Title:

# PIMA-INDIANS DIABETES PREDICTION MODEL

Name-Ganesh N Anil
Email-ganesh.n.anil007@gmail.com
Ph-7034007173

## Predicting the onset of diabetes.

*Abstract*-The diabetes dataset is a binary classification problem where it needs to be analysed whether a patient is suffering from the disease or not on the basis of many available features in the dataset. Different methods and procedures of cleaning the data, feature extraction, feature engineering and algorithms to predict the onset of diabetes are used based for diagnostic measure on Pima Indians Diabetes Dataset.

**Dataset details:**

The data sets comprise several variables of the medical predictor, and one objective variable. The forecasting variables include the patient's number of pregnancies, BMI levels, insulin levels, age, etc.

1: Pregnancies: Number of times pregnant

2: Glucose: Plasma glucose concentration 2 hours in an oral glucose tolerance test.

3: Blood Pressure: Diastolic blood pressure (mm Hg)

4: Skin Thickness: Triceps skinfold thickness (mm)

5: Insulin: 2-Hour serum insulin (mu U/ml)

6: BMI: Body mass index (weight in kg/ (height in m²))

7: Diabetes Pedigree Function: Diabetes pedigree function

8: Age: Age (years)

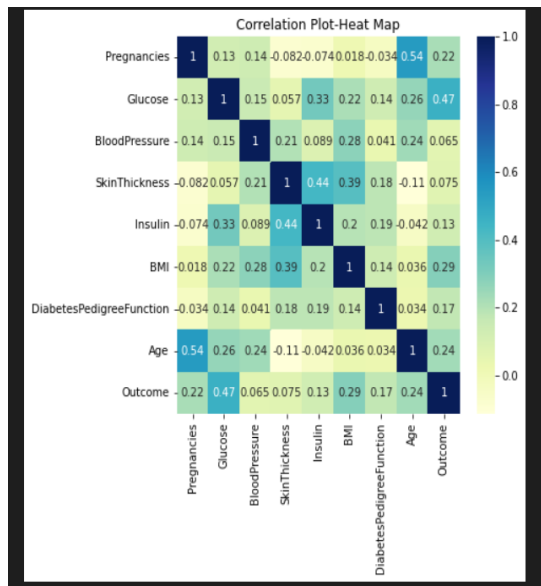9: Outcome: Class variable (0 or 1) 268 of 768 are 1, the others are 0

**Importing the required libraries:**

```
1   #Import the required libraries
2   import pandas as pd
3   import numpy as np
4   import seaborn as sns
5   import matplotlib.pyplot as plt
6   import warnings
7   warnings.filterwarnings('ignore', category=FutureWarning)
8   warnings.filterwarnings('ignore')
9
10  from sklearn.preprocessing import StandardScaler
11  from sklearn.model_selection import train_test_split
12  from sklearn.linear_model import LogisticRegression
13  from sklearn.tree import DecisionTreeClassifier
14  from sklearn.ensemble import RandomForestClassifier
15  from sklearn.ensemble import AdaBoostClassifier
16  from sklearn.ensemble import GradientBoostingClassifier
17  from sklearn.svm import SVC
18  from sklearn.ensemble import VotingClassifier
19  from sklearn.metrics import accuracy_score
```

**EDA (Exploratory Data Analysis):**

- Checking the unique values in target variable and plotting the target variable.
- Value count of output variable was found out and outcome is dataset is imbalance.
- Checking the basic info of the dataset such as shape, size, info etc.
- Finding out the missing values in each of the features.
- Plotting the heatmap and pairplot.

Correlation Plot-Heat Map

**Feature Engineering:**

- Handling the missing values using the function.
- Handling the outlier treatment using standard deviation method.

```python
1  #Replacing 0 with NaN values
2  df[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']] = df[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']].replace(0, np.NaN)
3
4  df.isnull().sum()
```

