

```

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

```

## Reading diabetes dataset

```

diabetes_df = pd.read_csv('/content/diabetes.csv')
diabetes_df.head()

```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeF
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	

```
diabetes_df.columns
```

```

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

```

```
diabetes_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
```

Describing dataset

```
diabetes_df.describe()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Di
count	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	
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	

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

```
Pregnancies      0
Glucose           0
BloodPressure     0
SkinThickness     0
Insulin           0
BMI               0
DiabetesPedigreeFunction  0
Age               0
Outcome           0
dtype: int64
```

checking null values

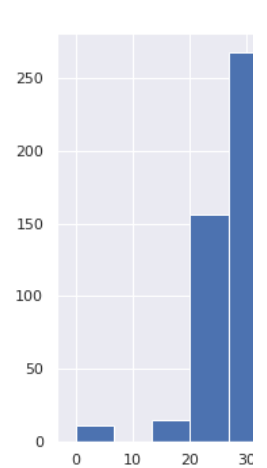
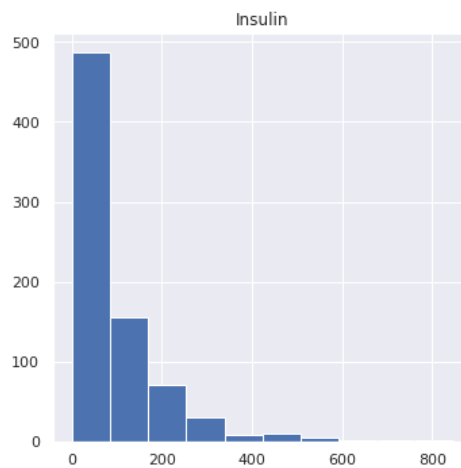
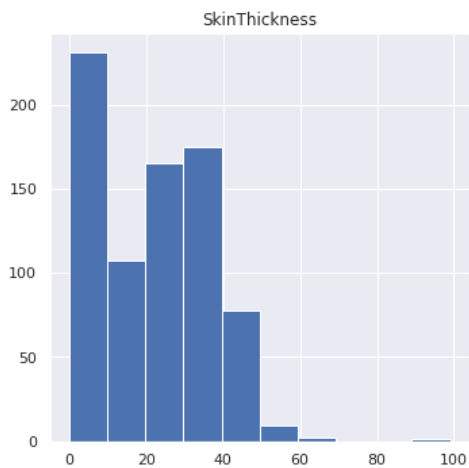
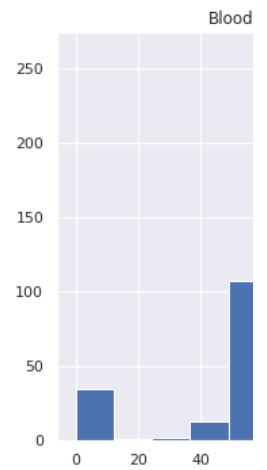
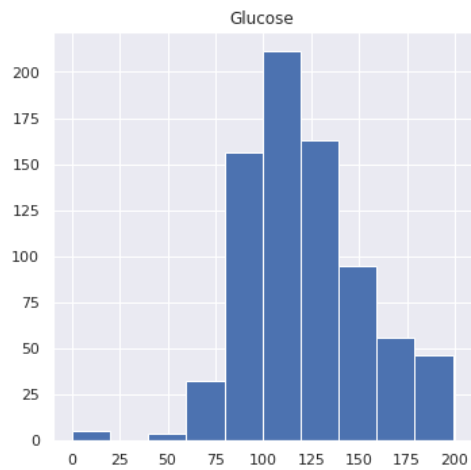
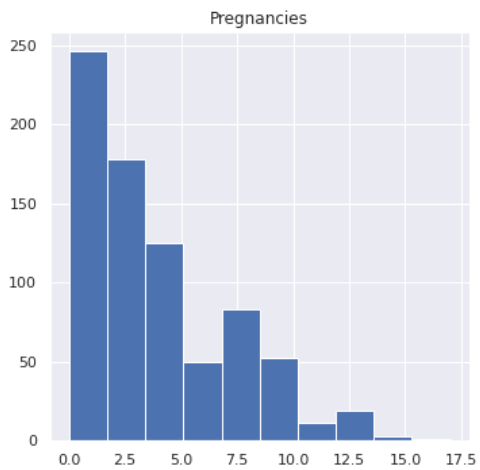
```
diabetes_df_copy = diabetes_df.copy(deep = True)
diabetes_df_copy[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']] = diabetes_df_c

