

# AI-BASED DIABETES PREDICTION SYSTEM

**TEAM MEMBER**

**510521205043 : SEENURAO**

**Phase 5 submission document**

**Project Title: Diabetes Prediction System**

**Phase 5: Project Documentation & Submission**

**Topic:** *In this section we will document the complete project and prepare it for submission.*



## INTRODUCTION:

- ❖ Diabetes is a chronic medical condition characterized by high levels of blood glucose (sugar) due to the body's inability to properly produce or use insulin. It is a significant global health concern, with the number of individuals diagnosed with diabetes steadily increasing.
- ❖ Early detection and management of diabetes are essential in preventing its complications and improving the quality of life for affected individuals. Machine learning, a subset of artificial intelligence, has emerged as a powerful tool for developing predictive models to identify individuals at risk of developing diabetes.

This technology has the potential to revolutionize healthcare by enabling proactive and personalized interventions.

- ❖ A Diabetes Prediction System using machine learning is a sophisticated application of data analysis and artificial intelligence techniques to predict the likelihood of an individual developing diabetes. It leverages historical patient data, including demographics, medical records, lifestyle factors, and genetic information, to build predictive models.
- ❖ These models can assist healthcare professionals in identifying high-risk individuals, enabling them to implement preventive measures and early interventions.

## TOOLS (In Predicted System)

Predicting diabetes using machine learning involves analyzing data to identify patterns and create models that can predict the likelihood of someone developing diabetes. *Several tools and libraries are commonly used for this purpose. Here are some of them:*

### 1) **Python:**

Python is a popular programming language for machine learning, and many libraries are available for data analysis, modeling, and visualization.

### 2) **Jupyter Notebooks:**

Jupyter is an interactive environment that allows you to write and execute code in a web-based notebook. It's commonly used for data exploration and model development.

### 3) **Pandas:**

Pandas is a Python library for data manipulation and analysis. It's often used to load, preprocess, and clean the dataset.

### 4) **NumPy:**

NumPy is a fundamental package for scientific computing in Python. It provides support for arrays and matrices, which are essential for numerical operations.

### 5) **Scikit-Learn:**

Scikit-Learn is a versatile machine learning library in Python. It provides a wide range of algorithms for classification, regression, clustering, and more. You can use it to build and evaluate diabetes prediction models.

### 6) **TensorFlow or PyTorch:**

These deep learning frameworks are useful if you want to build neural networks for diabetes prediction, especially if your dataset is large and complex.

## 7) **Keras:**

Keras is a high-level neural networks API that runs on top of TensorFlow or other backends. It's known for its user-friendliness and is suitable for building deep learning models quickly.

## 8) **XGBoost and LightGBM:**

These are gradient boosting libraries often used for tabular data like diabetes datasets. They are efficient and can provide high predictive accuracy.

## 9) **Matplotlib and Seaborn:**

These libraries are used for data visualization, which is important for understanding the data and model results.

## 10) **Data Preprocessing Tools:**

Standardization, normalization, handling missing values, and feature selection are crucial in data preprocessing. Scikit-Learn provides tools for these tasks.

## 11) **Cross-Validation Tools:**

Cross-validation techniques are essential for assessing the model's generalization performance. Scikit-Learn has built-in support for this.

## 12) **Hyperparameter Tuning Tools:**

Tools like GridSearchCV and RandomizedSearchCV in Scikit-Learn help find the best hyperparameters for your model.

## 13) **MLflow or TensorBoard:**

These tools can be used for experiment tracking, model versioning, and performance monitoring.

#### 14) **SQL or NoSQL Databases:**

Depending on the data source, you may need to interact with data bases to extract, store, or manage the data.

#### 15) **Web Frameworks:**

If you plan to deploy your diabetes prediction model as a web application, you might use Flask or Django for building the web interface.

#### 16) **Cloud Services:**

Cloud platforms like AWS, Google Cloud, or Microsoft Azure can be used for scalable and reliable hosting of your machine learning models and datasets.

*The choice of tools depends on the specifics of your project, including the dataset, the complexity of the model, and the deployment environment.*

### **Key Components of a Diabetes Prediction System:**

**Data Collection:** *The foundation of any machine l*

earning-based diabetes prediction system is the collection of diverse and comprehensive data. This data can include patient records, such as age, gender, family history, BMI, blood pressure, glucose levels, and lifestyle choices like diet and exercise habits. Genetic information and environmental factors may also be considered to enhance the accuracy of predictions.

#### **Data Preprocessing:**

Raw data needs to be cleaned, transformed, and standardized for effective analysis. Data preprocessing techniques include handling missing values, normalization, feature selection, and dimensionality reduction.

## **Feature Engineering:**

Feature engineering is a crucial step in selecting the most relevant variables or features for the prediction model. It involves choosing the factors that have the most significant impact on diabetes risk, thereby improving model accuracy.

## **Model Selection:**

Machine learning algorithms such as logistic regression, decision trees, random forests, support vector machines, and neural networks are commonly used for diabetes prediction. The choice of the algorithm depends on the complexity of the data and the desired level of prediction accuracy.

## **Training and Validation:**

The selected model is trained on a portion of the dataset and validated to assess its performance. Cross-validation techniques help ensure the model generalizes well to new, unseen data.

## **Prediction and Interpretation:**

Once the model is trained and validated, it can be used to predict the risk of diabetes for new patients. Interpretability of the model's predictions is crucial for healthcare professionals and patients to understand the factors contributing to the risk.

## **Deployment:**

The predictive model can be integrated into a healthcare system or application, making it accessible to physicians, clinicians, and patients. Real-time predictions can be generated, and appropriate interventions can be recommended.

## Continuous Improvement:

Machine learning models require ongoing monitoring and updates to maintain their accuracy. As new data becomes available and the model's performance is evaluated, adjustments may be necessary to ensure its effectiveness.

Data Set Link: ( <https://www.kaggle.com/datasets/mathchi/diabetes-data-set> )

Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
6	148	72	35	0	33.6	0.627	50	1
1	85	66	29	0	26.6	0.351	31	0
8	183	64	0	0	23.3	0.672	32	1
1	89	66	23	94	28.1	0.167	21	0
0	137	40	35	168	43.1	2.288	33	1
5	116	74	0	0	25.6	0.201	30	0
3	78	50	32	88	31	0.248	26	1
10	115	0	0	0	35.3	0.134	29	0
2	197	70	45	543	30.5	0.158	53	1
8	125	96	0	0	0	0.232	54	1
4	110	92	0	0	37.6	0.191	30	0
10	168	74	0	0	38	0.537	34	1
10	139	80	0	0	27.1	1.441	57	0
1	189	60	23	846	30.1	0.398	59	1
5	166	72	19	175	25.8	0.587	51	1
7	100	0	0	0	30	0.484	32	1
0	118	84	47	230	45.8	0.551	31	1
7	107	74	0	0	29.6	0.254	31	1
1	103	30	38	83	43.3	0.183	33	0
1	115	70	30	96	34.6	0.529	32	1
3	126	88	41	235	39.3	0.704	27	0

## ALL TECHNIQUES WITH CODE FOR DIABETES

### PREDICT ON SYSTEM

```
import numpy as np
```

```
# linear algebra
```

```
import pandas as pd
```

```
# data processing, CSV file I/O (e.g. pd.read_csv)
```

*# Input data files are available in the read-only "../input/" directory# For example, running this (by clicking run or pressing Shift+Enter) will list all files under the input directory*

```
import osfor dirname, _, filenames in os.walk('/kaggle/input'):
```

```
    for filename in filenames:
```

```
        print(os.path.join(dirname, filename))
```

*# You can write up to 20GB to the current directory (/kaggle/working/) that gets preserved as output when you create a version using "Save & Run All" # You can also write temporary files to /kaggle/temp/, but they won't be saved outside of the current session*

D:/diabetes-data-set/diabetes.csv

In [2]:

```
#### Data Visualisers
```

```
import numpy as npimport pandas as pd
```

```
#### Neural Network Implementors
```

```
import tensorflow as tf
```

```
from tensorflow
```

```
import keras
```

```
from keras.models
```

```
import Sequential
```

```
from keras.layers
```

```
import Dense
```

```
#### Random Forest Implementors
```



```
from sklearn.ensemble
```

```
import RandomForestClassifier
```

```
from sklearn.metrics
```

```
import accuracy_score, confusion_matrix, precision_score, recall_score, Confusion  
MatrixDisplay
```

```
from sklearn.model_selection
```

```
import RandomizedSearchCV, train_test_split
```

In [3]:

```
df = pd.read_csv("D:/diabetes-data-set/diabetes.csv")
```

```
### reading the csv file
```

In [4]:

```
print("Size of Dataframe : ", df.shape) """Size of Dataframe : (768,9)"""
```

Size of Dataframe : (768, 9)

Out[4]:

```
'\nSize of Dataframe : (768,9)\n'
```

In [5]:

```
df.head()
```

Out[5]:

	Pregna ncies	Gluc ose	BloodPre ssure	SkinThic kness	Insu lin	B MI	DiabetesPedigre eFunction	A ge	Outc ome
0	6	148	72	35	0	33. 6	0.627	50	1
1	1	85	66	29	0	26. 6	0.351	31	0

	Pregna ncies	Gluc ose	BloodPre ssure	SkinThic kness	Insu lin	B MI	DiabetesPedigre eFunction	A ge	Outc ome
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 [6]:

```
df.describe()
```

