

IBM NAAN MUDHALVAN  
ARTIFICIAL INTELLIGENCE  
PROJECT

AI-BASED DIABETES PREDICTION SYSTEM PHASE 3

**INTRODUCTION:**

An AI-based diabetes prediction system is to develop a model that can accurately predict the likelihood of an individual developed diabetes based on demographic, lifestyle, and health factors. The system should be able to analyse large amount of data, including medical records, genetic information, and lifestyle choices, to provide personalized predictions. The goal is assist to healthcare professionals in identifying individuals who are high risk of developing diabetes, allowing for early intervention and prevention strategies. The system should also user-friendly and easily accessible to both healthcare professionals and individual looking to access their own risk of diabetes.

**OVERVIEW:**

In this article, we will be predicting that whether the patient has diabetes or not on the basis of the features we will provide to our machine learning model, and for that, we will be using the famous Pima Indians Diabetes Database.

- 1.Data analysis: Here one will get to know about how the data analysis part is done in a data science life cycle.
- 2.Exploratory data analysis: EDA is one of the most important steps in the data science project life cycle and here one will need to know that how to make inferences from the visualizations and data analysis

3.Model building: Here we will be using 4 ML models and then we will choose the best performing model.

4.Saving model: Saving the best model using pickle to make the prediction from real data.

**DATASET:**

Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI
6	148	72	35	0	33.6
1	85	66	29	0	26.6
8	183	64	0	0	23.3
1	89	66	23	94	28.1
0	137	40	35	168	43.1
5	116	74	0	0	25.6
3	78	50	32	88	31
10	115	0	0	0	35.3
2	197	70	45	543	30.5
8	125	96	0	0	0
4	110	92	0	0	37.6
10	168	74	0	0	38
10	139	80	0	0	27.1
1	189	60	23	846	30.1
5	166	72	19	175	25.8
7	100	0	0	0	30
0	118	84	47	230	45.8
7	107	74	0	0	29.6
1	103	30	38	83	43.3
1	115	70	30	96	34.6
3	126	88	41	235	39.3
8	99	84	0	0	35.4
7	196	90	0	0	39.8
9	119	80	35	0	29
11	143	94	33	146	36.6
10	125	70	26	115	31.1
7	147	76	0	0	39.4
1	97	66	15	140	23.2
13	145	82	19	110	22.2
5	117	92	0	0	34.1
5	109	75	26	0	36
3	158	76	36	245	31.6
3	88	58	11	54	24.8
6	92	92	0	0	19.9
10	122	78	31	0	27.6
4	103	60	33	192	24
11	138	76	0	0	33.2
9	102	76	37	0	32.9
2	90	68	42	0	38.2
4	111	72	47	207	37.1
3	180	64	25	70	34
7	133	84	0	0	40.2
7	106	92	18	0	22.7
9	171	110	24	240	45.4
7	159	64	0	0	27.4

769 rows\*8 columns

## IMPORT LIBRARIES:

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

sns.set()

from mlxtend.plotting import plot_decision_regions
import missingno as msno
from pandas.plotting import scatter_matrix
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier

from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import classification_report
import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
```

Here we will be reading the dataset which is in the CSV format

```
diabetes_df = pd.read_csv('diabetes.csv')
diabetes_df.head()
```

## Output:

	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

## Exploratory Data Analysis (EDA)

Now let's see that what are columns available in our dataset.

```
diabetes_df.columns
```

## Output:

```
Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness',  
      'Insulin',  
      'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],  
      dtype='object')
```

## Information about the dataset

```
diabetes_df.info()
```

## Output:

```
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
```

## To know more about the dataset

```
diabetes_df.describe()
```

## Output:

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

## To know more about the dataset with transpose – here T is for the transpose

