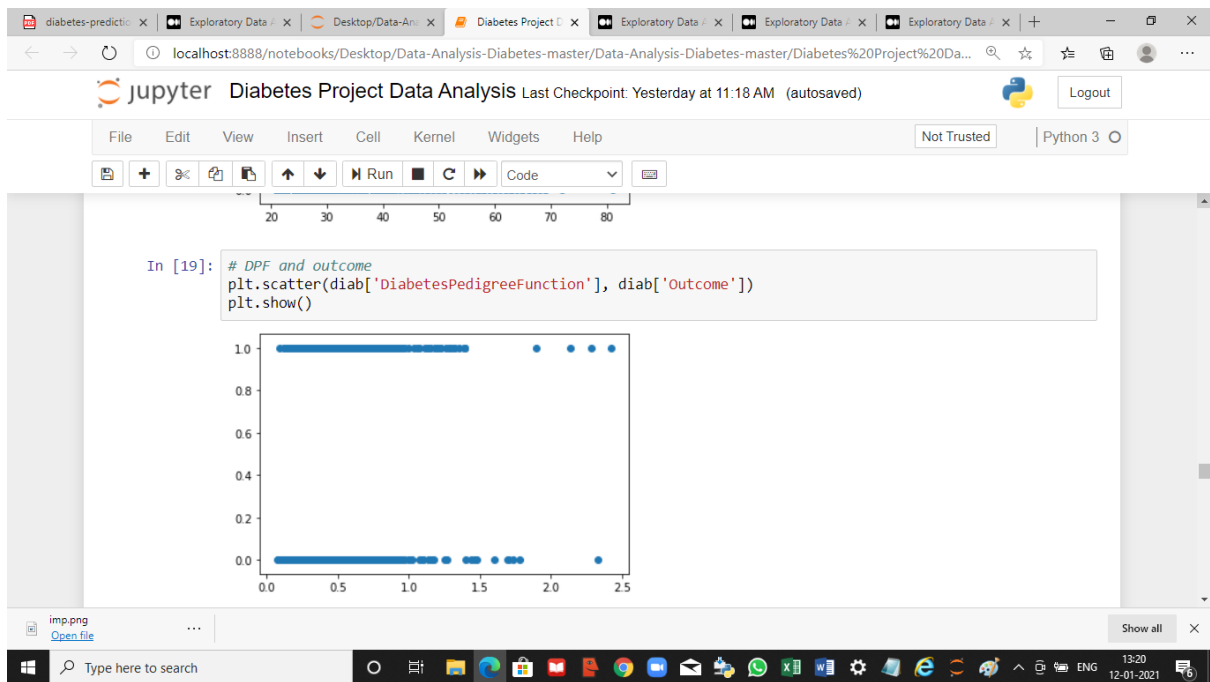
“Uni” means “one”, so in other words your data has only one variable. It doesn’t deal with causes or relationships. It’s major purpose is to describe; it takes data, summarizes that data and finds patterns in the data. This means here we deal and with only one variable or column of the data and try to find out it’s nature

