

```
In [1]: #import Library
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
%matplotlib inline
```

Importing important libraries for prediction

```
In [2]: from sklearn.metrics import accuracy_score
from pandas.plotting import scatter_matrix
```

Loading Dataset

```
In [3]: df=pd.read_csv('diabetes.csv')
df

Out[3]:
   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        180  32.9              0.171    63         0
764             2         122             70             27         0  26.6              0.340    27         0
765             5         121             72             23        112  26.2              0.245    30         0
766             1         126             60              0         0  30.1              0.349    47         1
767             1          93             70             31         0  30.4              0.315    23         0

768 rows x 9 columns

In [4]: # Top 5 Rows
df.head()

Out[4]:
   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

In [5]: # Bottom 5 Rows
df.tail()

Out[5]:
   Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin  BMI  DiabetesPedigreeFunction  Age  Outcome
763            10         101             76             48        180  32.9              0.171    63         0
764             2         122             70             27         0  26.6              0.340    27         0
765             5         121             72             23        112  26.2              0.245    30         0
766             1         126             60              0         0  30.1              0.349    47         1
767             1          93             70             31         0  30.4              0.315    23         0

In [6]: # No of rows and column
df.shape

Out[6]:
(768, 9)

In [7]: # summary of dataset
df.info()

Out[7]:
<bound method DataFrame.info of
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        180  32.9              0.171    63         0
764              2         122             70             27         0  26.6              0.340    27         0
765              5         121             72             23        112  26.2              0.245    30         0
766              1         126             60              0         0  30.1              0.349    47         1
767              1          93             70             31         0  30.4              0.315    23         0

DiabetesPedigreeFunction      Age      Outcome
0              0.627    50         1
1              0.351    31         0
2              0.672    32         1
3              0.167    21         0
4              2.288    33         1
...         ...         ...         ...
763             0.171    63         0
764             0.340    27         0
765             0.245    30         0
766             0.349    47         1
767             0.315    23         0

[768 rows x 9 columns]>

In [8]: # No of columns in dataset
df.columns

Out[8]:
Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
       'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
      dtype='object')

In [9]: # Datatypes
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      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

In [10]: df.describe()

Out[10]:
   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  946.000000  67.100000   2.420000  81.000000   1.000000
```

```
In [11]: # dataset with transpose
df.describe().T

Out[11]:
   count  mean  std  min  25%  50%  75%  max
Pregnancies  768.0  3.845052  3.369578  0.000  1.000000  3.0000  6.000000  17.00
Glucose      768.0  120.894531  31.972618  0.000  99.000000  117.0000  140.250000  199.00
BloodPressure  768.0  69.105469  19.355807  0.000  62.000000  72.0000  80.000000  122.00
SkinThickness  768.0  20.536458  15.952218  0.000  0.000000  23.0000  32.000000  99.00
Insulin      768.0  79.799479  115.244002  0.000  0.000000  30.5000  127.250000  946.00
BMI          768.0  31.992578  7.884160  0.000  27.300000  32.0000  36.600000  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  24.000000  29.0000  41.000000  81.00
Outcome      768.0  0.348958  0.476951  0.000  0.000000  0.0000  1.000000  1.00
```