```
diabetes_df.describe().T
```

## Output:

	count	mean	std	min	25%	50%	75%	max
Pregnancies	768.0	3.845052	3.369578	0.000	1.00000	3.0000	6.00000	17.00
Glucose	768.0	120.894531	31.972618	0.000	99.00000	117.0000	140.25000	199.00
BloodPressure	768.0	69.105469	19.355807	0.000	62.00000	72.0000	80.00000	122.00
SkinThickness	768.0	20.536458	15.952218	0.000	0.00000	23.0000	32.00000	99.00
Insulin	768.0	79.799479	115.244002	0.000	0.00000	30.5000	127.25000	846.00
BMI	768.0	31.992578	7.884160	0.000	27.30000	32.0000	36.60000	67.10
DiabetesPedigreeFunction	768.0	0.471876	0.331329	0.078	0.24375	0.3725	0.62625	2.42
Age	768.0	33.240885	11.760232	21.000	24.00000	29.0000	41.00000	81.00
Outcome	768.0	0.348958	0.476951	0.000	0.00000	0.0000	1.00000	1.00

Now let's check that if our dataset have null values or not

```
diabetes_df.isnull().head(10)
```

Output:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False
5	False	False	False	False	False	False	False	False	False
6	False	False	False	False	False	False	False	False	False
7	False	False	False	False	False	False	False	False	False
8	False	False	False	False	False	False	False	False	False
9	False	False	False	False	False	False	False	False	False

Now let's check the number of null values our dataset has.

```
diabetes_df.isnull().sum()
```

Output:

```
Pregnancies      0
Glucose           0
BloodPressure     0
SkinThickness     0
Insulin           0
BMI               0
```

```
DiabetesPedigreeFunction    0
Age                         0
Outcome                     0
dtype: int64
```

Here from the above code we first checked that is there any null values from the `IsNull()` function then we are going to take the sum of all those missing values from the `sum()` function and the inference we now get is that there are no missing values but that is actually not a true story as in **this particular dataset all the missing values were given the 0 as a value which is not good for the authenticity of the dataset.** Hence we will first **replace the 0 value with the NAN** value then start the imputation process.

```
diabetes_df_copy = diabetes_df.copy(deep = True)
diabetes_df_copy[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']] =
diabetes_df_copy[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']].replace(0,np.NaN)
```

**# Showing the Count of NANs**

```
print(diabetes_df_copy.isnull().sum())
```

**Output:**

```
Pregnancies                0
Glucose                    5
BloodPressure              35
SkinThickness              227
Insulin                   374
BMI                       11
DiabetesPedigreeFunction   0
Age                       0
```



```
Outcome          0
dtype: int64
```

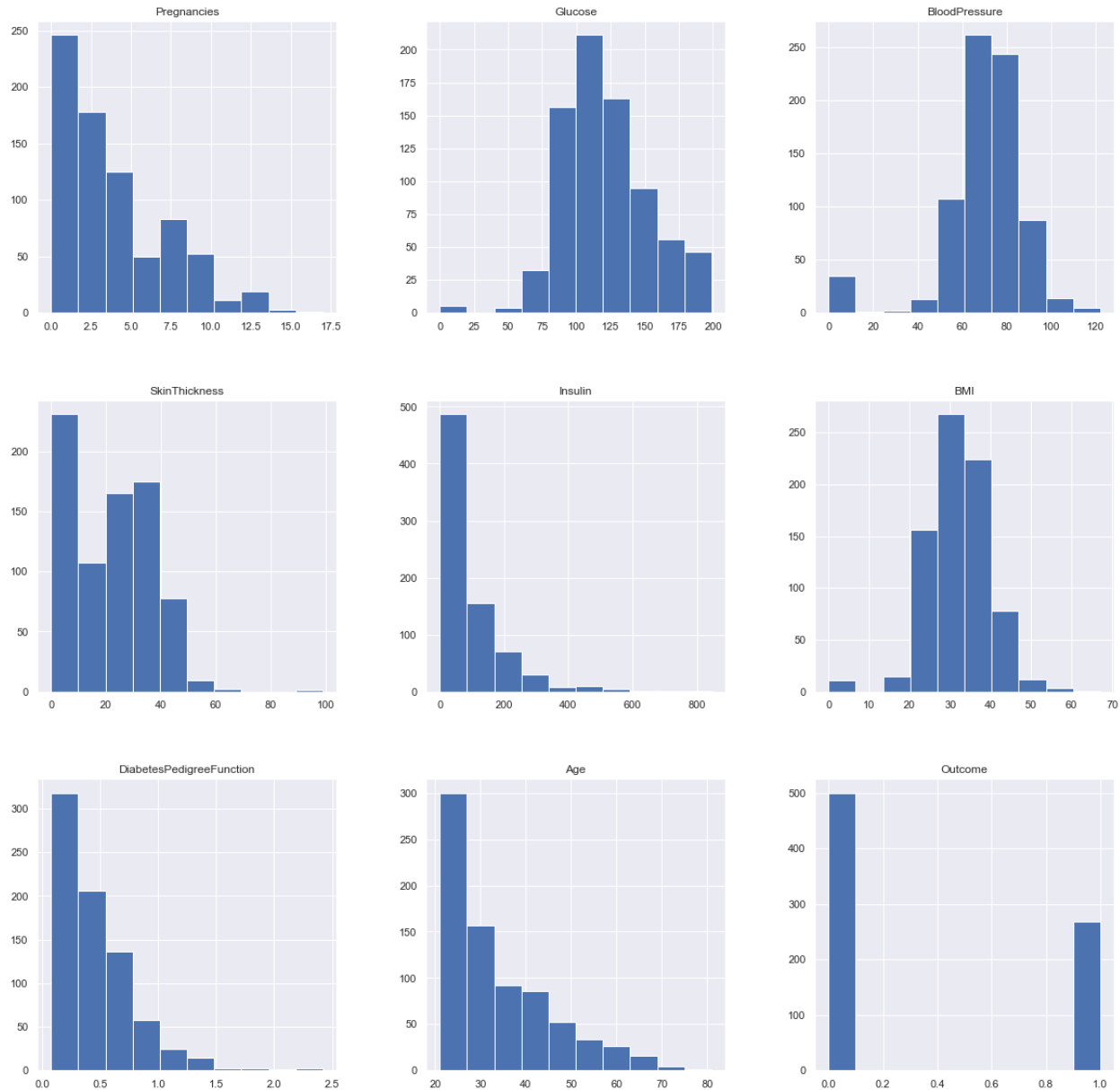
As mentioned above that now **we will be replacing the zeros with the NAN values** so that we can impute it later to maintain the authenticity of the dataset as well as trying to have a better Imputation approach i.e **to apply mean values of each column to the null values of the respective columns.**

## Data Visualization

**Plotting the data distribution plots before removing null values**

