

# Diabetes Project

November 27 2023

## DIABETIC PREDICTION USING CLASSIFICATION ALGORITHMS

```
In [48]: #IMPORT LIBRARIES
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

```
In [49]: df=pd.read_csv('diabetes (Task1).csv')
df
```

```
Out[49]:
```

	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	36.8	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 × 9 columns

```
In [50]: #generating the first five rows of the data
df.head()
```

```
Out[50]:
```

	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 [51]: #generating the last five rows of the data
df.tail()
```

```
Out[51]:
```

	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	36.8	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 [52]: #Checking the columns
df.describe()
```

Out[52]:

	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

In [53]:

df.dtypes

Out[53]:

Pregnancies

int64

Glucose

int64

BloodPressure

int64

SkinThickness

int64

Insulin

int64

BMI

float64

DiabetesPedigreeFunction

float64

Age

int64

Outcome

int64

dtype: object

In [54]:

df.describe

Out[54]:

<bound method NDFrame.describe of

Pregnancies

Glucose

BloodPressure

SkinThickness

Insulin

BMI

\

0

6

148

72

35

0

33.6

1

1

85

66

29

0

26.6

2

8

183

64

0

0

23.3

3

1

89

66

23

94

28.1

4

0

137

40

35

168

43.1

..

...

...

...

...

...

...

763

10

101

76

48

180

32.9

764

2

122

70

27

0

36.8

765

5

121

72

23

112

26.2

766

1

126

60

0

0

30.1

767

1

93

70

31

0

30.4

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 [55]:

# Information about all the columns

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 [56]:

#To check null values

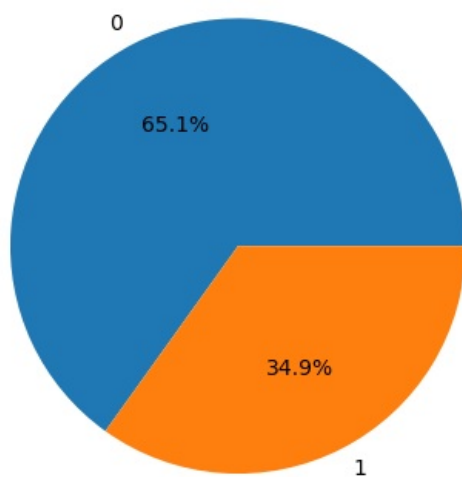
df.isna().sum()

```
Out[56]: Pregnancies      0
         Glucose          0
         BloodPressure    0
         SkinThickness    0
         Insulin          0
         BMI              0
         DiabetesPedigreeFunction  0
         Age              0
         Outcome          0
         dtype: int64
```

```
In [57]: #check the dataset is balanced or not
         df["Outcome"].value_counts()
```

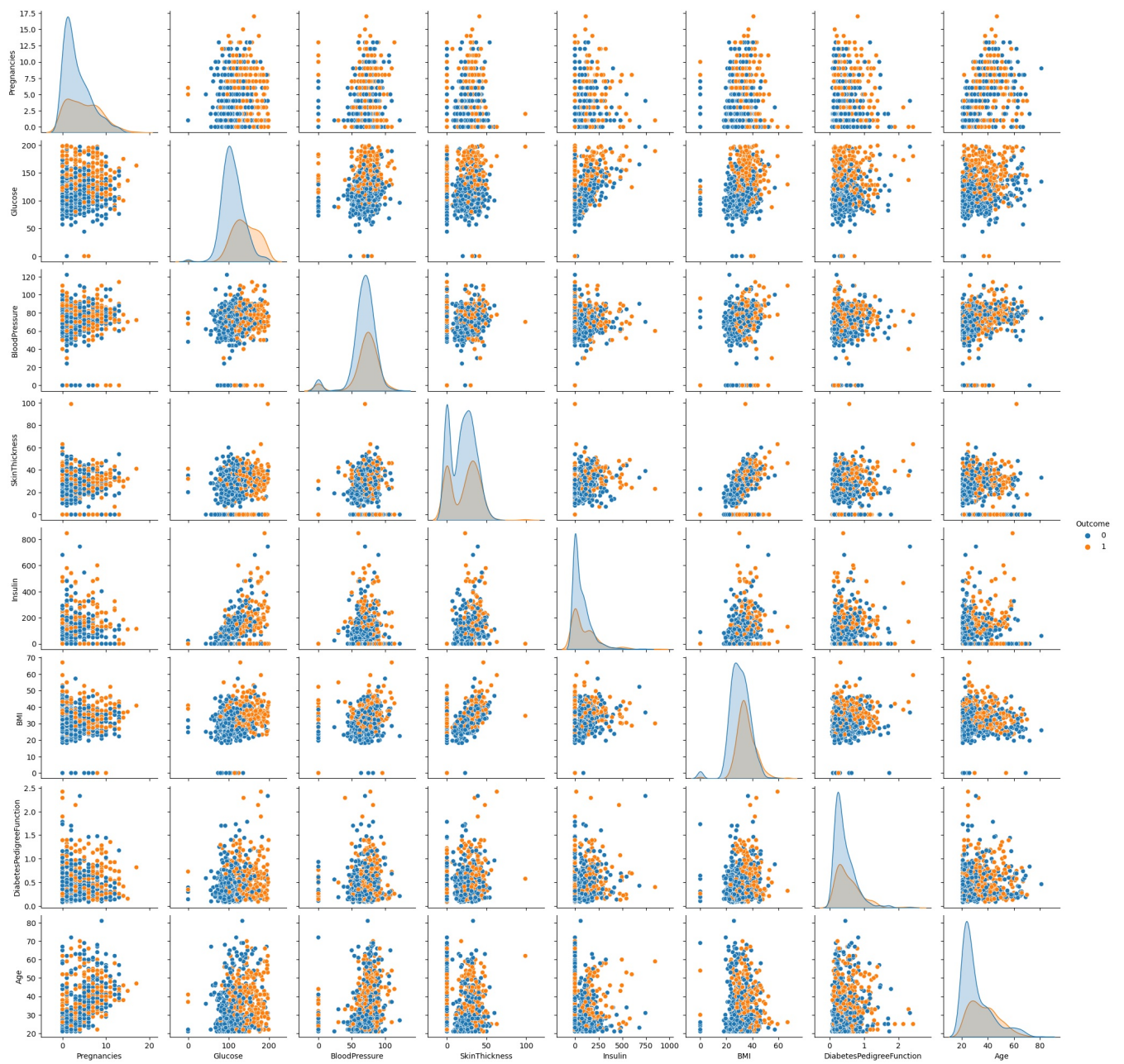
```
Out[57]: Outcome
         0      500
         1      268
         Name: count, dtype: int64
```