Out[6]:

	Pregn ancie s	Gluc ose	Blood Pressur e	SkinTh ickness	Insuli n	BMI	DiabetesPedi greeFunction	Age	Outc ome
co un t	768.0 0000 0	768.0 0000 0	768.00 0000	768.00 0000	768.0 0000 0	768.0 0000 0	768.000000	768.0 0000 0	768.0 0000 0
m ea n	3.845 052	120.8 9453 1	69.105 469	20.536 458	79.79 9479	31.99 2578	0.471876	33.24 0885	0.348 958
st d	3.369 578	31.97 2618	19.355 807	15.952 218	115.2 4400 2	7.884 160	0.331329	11.76 0232	0.476 951
mi n	0.000 000	0.000 000	0.0000 00	0.0000 00	0.000 000	0.000 000	0.078000	21.00 0000	0.000 000
25 %	1.000 000	99.00 0000	62.000 000	0.0000 00	0.000 000	27.30 0000	0.243750	24.00 0000	0.000 000
50 %	3.000 000	117.0 0000 0	72.000 000	23.000 000	30.50 0000	32.00 0000	0.372500	29.00 0000	0.000 000
75	6.000	140.2	80.000	32.000	127.2	36.60	0.626250	41.00	1.000

	Pregnancies	Glucose	Blood Pressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
%	000	50000	000	000	50000	0000		0000	000
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000	1.000000

In [7]:

```
df.isna().sum() """No Null Value"""
```

Out[7]:

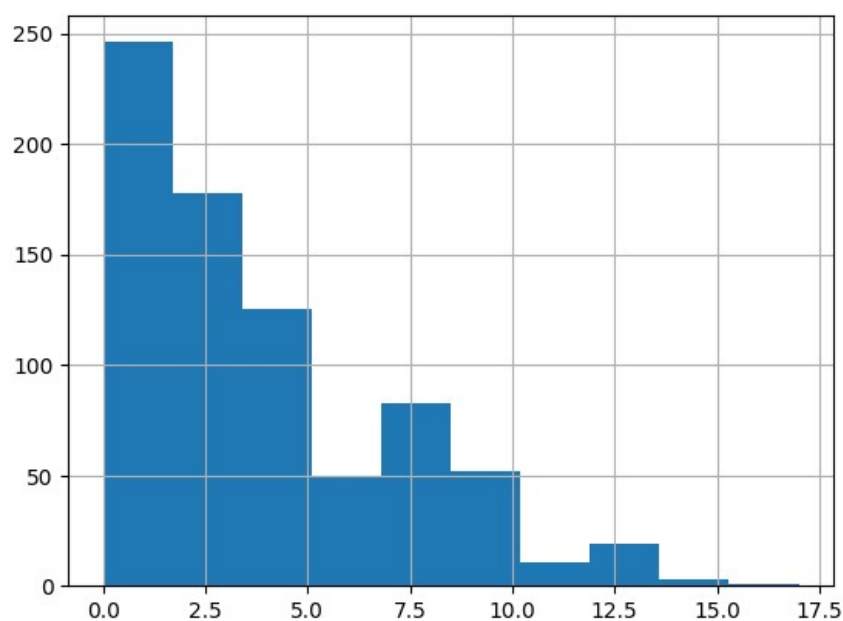
```
'\nNo Null Value\n'
```

In [8]:

```
"""The Values are quite apart from each other in every column. So, we scale it."""df.Pregnancies.hist()
```

Out[8]:

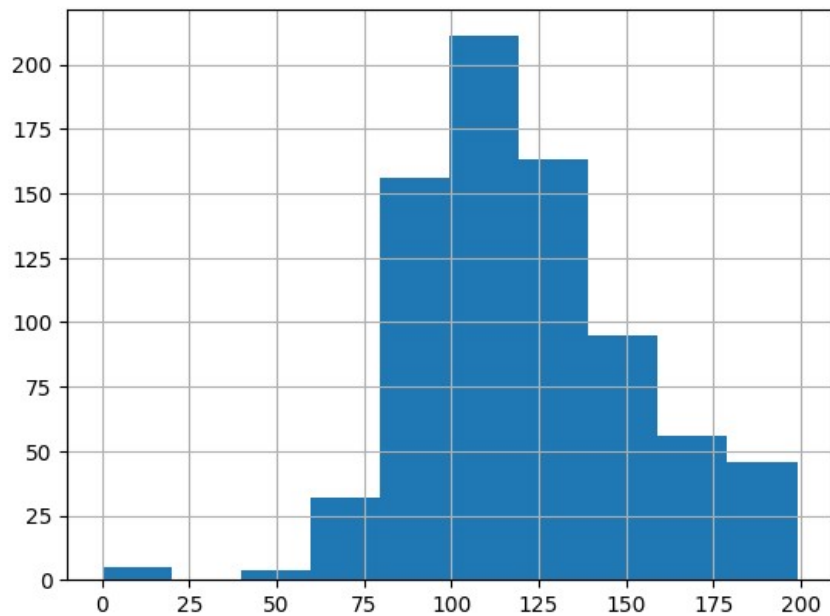
<Axes: >



In [9]:

```
df.Glucose.hist()
```

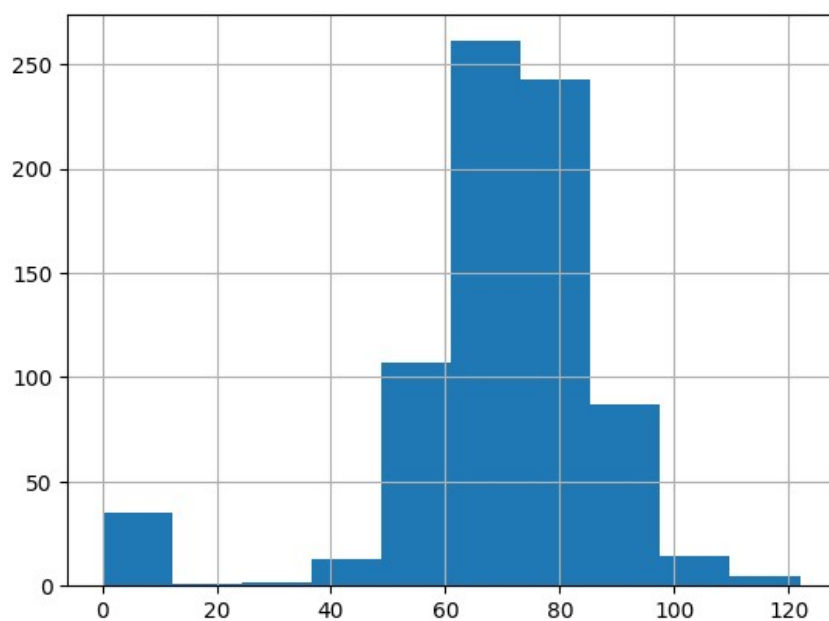
Out[9]: <Axes: >



In [10]:

```
df.BloodPressure.hist()
```

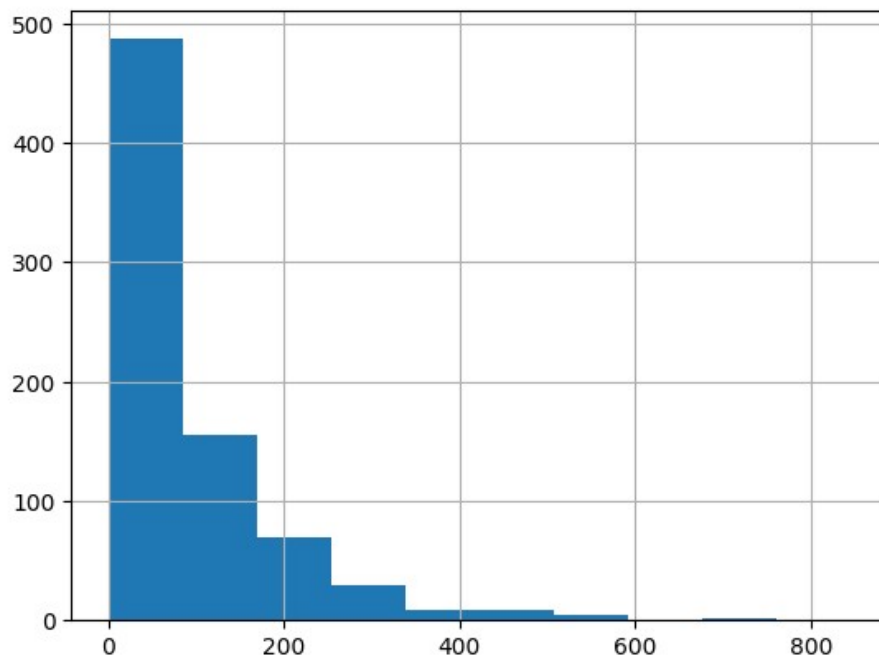
Out[10]: <Axes: >



In [11]:

```
df.Insulin.hist()
```

Out[11]: <Axes: >



In [12]:

```
df.drop_duplicates(inplace=True)
```

In [13]:

```
"""We are going to use MinMaxScaler 'cause the data is not mainly normalised. """
```

```
from sklearn.preprocessing import MinMaxScaler scaler = MinMaxScaler()
```

In [14]:

```
df.head(3)
```

Out[14]:

	Pregna ncies	Gluc ose	BloodPre ssure	SkinThic kness	Insu lin	B MI	DiabetesPedigre eFunction	A ge	Outc ome
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

In [15]:

```
"""X ---> inputy ---> output"""
```

```
X=df.loc[:, 'Pregnancies': 'Age']y = df.Outcome
```

In [16]:

```
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.3) ### training and testing sets
```

In [17]:

```
X_train.head(3)
```

Out[17]:

	Pregnan cies	Gluc ose	BloodPres sure	SkinThick ness	Insul in	B MI	DiabetesPedigreeF unction	A ge
31	3	158	76	36	245	31. 6	0.851	28
56 7	6	92	62	32	126	32. 0	0.085	46
99	1	122	90	51	220	49. 7	0.325	31

In [18]:

```
rf = RandomForestClassifier()rf.fit(X_train, y_train)
```

Out[18]:

**RandomForestClassifier:**

RandomForestClassifier()

In [19]:

```
y_pred = rf.predict(X_test)accuracy = accuracy_score(y_test, y_pred)print("Accuracy: ", accuracy)
```

Accuracy: 0.7056277056277056

In [20]:

```
X_scaled = pd.DataFrame(scaler.fit_transform(X))
```

In [21]:

```
X_scaled.head(3)
```

Out[21]:

	0	1	2	3	4	5	6	7
0	0.352941	0.743719	0.590164	0.353535	0.0	0.500745	0.234415	0.483333
1	0.058824	0.427136	0.540984	0.292929	0.0	0.396423	0.116567	0.166667
2	0.470588	0.919598	0.524590	0.000000	0.0	0.347243	0.253629	0.183333

In [22]:

```
X.head(3)
```

Out[22]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
0	6	148	72	35	0	33.6	0.627	50
1	1	85	66	29	0	26.6	0.351	31
2	8	183	64	0	0	23.3	0.672	32

In [23]:

```
np.array(y)
```

Out[23]:

```
array([1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 0, 0,
       1, 1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 1,
       0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0,
       1, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0,
       1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1,
       1, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 1, 1, 1, 1,
       1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
       1, 1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1,
       0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1,
```



1, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1, 1,  
1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 0,  
1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 0,  
1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 0,  
0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0,  
1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0,  
0, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0,  
0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 0,  
0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 1, 0,  
0, 1, 0, 1, 1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 1,  
0, 1, 0, 0, 0, 0, 1, 1, 0, 1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0,  
1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0,  
0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,  
1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,  
1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0,  
0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 0, 1, 1, 0, 0, 0,  
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0,  
0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 0,  
0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 1, 0,  
0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0,  
1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1,  
0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 0, 0, 1,

```
0, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0,
0, 0, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0,
0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0,
1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0, 1, 0])
```

In [24]:

```
target_count = df['Outcome'].value_counts()

print('Class 0:', target_count[0])

print('Class 1:', target_count[1])

print('Proportion:', round(target_count[0] / target_count[1], 2), ': 1')
```

Class 0: 500

Class 1: 268

Proportion: 1.87 : 1

In [25]:

```
X_scaled_shaped = np.array(X_scaled).reshape(768*8,)
```

In [26]:

```
X_scaled
```

Out[26]:

	0	1	2	3	4	5	6	7
0	0.35294 1	0.74371 9	0.59016 4	0.35353 5	0.00000 0	0.50074 5	0.23441 5	0.48333 3
1	0.05882 4	0.42713 6	0.54098 4	0.29292 9	0.00000 0	0.39642 3	0.11656 7	0.16666 7
2	0.47058	0.91959	0.52459	0.00000	0.00000	0.34724	0.25362	0.18333

	0	1	2	3	4	5	6	7
	8	8	0	0	0	3	9	3
3	0.05882 4	0.44723 6	0.54098 4	0.23232 3	0.11111 1	0.41877 8	0.03800 2	0.00000 0
4	0.00000 0	0.68844 2	0.32786 9	0.35353 5	0.19858 2	0.64232 5	0.94363 8	0.20000 0
...	...	...	...	...	...	...	...	...
76 3	0.58823 5	0.50753 8	0.62295 1	0.48484 8	0.21276 6	0.49031 3	0.03971 0	0.70000 0
76 4	0.11764 7	0.61306 5	0.57377 0	0.27272 7	0.00000 0	0.54843 5	0.11187 0	0.10000 0
76 5	0.29411 8	0.60804 0	0.59016 4	0.23232 3	0.13238 8	0.39046 2	0.07130 7	0.15000 0
76 6	0.05882 4	0.63316 6	0.49180 3	0.00000 0	0.00000 0	0.44858 4	0.11571 3	0.43333 3
76 7	0.05882 4	0.46733 7	0.57377 0	0.31313 1	0.00000 0	0.45305 5	0.10119 6	0.03333 3

*768 rows × 8 columns*

In [27]:

```
X_scaled_shaped
```

Out[27]:

```
array([0.35294118, 0.74371859, 0.59016393, ..., 0.45305514, 0.10119556,
       0.03333333])
```

In [28]:

```
"""Neural Network Implementation"""
```

```
model = Sequential([
```

```
Dense(20,activation='relu'),
```

```
Dense(10,activation='relu'),
```

```
Dense(5,activation='relu'),
```

```
Dense(1,activation='sigmoid'))])
```

In [29]:

```
X_train.head(3)
```

Out[29]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
31	3	158	76	36	245	31.6	0.851	28
567	6	92	62	32	126	32.0	0.085	46
99	1	122	90	51	220	49.7	0.325	31

I

In [31]:

```
X_train
```

Out[31]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age
31	3	158	76	36	245	31.6	0.851	28
567	6	92	62	32	126	32.0	0.085	46
99	1	122	90	51	220	49.7	0.325	31

	Pregnan cies	Gluc ose	BloodPres sure	SkinThick ness	Insul in	B MI	DiabetesPedigreeF unction	A ge
74 9	6	162	62	0	0	24. 3	0.178	50
41 9	3	129	64	29	115	26. 4	0.219	28
...	...	...	...	...	...	...	...	...
61 8	9	112	82	24	0	28. 2	1.282	50
48 0	3	158	70	30	328	35. 5	0.344	35
16 1	7	102	74	40	105	37. 2	0.204	45
29	5	117	92	0	0	34. 1	0.337	38
12	10	139	80	0	0	27. 1	1.441	57

*537 rows × 8 columns*

In [32]:

```
X_train_scaled = pd.DataFrame(scaler.fit_transform(X_train))
```

In [33]:

```
X_train_scaled
```

Out[33]:

	0	1	2	3	4	5	6	7
0	0.17647 1	0.79798 0	0.66666 7	0.36363 6	0.28959 8	0.59398 5	0.34340 3	0.13725 5
1	0.35294 1	0.46464 6	0.54386 0	0.32323 2	0.14893 6	0.60150 4	0.00311 0	0.49019 6

	0	1	2	3	4	5	6	7
2	0.05882 4	0.61616 2	0.78947 4	0.51515 2	0.26004 7	0.93421 1	0.10972 9	0.19607 8
3	0.35294 1	0.81818 2	0.54386 0	0.00000 0	0.00000 0	0.45676 7	0.04442 5	0.56862 7
4	0.17647 1	0.65151 5	0.56140 4	0.29292 9	0.13593 4	0.49624 1	0.06263 9	0.13725 5
...	...	...	...	...	...	...	...	...
53 2	0.52941 2	0.56565 7	0.71929 8	0.24242 4	0.00000 0	0.53007 5	0.53487 3	0.56862 7
53 3	0.17647 1	0.79798 0	0.61403 5	0.30303 0	0.38770 7	0.66729 3	0.11817 0	0.27451 0
53 4	0.41176 5	0.51515 2	0.64912 3	0.40404 0	0.12411 3	0.69924 8	0.05597 5	0.47058 8
53 5	0.29411 8	0.59090 9	0.80701 8	0.00000 0	0.00000 0	0.64097 7	0.11506 0	0.33333 3
53 6	0.58823 5	0.70202 0	0.70175 4	0.00000 0	0.00000 0	0.50939 8	0.60550 9	0.70588 2

*537 rows × 8 columns*

In [34]:

```
df.head(3)
```

Out[35]:

	Pregna ncies	Gluc ose	BloodPre ssure	SkinThic kness	Insu lin	B MI	DiabetesPedigre eFunction	A ge	Outc ome
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

In [36]:

```
df.head(3)
```

Out[36]:

	Pregna ncies	Gluc ose	BloodPre ssure	SkinThic kness	Insu lin	B MI	DiabetesPedigre eFunction	A ge	Outc ome
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

In [37]:

```
from sklearn.preprocessing import StandardScaler
```

In [38]:

```
scaler = StandardScaler()
```

In [39]:

```
X.head(3)
```

Out[39]:

	Pregnan cies	Gluco se	BloodPres sure	SkinThick ness	Insul in	B MI	DiabetesPedigreeF unction	Ag e
0	6	148	72	35	0	33. 6	0.627	50
1	1	85	66	29	0	26. 6	0.351	31

	Pregnan cies	Glucose	BloodPres sure	SkinThick ness	Insul in	B MI	DiabetesPedigreeF unction	Age
2	8	183	64	0	0	23. 3	0.672	32

In [40]:

```
y
```

Out[40]:

```
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*

In [41]:

```
X_scaled = pd.DataFrame(scaler.fit_transform(X))
```

In [42]:



```
X_scaled.head()
```

Out[42]:

	0	1	2	3	4	5	6	7
0	0.639947	0.848324	0.149641	0.907270	-0.692891	0.204013	0.468492	1.425995
1	-0.844885	-1.123396	-0.160546	0.530902	-0.692891	-0.684422	-0.365061	-0.190672
2	1.233880	1.943724	-0.263941	-1.288212	-0.692891	-1.103255	0.604397	-0.105584
3	-0.844885	-0.998208	-0.160546	0.154533	0.123302	-0.494043	-0.920763	-1.041549
4	-1.141852	0.504055	-1.504687	0.907270	0.765836	1.409746	5.484909	-0.020496

In [43]:

```
model
```

Out[43]:

<keras.engine.sequential.Sequential at 0x7896e8fce3b0>

In [44]:

```
X_train_scaled,X_test_scaled,y_train,y_test = train_test_split(X_scaled,y,test_size=0.3)
```

In [45]:

```
X_train_scaled.head(3)
```

Out[45]:

	0	1	2	3	4	5	6	7
140	-0.250952	0.222381	0.459827	-1.288212	-0.692891	-1.382478	-0.615731	1.851434
105	-0.844885	0.159787	-0.677523	0.530902	0.626910	-0.417892	0.993993	-1.041549
601	0.639947	-0.779128	-3.572597	-1.288212	-0.692891	-1.052488	-0.851300	-0.445935

In [48]:

```
df.head(3)
```

Out[48]:

	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

In [49]:

```
# Get the weights from the model's layers
```

```
weights = model.get_weights()
```

```
# Calculate feature importance as the sum of absolute weights for each feature
```

```
feature_importance = np.sum(np.abs(weights[0]), axis=1)
```

In [50]:

```
feature_importance
```

Out[50]:

```
array([5.448598 , 4.0084224, 5.4832244, 3.5369577, 4.5502095, 5.109804 ,  
       3.9921277, 6.2405443], dtype=float32)
```

In [51]:

```
# Compute gradients of the output with respect to the input
```

```
input_tensor = tf.convert_to_tensor(X)
```

```
with tf.GradientTape() as tape:
```

```
    tape.watch(input_tensor)
```

```
    output = model(input_tensor)
```

```
gradients = tape.gradient(output, input_tensor)
```

```
# Calculate feature importance as the absolute mean of the gradients
```

```
feature_importance = np.mean(np.abs(gradients.numpy()), axis=0)
```

In [52]:

feature\_importance

Out[52]:

```
array([2.40560107e-06, 2.09994840e-05, 4.16881339e-06, 3.96728967e-06,  
       9.45686434e-06, 7.34819230e-06, 5.94307293e-06, 7.19741934e-06])
```

In [53]:

```
X_train_scaled.drop(0,axis=1)
```

Out[53]:

	1	2	3	4	5	6	7
140	0.222381	0.459827	-1.288212	-0.692891	-1.382478	-0.615731	1.851434
105	0.159787	-0.677523	0.530902	0.626910	-0.417892	0.993993	-1.041549
601	-0.779128	-3.572597	-1.288212	-0.692891	-1.052488	-0.851300	-0.445935
645	1.129998	0.253036	0.907270	3.127584	0.940144	-1.020427	-0.275760
204	-0.560048	0.149641	0.719086	0.956860	0.724382	-0.446604	1.851434
...	...	...	...	...	...	...	...
310	-1.279882	-0.160546	0.593630	-0.692891	-0.735190	-0.479825	0.660206
558	-0.560048	-0.057150	1.220910	-0.692891	1.803195	-1.044587	0.745293
451	0.410164	0.046245	-1.288212	-0.692891	-0.392508	0.211782	-0.871374
231	0.410164	0.563223	1.032726	2.519781	1.803195	-0.706334	1.085644
48	-0.560048	-0.160546	0.719086	-0.692891	0.902069	-0.386202	-0.190672

537 rows × 7 columns

In [54]:

```
X_train_scaled.head(3)
```

Out[54]:

	0	1	2	3	4	5	6	7
140	-0.250952	0.222381	0.459827	-1.288212	-0.692891	-1.382478	-0.615731	1.851434
105	-0.844885	0.159787	-0.677523	0.530902	0.626910	-0.417892	0.993993	-1.041549
601	0.639947	-0.779128	-3.572597	-1.288212	-0.692891	-1.052488	-0.851300	-0.445935

In [55]:

```
X_train_scaled.drop([0,3],axis=1)
```

Out[55]:

	1	2	4	5	6	7
140	0.222381	0.459827	-0.692891	-1.382478	-0.615731	1.851434
105	0.159787	-0.677523	0.626910	-0.417892	0.993993	-1.041549
601	-0.779128	-3.572597	-0.692891	-1.052488	-0.851300	-0.445935
645	1.129998	0.253036	3.127584	0.940144	-1.020427	-0.275760
204	-0.560048	0.149641	0.956860	0.724382	-0.446604	1.851434
...	...	...	...	...	...	...
310	-1.279882	-0.160546	-0.692891	-0.735190	-0.479825	0.660206
558	-0.560048	-0.057150	-0.692891	1.803195	-1.044587	0.745293
451	0.410164	0.046245	-0.692891	-0.392508	0.211782	-0.871374
231	0.410164	0.563223	2.519781	1.803195	-0.706334	1.085644

	1	2	4	5	6	7
48	-0.560048	-0.160546	-0.692891	0.902069	-0.386202	-0.190672

537 rows  $\times$  6 columns

## Python Programming DIABETES PREDICTION IN PYTHON

@askpython.com



### **DATA VISUALIZATION TECHNIQUES**

*# Step 1: Import necessary libraries*

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings("ignore", category=UserWarning)
```

*# Step 2: Load the dataset*

```
df = pd.read_csv('D:/diabetes.csv')
```

*# Step 3: Data Cleaning*

*# Check for Missing Values*

```
missing_values = df.isnull().sum()
```

```
print("Missing Values:")
```

```
print(missing_values)
```

*# Handle missing values (if any)*

```
mean_fill = df.mean()df.fillna(mean_fill, inplace=True)
```

*# Check for Duplicate Rows*

```
duplicate_rows = df[df.duplicated()]
```

```
print("\nDuplicate Rows:")
```

```
print(duplicate_rows)
```

*# Handle duplicate rows (if any)*

```
df.drop_duplicates(inplace=True)
```

*# Step 4: Data Analysis*

*# Summary Statistics*

```
summary_stats = df.describe()
```

```
print("\nSummary Statistics:")
```

```
print(summary_stats)
```

*# Class Distribution*

```
class_distribution = df['Outcome'].value_counts()
```

```
print("\nClass Distribution:")print(class_distribution)
```

```
# Step 5: Data Visualization
```

```
sns.pairplot(df, hue='Outcome')  
plt.show()
```

Missing Values:

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

**Duplicate Rows:**

Empty DataFrame

Columns: [Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DiabetesPedigreeFunction, Age, Outcome]

Index: []

**Summary Statistics:**

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin \
count	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	20.536458	79.799479
std	3.369578	31.972618	19.355807	15.952218	115.244002
min	0.000000	0.000000	0.000000	0.000000	0.000000



25%	1.000000	99.000000	62.000000	0.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000
75%	6.000000	140.250000	80.000000	32.000000	127.250000
max	17.000000	199.000000	122.000000	99.000000	846.000000

	BMI	Diabetes	PedigreeFunction	Age	Outcome
count	768.000000		768.000000	768.000000	768.000000
mean	31.992578		0.471876	33.240885	0.348958
std	7.884160		0.331329	11.760232	0.476951
min	0.000000		0.078000	21.000000	0.000000
25%	27.300000		0.243750	24.000000	0.000000
50%	32.000000		0.372500	29.000000	0.000000
75%	36.600000		0.626250	41.000000	1.000000
max	67.100000		2.420000	81.000000	1.000000

Class Distribution:

Outcome

0 500

1 268

Name: count, dtype: int64



# DATA VISUALIZATION USING DIFFERENT PLOTTING TECHNIQUES

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
from sklearn.model_selection
import train_test_split
from sklearn.preprocessing
import StandardScalerfrom sklearn
import svmfrom sklearn.metrics
import classification_report
from sklearn.metrics
import confusion_matrix
from sklearn.metrics
import ConfusionMatrixDisplay
warnings.filterwarnings("ignore", category=UserWarning)
RED = "\033[91m"GREEN = "\033[92m"YELLOW = "\033[93m"BLUE = "\033[94
m"RESET = "\033[0m"
df = pd.read_csv("/kaggle/input/diabetes-data-set/diabetes.csv")
```

*# DATA CLEANING*

```
print(BLUE + "\nDATA CLEANING" + RESET)# --- Check for missing valuesmiss
ing_values = df.isnull().sum()
print(GREEN + "Missing Values : " + RESET)
print(missing_values)
```

*# --- Handle missing values*

```
mean_fill = df.fillna(df.mean())
df.fillna(mean_fill, inplace=True)
```

```
# --- Check for duplicate values
```

```
duplicate_values = df.duplicated().sum()  
print(GREEN + "Duplicate Values : " + RESET)  
print(duplicate_values)
```