```
p = diabetes_df.hist(figsize = (20,20))
```

**Output:**



**Inference:** So here we have seen the distribution of each features whether it is dependent data or independent data and one thing which could always strike that **why do we need to see the distribution of data?** So the answer is simple it is the best way to start the analysis of the dataset as **it shows the occurrence of every kind of value in the graphical structure which in turn lets us know the range of the data.**

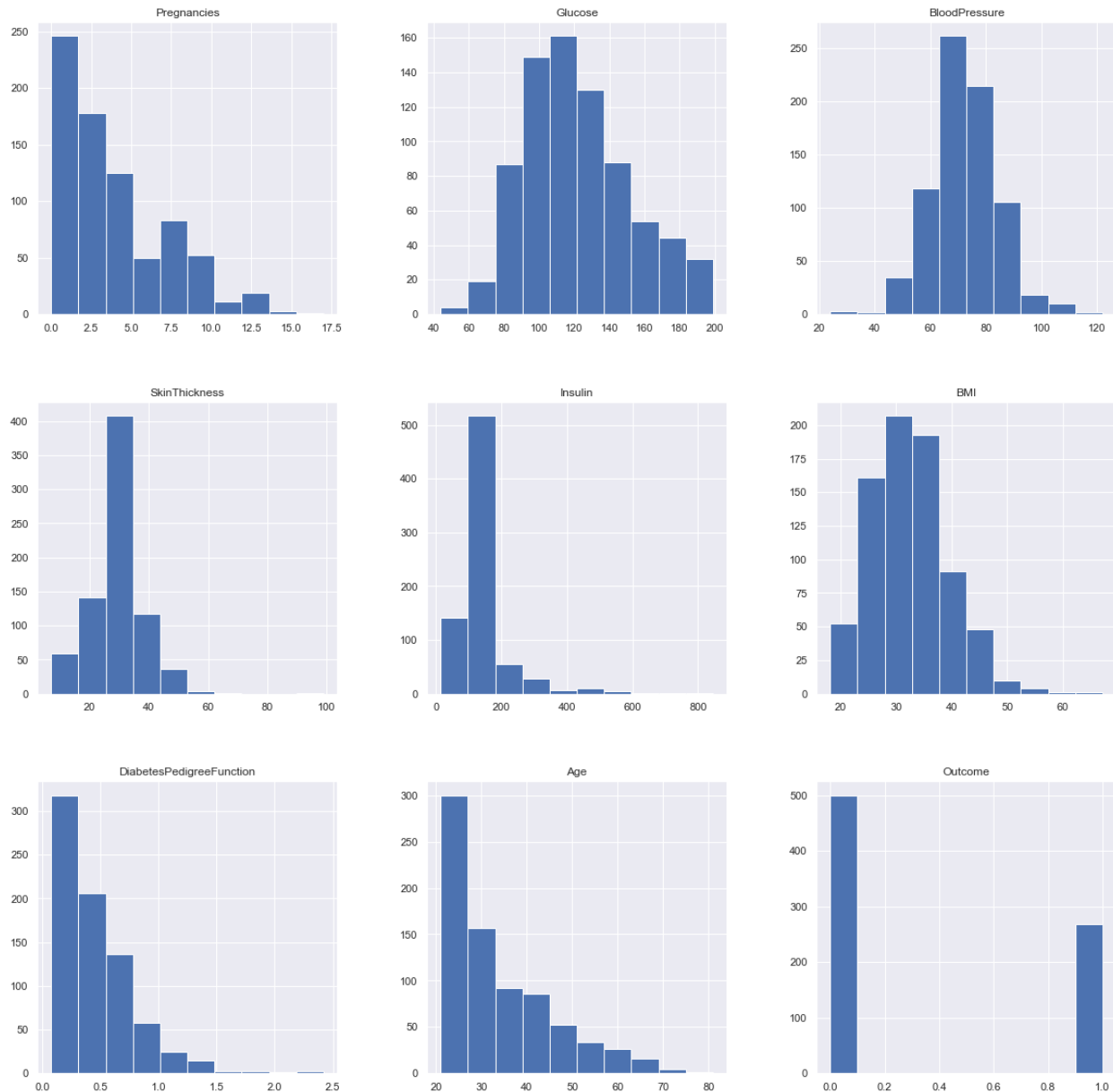
**Now we will be imputing the mean value of the column to each missing value of that particular column.**