# Showing the Count of NANS
print(diabetes_df_copy.isnull().sum())
```

Pregnancies	0
Glucose	5
BloodPressure	35
SkinThickness	227
Insulin	374
BMI	11
DiabetesPedigreeFunction	0
Age	0
Outcome	0
dtype: int64	

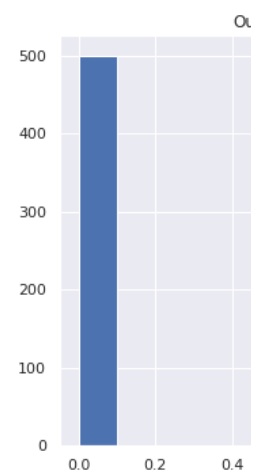
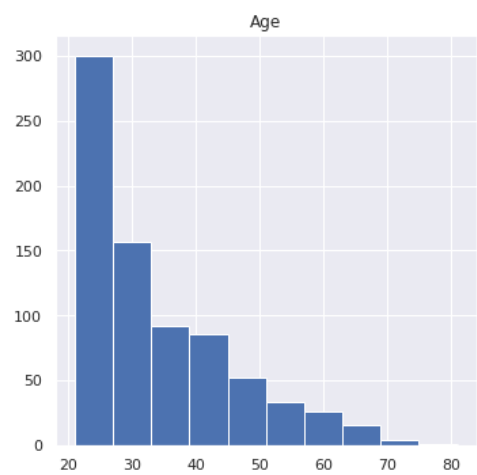
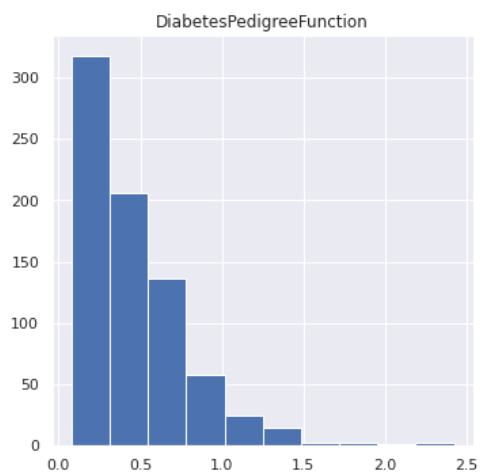
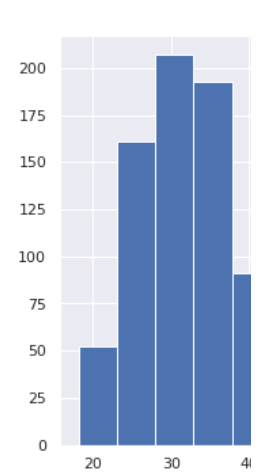