### Analysis of ‘Pregnancies’ parameter



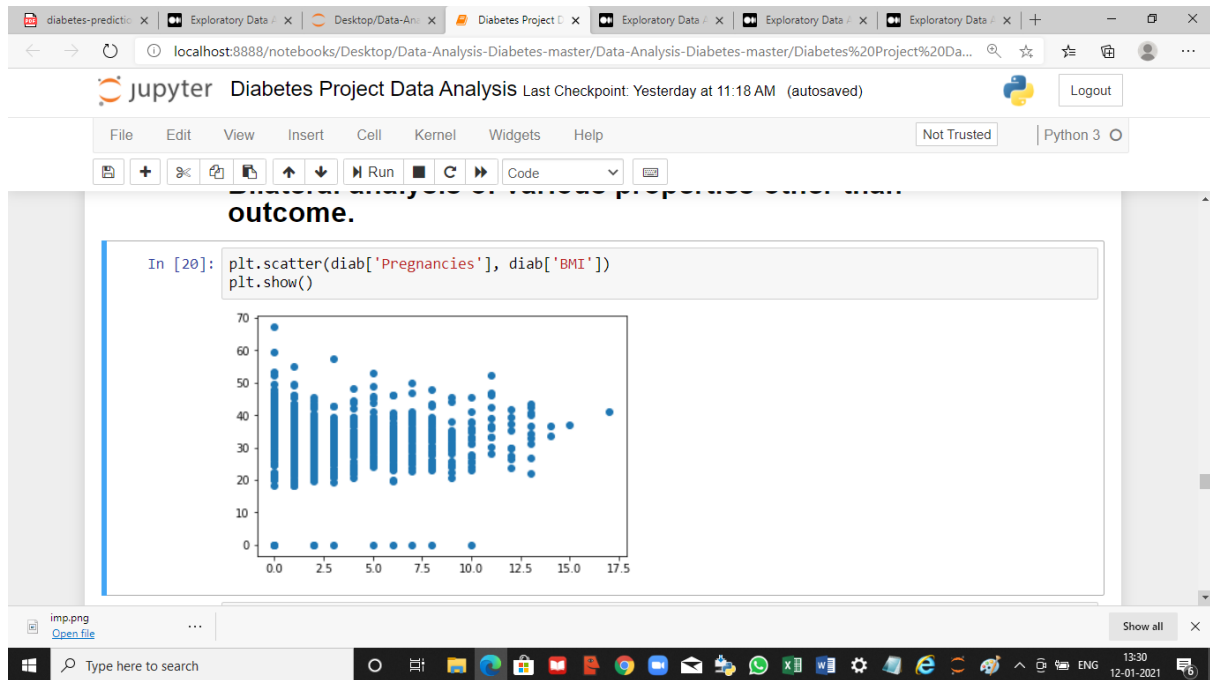