```
diabetes_df_copy['Glucose'].fillna(diabetes_df_copy['Glucose'].mean()  
( ), inplace = True)  
  
diabetes_df_copy['BloodPressure'].fillna(diabetes_df_copy['BloodPre  
ssure'].mean(), inplace = True)  
  
diabetes_df_copy['SkinThickness'].fillna(diabetes_df_copy['SkinThic  
kness'].median(), inplace = True)  
  
diabetes_df_copy['Insulin'].fillna(diabetes_df_copy['Insulin'].medi  
an(), inplace = True)  
  
diabetes_df_copy['BMI'].fillna(diabetes_df_copy['BMI'].median(),  
inplace = True)
```

**Plotting the distributions after removing the NAN values.**

```
p = diabetes_df_copy.hist(figsize = (20,20))
```

**Output:**

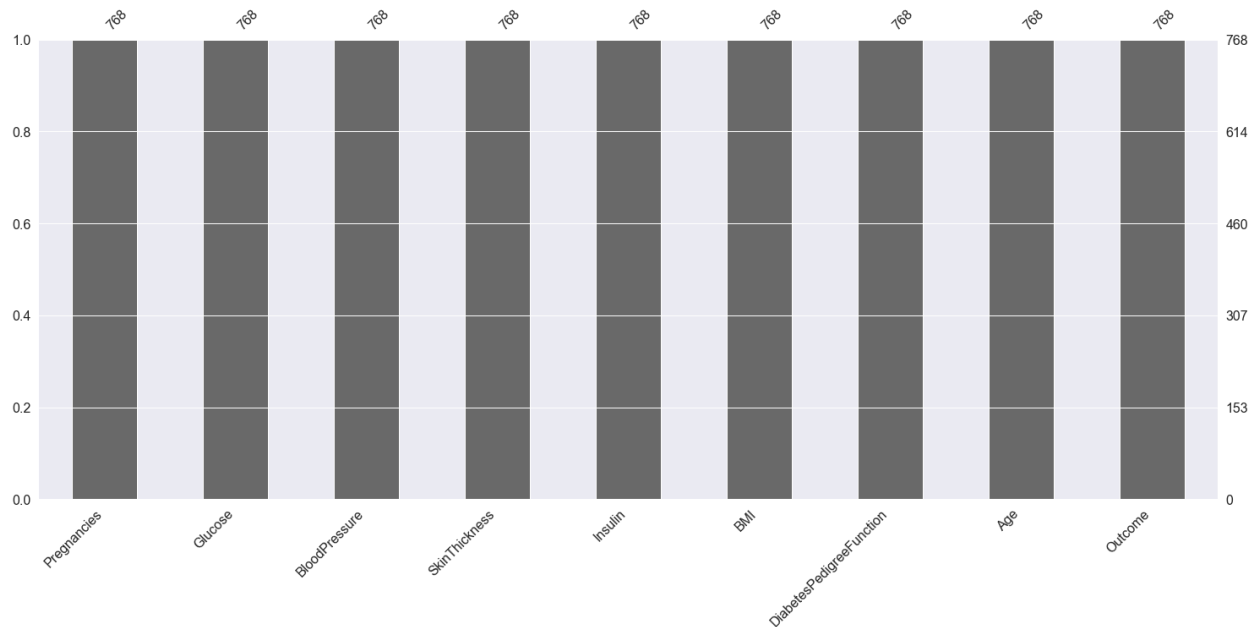


**Inference:** Here we are again using the hist plot to **see the distribution of the dataset** but this time we are using this visualization to see the changes that we can see after those null values are removed from the dataset and we can clearly see the difference **for example** – In age column after removal of the null values, we can see that there is a spike at the range of 50 to 100 which is quite logical as well.

