

PREDICTING DIABETES SYSTEM USING MACHINE LEARNING

TEAM LEADER: 922321106033

X.SHOBIA ROSALIN MARRIE

Phase - 5 submission document

Project title: **Diabetes Prediction System**

Phase-4: **Project Documentation and Submission**

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



INTRODUCTION:

An AI-based diabetes prediction system using machine learning is a system that uses machine learning algorithms to analyze data from patients and predict whether they are at risk of developing diabetes. This type of system can be used to help people identify their risk factors for diabetes and take steps to prevent it, or to help healthcare providers diagnose diabetes early.



Machine learning algorithms work by learning from data. In the case of a diabetes prediction system, the algorithm would be trained on a dataset of patient data that includes information such as age, gender, weight, height, blood sugar levels, and other health factors. Once the algorithm is trained, it can be used to predict the risk of diabetes for new patients.

AI-based diabetes prediction systems have the potential to improve the early detection and prevention of diabetes. They can also help to reduce the healthcare costs associated with diabetes.

Here is an example of how an AI-based diabetes prediction system might work:

1. The patient provides information about their health history, family history, and lifestyle factors.

2. The system uses this information to calculate the patient's risk of developing diabetes.
3. The system then provides the patient with personalized recommendations for reducing their risk of diabetes, or for managing their diabetes if they have already been diagnosed.

AI-based diabetes prediction systems are still under development, but they have the potential to be a valuable tool for improving the health and well-being of people with diabetes.

Benefits of AI-based diabetes prediction systems

- Early detection and prevention of diabetes
- Reduced healthcare costs
- Personalized recommendations for reducing risk or managing diabetes
- Improved quality of life for people with diabetes

Challenges of AI-based diabetes prediction systems

- Ensuring the accuracy and reliability of the system
- Accessibility and affordability of the system
- Ethical considerations, such as the potential for bias in the system.

An AI-based diabetes prediction system using machine learning is a technology that leverages algorithms and data analysis to predict the likelihood of an individual developing diabetes. It does so by utilizing historical patient data, such as medical records, lifestyle information, and genetic factors, to create predictive models. These models can assess the risk of diabetes for a person

and provide early warnings or recommendations for preventive measures. Machine learning techniques, like decision trees, support vector machines, or neural networks, are commonly employed to build these predictive models. The goal is to enable proactive healthcare interventions and personalized strategies for managing and preventing diabetes, ultimately improving the overall health and well-being of individuals at risk.

Given Dataset:

| | A | B | C | D | E | F | G | H | I |
|----|-----------|---------|-----------|------------|---------|------|-----------|-----|---------|
| 1 | Pregnanci | Glucose | BloodPres | SkinThickn | Insulin | BMI | DiabetesF | Age | Outcome |
| 2 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.627 | 50 | 1 |
| 3 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.351 | 31 | 0 |
| 4 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| 5 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 6 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |
| 7 | 5 | 116 | 74 | 0 | 0 | 25.6 | 0.201 | 30 | 0 |
| 8 | 3 | 78 | 50 | 32 | 88 | 31 | 0.248 | 26 | 1 |
| 9 | 10 | 115 | 0 | 0 | 0 | 35.3 | 0.134 | 29 | 0 |
| 10 | 2 | 197 | 70 | 45 | 543 | 30.5 | 0.158 | 53 | 1 |
| 11 | 8 | 125 | 96 | 0 | 0 | 0 | 0.232 | 54 | 1 |
| 12 | 4 | 110 | 92 | 0 | 0 | 37.6 | 0.191 | 30 | 0 |
| 13 | 10 | 168 | 74 | 0 | 0 | 38 | 0.537 | 34 | 1 |
| 14 | 10 | 139 | 80 | 0 | 0 | 27.1 | 1.441 | 57 | 0 |
| 15 | 1 | 189 | 60 | 23 | 846 | 30.1 | 0.398 | 59 | 1 |
| 16 | 5 | 166 | 72 | 19 | 175 | 25.8 | 0.587 | 51 | 1 |
| 17 | 7 | 100 | 0 | 0 | 0 | 30 | 0.484 | 32 | 1 |
| 18 | 0 | 118 | 84 | 47 | 230 | 45.8 | 0.551 | 31 | 1 |
| 19 | 7 | 107 | 74 | 0 | 0 | 29.6 | 0.254 | 31 | 1 |
| 20 | 1 | 103 | 30 | 38 | 83 | 43.3 | 0.183 | 33 | 0 |
| 21 | 1 | 115 | 70 | 30 | 96 | 34.6 | 0.529 | 32 | 1 |
| 22 | 3 | 126 | 88 | 41 | 235 | 39.3 | 0.704 | 27 | 0 |
| 23 | 8 | 99 | 84 | 0 | 0 | 35.4 | 0.388 | 50 | 0 |
| 24 | 7 | 196 | 90 | 0 | 0 | 39.8 | 0.451 | 41 | 1 |

769 Rows *9 Columns

DATASET LINK:

<https://www.kaggle.com/datasets/mathchi/diabetes-data-set>

Here's is the list of tools and software commonly used in the process:

The following are some commonly used tools and software in the process of AI-based diabetes prediction system using machine learning:

- **Programming languages:** Python, R, and Julia are popular programming languages for machine learning and data science.
- **Machine learning libraries:** Scikit-learn, TensorFlow, and PY-Torch are popular machine learning libraries that provide a variety of algorithms and tools for building and deploying machine learning models.
- **Cloud computing platforms:** Google Cloud Platform, Amazon Web Services, and Microsoft Azure are cloud computing platforms that provide scalable and affordable infrastructure for training and deploying machine learning models.
- **Data visualization tools:** Tableau, Power BI, and Vega-Lite are data visualization tools that can be used to explore and visualize data, and to create interactive dashboards and reports.

In addition to these general-purpose tools and software, there are also some specialized tools and software that can be used for AI-based diabetes prediction. For example, the following tools can be used to process and prepare medical data:

- **Medical imaging software:** Medical imaging software can be used to process and analyze medical images, such as MRI scans and CT scans.
- **Electronic health record (EHR) systems:** EHR systems can be used to access and extract patient data from medical records.

Once the data has been prepared, the following tools can be used to build and deploy machine learning models for diabetes prediction:

- **Machine learning frameworks:** Machine learning frameworks such as TensorFlow and PY Torch provide a variety of tools and libraries for building and deploying machine learning models.
- **Auto-ML tools:** Auto-ML tools can automate the process of building and tuning machine learning models, which can be helpful for non-experts.

Once the model has been trained and deployed, it can be used to predict the risk of diabetes for new patients. The following tools can be used to deploy and manage machine learning models at scale:

- **Model serving tools:** Model serving tools can be used to deploy machine learning models in production, so that they can be used to make predictions on new data.
- **Model monitoring tools:** Model monitoring tools can be used to monitor the performance of machine learning models in production, and to detect any changes in the model's performance.

This is not an exhaustive list of all the tools and software that can be used for AI-based diabetes prediction, but it covers some of the most commonly used tools and software.

Example workflow for building and deploying an AI-based diabetes prediction system using machine learning

The following is an example workflow for building and deploying an AI-based diabetes prediction system using machine learning:

1. **Collect and prepare the data.** This may involve collecting data from medical records, medical imaging software, and other sources. The data may need to be cleaned and preprocessed before it can be used to train a machine learning model.
2. **Choose a machine learning algorithm.** There are many different machine learning algorithms that can be used for diabetes prediction. Some popular choices include support vector machines (SVMs), random forests, and logistic regression.
3. **Train the machine learning model.** Once the data has been prepared and the algorithm has been chosen, the machine

learning model can be trained. This involves feeding the model the prepared data and allowing it to learn the patterns in the data.

4. Evaluate the machine learning model. Once the model has been trained, it is important to evaluate its performance on a held-out test set. This will help to ensure that the model is



generalizing well and is not overfitting the training data.

5. Deploy the machine learning model. Once the model has been evaluated and is performing well, it can be deployed to production. This may involve deploying the model to a cloud computing platform or to a local server.

Once the model has been deployed, it can be used to predict the risk of diabetes for new patients. For example, a doctor could enter a patient's medical data into the model and the model would predict the patient's risk of developing diabetes. This information can be used to help the doctor make decisions about the patient's care, such as whether to recommend further testing or lifestyle changes.

DESIGN THINKING AND PRESENT IT IN THE FORM OF DOCUMENT

Design thinking is a non-linear, iterative process that teams use to understand users, challenge assumptions, redefine problems, and create innovative solutions to prototype and test. In healthcare, design thinking can be used to develop new technologies and services that improve the patient experience and outcomes.

Here is a design thinking process for developing a diabetes prediction system using machine learning:

Empathize

The first step in design thinking is to empathize with the users. This means understanding their needs, wants, and pain points. In the context of a diabetes prediction system, this would involve talking to people with diabetes about their experiences, as well as interviewing healthcare professionals.

Some questions to consider include:

1. What are the biggest challenges people with diabetes face?
2. What would help them manage their diabetes better?
3. What are their concerns about using a machine learning-based prediction system?

Define

Once you have a good understanding of the users, you can start to define the problem you are trying to solve. In the case of a diabetes prediction system, this might be something like:

To develop a machine learning-based system that can accurately predict whether or not someone is at risk of developing diabetes, in order to help them take preventive measures.

Ideate

Once you have a clear definition of the problem, you can start to brainstorm ideas for solutions. This is where design thinking gets creative! Here are some ideas to get you started:

- Develop a mobile app that allows users to enter their personal health data and receive a personalized risk assessment.
- Create a web-based platform where users can connect with other people with diabetes and share their experiences.
- Develop a machine learning model that can be used by healthcare professionals to identify patients who are at high risk of developing diabetes, so that they can be targeted with preventive care.