```
In [58]: #check the dataset is balanced or not using pie chart
         plt.pie(df['Outcome'].value_counts(), labels=['0', '1'], autopct='%1.1f%%')
         plt.show()
```

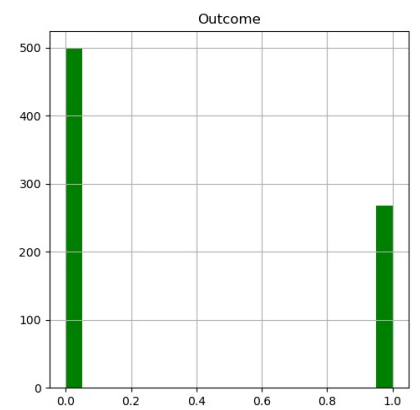
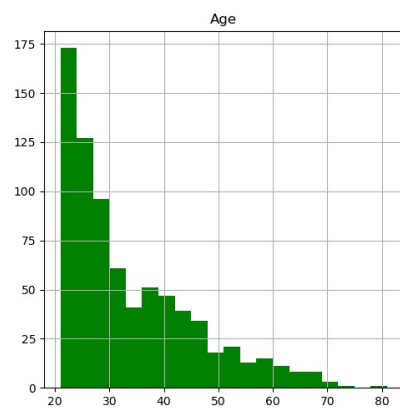
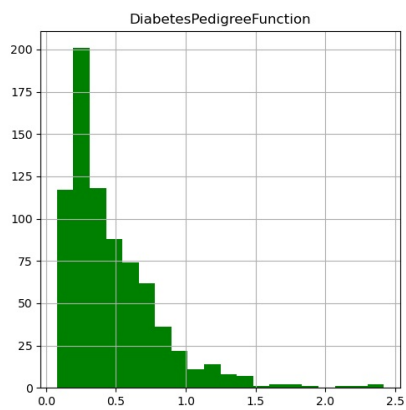
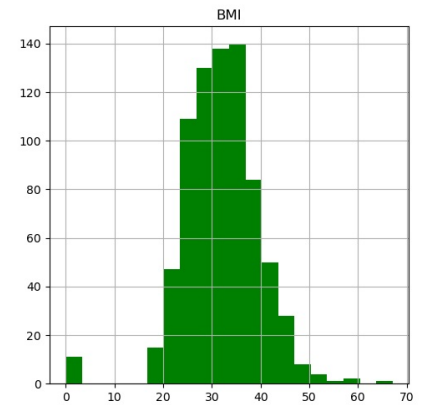
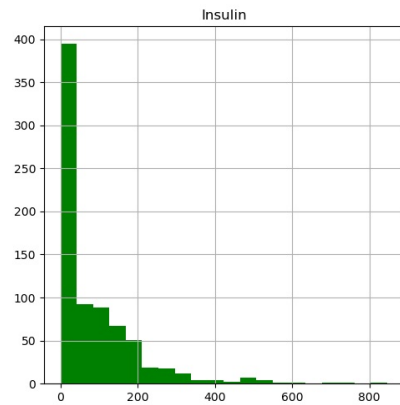
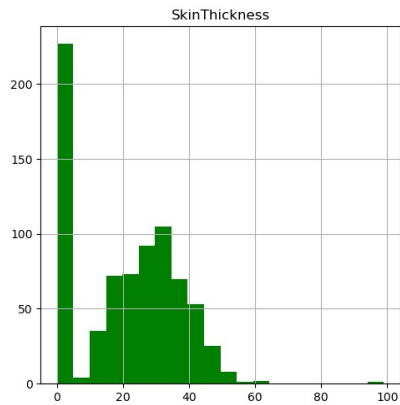
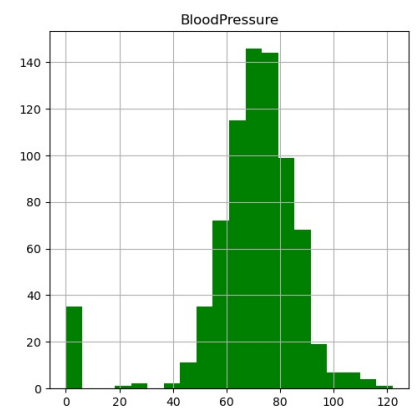
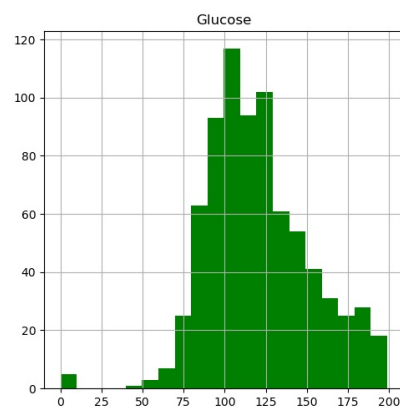
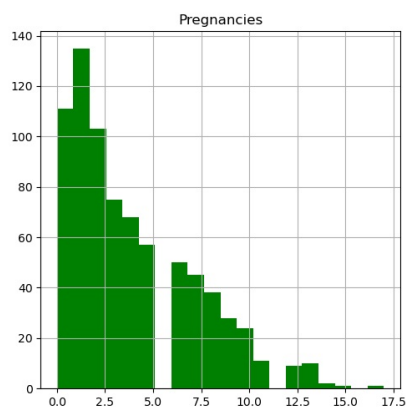