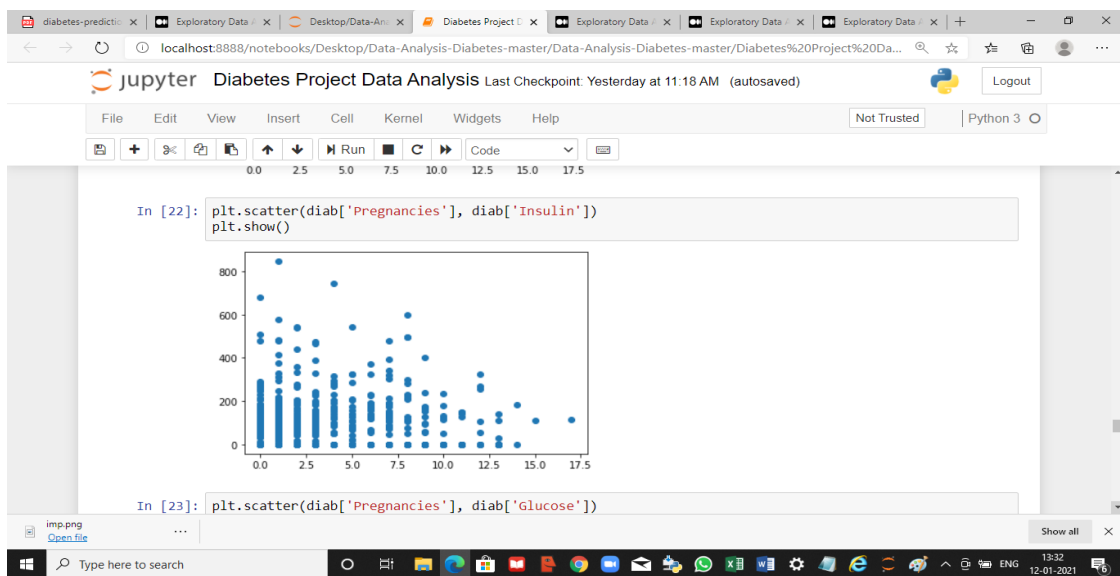
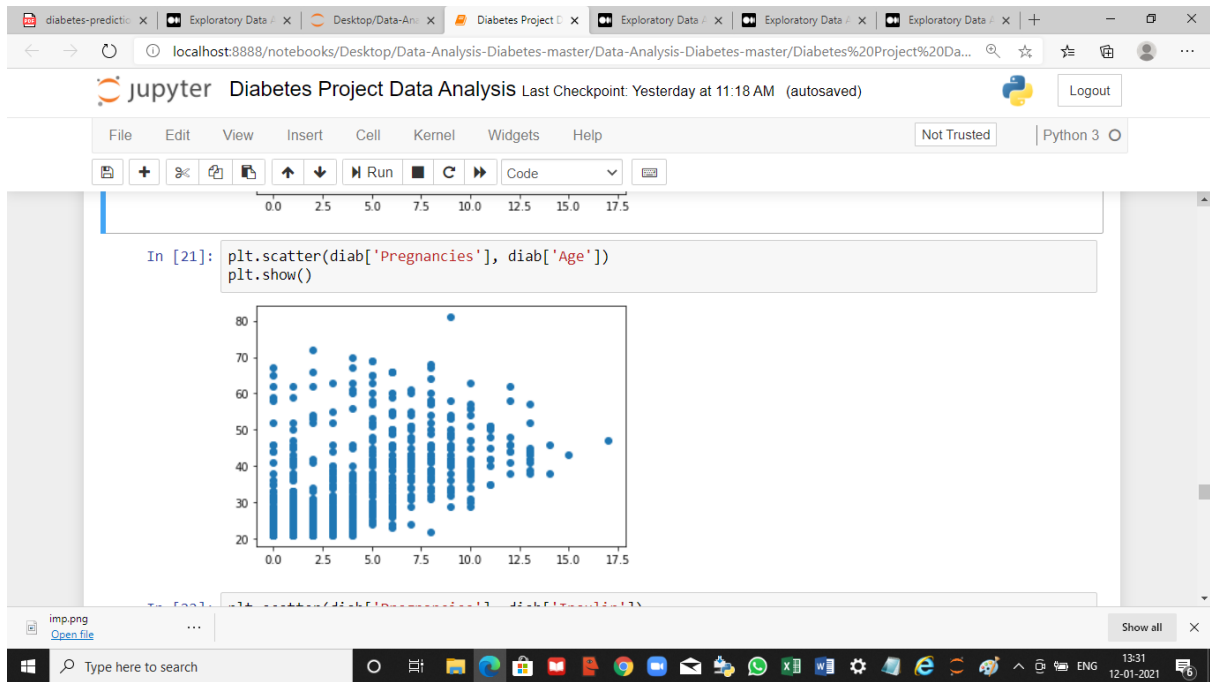