Prototype

Once you have some ideas, it's time to start prototyping. This means creating a working model of your solution so that you can test it with users and get feedback. In the context of a diabetes prediction system, this might involve developing a simple mobile app or web-based platform.

Test

Once you have a prototype, you can start to test it with users. This is an important step in the design thinking process, as it allows you to identify any problems with your solution and make necessary changes.

When testing your diabetes prediction system, you should focus on the following:

1. Is the system easy to use?
2. Is the system accurate?
3. Do users find the system helpful?

Deploy

Once you have tested your prototype and made any necessary changes, you can deploy your solution to the real world. This might involve launching your mobile app or web-based platform, or making your machine learning model available to healthcare professionals.

Monitor

Even after you have deployed your solution, it is important to continue to monitor it and collect feedback from users. This will help you to identify any areas where you can improve.

Design thinking is a powerful tool that can be used to develop innovative solutions to complex problems. By following the design thinking process, you can develop a diabetes prediction system that is accurate, user-friendly, and helpful.

Here are some additional considerations for designing a diabetes prediction system using machine learning:

- **Data collection:** Machine learning models need to be trained on large amounts of data. This data should be representative of the population that you are targeting with your system.
- **Model selection:** There are many different machine learning algorithms that can be used for diabetes prediction.

It is important to select the right algorithm for your specific dataset and target population.

- **Model evaluation:** Once you have trained a machine learning model, it is important to evaluate its performance on a held-out test set. This will give you an idea of how well the model will generalize to new data.
- **Model interpretation:** It is important to be able to interpret the results of your machine learning model. This will help you to understand why the model is making the predictions that it is making.
- **Model deployment:** Once you have trained and evaluated a machine learning model, you need to deploy it in a way that makes it accessible to users. This might involve developing a mobile app, web-based platform, or integrating the model into an existing healthcare system.

By following these considerations, you can develop a diabetes prediction system that is both accurate and useful.

DESIGN INTO INNOVATION IN DIABETES PREDICTION SYSTEM USING MACHINE LEARNING

Design Principles

The following design principles can be used to develop an innovative diabetes prediction system using machine learning:

- ❖ **Data-driven:** The system should be trained on a large and diverse dataset of medical records, including patient demographics, clinical data, and lifestyle factors. This will

allow the system to learn complex patterns in the data and make accurate predictions.

- ❖ **Explainable:** The system should be able to explain its predictions to users in a clear and concise way. This is important for building trust and enabling users to take informed decisions about their health.
- ❖ **Personalized:** The system should be able to personalize its predictions based on individual user characteristics. This will make the system more accurate and relevant for each user.
- ❖ **Scalable:** The system should be designed to be scalable to large numbers of users and data. This will make it feasible to deploy the system in real-world settings.

Innovative Features

In addition to the design principles above, the following innovative features can be incorporated into a diabetes prediction system using machine learning:

- **Multimodal data fusion:** The system can be trained on multiple types of data, such as medical records, wearable device data, and genetic data. This will improve the accuracy of the system's predictions.
- **Real-time monitoring:** The system can be used to monitor users' health data in real time and provide timely alerts if their risk of developing diabetes increases.
- **Personalized interventions:** The system can be used to develop personalized interventions for users at risk of developing diabetes. This could include lifestyle recommendations, dietary advice, or medication plans.

Implementation Considerations

The following implementation considerations should be taken into account when developing a diabetes prediction system using machine learning:

1. **Data security and privacy:** The system should be designed to protect user data from unauthorized access and use. This is important for maintaining user trust and compliance with data privacy regulations.
2. **Clinical validation:** The system should be clinically validated to ensure that its predictions are accurate and reliable.
3. **Integration with existing healthcare systems:** The system should be designed to be integrated with existing healthcare systems, such as electronic health records (EHRs) and clinical decision support systems (CDSSs). This will make it easier for healthcare providers to use the system to improve patient care.

Examples of Innovation

The following are some examples of innovative diabetes prediction systems that are being developed using machine learning:

1. **DeepMind's Streams:** This system uses machine learning to analyze data from wearable devices, such as glucose monitors and fitness trackers, to predict a person's risk of developing diabetes.

2. GLYTEC Glu-commander: This system uses machine learning to analyze data from insulin pumps and continuous glucose monitors to help people with diabetes manage their blood sugar levels.
3. Viz.ai: This system uses machine learning to analyze medical images, such as CT scans and MRI scans, to detect early signs of diabetes complications, such as heart disease and kidney disease.

Machine learning has the potential to revolutionize diabetes prediction and management. By incorporating the design principles and innovative features discussed above, developers can create diabetes prediction systems that are more accurate, personalized, and scalable than ever before. These systems can help people to identify their risk of developing diabetes early on and take steps to prevent or delay the onset of the disease. They can also help people with diabetes to better manage their condition and avoid complications.

BUILD LOADING AND PRE-PROCESSING THE DATASET

Loading and preprocessing data for a diabetes prediction system can be challenging due to various factors, including the nature of healthcare data and the complexity of the prediction task. Here are some specific challenges you might encounter:

1)Data Quality and Completeness:

Healthcare data can be incomplete or contain missing values, which can affect the accuracy of predictions. Ensuring data quality and completeness is crucial.

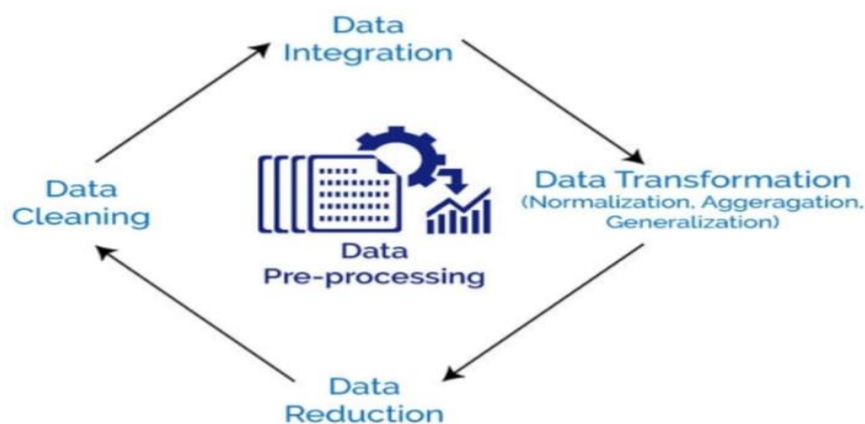
2) Imbalanced Data:

In many healthcare datasets, there is an imbalance between diabetes-positive and diabetes-negative cases. The model may be biased toward the majority class without proper handling

3) Data Privacy and Security:

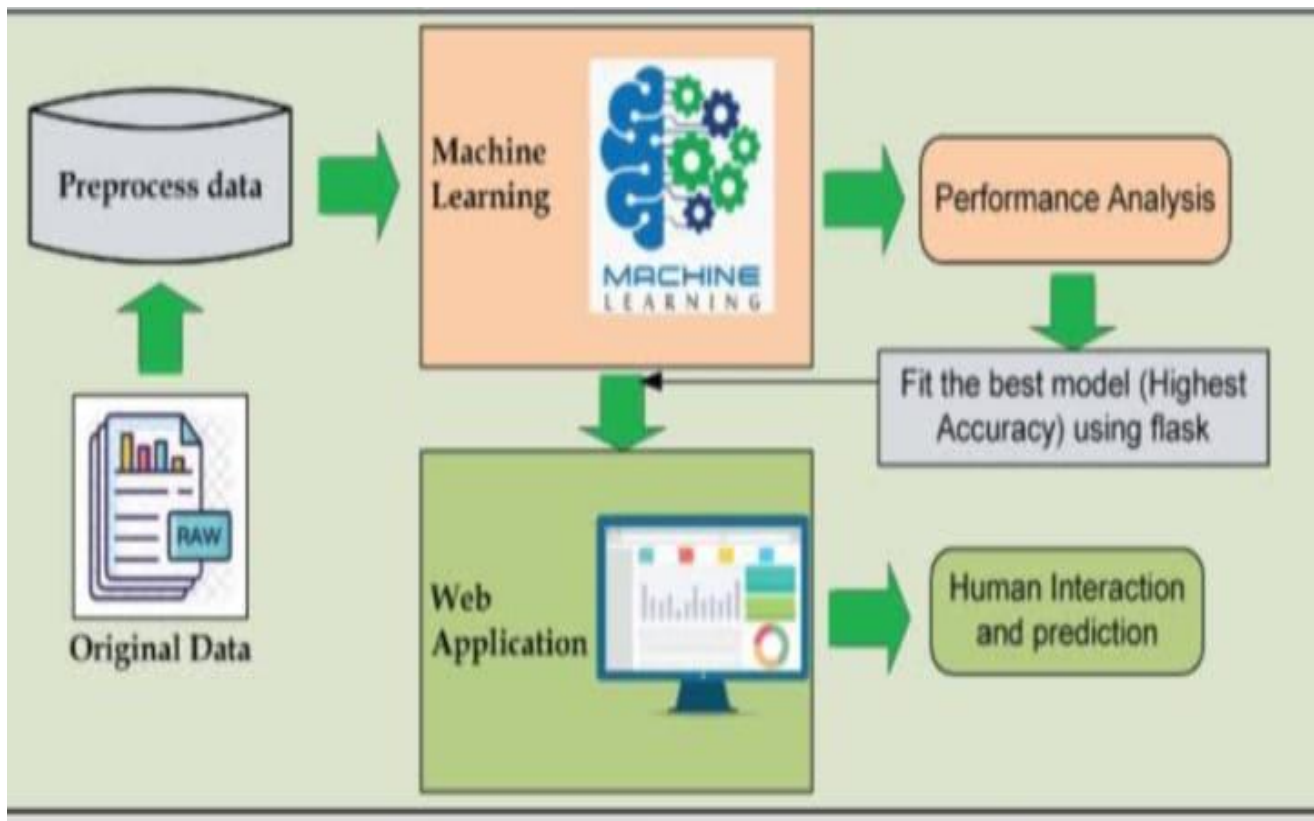
Healthcare data is highly sensitive, and there are strict regulations like HIPAA in the United States and GDPR in Europe. Handling and securing patient data while preserving privacy is challenging.