Checking of Missing Values

```
In [12]: # checking for total null values
df.isnull().sum()

Out[12]:
Pregnancies      0
Glucose          0
BloodPressure     0
SkinThickness     0
Insulin          0
BMI              0
DiabetesPedigreeFunction  0
Age              0
Outcome          0

DiabetesPedigreeFunction      Age      Outcome
0              False      False      False
1              False      False      False
2              False      False      False
3              False      False      False
4              False      False      False
...         ...         ...         ...
763             False      False      False
764             False      False      False
765             False      False      False
766             False      False      False
767             False      False      False

[768 rows x 9 columns]>

In [13]: sns.heatmap(df.isnull())

Out[13]:
<Axes: >

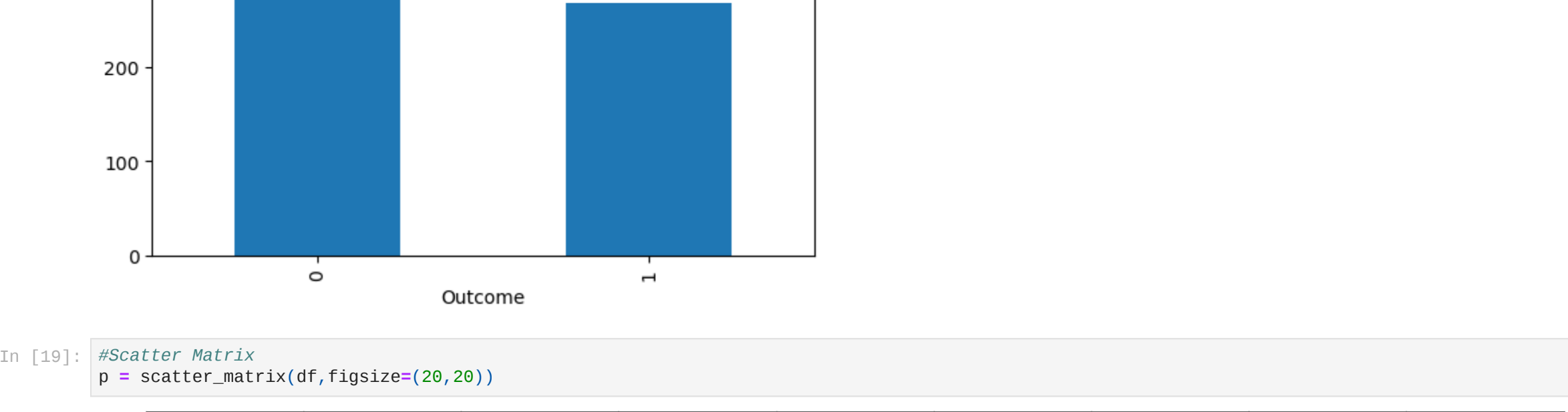