```
In [59]: sns.pairplot(df, hue = 'Outcome')
```

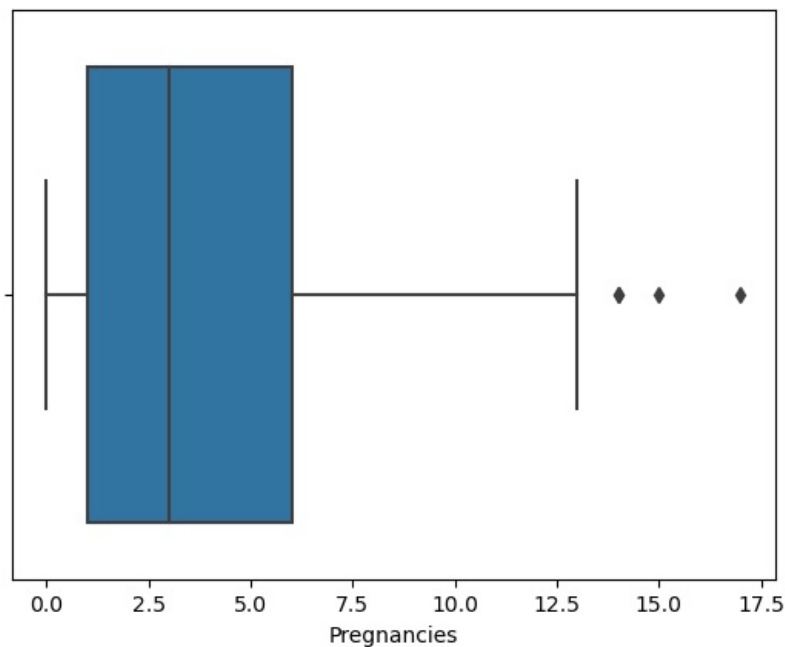
```
Out[59]: <seaborn.axisgrid.PairGrid at 0x1a6fc84ae10>
```

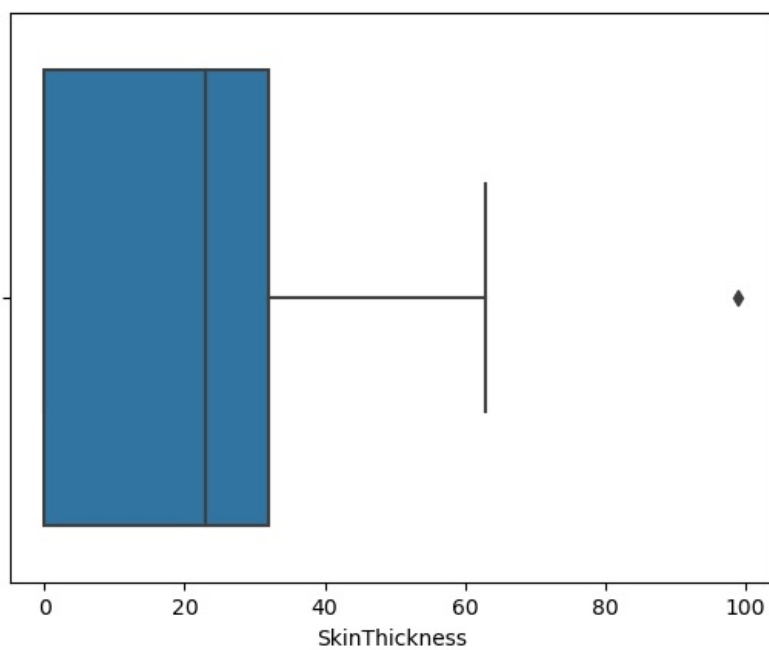