4) Feature Engineering: Selecting the most relevant features or variables for prediction is often a complex task. Medical data may have numerous features, and determining which ones are informative can be challenging.



5)Categorical Data:

Healthcare data often includes categorical variables such as gender, medication types, and medical procedures. Converting these into a numerical format that machine learning models can understand is a challenge.



PYTHON PROGRAM:

DATA COLLECTION

Patient Records: Gather historical health records, including blood glucose measurements, medical history, and medication usage.

Wearable Devices: Collect real-time data from wearable devices like continuous glucose monitors (CGMs), fitness trackers, and smartwatches.

Surveys and Questionnaires: Administer surveys to collect lifestyle information, dietary habits, physical activity, and family history.

Genetics: Incorporate genetic data, if available, to assess the genetic predisposition to diabetes.

DATA PRE-PROCESSING

Data Cleaning: Remove duplicates, correct errors, and handle missing values in the collected data.

Feature Selection: Identify relevant features (variables) for prediction, such as glucose levels, BMI, age, and family history.

Data Transformation: Normalize or scale numerical features to bring them to a similar range and encode categorical data.

Time Series Data: If using time series data (e.g., CGM data), handle irregular sampling rates, and possibly apply smoothing techniques.

In [1]:

```
#Installation of required libraries
```

```
import numpy as np
```

```
import pandas as pd
```

```
import statsmodels.api as sm
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.preprocessing import scale, StandardScaler
from sklearn.model_selection import train_test_split,
GridSearchCV, cross_val_score
from sklearn.metrics import confusion_matrix, accuracy_score,
mean_squared_error, r2_score, roc_auc_score, roc_curve,
classification_report
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.neural_network import MLPClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.ensemble import GradientBoostingClassifier
from lightgbm import LGBMClassifier
from sklearn.model_selection import KFold
import warnings
warnings.simplefilter(action = "ignore")
In [2]:
#Reading the dataset
```

```
df = pd.read_csv("../input/pima-indians-diabetes-
database/diabetes.csv")
```

In [3]:

```
# The first 5 observation units of the data set were accessed.
```

```
df.head()
```

Out[3]:

| | 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 [4]:

```
# The size of the data set was examined. It consists of 768
observation units and 9 variables.
```

```
df.shape
```

Out[4]:

```
(768, 9)
```

In [5]:

```
#Feature information
```

```
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 [6]:

Descriptive statistics of the data set accessed.

df.describe([0.10,0.25,0.50,0.75,0.90,0.95,0.99]).T

Out[6]:

| count | mean | std | min | 10% | 25% | 50% | 75% | 90% | 95% | 99% max |
|-------|------|-----|-----|-----|-----|-----|-----|-----|-----|---------|
|-------|------|-----|-----|-----|-----|-----|-----|-----|-----|---------|

| | | | | | |
|--------------------------|-----------|------------|------------|-----------|-------------------|
| Pregnancies | 768.0 | 3.845052 | 3.369578 | 0.000 | 0.000 |
| | 1.00000 | 3.0000 | 6.00000 | 9.0000 | 10.00000 13.00000 |
| | 17.00 | | | | |
| Glucose | 768.0 | 120.894531 | 31.972618 | 0.000 | |
| | 85.000 | 99.00000 | 117.0000 | 140.25000 | 167.0000 |
| | 181.00000 | 196.00000 | 199.00 | | |
| BloodPressure | 768.0 | 69.105469 | 19.355807 | 0.000 | |
| | 54.000 | 62.00000 | 72.0000 | 80.00000 | 88.0000 90.00000 |
| | 106.00000 | 122.00 | | | |
| SkinThickness | 768.0 | 20.536458 | 15.952218 | 0.000 | |
| | 0.000 | 0.00000 | 23.0000 | 32.00000 | 40.0000 44.00000 |
| | 51.33000 | 99.00 | | | |
| Insulin | 768.0 | 79.799479 | 115.244002 | 0.000 | |
| | 0.000 | 0.00000 | 30.5000 | 127.25000 | 210.0000 |
| | 293.00000 | 519.90000 | 846.00 | | |
| BMI | 768.0 | 31.992578 | 7.884160 | 0.000 | 23.600 |
| | 27.30000 | 32.0000 | 36.60000 | 41.5000 | 44.39500 50.75900 |
| | 67.10 | | | | |
| DiabetesPedigreeFunction | 768.0 | 0.471876 | 0.331329 | | |
| | 0.078 | 0.165 | 0.24375 | 0.3725 | 0.62625 0.8786 |
| | 1.13285 | 1.69833 | 2.42 | | |
| Age | 768.0 | 33.240885 | 11.760232 | 21.000 | 22.000 |
| | 24.00000 | 29.0000 | 41.00000 | 51.0000 | 58.00000 67.00000 |
| | 81.00 | | | | |
| Outcome | 768.0 | 0.348958 | 0.476951 | 0.000 | 0.000 |
| | 0.00000 | 0.0000 | 1.00000 | 1.0000 | 1.00000 1.00000 |
| | 1.00 | | | | |

In [7]:

```
# The distribution of the Outcome variable was examined.
```

```
df["Outcome"].value_counts()*100/len(df)
```

Out[7]:

```
0    65.104167
```

```
1    34.895833
```

Name: Outcome, dtype: float64

In [8]:

```
# The classes of the outcome variable were examined.
```

```
df.Outcome.value_counts()
```

Out[8]:

```
0    500
```

```
1    268
```

Name: Outcome, dtype: int64

In [9]:

```
# The histogram of the Age variable was reached.
```

```
df["Age"].hist(edgecolor = "black");
```

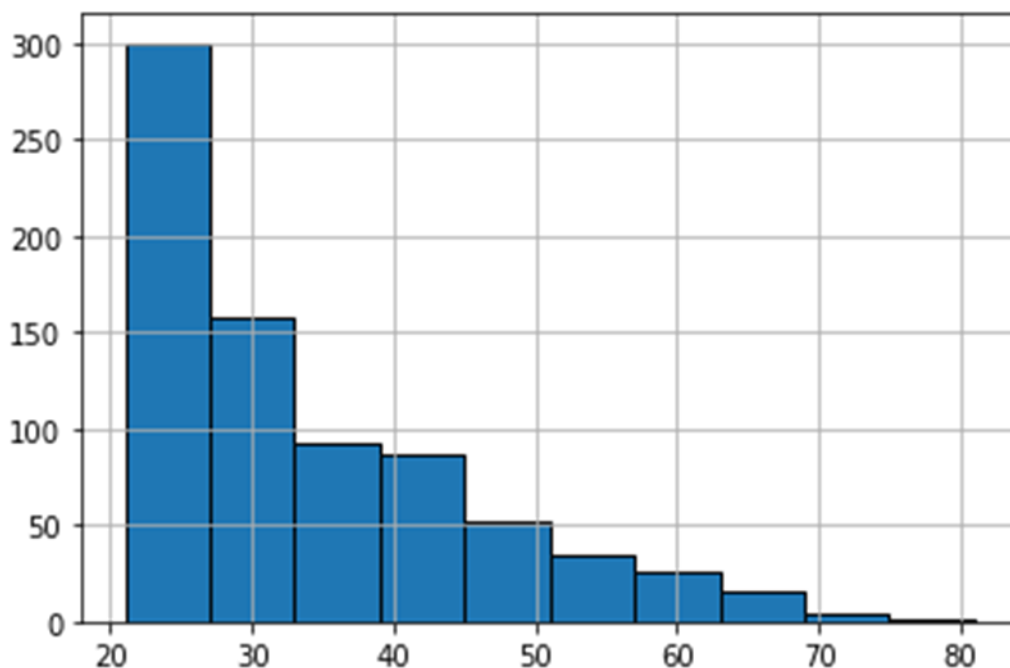
In[10]:

```
print("Max Age: " + str(df["Age"].max()) + " Min Age: " +  
str(df["Age"].min()))
```

Max Age: 81 Min Age: 21

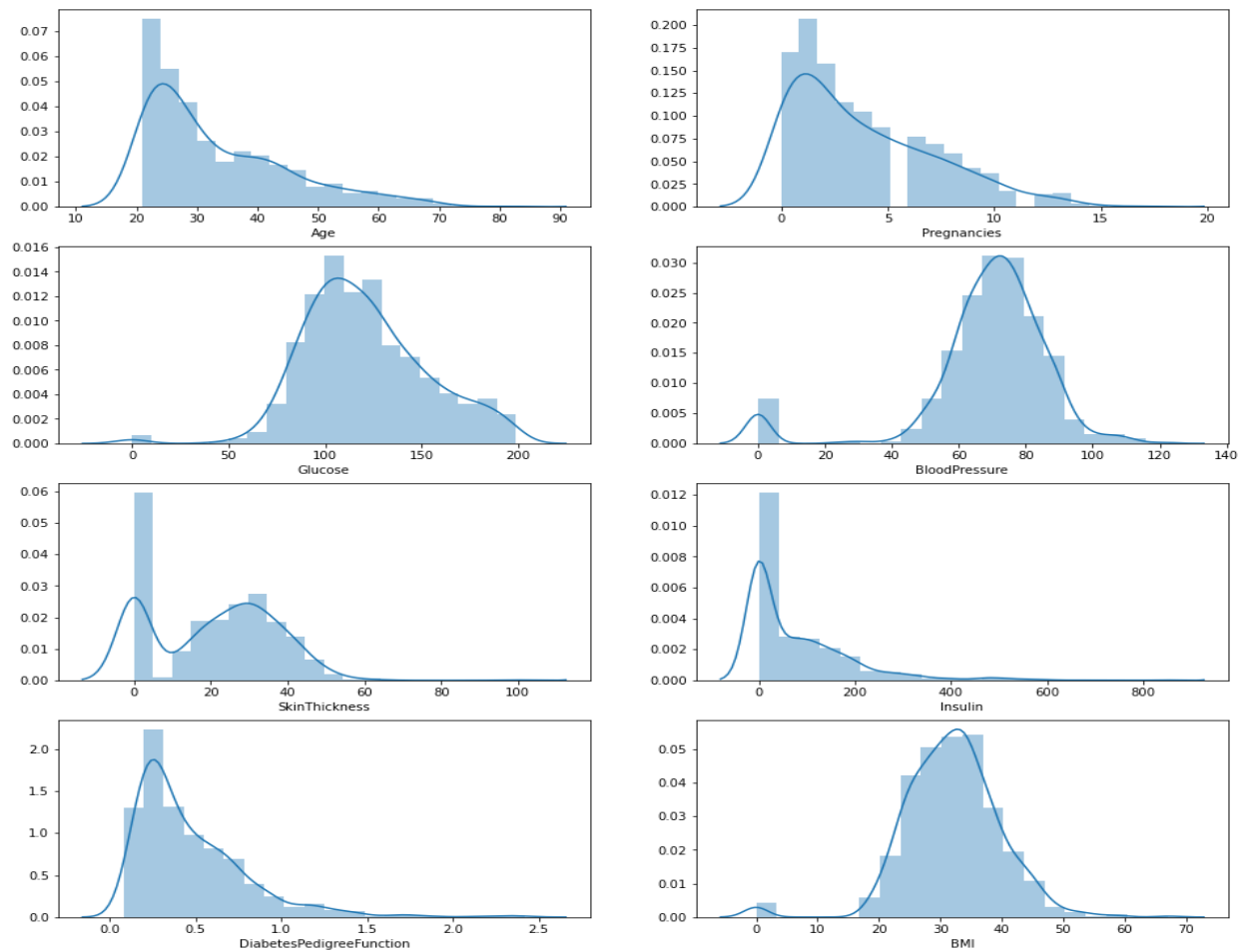
In [11]:

```
# Histogram and density graphs of all variables were accessed.  
fig, ax = plt.subplots(4,2, figsize=(16,16))  
sns.distplot(df.Age, bins = 20, ax=ax[0,0])  
sns.distplot(df.Pregnancies, bins = 20, ax=ax[0,1])  
sns.distplot(df.Glucose, bins = 20, ax=ax[1,0])  
sns.distplot(df.BloodPressure, bins = 20, ax=ax[1,1])  
sns.distplot(df.SkinThickness, bins = 20, ax=ax[2,0])  
sns.distplot(df.Insulin, bins = 20, ax=ax[2,1])  
sns.distplot(df.DiabetesPedigreeFunction, bins = 20, ax=ax[3,0])  
sns.distplot(df.BMI, bins = 20, ax=ax[3,1])
```



Out[11]:

<matplotlib.axes._subplots.AxesSubplot at 0x7f77b83d5950>



In[12]:

```
df.groupby("Outcome").agg({"Pregnancies":"mean"})
```

Out[12]:

Pregnancies

Outcome

0 3.298000

1 4.865672

In[13]:

```
df.groupby("Outcome").agg({"Age": "mean"})
```

Out[13]:

| | Age |
|---------|-----------|
| Outcome | |
| 0 | 31.190000 |
| 1 | 37.067164 |

In[14]:

```
df.groupby("Outcome").agg({"Age": "max"})
```

Out[14]:

| | Age |
|---------|-----|
| Outcome | |
| 0 | 81 |
| 1 | 70 |

In[15]:

```
df.groupby("Outcome").agg({"Insulin": "mean"})
```

Out[15]:

| | Insulin |
|---------|------------|
| Outcome | |
| 0 | 68.792000 |
| 1 | 100.335821 |

In[16]:

```
df.groupby("Outcome").agg({"Insulin": "max"})
```

Out[16]:

| | Insulin |
|---------|---------|
| Outcome | |
| 0 | 744 |
| 1 | 846 |

In[17]:

```
df.groupby("Outcome").agg({"Glucose": "mean"})
```

Out[17]:

| | Glucose |
|---------|------------|
| Outcome | |
| 0 | 109.980000 |
| 1 | 141.257463 |

In[18]:

```
df.groupby("Outcome").agg({"Glucose": "max"})
```

Out[18]:

| | Glucose |
|---------|---------|
| Outcome | |
| 0 | 197 |

```
1    199
```

```
In [19]:
```

```
df.groupby("Outcome").agg({"BMI": "mean"})
```

```
Out[19]:
```

```
      BMI
Outcome
0    30.304200
1    35.142537
```

```
In[20]:
```

```
# The distribution of the outcome variable in the data was
examined and visualized.
```

```
f,ax=plt.subplots(1,2,figsize=(18,8))
```

```
df['Outcome'].value_counts().plot.pie(explode=[0,0.1],autopct='
%1.1f%%',ax=ax[0],shadow=True)
```

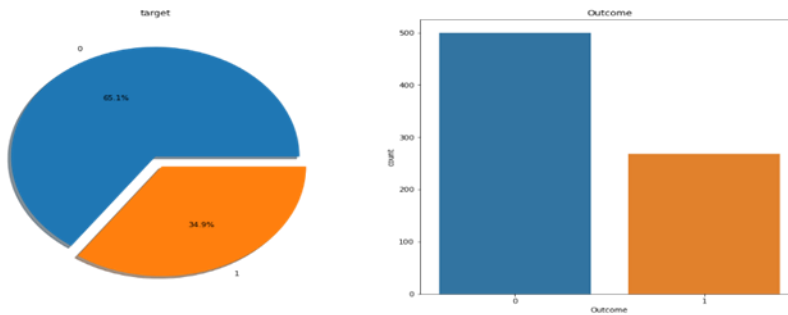
```
ax[0].set_title('target')
```

```
ax[0].set_ylabel("")
```

```
sns.countplot('Outcome',data=df,ax=ax[1])
```

```
ax[1].set_title('Outcome')
```

```
plt.show()
```



In[21]:

Access to the correlation of the data set was provided. What kind of relationship is examined between the variables.

If the correlation value is > 0 , there is a positive correlation. While the value of one variable increases, the value of the other variable also increases.

Correlation = 0 means no correlation.

If the correlation is < 0 , there is a negative correlation. While one variable increases, the other variable decreases.

When the correlations are examined, there are 2 variables that act as a positive correlation to the Salary dependent variable.

These variables are Glucose. As these increase, Outcome variable increases.

`df.corr()`

Out [21]:

| | | | | |
|-------------|--------------------------|---------------|---------------|---------|
| Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin |
| BMI | DiabetesPedigreeFunction | Age | Outcome | |

```

Pregnancies 1.000000 0.129459 0.141282 -0.081672 -
0.073535 0.017683 -0.033523 0.544341 0.221898
Glucose 0.129459 1.000000 0.152590 0.057328 0.331357
0.221071 0.137337 0.263514 0.466581
BloodPressure 0.141282 0.152590 1.000000 0.207371 0.088933
0.281805 0.041265 0.239528 0.065068
SkinThickness -0.081672 0.057328 0.207371 1.000000
0.436783 0.392573 0.183928 -0.113970 0.074752
Insulin -0.073535 0.331357 0.088933 0.436783 1.000000
0.197859 0.185071 -0.042163 0.130548
BMI 0.017683 0.221071 0.281805 0.392573 0.197859 1.000000
0.140647 0.036242 0.292695
DiabetesPedigreeFunction -0.033523 0.137337 0.041265
0.183928 0.185071 0.140647 1.000000 0.033561 0.173844
Age 0.544341 0.263514 0.239528 -0.113970 -0.042163
0.036242 0.033561 1.000000 0.238356
Outcome 0.221898 0.466581 0.065068 0.074752 0.130548
0.292695 0.173844 0.238356 1.000000

```

In[22]:

```

# Correlation matrix graph of the data set
f, ax = plt.subplots(figsize= [20,15])
sns.heatmap(df.corr(), annot=True, fmt=".2f", ax=ax, cmap =
"magma" )
ax.set_title("Correlation Matrix", fontsize=20)

```