```

Data Visualization



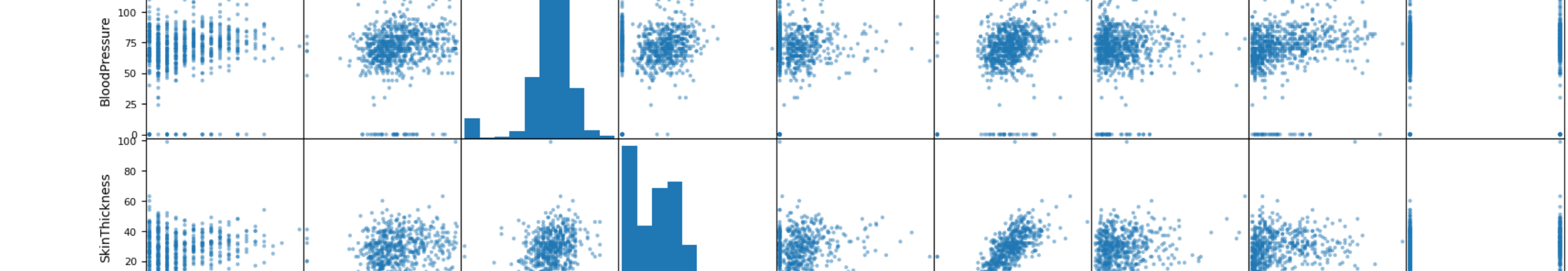
```
In [15]: #aiming to input NAN values
df['Glucose'].fillna(df['Glucose'].mean(),inplace=True)
df['BloodPressure'].fillna(df['BloodPressure'].mean(),inplace=True)
df['SkinThickness'].fillna(df['SkinThickness'].median(),inplace=True)
df['Insulin'].fillna(df['Insulin'].median(),inplace=True)
df['BMI'].fillna(df['BMI'].median(),inplace=True)

In [16]: #Distribution after removing NAN values
p=df.hist(figsize=(20,20))
```



```
In [18]: color_wheel=plt.cm.get_cmap('hsv',256)
color= df['Outcome'].map(lambda x:color_wheel.get(x+1))
print(df.Outcome.value_counts())
p = df.Outcome.value_counts().plot(kind="bar")

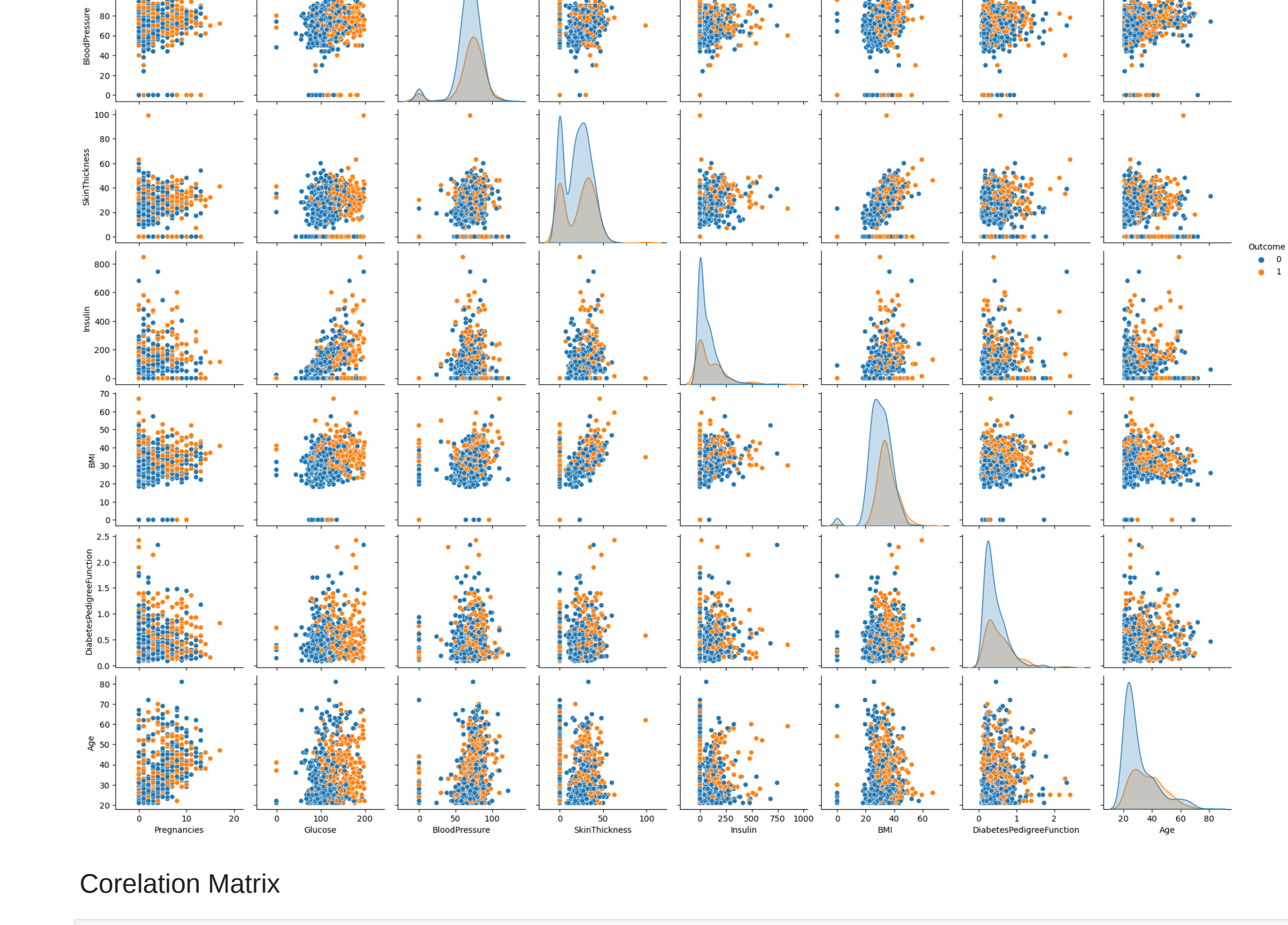
Outcome
0      580
1      268
Name: count, dtype: int64
```



```
In [19]: #Scatter Matrix
p = scatter_matrix(df,figsize=(20,20))
```



```
In [22]: # pairs plot for data
sns.pairplot(df,hue='Outcome')
plt.show()
```