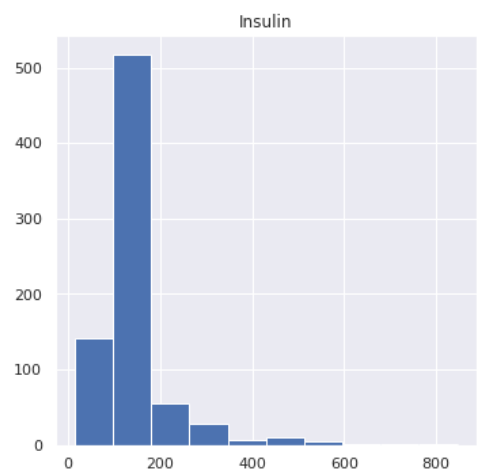
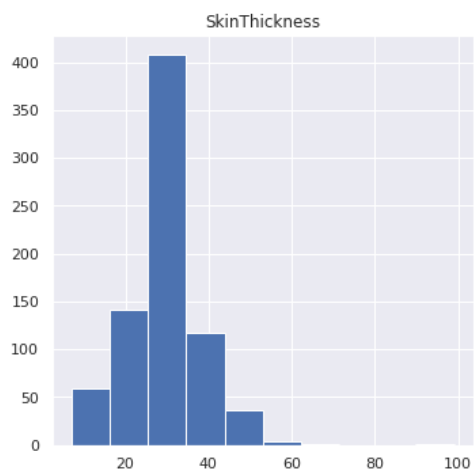
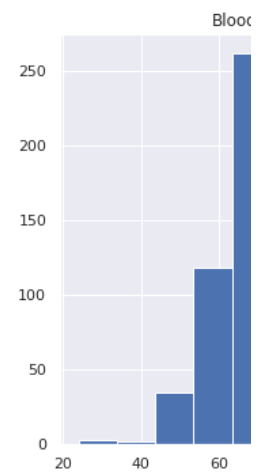
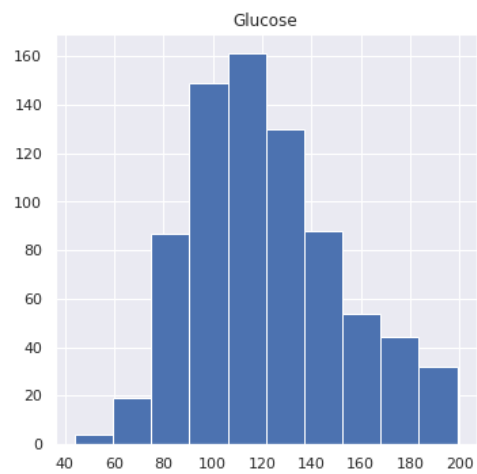
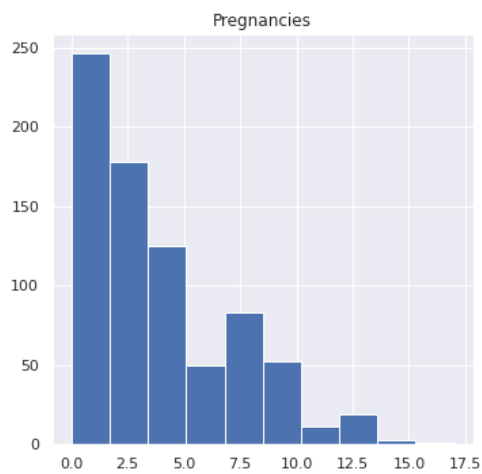
### plotting graphs on diabetes dataset

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

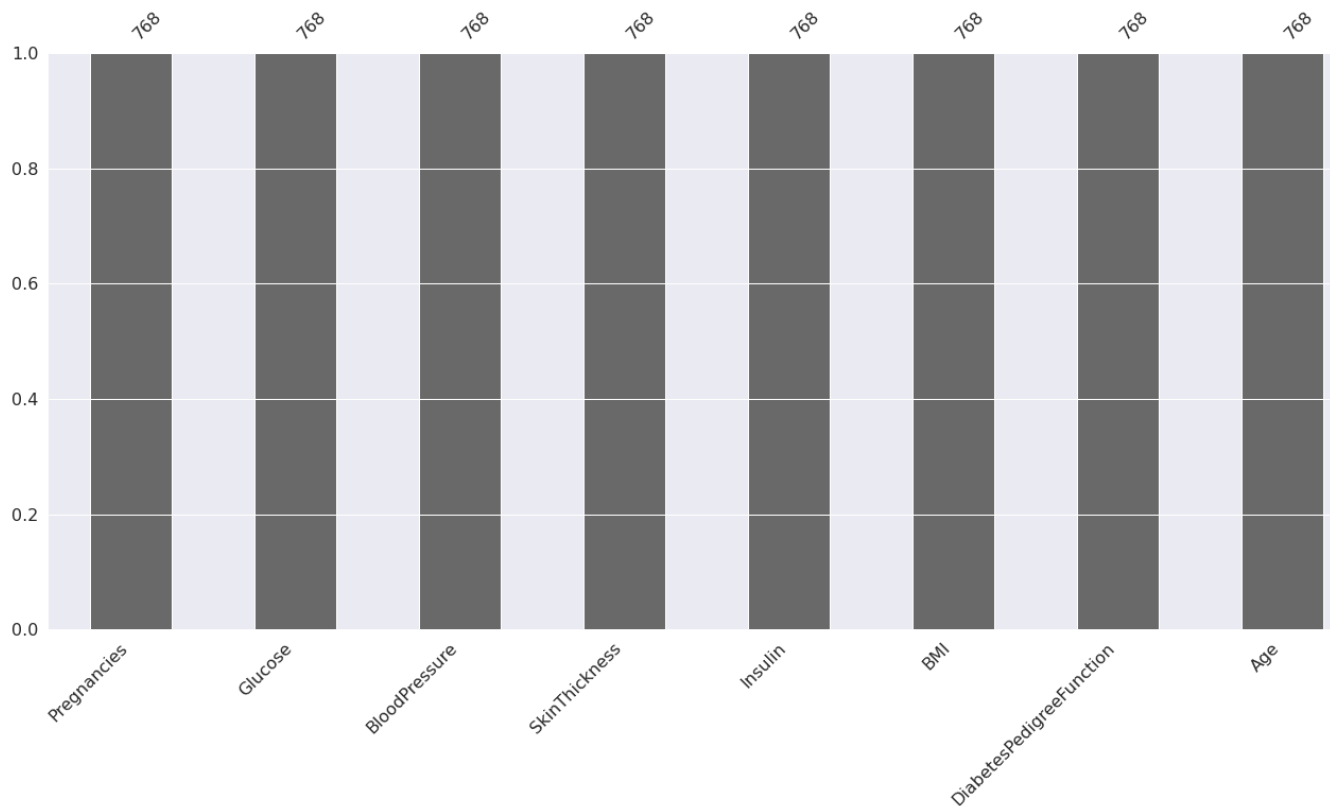