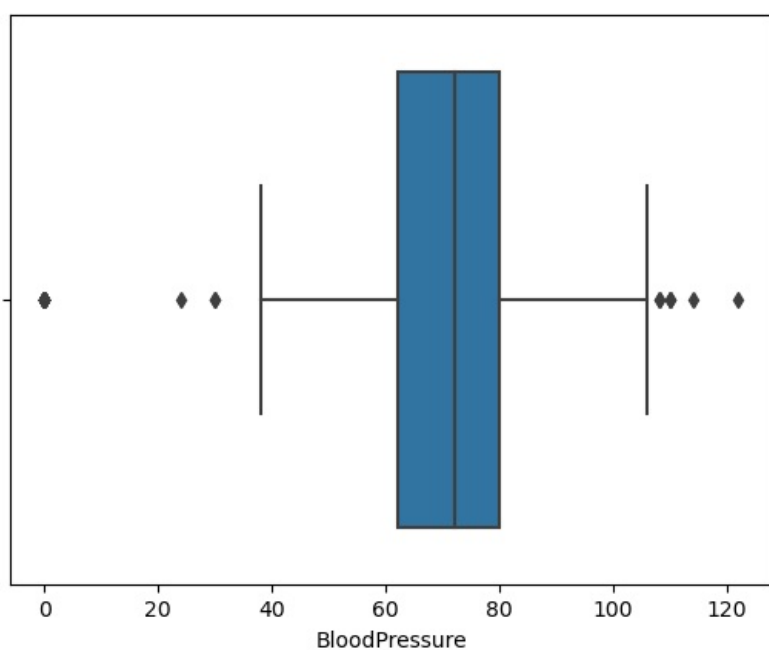
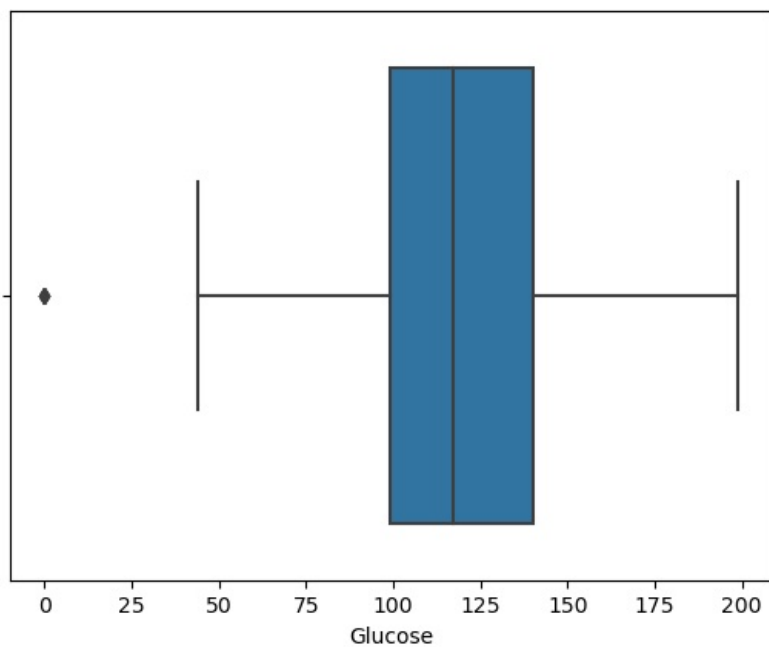


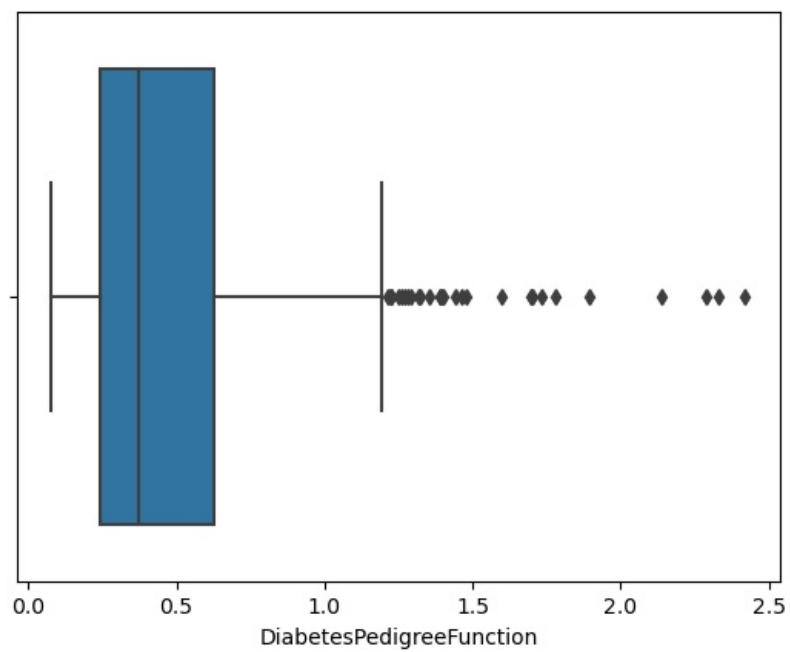
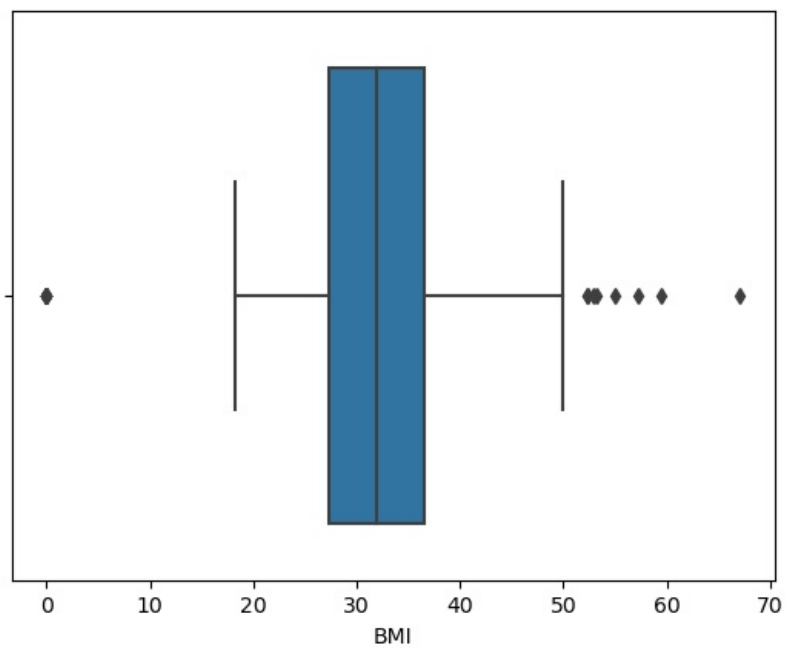