```
# --- Drop duplicate values
```

```
df.drop_duplicates(inplace=True)
```

```
# DATA ANALYSIS
```

```
print(BLUE + "\nDATA ANALYSIS" + RESET)# --- Summary Statistics  
summary_stats = df.describe()  
print(GREEN + "Summary Statistics : " + RESET)  
print(summary_stats)
```

```
# --- Class Distribution
```

```
class_distribution = df["Outcome"].value_counts()  
print(GREEN + "Class Distribution : " + RESET)  
print(class_distribution)
```

```
# Support Vector Machine Modelling
```

```
# --- Separate features and target variable
```

```
print(BLUE + "\nMODELLING" + RESET)X = df.drop("Outcome", axis=1)y = df["Outcome"]
```

```
# --- Splitting the data into training and testing sets  
X_train, X_test, y_train, y_test = train_test_split(  
    X, y, test_size=0.2, random_state=42)
```

```
# --- Standardize Features
```

```
scaler = StandardScaler()
```

```
X_train = scaler.fit_transform(X_train)X_test = scaler.transform(X_test)
```

```
# --- init and train SVM model
```

```
model = svm.SVC(kernel="linear")model.fit(X_train, y_train)
```

```
# --- Predict on test data
```

```
y_pred = model.predict(X_test)
```

```
# --- Evaluate model performance
```

```
accuracy = model.score(X_test, y_test)
```

```
print(GREEN + "Model Accuracy : " + RESET)
```

```
print(accuracy)
```

```
# --- Classification Report and Confusion Matrix
```

```
print(GREEN + "Classification Report : " + RESET)
```

```
print(classification_report(y_test, y_pred))
```

```
print(GREEN + "Confusion Matrix : " + RESET)
```

```
cm = ConfusionMatrixDisplay.from_predictions(y_test, y_pred)
```

```
sns.heatmap(cm.confusion_matrix, annot=True, cmap="Blues")
```

```
plt.show()
```

```
print("Displayed")
```

```
# DATA VISUALIZATION
```

```
print(BLUE + "\nDATA VISUALIZATION" + RESET)
```

```
# --- Pair Plot
```

```
print(GREEN + "PairPlot : " + RESET)
```

```
sns.pairplot(df, hue='Outcome',diag_kind='kde',palette = "Blues")
```

```
plt.title("Pairwise Relationships")
```

```
plt.show()
```

```
# --- Histogram for age distribution
```

```
print(GREEN + "Histogram : " + RESET)
```

```
sns.histplot(df["Age"], bins=10, kde=True,palette = "Blues")
```

```
plt.xlabel("Age")
```

```
plt.ylabel("Count")
```

```
plt.title("Age Distribution")
```

```
plt.show()
```

```
# --- Box plot to visualize glucose levels by outcome
```

```
print(GREEN + "BoxPlot : " + RESET)
```

```
sns.boxplot(x="Outcome", y="Glucose", data=df,palette = "Blues")
```

```
plt.xlabel("Diabetes Outcome (0: No, 1: Yes)")
```

```
plt.ylabel("Glucose Level")
```

```
plt.title("Glucose Levels by Diabetes Outcome")
```

```
plt.show()
```

```
# --- Correlation heatmap
```

```
print(GREEN + "Correlation Heatmap : " + RESET)correlation_matrix = df.corr()sns.heatmap(correlation_matrix, annot=True, cmap="Blues")
```

```
plt.title("Correlation Heatmap")
plt.show()
```

*# SAVING THE FILE*

```
df.to_csv("/kaggle/working/cleaned_diabetes.csv", index=False)
print(BLUE + "\nDATA SAVING" + RESET)
print(GREEN + "Data Cleaned and Saved !" + RESET)
print("\n")
```

### DATA CLEANING Missing Values :

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

### Duplicate Values :