```
diabetes_df_copy['Glucose'].fillna(diabetes_df_copy['Glucose'].mean(), inplace = True)
diabetes_df_copy['BloodPressure'].fillna(diabetes_df_copy['BloodPressure'].mean(), inplace = True)
diabetes_df_copy['SkinThickness'].fillna(diabetes_df_copy['SkinThickness'].median(), inplace = True)
diabetes_df_copy['Insulin'].fillna(diabetes_df_copy['Insulin'].median(), inplace = True)
diabetes_df_copy['BMI'].fillna(diabetes_df_copy['BMI'].median(), inplace = True)
```

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



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



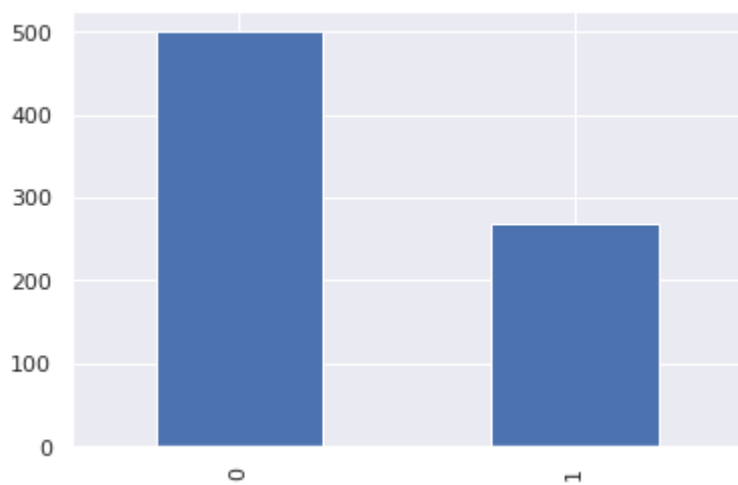
### count plot on diabetes dataset