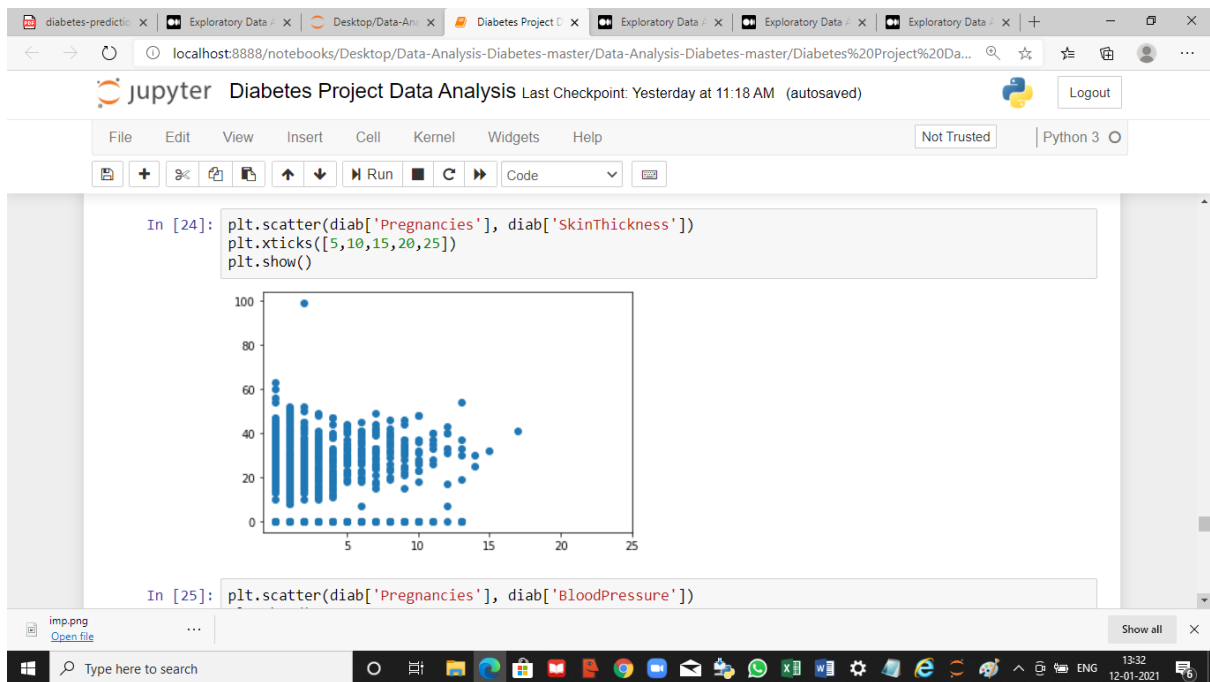
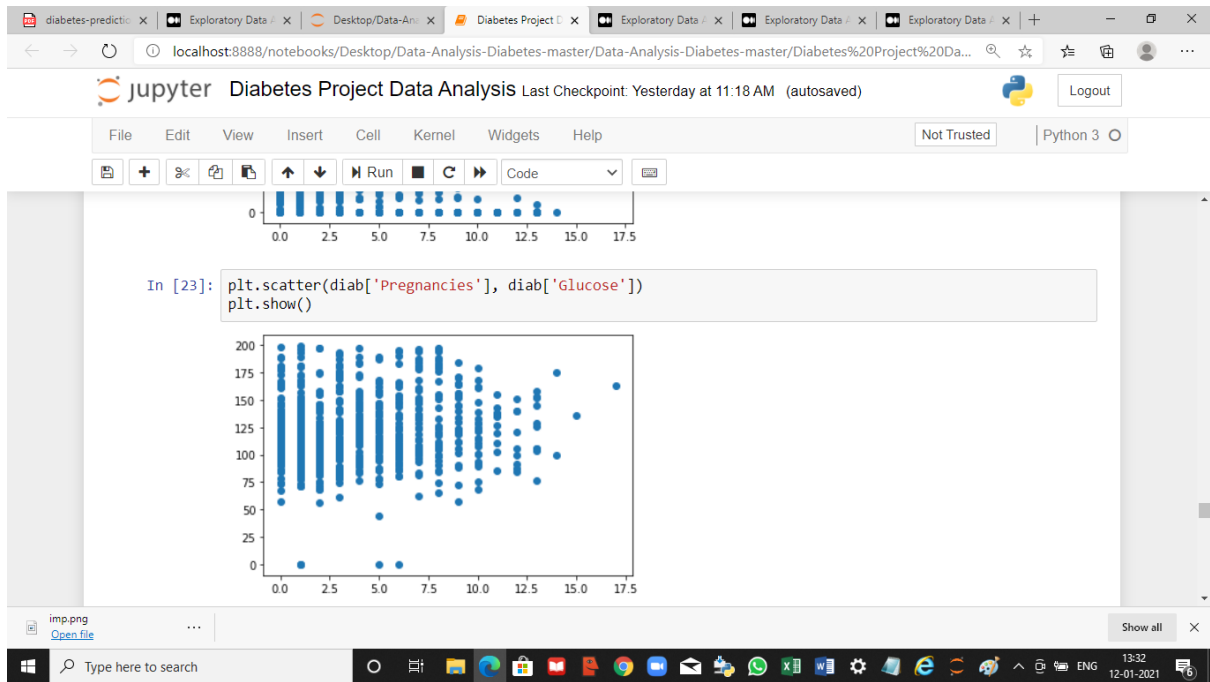
1. Pregnancies and age have some kind of a linear line.
2. BloodPressure and age have little relation. Most of the aged people have BloodPressure.
3. Insulin and Glucose have some relation.