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

```python
1   #Filling the missing values using function
2   def fill_missing_values(data,fill_value,fill_type,columns,dataframe):
3       print("Missing values before removal in",dataframe,"data")
4       display(data.isnull().sum())
5       for column in columns:
6           if "Random Sample Imputation" in fill_type:
7               data[column+"_random"]=data[column]
8               random_sample = data[column].dropna().sample(data[column].isnull().sum(),random_state=0)
9               random_sample.index=data[data[column].isnull()].index
10              data.loc[data[column].isnull(),column+'_random']=random_sample
11              data[column]=data[column+"_random"]
12              data.drop([column+"_random"],axis=1,inplace=True)
13          if "New_Feature_Importance" in fill_type:
14              data[column] = np.where(data[column].isnull(),1,0)
15              data[column].fillna(data.column.median(),inplace=True)
16          if "Median Fill" in fill_type:
17              data[column].fillna(data[column].median(),inplace=True)
18          if "Mode Fill" in fill_type:
19              data[column].fillna(data[column].mode()[0],inplace=True)
20          if "Value_Fill" in fill_type:
21              data[column].fillna(fill_value,inplace=True)
22          if "Forward_Fill" in fill_type:
23              data[column].ffill(axis = 0, inplace=True)
24          if "Backward_Fill" in fill_type:
25              data[column] = data[column].bfill(axis=0)
26      print("Missing Values AFTER REMOVAL in ",dataframe," data")
27      display(data.isnull().sum())
28      return(data)
```

```
1  fill_missing_values(df,fill_value=0,fill_type=["Random Sample Imputation"],columns=["Glucose","BloodPressure","SkinThickness","Insulin","BMI"],dataframe='Pima Indians Diabetes')
```
Python

```
Missing values before removal in Pima Indians Diabetes data

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

Missing Values AFTER REMOVAL in  Pima Indians Diabetes  data

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

**Train Test Split:**

- Splitting into X_train, X_test, y_train, y_test.

```
1  X=df.drop("Outcome",axis=1)
2  y=df["Outcome"]
3  X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.3,random_state=42)
```
[24]

**Scaling:**

- Scale down using Standardscaler function.
- Train using fit_transform and test using transform.
- Created pickle file using joblib.

# Data Modelling.

### Applying All Algorithms:

- Applying all classification algorithms such as logistic regression, decision tree classifier, random forest classifier, adaboost classifier, gradient boosting classifier and svc classifier.

### Hyperparameter Tuning:

- Train accuracy and test accuracy for each classifier are found out using hyperparameter tuning.
- Best parameter for each classifier is found out.
- The best model was selected by comparing the train and test accuracy.

- In my case, I selected Random Forest Classifier as it is the best generalised model.

```
The Classifier is LogisticRegression() and its hyper params are [{'penalty': ['l1', 'l2'], 'solver': ['saga']}]
The Train accuracy for the LogisticRegression() is 0.7728119180633147
The Test accuracy for the LogisticRegression() is 0.7402597402597403
The Best param for the LogisticRegression() is {'penalty': 'l1', 'solver': 'saga'}
==================

The Classifier is DecisionTreeClassifier() and its hyper params are [{'criterion': ['gini', 'entropy'], 'splitter': ['best', 'random'], 'max_depth': [3, 4, 5], 'min_samples_split':
[2, 3, 4], 'max_features': ['auto', 'sqrt', 'log2']}]
The Train accuracy for the DecisionTreeClassifier() is 0.8156424581005587
The Test accuracy for the DecisionTreeClassifier() is 0.7012987012987013
The Best param for the DecisionTreeClassifier() is {'criterion': 'entropy', 'max_depth': 5, 'max_features': 'log2', 'min_samples_split': 2, 'splitter': 'best'}
==================

The Classifier is RandomForestClassifier() and its hyper params are [{'n_estimators': [4, 6, 9], 'max_features': ['log2', 'sqrt', 'auto'], 'criterion': ['entropy', 'gini'],
'max_depth': [2, 3, 5, 10]}]
The Train accuracy for the RandomForestClassifier() is 0.8677839851024208
The Test accuracy for the RandomForestClassifier() is 0.7359307359307359
The Best param for the RandomForestClassifier() is {'criterion': 'entropy', 'max_depth': 5, 'max_features': 'log2', 'n_estimators': 9}
==================