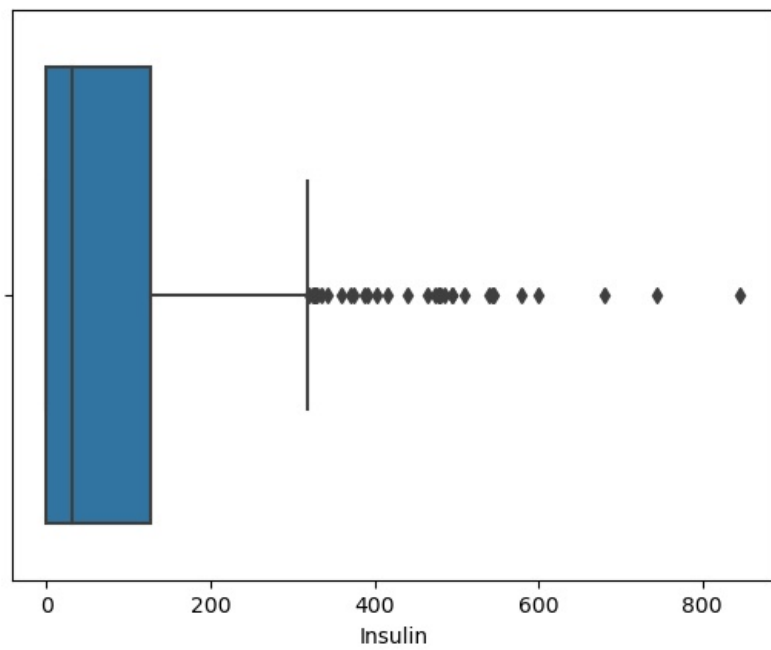
```
In [60]: df.hist(bins=20, figsize=(20, 20), color='green')
plt.show()
```

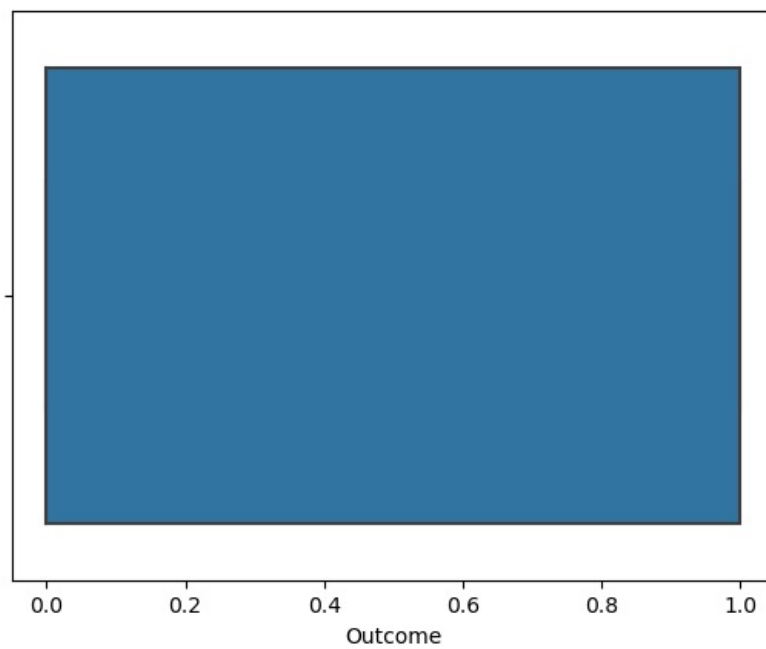
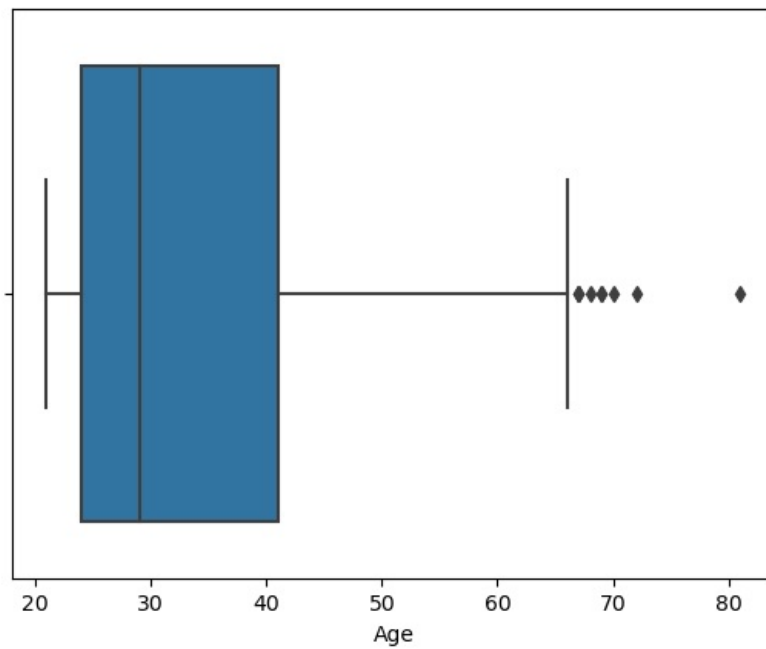