```
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")
```

```
0    500
```

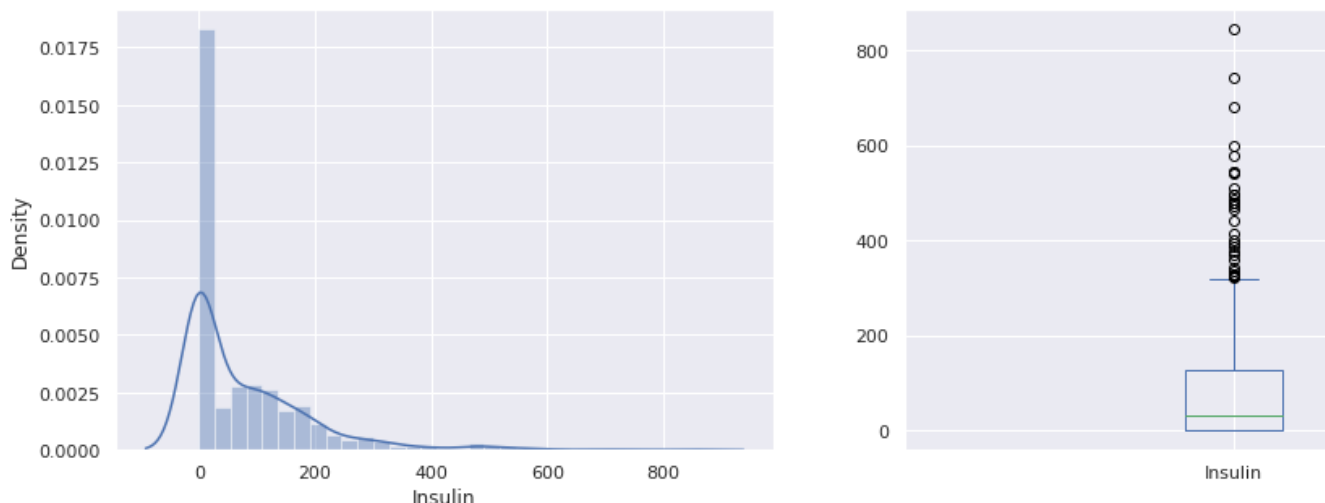
```
1    268
```

```
Name: Outcome, dtype: int64
```



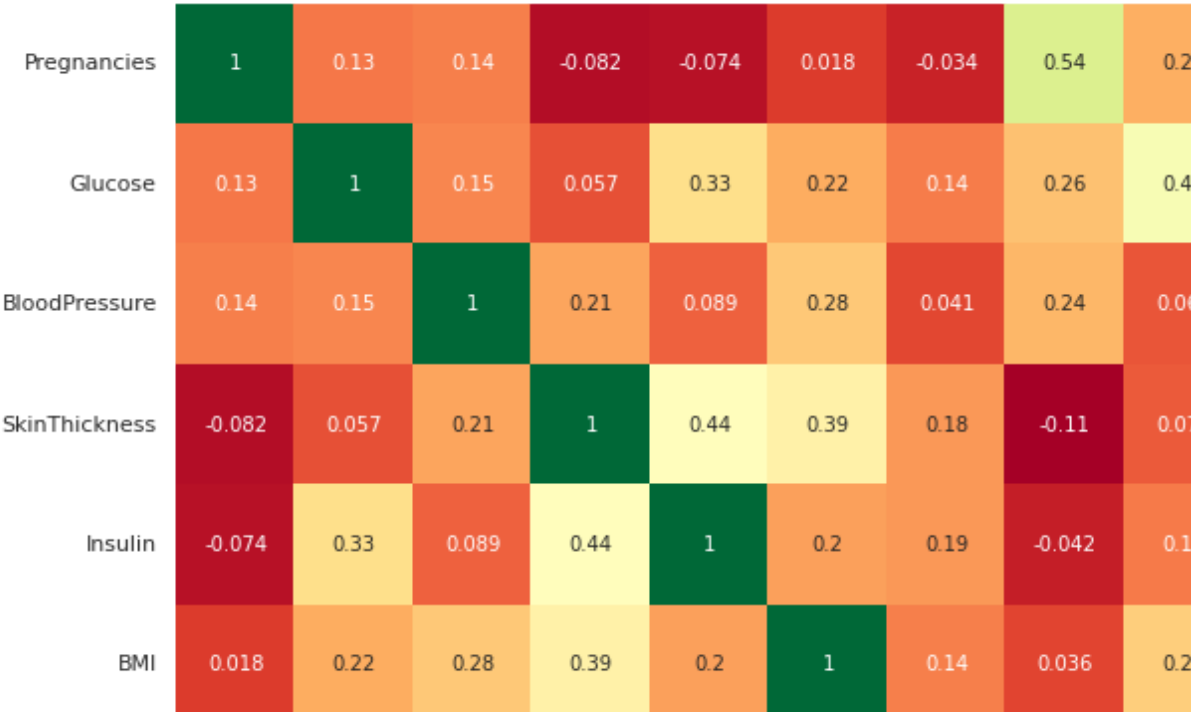
## checking insulin

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



## Feature selection using heat map

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