The Classifier is AdaBoostClassifier() and its hyper params are [{'n_estimators': [10, 50, 250, 1000], 'learning_rate': [0.01, 0.1]}]
The Train accuracy for the AdaBoostClassifier() is 0.7914338919925512
The Test accuracy for the AdaBoostClassifier() is 0.7575757575757576
The Best param for the AdaBoostClassifier() is {'learning_rate': 0.01, 'n_estimators': 250}
==================
```

```
The Classifier is GradientBoostingClassifier() and its hyper params are [{'loss': ['log_loss', 'deviance', 'exponential'], 'learning_rate': [1, 7, 9], 'criterion': ['friedman_mse',
'mse']}]
The Train accuracy for the GradientBoostingClassifier() is 1.0
The Test accuracy for the GradientBoostingClassifier() is 0.7229437229437229
The Best param for the GradientBoostingClassifier() is {'criterion': 'friedman_mse', 'learning_rate': 1, 'loss': 'deviance'}
==================

The Classifier is SVC() and its hyper params are [{'kernel': ['linear', 'poly', 'rbf'], 'degree': [3, 4, 5]}]
The Train accuracy for the SVC() is 0.7783985102420856
The Test accuracy for the SVC() is 0.7272727272727273
The Best param for the SVC() is {'degree': 3, 'kernel': 'linear'}
==================
```

- Created pickle file using joblib.

## Web Framework.

- Using the flask backend api was set
- UI (User Interface) was set using html code.

```python
#import relevant libraries for flask, html rendering and loading the ML model
from distutils.log import debug
from re import A
from flask import Flask,request, url_for, redirect, render_template
import pickle
import pandas as pd
import joblib

app = Flask(__name__)


# model = pickle.load(open("model.pkl","rb"))
model=joblib.load("model.pkl")
# scale = pickle.load(open("scale.pkl","rb"))
scale=joblib.load("scale.pkl")


@app.route("/")
def landingPage():
    return render_template('index.html')


@app.route("/predict",methods=['POST','GET'])
def predict():
    pregnancies = request.form['1']
    glucose = request.form['2']
    bloodPressure = request.form['3']
    skinThickness = request.form['4']
    insulin = request.form['5']
    bmi = request.form['6']
    dpf = request.form['7']
    age = request.form['8']
```

```python
    rowDF= pd.DataFrame([pd.Series([pregnancies,glucose,bloodPressure,skinThickness,insulin,bmi,dpf,age])])
    rowDF_new = pd.DataFrame(scale.transform(rowDF))

    print(rowDF_new)

#  model prediction
    prediction= model.predict_proba(rowDF_new)
    print(f"The  Predicted values is :{prediction[0][1]}")

    if prediction[0][1] >= 0.5:
        valPred = round(prediction[0][1],3)
        print(f"The Round val {valPred*100}%")
        return render_template('result.html',pred=f'You have a chance of having diabetes.\n\nProbability of you being a diabetic is {valPred*100}%.\n\nAdvice : Exercise R
    else:
        valPred = round(prediction[0][0],3)
        return render_template('result.html',pred=f'Congratulations!!!, You are in a Safe Zone.\n\n Probability of you being a non-diabetic is {valPred*100}%.\n\n Advice


if __name__ == "__main__":
    app.run(debug=True)
```

## User Interface.



Are you worried that you might be Diabetic?    Return to home

# Predict diabetes (using AI in the background)

Predict the probability that you may be diabetic

Pregnancies (0-15)       Glucose (40-250)        BloodPressure (20-140)
No. of Pregnancies       Glucose level in sugar   BloodPressure

SkinThickness (5-80)     Insulin (0-1000)        BMI(10-100)
SkinThickness            Insulin level           Body Mass Index

DiabetesPedigreeFunction (0-2.5)   Age (10-120)
DiabetesPedigreeFunction            Age

SUBMIT AND PREDICT PROBABILITY

Diabetes Prediction     Home

Congratulations!!!, You are in a Safe Zone. Probability of you being a non-diabetic is 97.8%. Advice : Exercise Regularly and maintain like this..!

© Ganesh N Anil

ganesh.n.anil007@gmail.com

## Model Deployment.

- I Have done the project in Google Collab and later pushed the file into my GitHub repository ( https://github.com/GaneshNAnil/Pima_Indians_Diabetes.git).
- I Have deployed the project through Heroku (https://diabetispred.herokuapp.com/)