```
In [61]: #To check the presence of outliers
for i in df:
    sns.boxplot(x=df[i])
    plt.show()
```





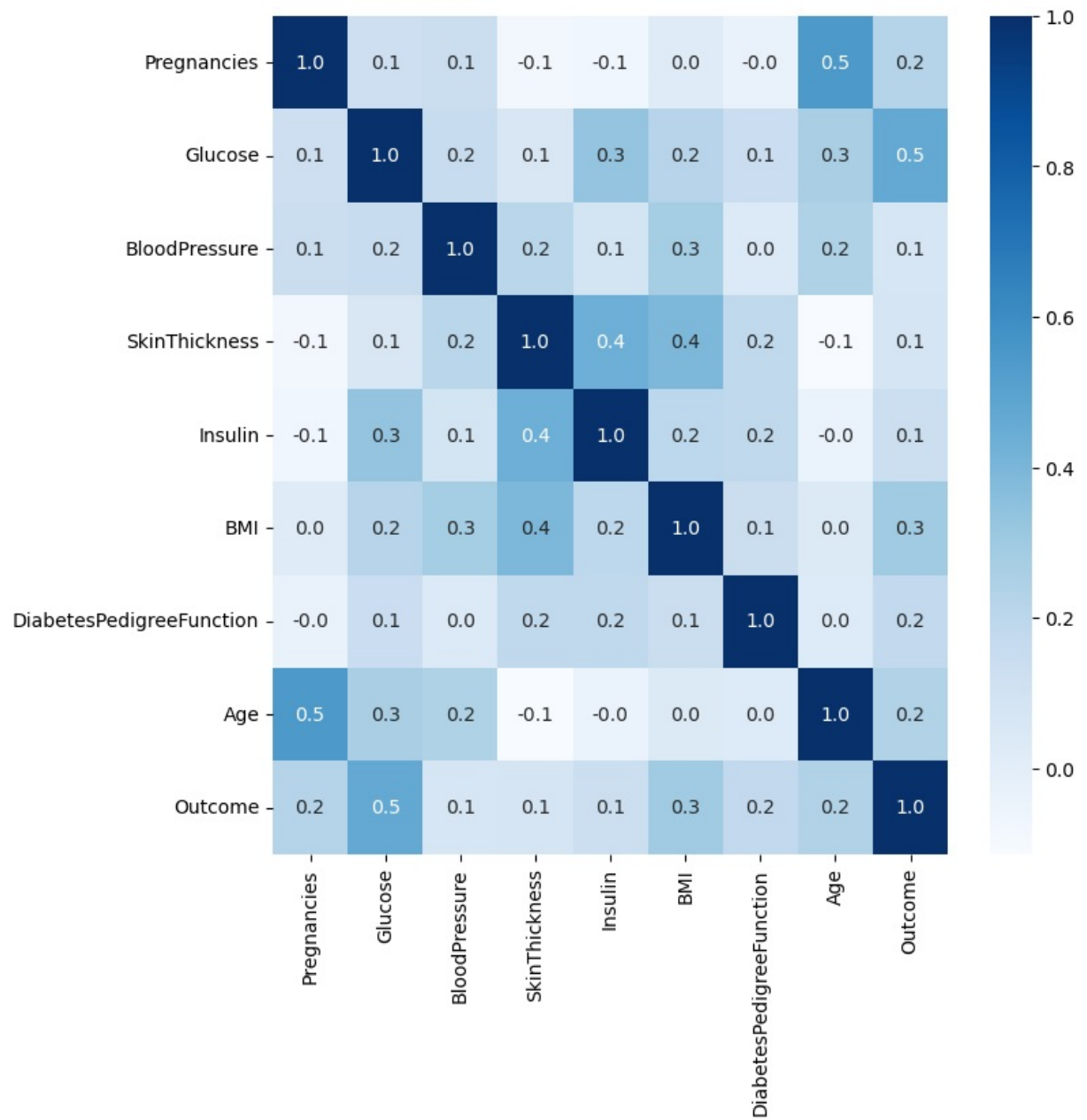