```
plt.show()
```

2) Data Preprocessing

2.1) Missing Observation Analysis

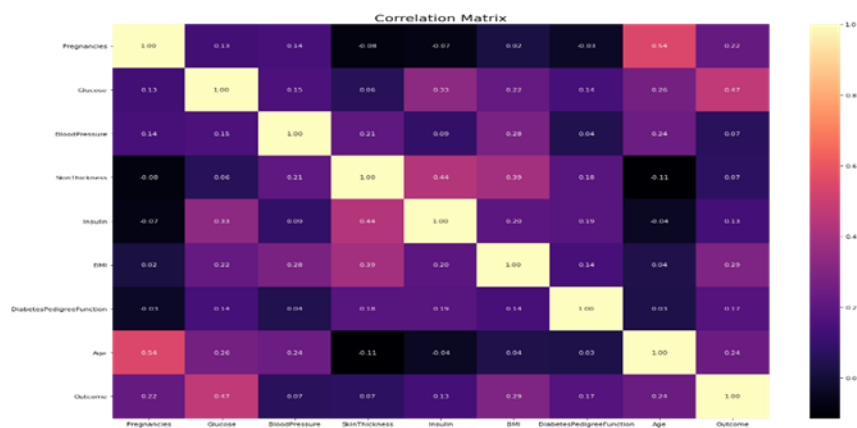
In [23]:

```
df[['Glucose','BloodPressure','SkinThickness','Insulin','BMI']] =  
df[['Glucose','BloodPressure','SkinThickness','Insulin','BMI']].re  
place(0,np.NaN)
```

In [24]:

```
df.head()
```

Out[24]:



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

```
0    6    148.0    72.0  35.0  NaN  33.6  0.627    50    1
```

```
1    1    85.0  66.0  29.0  NaN  26.6  0.351    31    0
```

```

2    8    183.0    64.0 NaN NaN 23.3 0.672    32    1
3    1    89.0 66.0 23.0 94.0 28.1 0.167    21    0
4    0    137.0    40.0 35.0 168.0    43.1 2.288    33    1

```

In[25]:

```
# Now, we can look at where are missing values
```

```
df.isnull().sum()
```

Out[25]:

```

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

```

```
dtype: int64
```

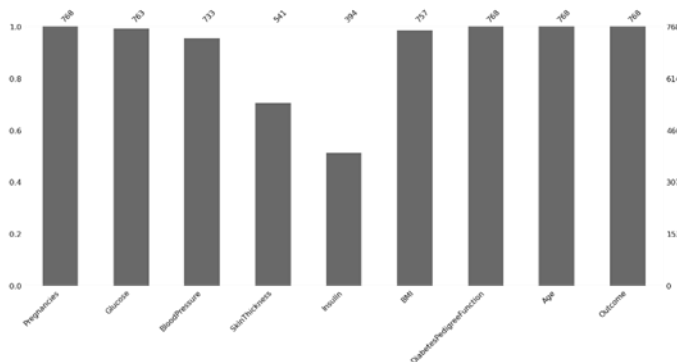
In [26]:

```
# Have been visualized using the missingno library for the
visualization of missing observations.
```


Plotting

```
import missingno as msno
```

```
msno.bar(df);
```



In[27]:

The missing values will be filled with the median values of each variable.

```
def median_target(var):
```

```
    temp = df[df[var].notnull()]
```

```
    temp = temp[[var,
'Outcome']].groupby(['Outcome'])[var].median().reset_index()
```

```
    return temp
```

In [28]:

The values to be given for incomplete observations are given the median value of people who are not sick and the median values of people who are sick.

```
columns = df.columns
```

```

columns = columns.drop("Outcome")
for i in columns:
    median_target(i)
    df.loc[(df['Outcome'] == 0 ) & (df[i].isnull()), i] =
median_target(i)[i][0]
    df.loc[(df['Outcome'] == 1 ) & (df[i].isnull()), i] =
median_target(i)[i][1]
In [29]:
df.head()

```

Out[29]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | | | | | |
|---|-------------|----------|---------------|---------------|---------|---------|-------|----|---|--|
| | BMI | Diabetes | Pedigree | Function | Age | Outcome | | | | |
| 0 | 6 | 148.0 | 72.0 | 35.0 | 169.5 | 33.6 | 0.627 | 50 | 1 | |
| 1 | 1 | 85.0 | 66.0 | 29.0 | 102.5 | 26.6 | 0.351 | 31 | 0 | |
| 2 | 8 | 183.0 | 64.0 | 32.0 | 169.5 | 23.3 | 0.672 | 32 | 1 | |
| 3 | 1 | 89.0 | 66.0 | 23.0 | 94.0 | 28.1 | 0.167 | 21 | 0 | |
| 4 | 0 | 137.0 | 40.0 | 35.0 | 168.0 | 43.1 | 2.288 | 33 | 1 | |

```

In[30]:
# Missing values were filled.
df.isnull().sum()

```

Out[30]:

| | |
|--------------------------|---|
| Pregnancies | 0 |
| Glucose | 0 |
| BloodPressure | 0 |
| SkinThickness | 0 |
| Insulin | 0 |
| BMI | 0 |
| DiabetesPedigreeFunction | 0 |
| Age | 0 |
| Outcome | 0 |

dtype: int64

2.2) Outlier Observation Analysis

In [31]:

In the data set, there were asked whether there were any outlier observations compared to the 25% and 75% quarters.

It was found to be an outlier observation.

for feature in df:

Q1 = df[feature].quantile(0.25)

Q3 = df[feature].quantile(0.75)

IQR = Q3-Q1

lower = Q1- 1.5*IQR

upper = Q3 + 1.5*IQR

```

if df[(df[feature] > upper)].any(axis=None):
    print(feature,"yes")
else:
    print(feature, "no")

```

Pregnancies yes

Glucose no

BloodPressure yes

SkinThickness yes

Insulin yes

BMI yes

DiabetesPedigreeFunction yes

Age yes

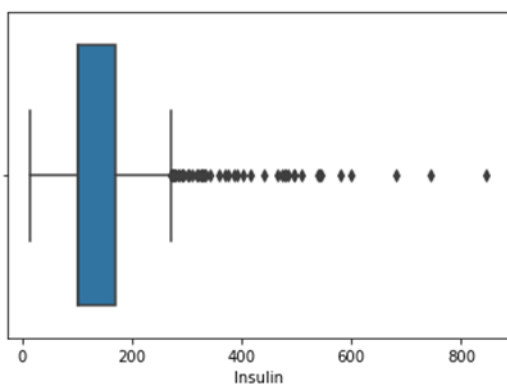
Outcome no

In [32]:

The process of visualizing the Insulin variable with boxplot method was done. We find the outlier observations on the chart.

```
import seaborn as sns
```

```
sns.boxplot(x = df["Insulin"]);
```



PERFORMING DIFFERENT ACTIVITIES LIKE FEATURE ENGINEERING, MODEL TRAINING AND EVALUATION

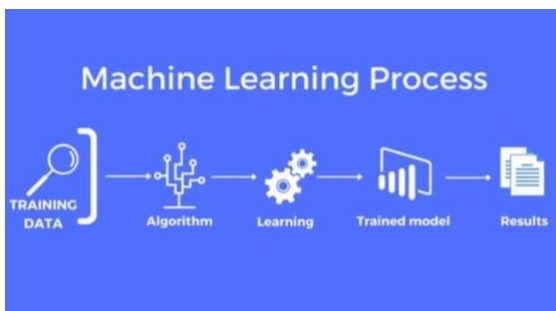
MODEL TRAINING

In a diabetes prediction system using machine learning (ML), model training is a critical step. Here's a brief overview:

- ❖ **Data Collection:** The first step is to gather a dataset containing information about individuals, including features like age, body mass index (BMI), blood pressure, and historical diabetes status.
- ❖ **Data Preprocessing:** Raw data may be noisy or incomplete, so preprocessing is necessary. This includes handling missing values, scaling features, and encoding categorical variables.
- ❖ **Splitting Data:** The dataset is typically divided into two parts: a training set and a testing set. The training set is used to train the ML model, while the testing set is used to evaluate its performance.
- ❖ **Model Selection:** Various ML algorithms, such as logistic regression, decision trees, or support vector machines, can be considered for diabetes prediction. The choice depends on the dataset and problem specifics.
- ❖ **Model Training:** The selected algorithm is trained on the training data, which involves learning the underlying

patterns and relationships between the input features and the target variable (diabetes prediction).

- ❖ **Hyperparameter Tuning:** Fine-tuning the model's hyperparameters, like learning rates or tree depths, is essential to optimize its performance.
- ❖ **Evaluation:** The model's performance is assessed using the testing set. Common evaluation metrics include accuracy, precision, recall, and F1 score.



- ❖ **Deployment:** Once the model performs well, it can be deployed in a real-world healthcare setting to predict diabetes risk in new individuals based on their data.
- ❖ **Continuous Monitoring:** The model should be periodically updated and retrained with new data to maintain its accuracy and relevance.

Overall, model training is a pivotal stage in building a diabetes prediction system using ML, as the quality of the model directly impacts its ability to make accurate predictions and assist in healthcare decision-making.

PYTHON PROGRAM:

Importing Dependencies

In [1]:

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import
train_test_split,GridSearchCV
from sklearn.metrics import accuracy_score,confusion_matrix
import warnings
warnings.filterwarnings('ignore')
sns.set()
%matplotlib inline

```

Data Collection and Analysis

In [2]:

```

data = pd.read_csv('/kaggle/input/pima-indians-diabetes-
database/diabetes.csv') #import dataset