```
diabetes_df_copy.head()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeF
0	6	148.0	72.0	35.0	125.0	33.6	
1	1	85.0	66.0	29.0	125.0	26.6	
2	8	183.0	64.0	29.0	125.0	23.3	
3	1	89.0	66.0	23.0	94.0	28.1	
4	0	137.0	40.0	35.0	168.0	43.1	

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

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPe
0	0.639947	0.865108	-0.033518	0.670643	-0.181541	0.166619	
1	-0.844885	-1.206162	-0.529859	-0.012301	-0.181541	-0.852200	
2	1.233880	2.015813	-0.695306	-0.012301	-0.181541	-1.332500	
3	-0.844885	-1.074652	-0.529859	-0.695245	-0.540642	-0.633881	
4	-1.141852	0.503458	-2.680669	0.670643	0.316566	1.549303	

target values



```
y = diabetes_df_copy.Outcome
y
```

```
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
```

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

## splitting data in train and test

```
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)
```

## Using Random Forest Classifier algorithm

```
from sklearn.ensemble import RandomForestClassifier

rfc = RandomForestClassifier(n_estimators=200)
rfc.fit(X_train, y_train)

RandomForestClassifier(n_estimators=200)

rfc_train = rfc.predict(X_train)
from sklearn import metrics

print("Accuracy_Score =", format(metrics.accuracy_score(y_train, rfc_train)))

Accuracy_Score = 1.0

from sklearn import metrics

rfc_predictions = rfc.predict(X_test)
print("Random Forest Accuracy Score =", format(metrics.accuracy_score(y_test, rfc_predictions)))
```

```
Random Forest Accuracy_Score = 0.7559055118110236
```

## Using Decision Tree Classifier algorithm

```
from sklearn.tree import DecisionTreeClassifier
```

```
dtree = DecisionTreeClassifier()  
dtree.fit(X_train, y_train)
```

```
DecisionTreeClassifier()
```

```
from sklearn import metrics
```

```
predictions = dtree.predict(X_test)  
print("Decission Tree Accuracy Score =", format(metrics.accuracy_score(y_test, predictions)))
```

```
Decission Tree Accuracy Score = 0.7086614173228346
```

## Using XGBoost Classifier algorithm

```
from xgboost import XGBClassifier
```

```
xgb_model = XGBClassifier(gamma=0)  
xgb_model.fit(X_train, y_train)
```

```
XGBClassifier()
```

```
from sklearn import metrics
```

```
xgb_pred = xgb_model.predict(X_test)  
print("XGBoost Accuracy Score =", format(metrics.accuracy_score(y_test, xgb_pred)))
```

```
XGBoost Accuracy Score = 0.7795275590551181
```

## Using SVM Classifier

```
from sklearn.svm import SVC
```

```
svc_model = SVC()  
svc_model.fit(X_train, y_train)
```

```
SVC()
```

```
svc_pred = svc_model.predict(X_test)
```

```
from sklearn import metrics

print("SVC Accuracy Score =", format(metrics.accuracy_score(y_test, svc_pred)))

SVC Accuracy Score = 0.7480314960629921
```

### Saving model in .h5 format

```
import joblib
joblib_file = "Random_Forest_Diabetes.h5"
joblib.dump(rfc, joblib_file)

['Random_Forest_Diabetes.h5']
```

---

✓ 0s completed at 5:15 PM