## Plotting Null Count Analysis Plot

```
p = msno.bar(diabetes_df)
```

**Output:**



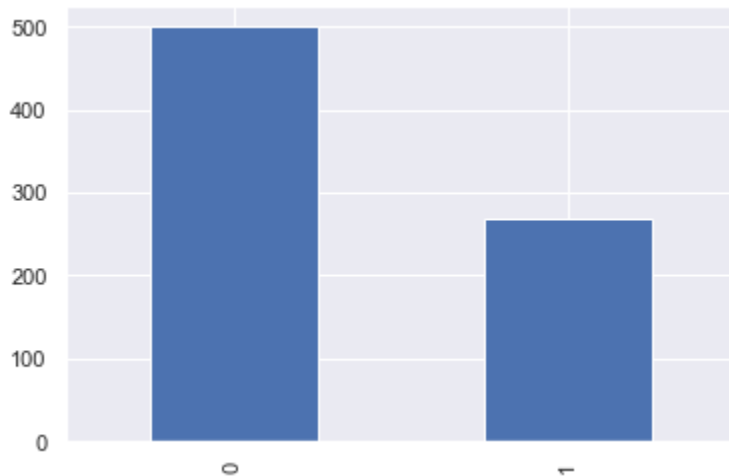
**Inference:** Now in the above graph also we can clearly see that there are **no null** values in the dataset.

**Now, let's check that how well our outcome column is balanced**

```
color_wheel = {1: "#0392cf", 2: "#7bc043"}
colors = diabetes_df["Outcome"].map(lambda x: color_wheel.get(x + 1))
print(diabetes_df.Outcome.value_counts())
p=diabetes_df.Outcome.value_counts().plot(kind="bar")
```

**Output:**

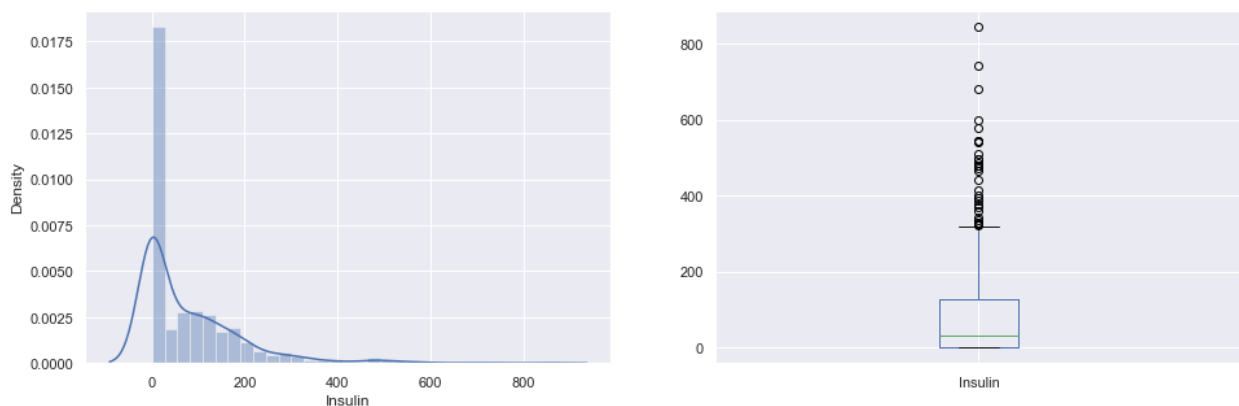
```
0    500
1    268
Name: Outcome, dtype: int64
```



**Inference:** Here from the above visualization it is clearly visible that our **dataset is completely imbalanced** in fact the number of patients who are **diabetic is half of the patients who are non-diabetic**.