## 7) Bilateral analysis of various properties other than outcome

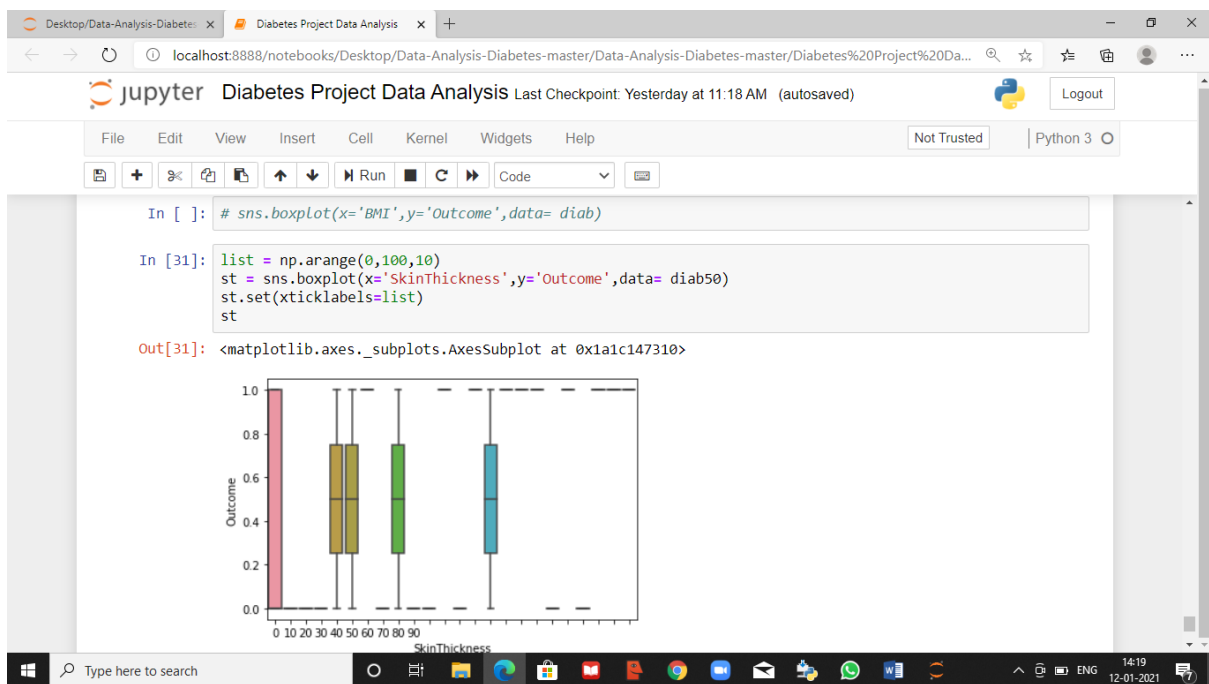
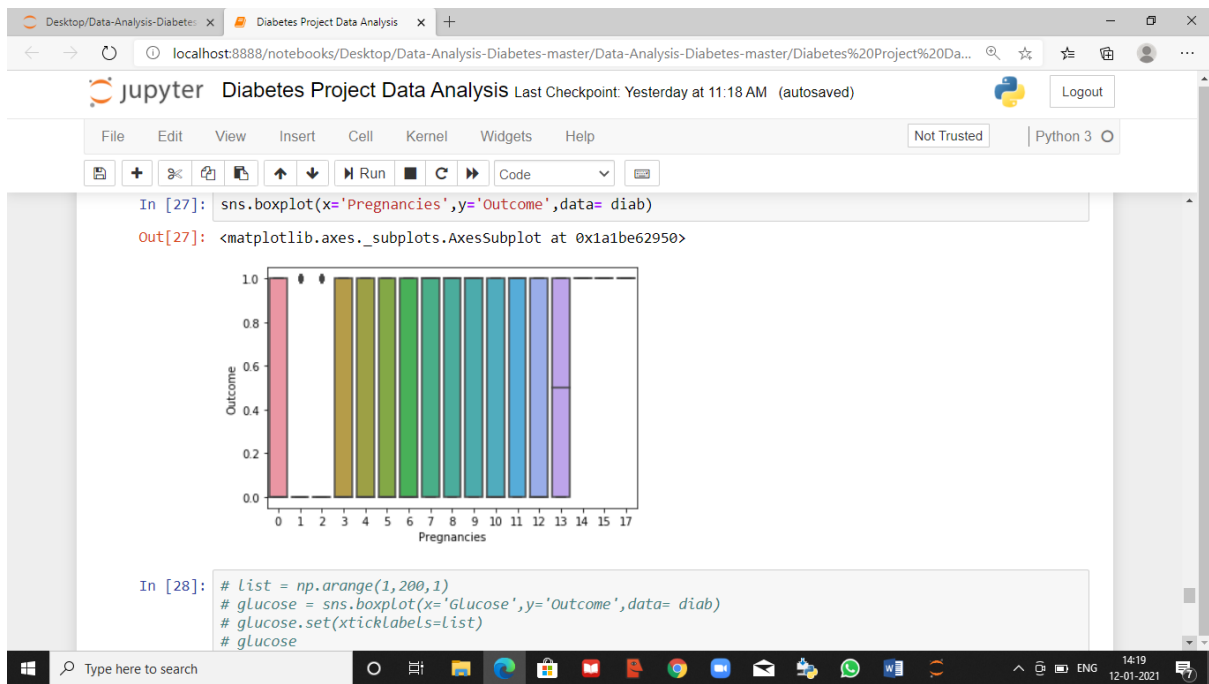












## Conclusion

Hence we can say from the above that EDA is very important as it helps us in achieving the following :-

1. It helps in detection of mistakes (like missing values and outliers)
2. It determines relationships between explanatory variables
3. Assessing the direction and rough size of relationships between explanatory and outcome variables.
4. It makes our data ready for machine learning algorithm