## ⌄ Diabeties Predictor Project

using Machine Learning Algorithm

```
1 # Importing pandas and basic libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
```

```
1 data= pd.read_csv("/content/DataSet/diabetes.csv")
2 data.head()
```

|   | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigre |
|---|---|---|---|---|---|---|---|
| 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 | |

Next steps:  | Generate code with `data` |   | ◉ View recommended plots |

Importing Machine Learning Libraries

```
1 from sklearn.preprocessing import StandardScaler  #Standard Scaler
2 from sklearn.linear_model  import LogisticRegression #Regression Model
3 from sklearn.model_selection import train_test_split #Train-Test-Split
4 from sklearn.metrics import accuracy_score, confusion_matrix #Accuracy parameters
```

```
1 data.describe()
```

|   | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI |
|---|---|---|---|---|---|---|
| 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 |

Any null values in the set?

```
1 data.isnull().sum()
```

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

```
Outcome                        0
dtype: int64
```
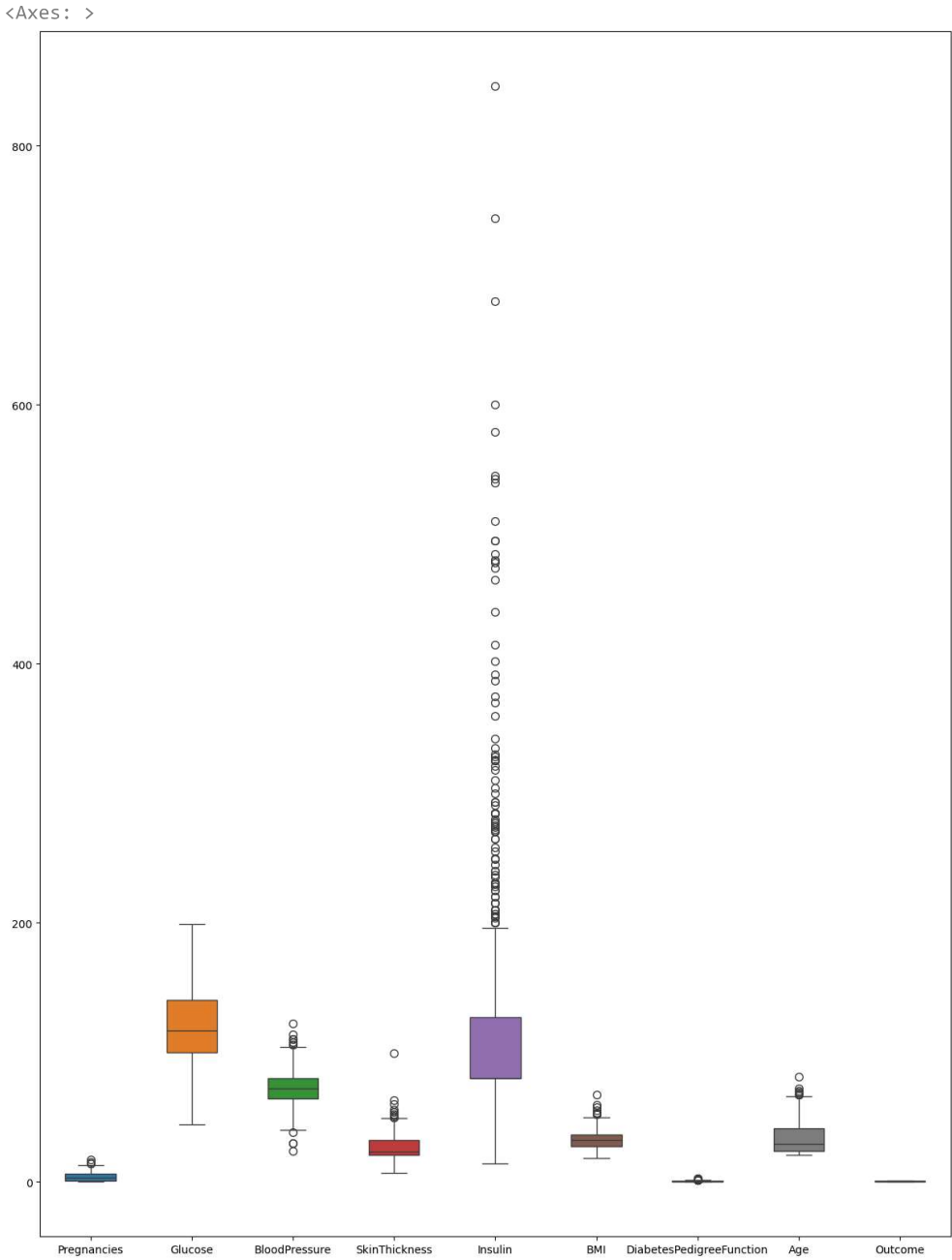
No null values

There are some data in the columns line `Insulin` which are 0, that is not practically possible. So, we will replace the 0 with the mean value of the columns.