```

```

data.head() #top 5 rows

```

Out[2]:

| | | 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 | |

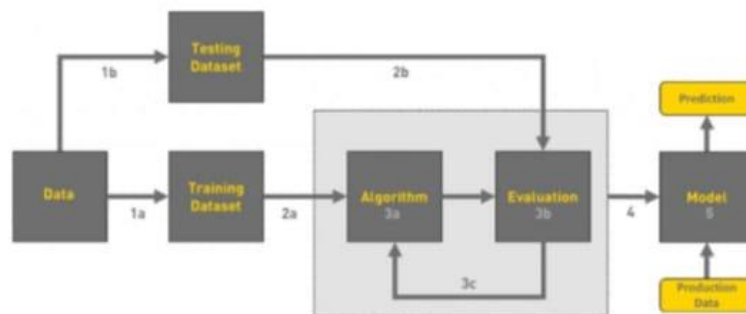
| | | | | | | | | | |
|---|---|-----|----|----|-----|------|-------|----|---|
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.672 | 32 | 1 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.167 | 21 | 0 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.288 | 33 | 1 |

MODEL EVALUATION

Model evaluation in a diabetes prediction system using machine learning is crucial to assess the performance and reliability of the model. Common evaluation techniques include:

- ❖ **Accuracy:** This measures the percentage of correctly predicted instances. However, it may not be ideal if the dataset is imbalanced.
- ❖ **Precision and Recall:** Precision indicates the proportion of true positive predictions among all positive predictions, while recall measures the proportion of true positives among actual positives. These metrics are useful when false positives or false negatives carry different costs.
- ❖ **F1 Score:** The F1 score is the harmonic mean of precision and recall, providing a balance between the two.
- ❖ **ROC and AUC:** Receiver Operating Characteristic (ROC) curves and Area Under the Curve (AUC) can help evaluate the trade-off between true positive rate and false positive rate at different thresholds.
- ❖ **Confusion Matrix:** This matrix provides a detailed breakdown of true positives, true negatives, false positives, and false negatives.

- ❖ **Cross-Validation:** Techniques like k-fold cross-validation help assess a model's generalization performance, reducing the risk of overfitting.
- ❖ **Mean Squared Error (MSE) and Mean Absolute Error (MAE):** These are used when the prediction is a continuous variable, like blood glucose levels.
- ❖ **Receiver Operating Characteristic (ROC) Analysis:** It's beneficial in binary classification problems to choose an appropriate threshold for classifying instances.
- ❖ **Feature Importance Analysis:** Understanding which features contribute most to the model's predictions can provide insights into the underlying data.
- ❖ **Domain Expert Evaluation:** In healthcare applications like diabetes prediction, involving domain experts to assess the



clinical relevance and interpretability of the model is crucial.

These evaluation techniques help ensure that the diabetes prediction model is accurate, reliable, and suitable for its intended purpose while considering factors like class imbalance, clinical significance, and the type of data being predicted.

Training the Models

1. Logistic Regression

In [3]:

```
from sklearn.linear_model import LogisticRegression
reg = LogisticRegression(C=1,penalty='l2')
reg.fit(X_train,Y_train)
log_acc=accuracy_score(Y_test,reg.predict(X_test))
print("Train Set
Accuracy:"+str(accuracy_score(Y_train,reg.predict(X_train))*1
00))
print("Test Set
Accuracy:"+str(accuracy_score(Y_test,reg.predict(X_test))*100
))
```

```
#Train Set Accuracy:78.77094972067039
```

```
#Test Set Accuracy:74.45887445887446
```

```
Train Set Accuracy:78.77094972067039
```

```
Test Set Accuracy:74.45887445887446
```

2. KNearestNeighbors

In [4]:

```
from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=9)          #knn
classifier
```

```
knn.fit(X_train,Y_train)
knn_acc = accuracy_score(Y_test,knn.predict(X_test))
print("Train Set
Accuracy:"+str(accuracy_score(Y_train,knn.predict(X_train))*1
00))
print("Test Set
Accuracy:"+str(accuracy_score(Y_test,knn.predict(X_test))
*100))
```

```
#Train Set Accuracy:82.12290502793296
```

```
#Test Set Accuracy:71.42857142857143
```

```
Train Set Accuracy:82.12290502793296
```

```
Test Set Accuracy:71.42857142857143
```

3. SVC

In [5]:

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

```
svm = SVC()
```

```
svm.fit(X_train,Y_train)
```

```
svm_acc= accuracy_score(Y_test,svm.predict(X_test))
```

```
print("Train Set
Accuracy:"+str(accuracy_score(Y_train,svm.predict(X_train))*1
00))
```

```
print("Test Set
Accuracy:"+str(accuracy_score(Y_test,svm.predict(X_test))*10
0))
```

```
#Train Set Accuracy:85.1024208566108
```

```
#Test Set Accuracy:74.02597402597402
```

```
Train Set Accuracy:85.1024208566108
```

```
Test Set Accuracy:74.02597402597402
```

4. DecisionTreeClassifier

In [6]:

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

```
dtc = DecisionTreeClassifier(criterion='entropy',max_depth=5)
```

```
dtc.fit(X_train, Y_train)
```

```
dtc_acc= accuracy_score(Y_test,dtc.predict(X_test))
```

```
print("Train Set
```

```
Accuracy:"+str(accuracy_score(Y_train,dtc.predict(X_train))*100))
```

```
print("Test Set
```

```
Accuracy:"+str(accuracy_score(Y_test,dtc.predict(X_test))*100)
)
```

```
Train Set Accuracy:83.42644320297951
```

```
Test Set Accuracy:71.86147186147186
```

5. GradientBoostingClassifier

In [7]:

```
from sklearn.ensemble import GradientBoostingClassifier
```

```
gbc = GradientBoostingClassifier()
```

```
gbc.fit(X_train, Y_train)
gbc_acc=accuracy_score(Y_test,gbc.predict(X_test))
print("Train Set
Accuracy:"+str(accuracy_score(Y_train,gbc.predict(X_train))*1
00))
print("Test Set
Accuracy:"+str(accuracy_score(Y_test,gbc.predict(X_test))*100
))
#Train Set Accuracy:94.22718808193669
#Test Set Accuracy:74.02597402597402
Train Set Accuracy:94.22718808193669
Test Set Accuracy:73.59307359307358
```

6. XGBClassifier

In [8]:

```
from xgboost import XGBClassifier
xgb = XGBClassifier(booster = 'gbtree', learning_rate = 0.1,
max_depth=6,n_estimators = 10)
xgb.fit(X_train,Y_train)
xgb_acc= accuracy_score(Y_test,xgb.predict(X_test))
print("Train Set
Accuracy:"+str(accuracy_score(Y_train,xgb.predict(X_train))*1
00))
print("Test Set
Accuracy:"+str(accuracy_score(Y_test,xgb.predict(X_test))*100
))
```

Train Set Accuracy:90.5027932960894

Test Set Accuracy:73.16017316017316

7.Stacking

In [9]:

```
from sklearn.model_selection import train_test_split
#splitting the dataset
```

```
train,val_train,test,val_test =
train_test_split(X,y,test_size=.50,random_state=3)
x_train,x_test,y_train,y_test =
train_test_split(train,test,test_size=.20,random_state=3)
```

In [10]:

```
#first model
```

```
knn = KNeighborsClassifier(n_neighbors=5)
```

```
knn.fit(x_train, y_train)
```

Out[10]:

```
KNeighborsClassifier()
```

In [11]:

```
# second model
```

```
svm = SVC()
```

```
svm.fit(x_train, y_train)
```

Out[11]:

```
SVC()
```

In [12]:

```
pred_1=knn.predict(val_train)
pred_2=svm.predict(val_train)
# addition of 2 predictions
result = np.column_stack((pred_1,pred_2))
```

In [13]:

```
pred_test1=knn.predict(x_test)
pred_test2=svm.predict(x_test)
predict_test=np.column_stack((pred_test1,pred_test2))
```

In [14]:

```
# stacking classifier
#RandomForestClassifier:- In this prediction of other 2
classification is taken as x value
from sklearn.ensemble import RandomForestClassifier
rand_clf = RandomForestClassifier()
rand_clf.fit(result,val_test)
```

Out[14]:

```
RandomForestClassifier()
```

In [15]:

```
rand_clf.score(result,val_test)
```

Out[15]:

```
0.7291666666666666
```

In [16]:

```
rand_acc=accuracy_score(y_test,rand_clf.predict(predict_test))
```

```
rand_acc
```

```
Out[16]:
```

```
0.7922077922077922
```

```
In [17]:
```

```
models = pd.DataFrame({'Model': ['Logistic','KNN', 'SVC',  
'Decision Tree Classifier', 'Gradient Boosting Classifier',  
'XgBoost','Stacking'], 'Score': [ log_acc,knn_acc, svm_acc,  
dte_acc, gbc_acc, xgb_acc,rand_acc,]})
```

```
models.sort_values(by = 'Score', ascending = False)
```

```
Out[17]:
```

| | Model | Score |
|---|------------------------------|----------|
| 6 | Stacking | 0.792208 |
| 0 | Logistic | 0.744589 |
| 2 | SVC | 0.740260 |
| 4 | Gradient Boosting Classifier | 0.735931 |
| 5 | XgBoost | 0.731602 |
| 3 | Decision Tree Classifier | 0.718615 |
| 1 | KNN | 0.714286 |