```
0
```

### DATA ANALYSIS Summary Statistics :

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin \
count	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845052	120.894531	69.105469	20.536458	79.799479
std	3.369578	31.972618	19.355807	15.952218	115.244002
min	0.000000	0.000000	0.000000	0.000000	0.000000
25%	1.000000	99.000000	62.000000	0.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.500000
75%	6.000000	140.250000	80.000000	32.000000	127.250000
max	17.000000	199.000000	122.000000	99.000000	846.000000

	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000
mean	31.992578	0.471876	33.240885	0.348958
std	7.884160	0.331329	11.760232	0.476951
min	0.000000	0.078000	21.000000	0.000000
25%	27.300000	0.243750	24.000000	0.000000
50%	32.000000	0.372500	29.000000	0.000000
75%	36.600000	0.626250	41.000000	1.000000
max	67.100000	2.420000	81.000000	1.000000

### Class Distribution :

Outcome

0 500

1 268

Name: count, dtype: int64

### MODELLING Model Accuracy :

0.7597402597402597Classification Report :

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.81	0.82	0.81	99
---	------	------	------	----

1	0.67	0.65	0.66	55
---	------	------	------	----

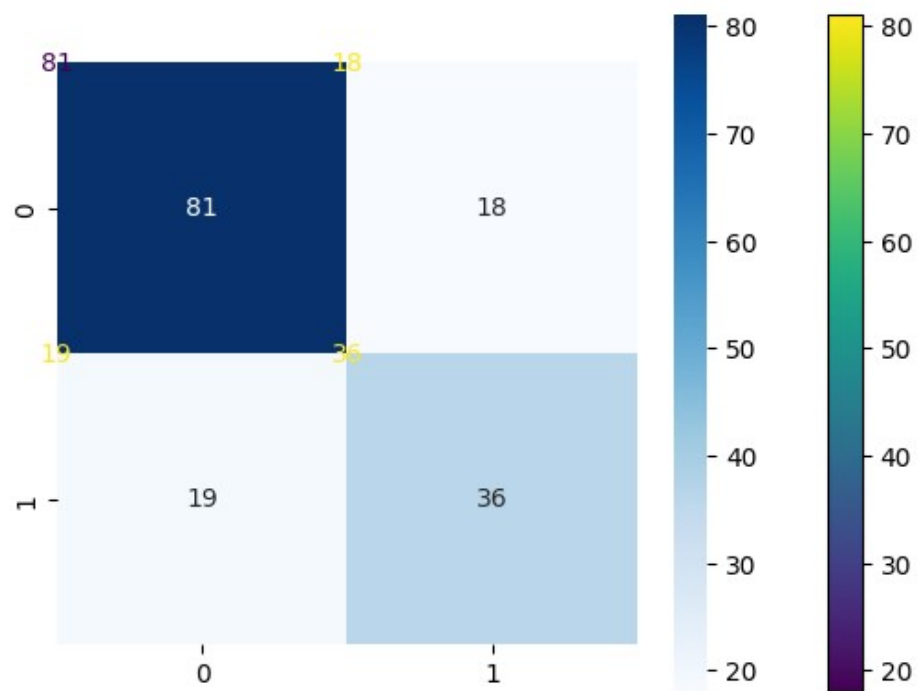
accuracy		0.76		154
----------	--	------	--	-----

macro avg	0.74	0.74	0.74	154
-----------	------	------	------	-----

weighted avg	0.76	0.76	0.76	154
--------------	------	------	------	-----

### Confusion Matrix :

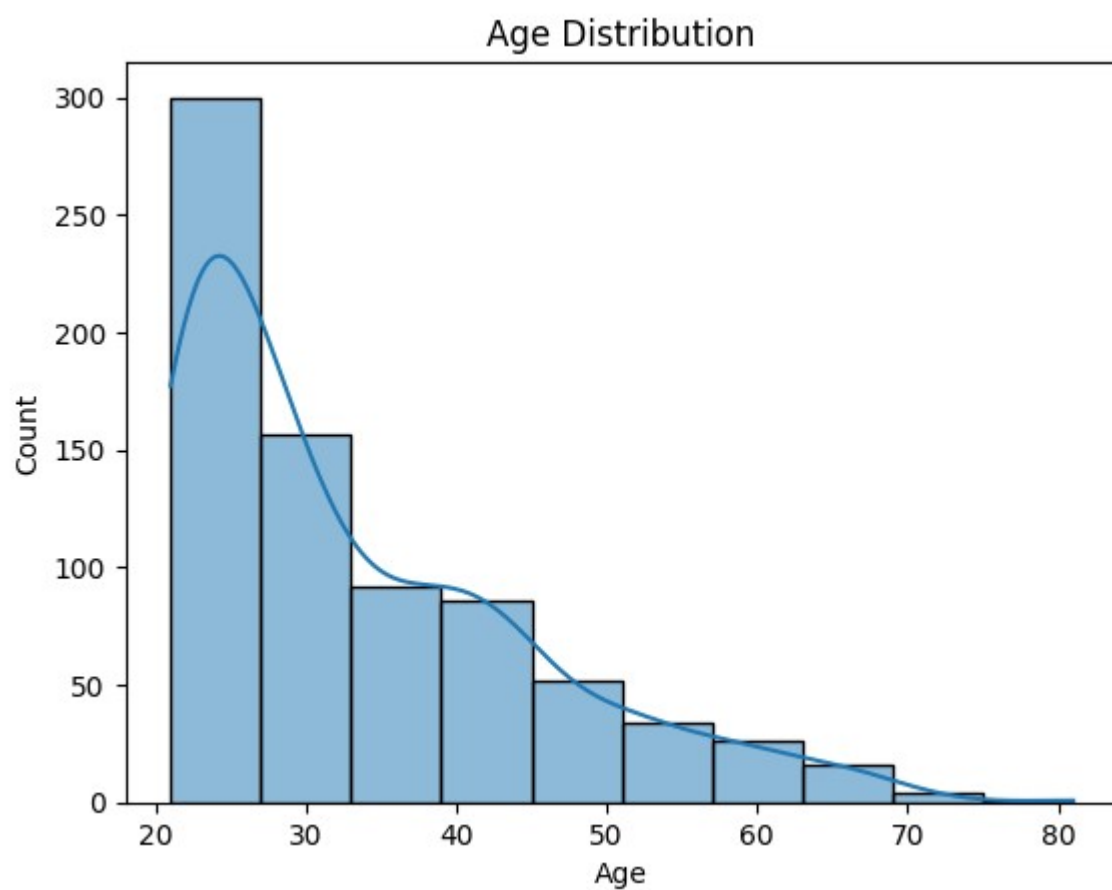




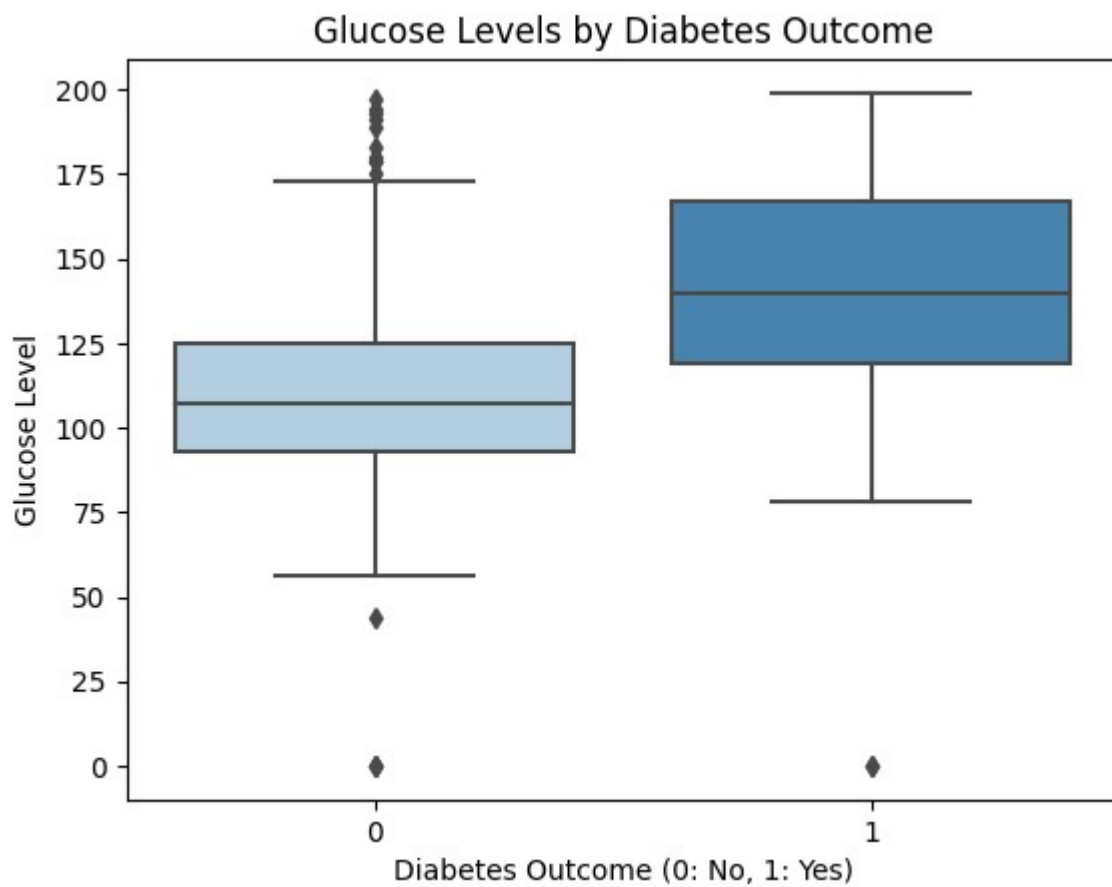
Displayed !

DATA VISUALIZATIONPairPlot :

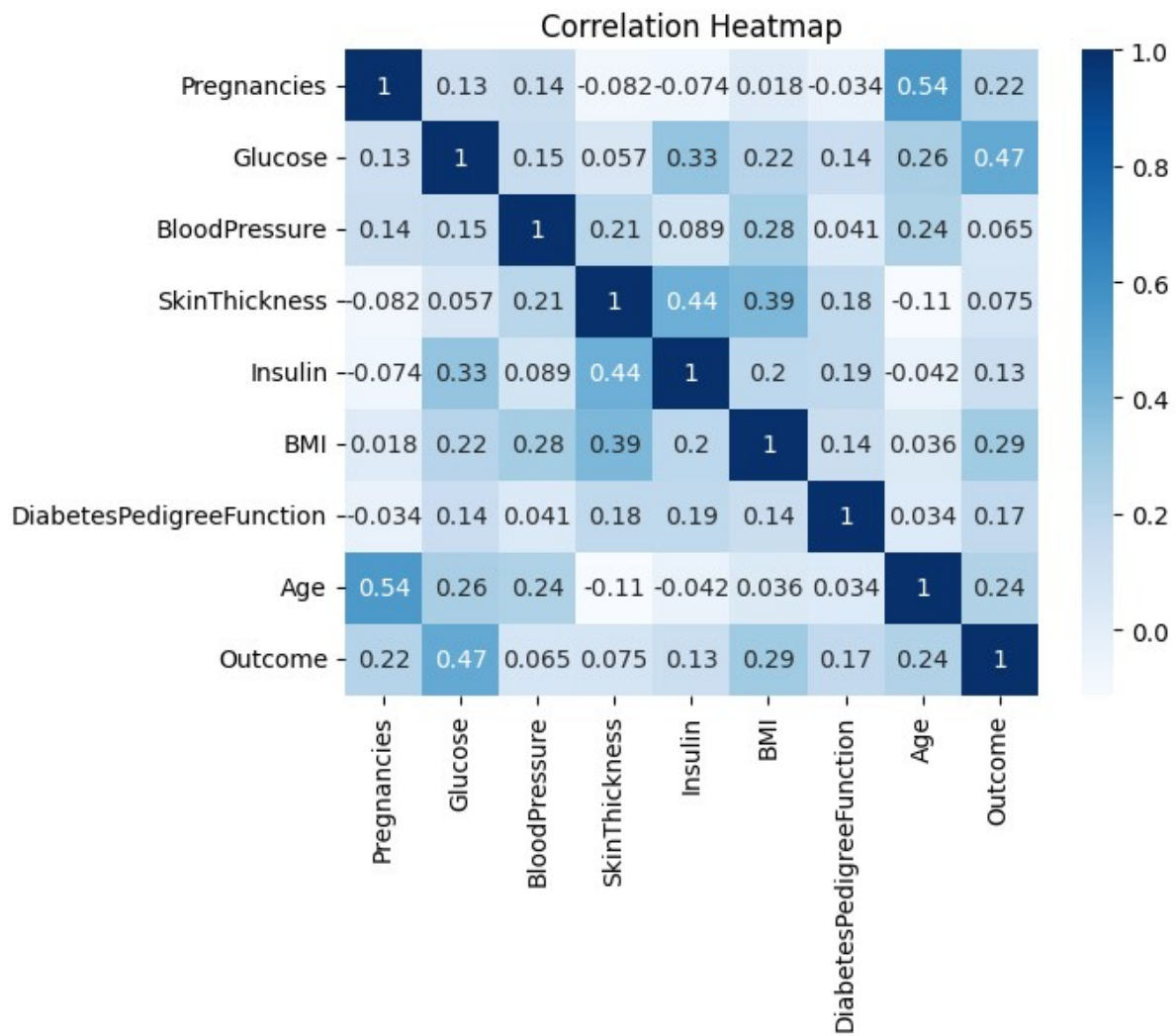
Histogram :



BoxPlot :



Correlation Heatmap :



## VISUALIZING USING PANDAS

DATA SAVING Data Cleaned and Saved !

### Pandas Profiling

In [1]:

```
pp.ProfileReport(df)
```

Summarize dataset: 100%

82/82 [00:18<00:00, 6.19it/s, Completed]

Generate report structure: 100%

1/1 [00:05<00:00, 5.09s/it]

Render HTML: 100%

1/1 [00:02<00:00, 2.82s/it]

Out[1]:

### DecisionTreeClassifier

In [2]:

```
# Input variables X = df.iloc[:, :-1] # Output variable y = df.iloc[:, -1:]
```

In [3]:

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.45, random_state=42)
```

In [4]:

```
X_train.head()
```

Out[4]:

	Pregnan cies	Gluc ose	BloodPres sure	SkinThick ness	Insul in	B MI	DiabetesPedigreeF unction	A ge
36 8	3	81	86	16	66	27. 5	0.306	22
67 2	10	68	106	23	49	35. 5	0.285	47
43 4	1	90	68	8	0	24. 5	1.138	36
27 8	5	114	74	0	0	24. 9	0.744	57
60 2	1	124	74	36	0	27. 8	0.100	30

In [5]:

```
y_train.head()
```

Out[5]:

	Outcome
368	0
672	0
434	0
278	0
602	0

In [6]:

```
clf = DecisionTreeClassifier(max_leaf_nodes=15,random_state=0)clf.fit(X_train,y_t  
rain)
```