```
plt.subplot(121), sns.distplot(diabetes_df['Insulin'])
plt.subplot(122), diabetes_df['Insulin'].plot.box(figsize=(16,5))
plt.show()
```

**Output:**



**Inference:** That's how **Distplot** can be helpful where one will be able to see the distribution of the data as well as with the help of **boxplot** one can see the

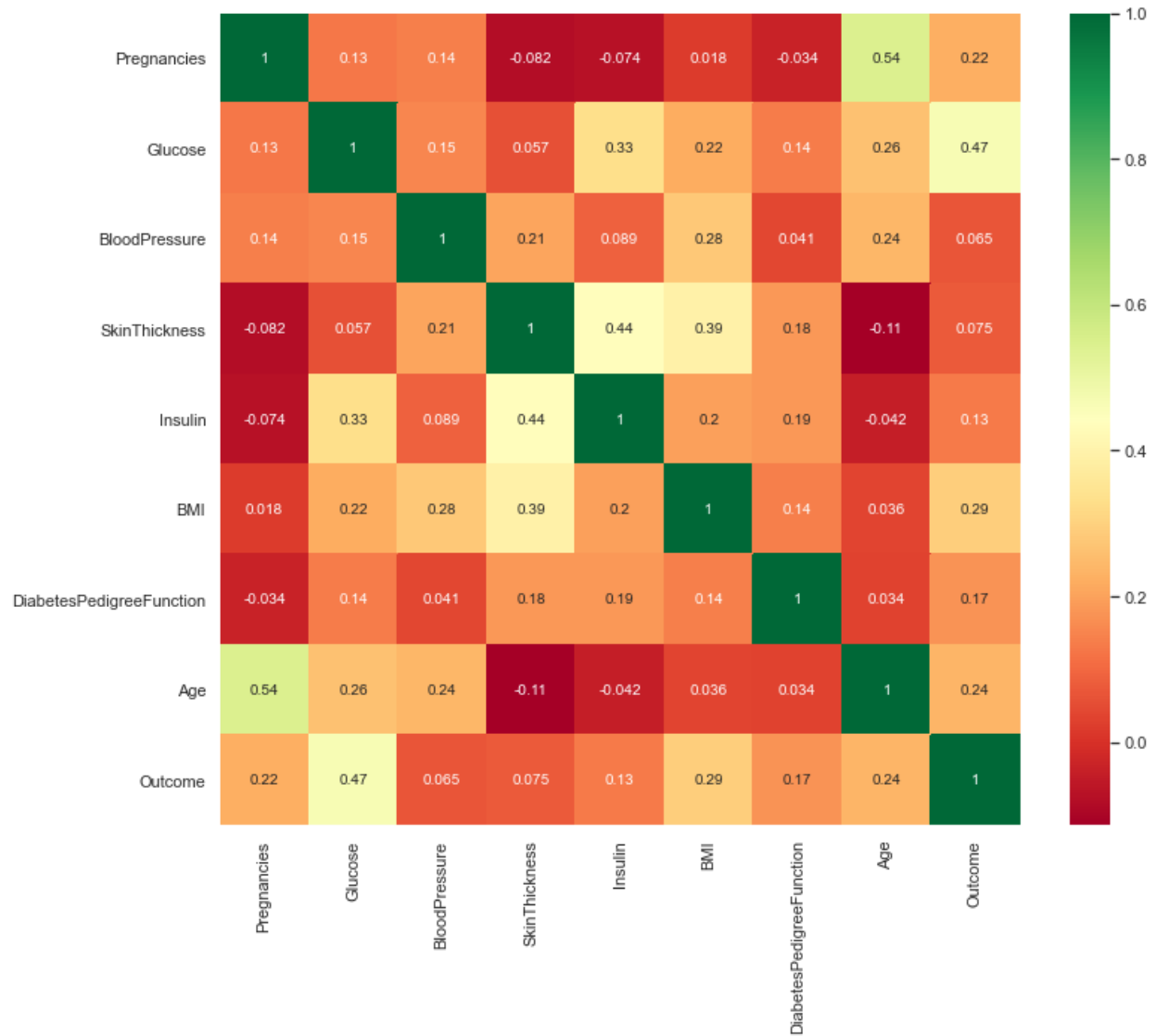
**outliers in that column** and other information too which can be derived by the **box and whiskers plot**.

## Correlation between all the features

### Correlation between all the features before cleaning