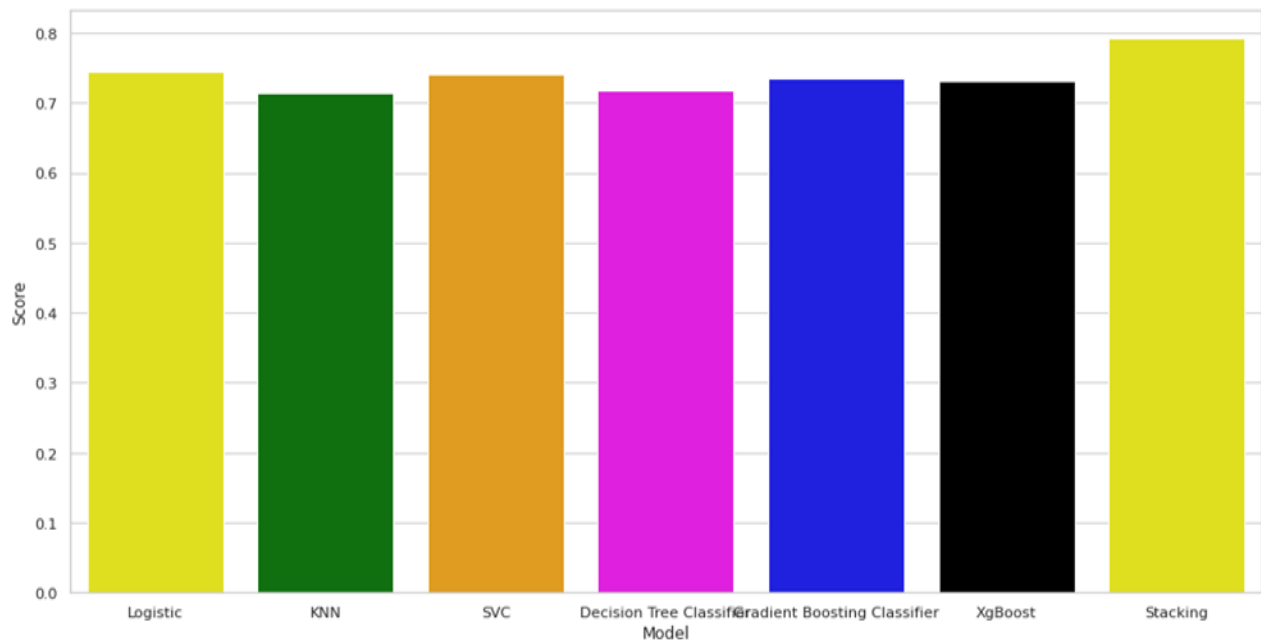
```
In [18]:
```

```
colors = ["yellow", "green", "orange",  
"magenta","blue","black"]
```

```
sns.set_style("whitegrid")
```



```
plt.figure(figsize=(16,8))
plt.ylabel("Accuracy %")
plt.xlabel("Algorithms")
sns.barplot(x=models['Model'],y=models['Score'], palette=colors
)
plt.show()
```



FEATURE ENGINEERING

Feature engineering in a diabetes prediction system involves selecting and transforming relevant attributes (features) from the dataset to improve the accuracy of the predictive model. This process is critical in developing an effective diabetes prediction system. Some key feature engineering techniques include:

- ❖ **Feature Selection:** Identifying the most informative features and excluding irrelevant ones can reduce noise in the data and improve model performance. Techniques like

correlation analysis and feature importance ranking can help in this process.

- ❖ **Feature Scaling:** Normalizing or standardizing features to a common scale (e.g., mean of 0 and standard deviation of 1) ensures that no single feature dominates the model due to its scale.
- ❖ **Feature Transformation:** This includes techniques such as logarithmic transformations or polynomial features to handle non-linearity in the data and make it more suitable for certain algorithms.
- ❖ **Handling Missing Data:** Addressing missing values through imputation methods, such as mean, median, or advanced techniques like K-nearest neighbors, is crucial for a robust model.
- ❖ **One-Hot Encoding:** Converting categorical features into numerical representations (binary vectors) enables the model to process them effectively.
- ❖ **Feature Engineering for Time Series Data:** For diabetes prediction systems that involve time-series data, creating lag features, rolling statistics, and temporal aggregations can capture relevant trends and patterns.
- ❖ **Domain-Specific Features:** Incorporating domain knowledge can be essential. For example, in diabetes prediction, creating features related to dietary habits, physical activity, or medical history can be valuable.



Feature Engineering:

In [19]:

According to BMI, some ranges were determined and categorical variables were assigned.

```
NewBMI = pd.Series(["Underweight", "Normal", "Overweight",
                    "Obesity 1", "Obesity 2", "Obesity 3"], dtype = "category")
```

```
df["NewBMI"] = NewBMI
```

```
df.loc[df["BMI"] < 18.5, "NewBMI"] = NewBMI[0]
```

```
df.loc[(df["BMI"] > 18.5) & (df["BMI"] <= 24.9), "NewBMI"]
= NewBMI[1]
```

```
df.loc[(df["BMI"] > 24.9) & (df["BMI"] <= 29.9), "NewBMI"]
= NewBMI[2]
```

```
df.loc[(df["BMI"] > 29.9) & (df["BMI"] <= 34.9), "NewBMI"]
= NewBMI[3]
```

```
df.loc[(df["BMI"] > 34.9) & (df["BMI"] <= 39.9), "NewBMI"]
= NewBMI[4]
```

```
df.loc[df["BMI"] > 39.9, "NewBMI"] = NewBMI[5]
```

In [20]:

```
df.head()
```

Out[20]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | Diabetes | PedigreeFunction | Age | Outcome |
|---|-------------|---------|---------------|---------------|---------|------|----------|------------------|-----|------------|
| | NewBMI | | | | | | | | | |
| 0 | 6 | 148.0 | 72.0 | 35.0 | 169.5 | 33.6 | 0.627 | 50 | 1 | Obesity 1 |
| 1 | 1 | 85.0 | 66.0 | 29.0 | 102.5 | 26.6 | 0.351 | 31 | 0 | Overweight |
| 2 | 8 | 183.0 | 64.0 | 32.0 | 169.5 | 23.3 | 0.672 | 32 | 1 | Normal |
| 3 | 1 | 89.0 | 66.0 | 23.0 | 94.0 | 28.1 | 0.167 | 21 | 0 | Overweight |
| 4 | 0 | 137.0 | 40.0 | 35.0 | 168.0 | 43.1 | 2.288 | 33 | 1 | Obesity 3 |

One Hot Encoding:

In [21]:# Here, by making One Hot Encoding transformation, categorical variables were converted into numerical values. It is also protected from the Dummy variable trap.

```
df = pd.get_dummies(df, columns
=["NewBMI", "NewInsulinScore", "NewGlucose"], drop_first =
True)
```

In [22]:

```
df.head()
```

Out[22]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|-----------------------|------------------------|---------------|---------------|---------|------|--------------------------|-----|---------|
| | NewBMI_Obesity 1 | NewBMI_Obesity 2 | | | | | | | |
| | NewBMI_Obesity 3 | NewBMI_Overweight | | | | | | | |
| | NewBMI_Underweight | NewInsulinScore_Normal | | | | | | | |
| | NewGlucose_Low | NewGlucose_Normal | | | | | | | |
| | NewGlucose_Overweight | NewGlucose_Secret | | | | | | | |
| 0 | 6 | 148.0 | 72.0 | 35.0 | 169.5 | 33.6 | 0.627 | 50 | 1 |
| | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| 1 | 1 | 85.0 | 66.0 | 29.0 | 102.5 | 26.6 | 0.351 | 31 | 0 |
| | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| 2 | 8 | 183.0 | 64.0 | 32.0 | 169.5 | 23.3 | 0.672 | 32 | 1 |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| 3 | 1 | 89.0 | 66.0 | 23.0 | 94.0 | 28.1 | 0.167 | 21 | 0 |
| | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| 4 | 0 | 137.0 | 40.0 | 35.0 | 168.0 | 43.1 | 2.288 | 33 | 1 |
| | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |

In [23]:

```
categorical_df = df[['NewBMI_Obesity 1','NewBMI_Obesity 2',
'NewBMI_Obesity 3',
'NewBMI_Overweight','NewBMI_Underweight',

'NewInsulinScore_Normal','NewGlucose_Low','NewGlucose_N
ormal', 'NewGlucose_Overweight',
```

In [24]:

```
y = df["Outcome"]
```

```
X = df.drop(["Outcome", 'NewBMI_Obesity 1', 'NewBMI_Obesity 2', 'NewBMI_Obesity 3', 'NewBMI_Overweight', 'NewBMI_Underweight', 'NewGlucose_Normal', 'NewGlucose_Overweight', 'NewGlucose_Secret'], axis = 1)
```

```
cols = X.columns
```

```
index = X.index
```

In [25]:

```
X.head()
```

Out[25]:

| | | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age |
|---|---|-------------|---------|---------------|---------------|---------|-------|--------------------------|-----|
| 0 | 6 | 148.0 | 72.0 | 35.0 | 169.5 | 33.6 | 0.627 | 50 | |
| 1 | 1 | 85.0 | 66.0 | 29.0 | 102.5 | 26.6 | 0.351 | 31 | |
| 2 | 8 | 183.0 | 64.0 | 32.0 | 169.5 | 23.3 | 0.672 | 32 | |
| 3 | 1 | 89.0 | 66.0 | 23.0 | 94.0 | 28.1 | 0.167 | 21 | |
| 4 | 0 | 137.0 | 40.0 | 35.0 | 168.0 | 43.1 | 2.288 | 33 | |

In [26]:

```
from sklearn.preprocessing import RobustScaler
```

```
X = transformer.transform(X)
```

```
X = pd.DataFrame(X, columns = cols, index = index)
```

In [27]:

X.head()