```
In [62]: #Removing outliers using z_scores technique
from scipy import stats
z_score_threshold = 3.0
z_scores = np.abs(stats.zscore(df))
outlier_mask = (z_scores > z_score_threshold).any(axis=1)
df1 = df[~outlier_mask]
```

```
In [63]: plt.figure(figsize=(8,8))
sns.heatmap(df.corr(), annot=True, cmap="Blues", fmt=".1f")
plt.show()
```





```
In [64]: x=df.iloc[:, :-1].values
          x
          y=df.iloc[:, -1].values
          y
```



```

from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
knn=KNeighborsClassifier(n_neighbors=7)
nb=GaussianNB()
svm=SVC()
tree_model=DecisionTreeClassifier(criterion='entropy')
forest_model=RandomForestClassifier(n_estimators=100)
lst=[knn,nb,svm,tree_model,forest_model]

```

```

In [71]: from sklearn.metrics import confusion_matrix,accuracy_score,classification_report,ConfusionMatrixDisplay
for i in lst:
    print(i)
    print('*****')
    i.fit(x_train,y_train)
    y_pred=i.predict(x_test)
    print(confusion_matrix(y_test,y_pred))
    print(classification_report(y_test,y_pred))
    labels=[0,1]
    result=confusion_matrix(y_test,y_pred)
    cmd=ConfusionMatrixDisplay(result,display_labels=labels)
    cmd.plot()

```

KNeighborsClassifier(n\_neighbors=7)

\*\*\*\*\*

```

[[121 30]
 [ 38 42]]

```

	precision	recall	f1-score	support
0	0.76	0.80	0.78	151
1	0.58	0.53	0.55	80
accuracy			0.71	231
macro avg	0.67	0.66	0.67	231
weighted avg	0.70	0.71	0.70	231

GaussianNB()

\*\*\*\*\*

```

[[122 29]
 [ 28 52]]

```

	precision	recall	f1-score	support
0	0.81	0.81	0.81	151
1	0.64	0.65	0.65	80
accuracy			0.75	231
macro avg	0.73	0.73	0.73	231
weighted avg	0.75	0.75	0.75	231

SVC()

\*\*\*\*\*

```

[[126 25]
 [ 32 48]]

```

	precision	recall	f1-score	support
0	0.80	0.83	0.82	151
1	0.66	0.60	0.63	80
accuracy			0.75	231
macro avg	0.73	0.72	0.72	231
weighted avg	0.75	0.75	0.75	231

DecisionTreeClassifier(criterion='entropy')

\*\*\*\*\*