```
plt.figure(figsize=(12,10))  
# seaborn has an easy method to showcase heatmap  
p = sns.heatmap(diabetes_df.corr(), annot=True, cmap = 'RdYlGn')
```

**Output:**



## Scaling the Data

Before scaling down the data let's have a look into it

```
diabetes_df_copy.head()
```

Output:



	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148.0	72.0	35.0	125.0	33.6	0.627	50	1
1	1	85.0	66.0	29.0	125.0	26.6	0.351	31	0
2	8	183.0	64.0	29.0	125.0	23.3	0.672	32	1
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	168.0	43.1	2.288	33	1

## After Standard scaling

```
sc_X = StandardScaler()
X =
pd.DataFrame(sc_X.fit_transform(diabetes_df_copy.drop(["Outcome"],a
xis = 1)), columns=['Pregnancies',
'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI',
'DiabetesPedigreeFunction', 'Age'])
X.head()
```

## Output:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	0.639947	0.865108	-0.033518	0.670643	-0.181541	0.166619	0.468492	1.425995
1	-0.844885	-1.206162	-0.529859	-0.012301	-0.181541	-0.852200	-0.365061	-0.190672
2	1.233880	2.015813	-0.695306	-0.012301	-0.181541	-1.332500	0.604397	-0.105584
3	-0.844885	-1.074652	-0.529859	-0.695245	-0.540642	-0.633881	-0.920763	-1.041549
4	-1.141852	0.503458	-2.680669	0.670643	0.316566	1.549303	5.484909	-0.020496

That's how our dataset will be looking like when it is scaled down or we can see every value now is on the same scale which will help our **ML model** to give a **better result**.

## Let's explore our target column

## Output:

```
0      1
1      0
2      1
```

```
3      0
4      1
      ..
763    0
764    0
765    0
766    1
767    0
```

Name: Outcome, Length: 768, dtype: int64

## Model Building

### Splitting the dataset

```
X = diabetes_df.drop('Outcome', axis=1)
y = diabetes_df['Outcome']
```

Now we will split the data into training and testing data using the **train\_test\_split** function

```
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y,
                                                    test_size=0.33, random_state=7)
```

### Importance of loading and processing dataset:

Loading and preprocessing the dataset is an important first step in building any machine learning model. However, it is especially important for house prediction models, as house price datasets are often complex and noisy.

By loading and preprocessing the dataset, we can ensure that the machine learning algorithm is able to learn from the data effectively and accurately.

### CONCLUSION:

In conclusion, the development of an AI-based diabetes prediction system hold significant potential for improving healthcare outcomes related to diabetes. By leveraging artificial intelligence techniques, such as machine learning and data analysis, the system can provide accurate risk assessments and enable early intervention and prevention strategies.

The AI-based diabetes prediction system offers several advantages. Firstly, it can assist healthcare professionals in identifying individuals at high risk of developing diabetes, allowing for targeted interventions and personalized care plans. Secondly, it can contribute to reducing healthcare costs by focusing resources on prevention and early detection, which can help mitigate the long-term complications associated with diabetes. Additionally, the system can empower individuals to take proactive steps towards managing their health and making informed lifestyle choices.