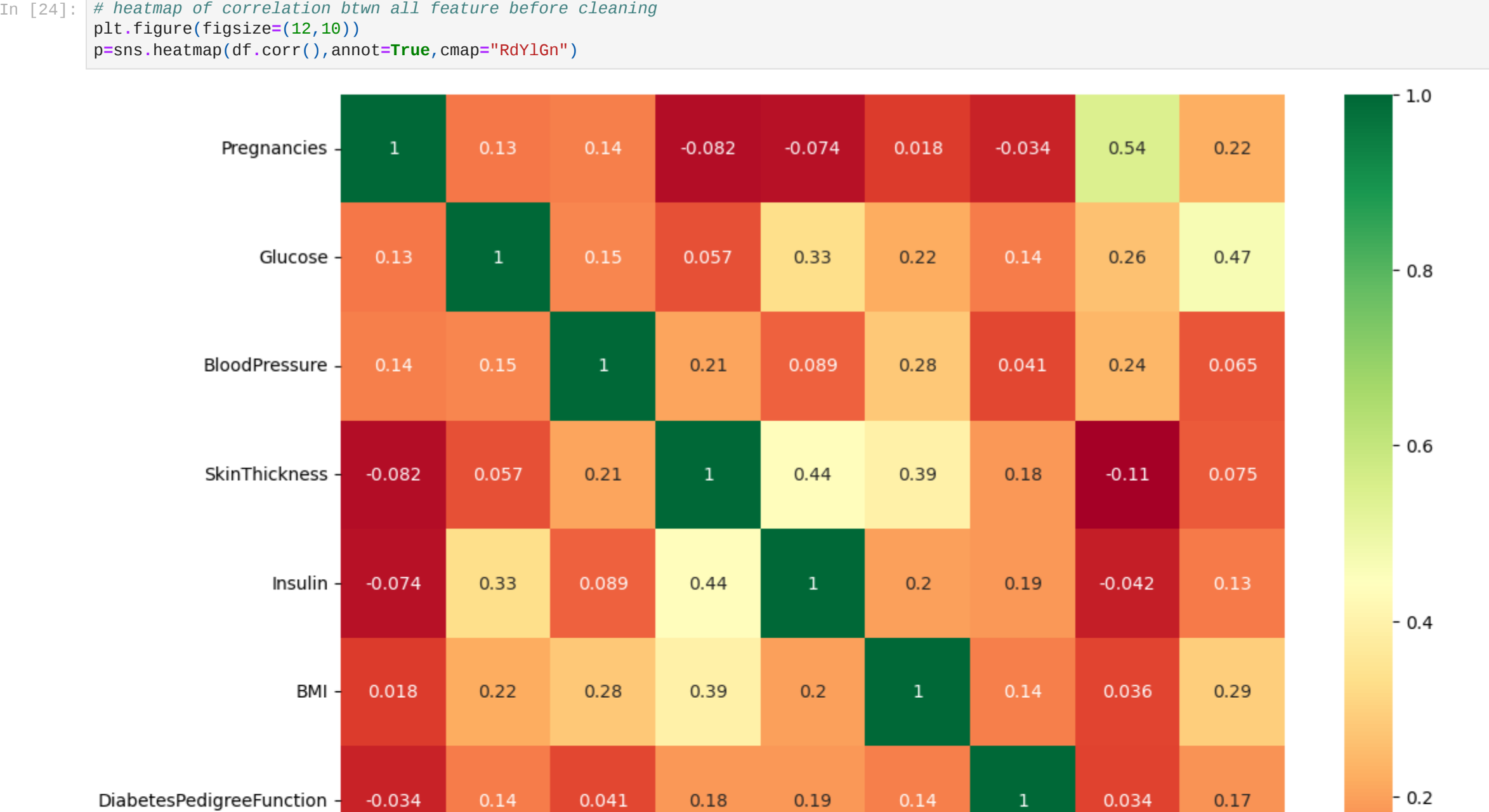


Correlation Matrix

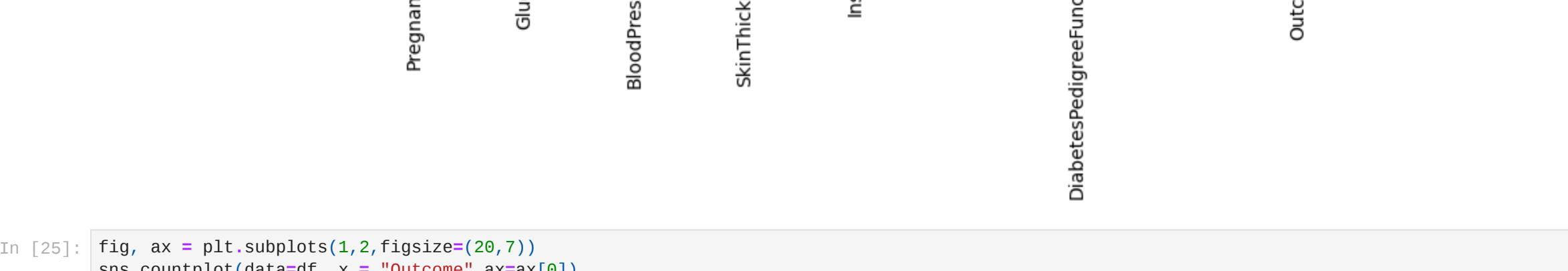
```
In [23]: correlation=df.corr()
correlation

Pregnancies      Glucose      BloodPressure      SkinThickness      BMI      DiabetesPedigreeFunction      Age      Outcome
Pregnancies      1.000000      0.129459      0.141282      -0.081672      0.073528      0.141282      0.073528
Glucose           0.129459      1.000000      0.152590      0.097328      0.267371      0.009090      0.436783
BloodPressure     0.141282      0.152590      1.000000      0.088933      0.436783      0.009090      0.392573
SkinThickness     0.081672      0.057328      0.267371      1.000000      0.183928      0.183928      0.183928
Insulin          0.073528      0.392573      0.436783      0.088933      1.000000      0.183928      0.183928
BMI              0.073528      0.392573      0.436783      0.088933      0.183928      1.000000      0.183928
DiabetesPedigreeFunction  0.009090      0.183928      0.183928      0.183928      0.183928      0.183928      1.000000
Age              0.141282      0.073528      0.073528      0.073528      0.073528      0.073528      0.183928
Outcome          0.436783      0.392573      0.392573      0.183928      0.183928      0.183928      0.183928

In [24]: # heatmap of correlation btwn all feature before cleaning
plt.figure(figsize=(12,10))
sns.heatmap(df.corr(),annot=True,cmap='RdYlGn')
```



```
In [25]: fig, ax = plt.subplots(1,2,figsize=(20,7))
sns.countplot(data=df, x = "Outcome",ax=ax[0])
df['Outcome'].value_counts().plot.pie(explode=[0.1,0],autopct='%1.1f%%',labels=["No","Yes"],shadow=True,ax=ax[1])
plt.show()
```



we observe from above plot that:

65.1% patients in the dataset do not have diabetes. 34.9% patients in the dataset has diabetes.

conclusion

1 It has a decent level of precision, indicating that when it predicts positive cases (diabetes), it's correct about 65% of the time. 2 out of the 768 patients, 268 have been diagnosed with diabetes. 3 patients with high blood pressure have greater chances of diabetes. 4 an increase in blood pressure BMI and skin thickness also increases. 5 Increasing level of glucose and insulin increases chances of diabetes.