```

[[112 39]
 [ 25 55]]

```

	precision	recall	f1-score	support
0	0.82	0.74	0.78	151
1	0.59	0.69	0.63	80
accuracy			0.72	231
macro avg	0.70	0.71	0.70	231
weighted avg	0.74	0.72	0.73	231

RandomForestClassifier()

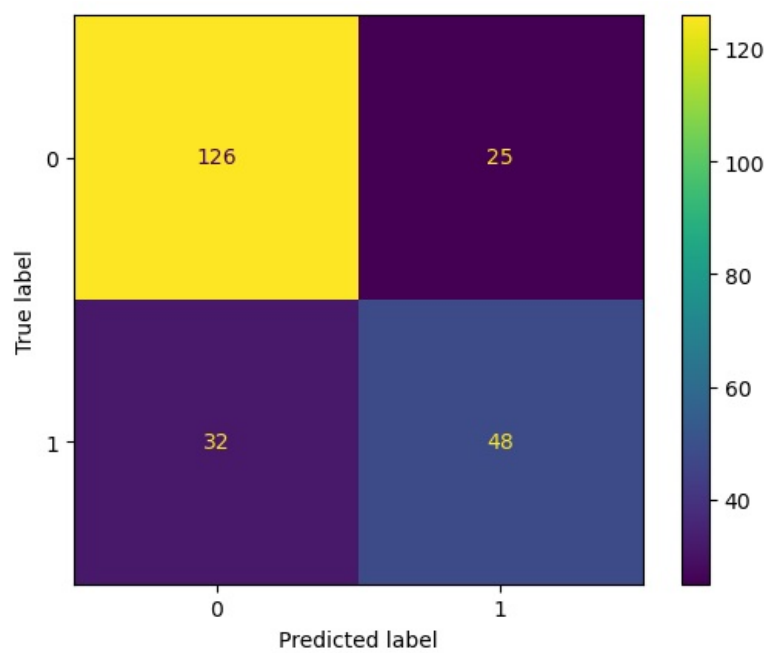
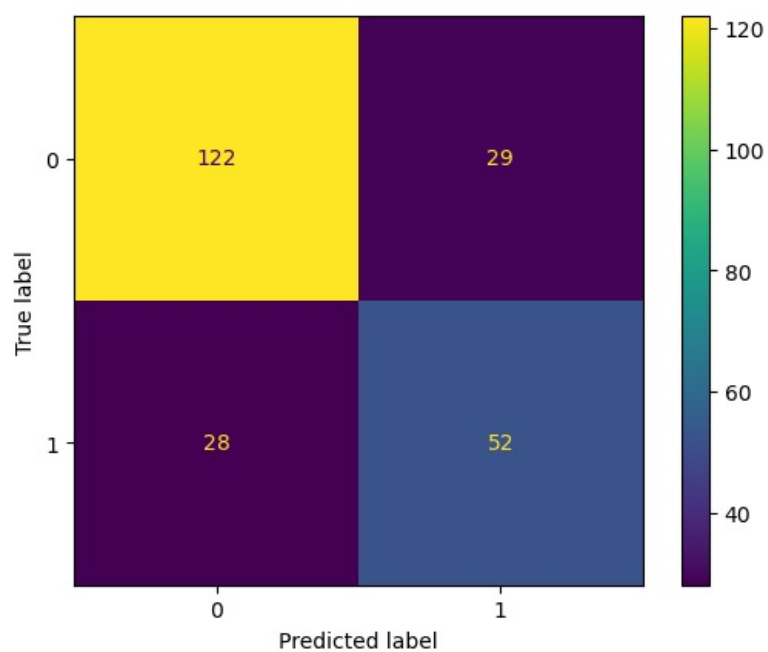
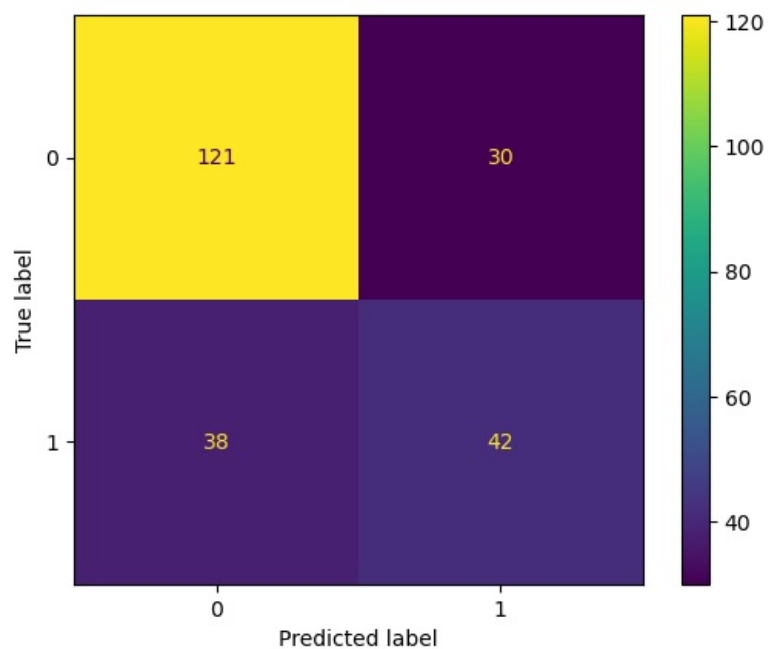
\*\*\*\*\*

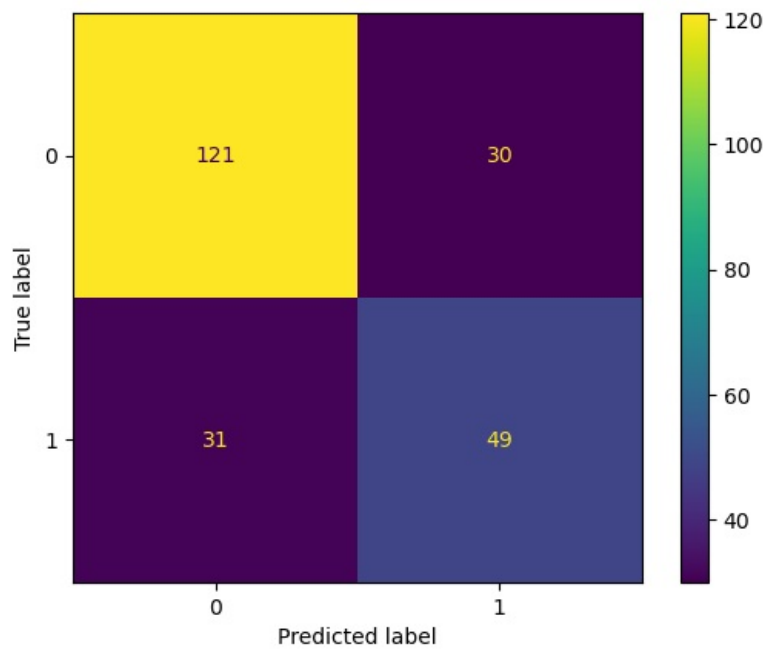
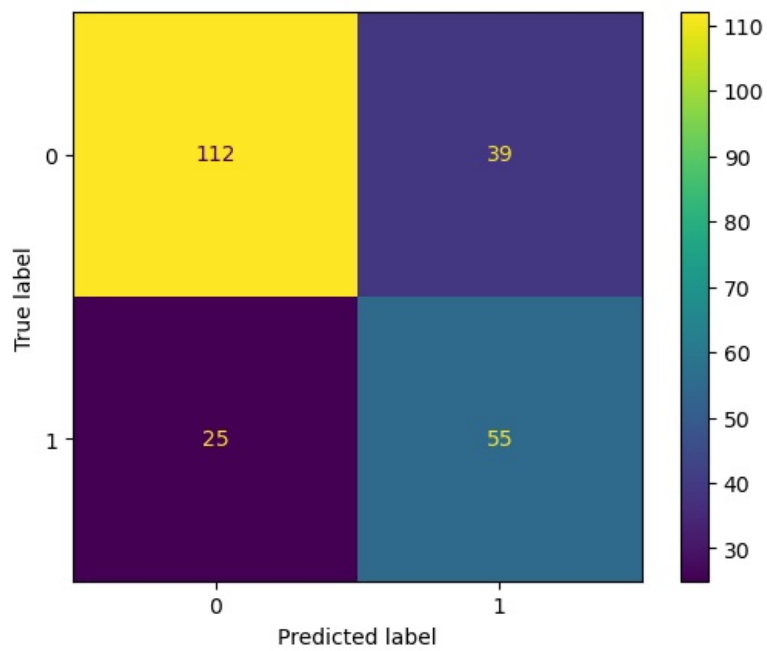
```

[[121 30]
 [ 31 49]]

```

	precision	recall	f1-score	support
0	0.80	0.80	0.80	151
1	0.62	0.61	0.62	80
accuracy			0.74	231
macro avg	0.71	0.71	0.71	231
weighted avg	0.74	0.74	0.74	231





```
In [72]: data=nb.predict(scaler.transform([[5,200,90,52,290,32,0.300,27]]))
if data==0:
    print('diabetic patient')
else:
    print('not diabetic')
```

not diabetic

## CONCLUSION

From above observation we concluded that in comparison to other models Naive Bayes works more better(accuracy-73%).