The columns name with 0 values are, `BMI, BloodPressure, Glucose, Insulin, SkinThikness`.

```
1 data['BMI'] = data['BMI'].replace(0,data['BMI'].mean())
2 data['BloodPressure'] = data['BloodPressure'].replace(0,data['BloodPressure'].mean())
3 data['Glucose'] = data['Glucose'].replace(0,data['Glucose'].mean())
4 data['Insulin'] = data['Insulin'].replace(0,data['Insulin'].mean())
5 data['SkinThickness'] = data['SkinThickness'].replace(0,data['SkinThickness'].mean())
```

Are there any outliers?

```
1 fig, ax=plt.subplots(figsize=(15,20))
2 sns.boxplot(data,width=0.5, ax=ax, fliersize=7)
```

`<Axes: >`



Separating Dependent and Independent(Prediction) data from the dataset.

```
1 x=data.drop(columns=['Outcome'])
2 y=data['Outcome'] #to be predicted
```

Separating the train, test data -

```
1 xtrain, xtest, ytrain, ytest = train_test_split(x,y, test_size=0.25, random_state=0)
```

```
1 xtrain.shape, xtest.shape
```

```
((576, 8), (192, 8))
```

Scaling the data and pikiling the scaler -

```
1 import pickle
```

```
1 def scaler_standard(xtrain,xtest):
2   scaler=StandardScaler()
3   xtrainscaled=scaler.fit_transform(xtrain)
4   xtestscaled=scaler.transform(xtest)
5   #save
6   file=open("/content/Model/standarscaler.pkl", 'wb')
7   pickle.dump(scaler,file)
8   file.close()
9   return xtrainscaled, xtestscaled
```

```
1 xtrain_scaled, xtest_scaled=scaler_standard(xtrain,xtest)
```

```
1 xtrain_scaled
```

```
array([[ 1.50755225, -1.09947934, -0.89942504, ..., -1.45561965,
        -0.98325882, -0.04863985],
       [-0.82986389, -0.1331471 , -1.23618124, ...,  0.09272955,
        -0.62493647, -0.88246592],
       [-1.12204091, -1.03283573,  0.61597784, ..., -0.03629955,
         0.39884168, -0.5489355 ],
       ...,
       [ 0.04666716, -0.93287033, -0.64685789, ..., -1.14021518,
        -0.96519215, -1.04923114],
       [ 2.09190629, -1.23276654,  0.11084355, ..., -0.36604058,
        -0.5075031 ,  0.11812536],
       [ 0.33884418,  0.46664532,  0.78435594, ..., -0.09470985,
         0.51627505,  2.953134  ]])
```

Logistic Regeression Model

```
1 logreg=LogisticRegression()
2 logreg.fit(xtrain_scaled, ytrain)
```
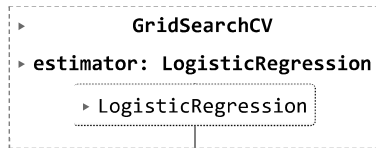
```
▾ LogisticRegression
LogisticRegression()
```

Any Hyperparameter Training?

```
1 # Grid Search CV
2 from sklearn.model_selection import GridSearchCV
3 import warnings
4 warnings.filterwarnings('ignore')
```

```
1 parameters={ 'penalty':['l1','l2'], 'C': np.logspace(-3,3,7), 'solver':['newton-cg', 'lbfgs', 'liblinear'],}
```

```
1 logreg=LogisticRegression()
```

```
1 c.fit(xtrain_scaled, ytrain)
```

```
▸              GridSearchCV
▸ estimator: LogisticRegression
        ▸ LogisticRegression
```

```
1 c.best_params_
```

```
{'C': 1.0, 'penalty': 'l2', 'solver': 'liblinear'}
```

```
1 c.best_score_
```

```
0.763793103448276
```

Let's do the prediction

```
1 ypred=c.predict(xtest_scaled)
```

```
1 accuracy=accuracy_score(ytest,ypred)
2 print("The accuracy of the model is : ", accuracy)
```

```
The accuracy of the model is :  0.796875
```

```
1 conmat=confusion_matrix(ytest,ypred)
2 conmat
```

```
array([[117,  13],
       [ 26,  36]])
```

```
1 tp = conmat[0][0]
2 fp = conmat[0][1]
3 fn = conmat[1][0]
4 tn = conmat[1][1]
```

```
1 accu=(tp+tn)/(tp+tn+fp+fn)
2 print("The accuracy of the model is : ", accu)
```