Out[6]:

```
DecisionTreeClassifier
```

```
DecisionTreeClassifier(max_leaf_nodes=15, random_state=0)
```

In [7]:

```
y_pred = clf.predict(X_test)
```

```
accuracy_score(y_test,y_pred)
```

Out[7]:

```
0.7485549132947977
```

In [8]:

```
# cross validation np.mean(cross_val_score(DecisionTreeClassifier(),X,y,cv=10,scoring='accuracy'))
```

Out[8]:

```
0.70311004784689
```

**EDA:**

**Input:**

```
#using heatmap to understand correlation better in dataset data
```

```
#Heatmap of correlation
```

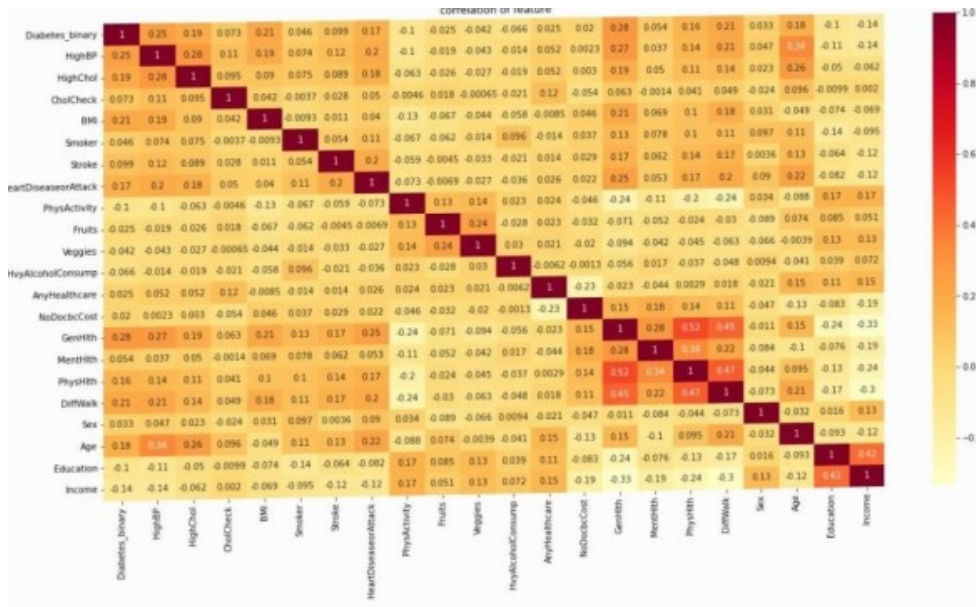
```
plt.figure(figsize = (20,10))
```

```
sns.heatmap(data.corr(),annot=True , cmap = 'YlOrRd' )
```

```
plt.title("correlation of feature")
```

## Output:

text(0.5, 1.0, 'correlation of feature')



## Correlation heatmap show relation between columns:

(GenHlth ,PhysHlth ),(PhysHlth , DiffWalk),(GenHlth ,DiffWalk )are highly correlated with each other => positive relation

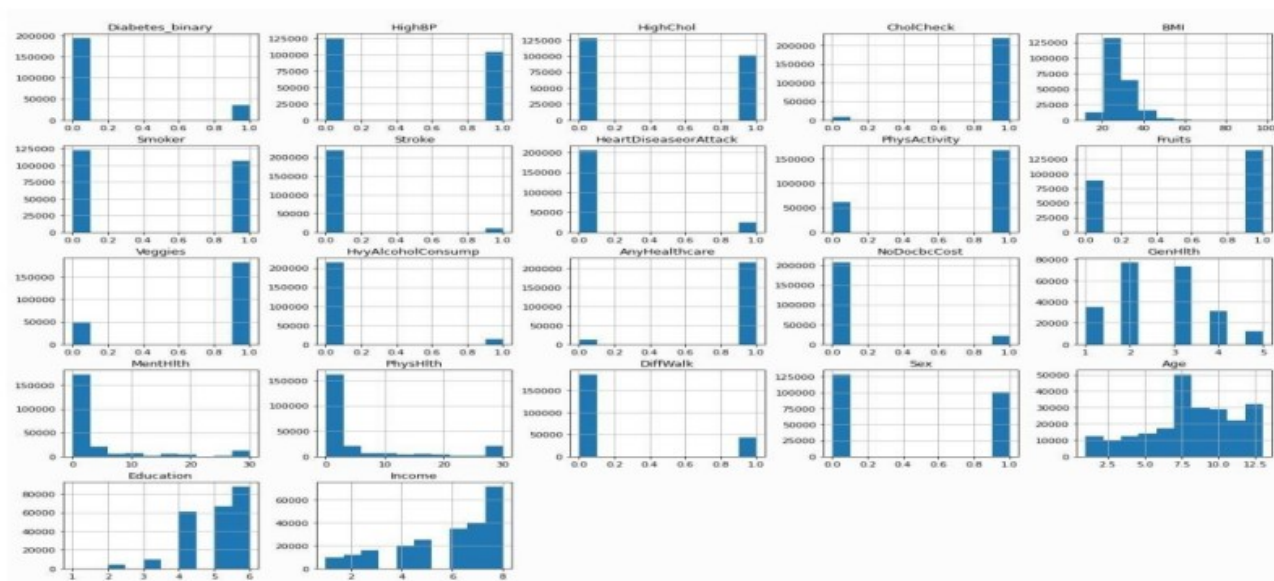
(GenHlth ,Income ) , (DiffWalk , Income) are highly correlated with each other => Negative relation

#using histogram to understand dataset data better

```
data.hist(figsize=(20,15));
```



## Output:



## Visualization Of [Yes – NO] Columns and their relation with the target:

### Input:

```
Cols = ['HighBP', 'HighChol', 'CholCheck', 'Smoker',
```

```
        'Stroke', 'HeartDiseaseorAttack', 'PhysActivity', 'Veggies',
```

```
        'HvyAlcoholConsump', 'AnyHealthcare', 'NoDocbcCost',
```

```
        'DiffWalk']
```

```
def create_plot_pivot(data2, x_column):
```

```
    """ Create a pivot table for satisfaction versus another rating for easy
    plotting. """
```

```
    _df_plot = data2.groupby([x_column, 'Diabetes_binary']).size() \
```

```
.reset_index().pivot(columns='Diabetes_binary', index=x_column,
values=0) return _df_plot
```

```
fig, ax = plt.subplots(3, 4, figsize=(20,20))
```

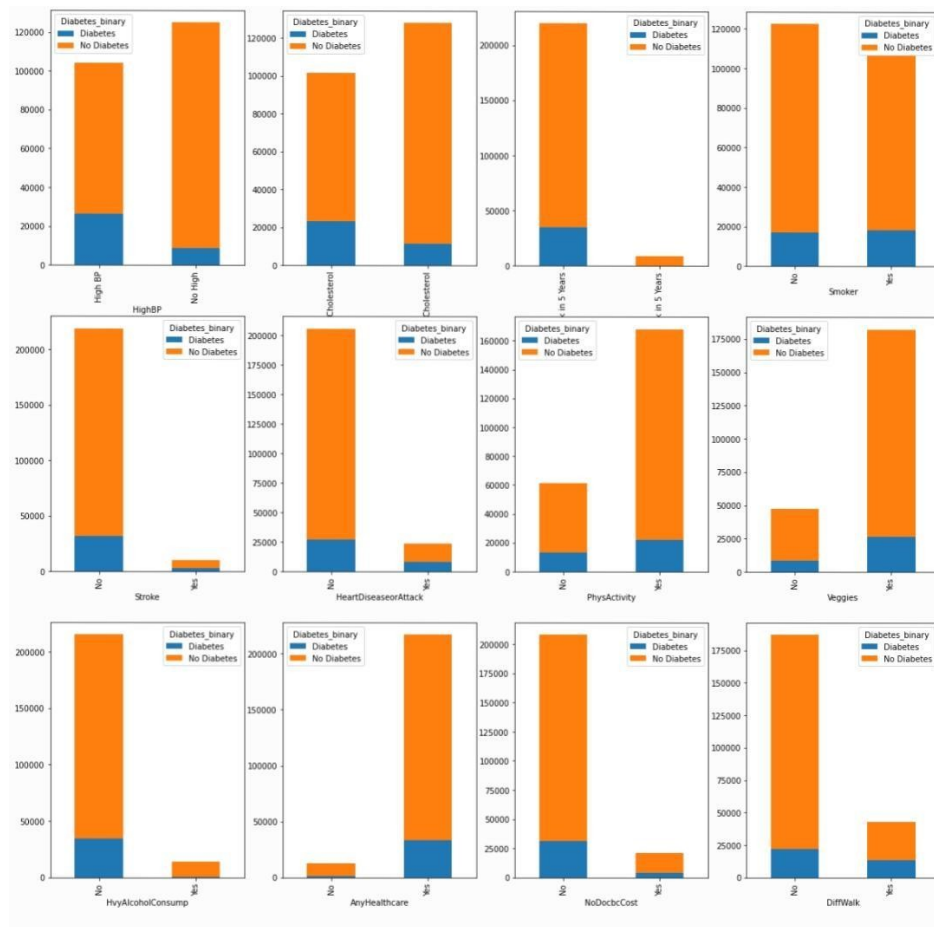
```
axe = ax.ravel() c = len(cols)
```

```
for i in range(c):
```

```
create_plot_pivot(data2, cols[i]).plot(kind='bar',stacked=True,
ax=axe[i])
```

```
axe[i].set_xlabel(cols[i])
```

## Output:



Let's view our target values ***"Diabetes\_binary"***

**Input:**

```
#average of column Daibetes_binary
```

```
# 0 for non-Diabetic person and 1 for Diabetic person
```

```
data2["Diabetes_binary"].value_counts()
```

**Output:**

```
No Diabetes    194377
```

```
Diabetes      35097
```

```
Name: Diabetes_binary, dtype: int64
```

**Input:**

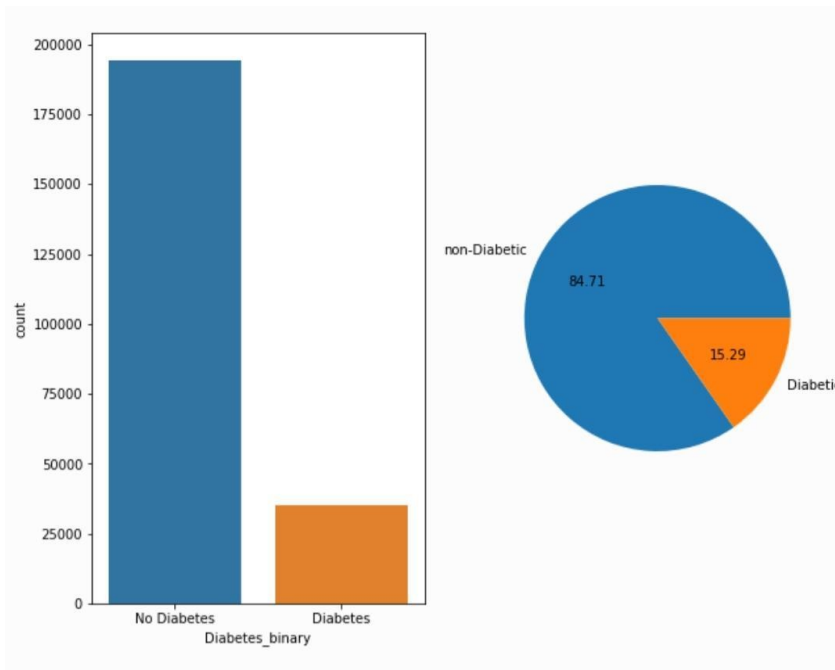
```
#checking the value count of Diabetes_binary_str by using countplot
```

```
figure1, plot1 = plt.subplots(1,2,figsize=(10,8))
```

```
sns.countplot(data2['Diabetes_binary'],ax=plot1[0]) #checking  
diabetic and non diabetic pepoles average by pie labels=["non-  
Diabetic","Diabetic"]
```

```
plt.pie(data2["Diabetes_binary"].value_counts() , labels  
=labels ,autopct='%0.02f' );
```

## Output:



## Benefits of a Diabetes Prediction System Using Machine Learning:

- ✓ **Early Detection:** Machine learning models can identify individuals at risk of diabetes before symptoms become apparent, enabling early intervention and prevention.
- ✓ **Personalized Care:** Predictive systems can tailor interventions to an individual's specific risk factors, optimizing the effectiveness of treatment and lifestyle modifications.
- ✓ **Improved Healthcare Resource Allocation:** Healthcare providers can allocate resources more efficiently by focusing on patients at higher risk, reducing the burden on the healthcare system.

- ✓ **Reduced Healthcare Costs:** Early intervention and prevention can lead to cost savings by reducing the need for expensive treatments and hospitalizations associated with diabetes complications.
- ✓ **Enhanced Patient Engagement:** Patients can be actively involved in managing their health, as they gain insight into their diabetes risk factors and are motivated to make healthier choices.

## **ADVANTAGES:**

- **Early Detection:** Machine learning models can analyze a wide range of patient data, including medical history, lifestyle factors, and genetic information, to identify patterns and risk factors that may not be evident to healthcare providers. This can lead to early detection of diabetes or prediabetes, allowing for timely intervention and management.
- **Personalized Care:** Machine learning models can provide personalized risk assessments and recommendations, taking into account an individual's unique characteristics. This allows for tailored treatment plans and lifestyle interventions, improving the effectiveness of diabetes management.
- **Improved Accuracy:** Machine learning algorithms can process vast amounts of data and extract intricate patterns and relationships that human clinicians might overlook. This results in more accurate predictions and risk assessments.
- **Cost-Effective:** Early detection and prevention can reduce the financial burden of managing diabetes. By identifying individuals at high risk of

developing diabetes, healthcare resources can be allocated more efficiently.

- **Accessibility:** Diabetes prediction systems can be deployed as mobile apps or web platforms, making them easily accessible to a broad population. This can help in reaching individuals who may not have regular access to healthcare facilities.
- **Continuous Monitoring:** Some machine learning models can offer continuous monitoring and real-time alerts, helping individuals manage their condition more effectively by making immediate adjustments to their lifestyle and medication.
- **Research and Data Analysis:** Diabetes prediction systems can generate valuable data for research, enabling healthcare professionals to better understand the disease, its progression, and the effectiveness of various interventions.

### **DIS-ADVANTAGES:**

- **Data Privacy and Security:** The collection and storage of sensitive medical data raise concerns about patient privacy and data security. Unauthorized access or data breaches can result in significant harm to individuals.
- **Data Quality:** The accuracy of predictions heavily depends on the quality and quantity of the data used for training. Inaccurate or incomplete data can lead to unreliable predictions.

- **Bias and Fairness:** Machine learning models may inherit biases present in the training data, which can lead to unfair predictions, especially when it comes to marginalized or underrepresented groups.
- **Lack of Medical Expertise:** Machine learning models do not possess the medical expertise and clinical judgment of healthcare professionals. They can provide predictions, but these should be used as a supplementary tool rather than a substitute for medical advice.
- **Overreliance on Technology:** Patients and healthcare providers may become overly reliant on the predictions of the system, potentially neglecting other important aspects of diabetes management, such as regular check-ups and lifestyle changes.
- **Cost of Implementation:** Developing and implementing a machine learning-based diabetes prediction system can be costly. This includes the cost of data collection, model development, software development, and ongoing maintenance.
- **Ethical Concerns:** The use of predictive models in healthcare raises ethical dilemmas, such as how to communicate predictions to patients and how to balance the benefits of early detection with the potential psychological burden of living with the knowledge of an increased risk.

## **CONCLUSION:**

The development of a Diabetes Prediction System using machine learning has shown significant promise in the field of healthcare. This system leverages the power of data analysis, pattern recognition, and predictive modeling to assist in the early detection and management of diabetes. The key findings and takeaways from such a system are as follows:

- ❖ **Early Detection:** Machine learning algorithms can effectively identify individuals at risk of developing diabetes even before the onset of symptoms. This early detection can enable healthcare professionals to intervene with preventive measures and lifestyle changes.
- ❖ **Personalized Medicine:** These systems can offer personalized recommendations and treatment plans based on an individual's unique health profile. This tailoring of care can lead to better outcomes and improved quality of life for diabetes patients.
- ❖ **Data-Driven Insights:** The system can analyze large datasets, which is essential for understanding the various risk factors and contributing variables associated with diabetes. This information can guide public health policies and individual health choices.
- ❖ **Reduced Healthcare Costs:** By preventing or delaying the onset of diabetes through early intervention, the system can potentially reduce the long-term healthcare costs associated with diabetes management and its complications.



❖ **Continuous Monitoring:** Machine learning can be integrated with wearable devices to provide real-time monitoring of glucose levels and other health parameters, offering a more comprehensive view of a patient's health.

*However, there are some significant challenges and issues associated with the development and deployment of a Diabetes Prediction System using machine learning:*

**Data Quality and Privacy:** The effectiveness of machine learning models heavily depends on the quality and quantity of data. Obtaining clean and representative data is often a challenge. Additionally, maintaining the privacy and security of sensitive health information is a paramount concern.

**Model Generalization:** Machine learning models must be trained on diverse and representative datasets to ensure they generalize well to different populations and demographics. Biases in the data can lead to inaccurate predictions.

**Interpretability:** Many machine learning models, especially deep learning models, are often considered black boxes, making it challenging for healthcare professionals to understand the reasoning behind predictions. Interpretable AI is an ongoing research area.

**Ethical Concerns:** The use of predictive models in healthcare raises ethical concerns, including potential discrimination, bias, and the need for transparency and accountability in decision-making.

**Regulatory and Legal Issues:** Compliance with healthcare regulations, such as HIPAA in the United States, and legal liability in case of inaccurate predictions or system failures must be addressed.

**User Acceptance:** Patients and healthcare providers may be skeptical of or resistant to relying on machine learning systems for critical healthcare decisions. Overcoming user acceptance issues is crucial for the adoption of such systems.

**Cost of Implementation:** Building and maintaining a robust machine learning-based Diabetes Prediction System can be costly, which may limit its availability, especially in resource-constrained healthcare settings.

In conclusion, a Diabetes Prediction System using machine learning holds great promise for early detection and personalized management of diabetes. However, addressing data quality, privacy, bias, and ethical concerns, along with ensuring regulatory compliance, user acceptance, and cost-effectiveness, are essential for its successful development and implementation in healthcare.