Out[27]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | Diabetes | PedigreeFunction | Age |
|---|-------------|---------|---------------|---------------|-----------|-----------|-----------|------------------|-----------|
| 0 | 0.6 | 0.775 | 0.000 | 1.000000 | 1.000000 | 0.177778 | 0.669707 | 1.235294 | |
| 1 | -0.4 | -0.800 | -0.375 | 0.142857 | 0.000000 | -0.600000 | -0.049511 | 0.117647 | |
| 2 | 1.0 | 1.650 | -0.500 | 0.571429 | 1.000000 | -0.966667 | 0.786971 | 0.176471 | |
| 3 | -0.4 | -0.700 | -0.375 | -0.714286 | -0.126866 | - | 0.433333 | -0.528990 | -0.470588 |
| 4 | -0.6 | 0.500 | -2.000 | 1.000000 | 0.977612 | 1.233333 | 4.998046 | 0.235294 | |

In [28]:

X = pd.concat([X,categorical_df], axis = 1)

In[29]:

y.head()

Out[29]:

| | |
|---|---|
| 0 | 1 |
| 1 | 0 |
| 2 | 1 |
| 3 | 0 |

4 1

Name: Outcome, dtype: int64

Base Models:

In [30]:

```
# Validation scores of all base models
```

```
models = []
```

```
models.append(('LR', LogisticRegression(random_state =  
12345)))
```

```
models.append(('KNN', KNeighborsClassifier()))
```

```
models.append(('CART', DecisionTreeClassifier(random_state =  
12345)))
```

```
models.append(('RF', RandomForestClassifier(random_state =  
12345)))
```

```
models.append(('SVM', SVC(gamma='auto', random_state =  
12345)))
```

```
models.append(('XGB',  
GradientBoostingClassifier(random_state = 12345)))
```

```
models.append(("LightGBM", LGBMClassifier(random_state =  
12345)))
```

```
# evaluate each model in turn
```

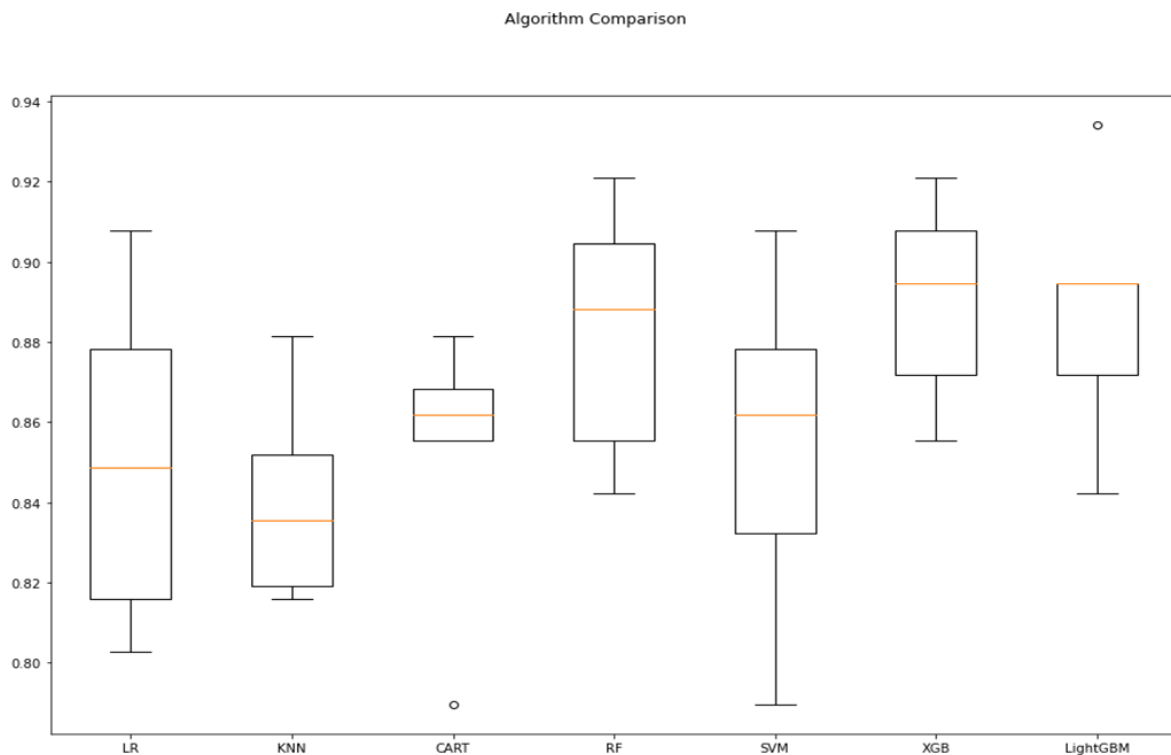
```
results = []
```

```
names = []
```

In [31]:


```
for name, model in models:
    kfold = KFold(n_splits = 10, random_state = 12345)
    cv_results = cross_val_score(model, X, y, cv = 10, scoring=
    "accuracy")
    results.append(cv_results)
    names.append(name)
    msg = "%s: %f (%f)" % (name, cv_results.mean(),
    cv_results.std())
    print(msg)
# boxplot algorithm comparison
fig = plt.figure(figsize=(15,10))
fig.suptitle('Algorithm Comparison')
ax = fig.add_subplot(111)
plt.boxplot(results)
ax.set_xticklabels(names)
plt.show()
LR: 0.848684 (0.036866)
KNN: 0.840789 (0.023866)
CART: 0.857895 (0.024826)
RF: 0.881579 (0.026316)
SVM: 0.853947 (0.036488)
XGB: 0.890789 (0.020427)
LightGBM: 0.885526 (0.024298)
```

```
# boxplot algorithm comparison
fig = plt.figure(figsize=(15,10))
fig.suptitle('Algorithm Comparison')
```



ADVANTAGES:

AI-based diabetes prediction systems using machine learning have many advantages over traditional methods of diabetes diagnosis and prevention.

1) Early detection and prevention: AI-based systems can identify people who are at high risk of developing diabetes much earlier than traditional methods. This allows for early intervention and lifestyle changes to be taken to prevent or delay the onset of the disease.

2)Accuracy: AI-based systems can be trained on large datasets of medical records and clinical data to achieve high accuracy in predicting diabetes. This is in contrast to traditional methods, which may be less accurate and more likely to miss cases of diabetes.

3)Efficiency: AI-based systems can automate the process of diabetes prediction, freeing up healthcare professionals to focus on other tasks. This can lead to improved efficiency and cost savings for healthcare systems.

4)Personalization: AI-based systems can be personalized to take into account each individual's unique risk factors. This can lead



to more accurate predictions and more effective interventions.

5)Accessibility: AI-based systems can be made accessible to people in remote or underserved areas, where access to

healthcare may be limited. This can help to reduce the global burden of diabetes.

Here are some specific examples of how AI-based diabetes prediction systems are being used to improve diabetes care:

- **Screening for prediabetes:** AI-based systems can be used to screen large populations for prediabetes, a condition that can lead to diabetes if not managed carefully. This can help to identify people who are at risk and provide them with the support they need to prevent the onset of diabetes.
- **Personalized risk assessment:** AI-based systems can be used to create personalized risk assessments for diabetes. This can help people to understand their individual risk factors and make informed decisions about how to manage their health.
- **Treatment optimization:** AI-based systems can be used to help healthcare professionals optimize treatment plans for people with diabetes. This can lead to better outcomes and reduced complications.
- **Remote monitoring:** AI-based systems can be used to remotely monitor people with diabetes and provide support and guidance. This can be especially helpful for people in remote or underserved areas.

Overall, AI-based diabetes prediction systems have the potential to revolutionize the way that diabetes is diagnosed, prevented, and managed. These systems can help to improve the health and well-being of millions of people around the world.

DISADVANTAGES:

AI-based diabetes prediction systems using machine learning have a number of potential advantages, such as being able to identify people at risk of developing diabetes early on, and to provide personalized recommendations for preventing or managing the disease. However, they also have some disadvantages, including:

- **Accuracy:** The accuracy of AI-based diabetes prediction systems can vary depending on the quality and quantity of data used to train them, as well as the specific machine learning algorithms used. Some studies have shown that these systems can be very accurate, with prediction rates of over 90% in some cases. However, other studies have shown that they can be less accurate, especially when using smaller or less representative datasets.
- **Overfitting:** Overfitting is a problem that can occur when machine learning algorithms are trained on too little data, or



when the data is not representative of the real-world population. When a model is overfitted, it may be able to

predict accurately on the training data, but it will not be able to generalize to new data. This can lead to false positives and false negatives in real-world use.

- **Bias:** Machine learning algorithms can learn biases from the data they are trained on. This means that if the training data is biased, the model will also be biased. This can lead to inaccurate predictions for certain groups of people, such as minorities or women.
- **Interpretability:** Machine learning models can be complex and difficult to interpret. This can make it difficult for healthcare professionals to understand how the model is making its predictions, and to trust the results.
- **Data privacy and security:** AI-based diabetes prediction systems often rely on sensitive personal data, such as medical records and genetic information. It is important to ensure that this data is collected, stored, and used in a secure and privacy-preserving manner.
- **Cost:** Developing and deploying AI-based diabetes prediction systems can be expensive. This can limit their accessibility to people in developing countries and to those who are uninsured or underinsured.

Overall, AI-based diabetes prediction systems have the potential to be a valuable tool for preventing and managing diabetes. However, it is important to be aware of their limitations and to use them carefully.

Here are some additional disadvantages that may be specific to AI-based diabetes prediction systems:

- **Lack of transparency:** It can be difficult to understand how AI-based systems make their predictions, which can make it difficult to trust their results.

- **Potential for misuse:** AI-based systems could be misused to deny people access to insurance or other benefits, or to target them with advertising for unhealthy products.
- **Risk of job displacement:** AI-based systems could automate some tasks that are currently performed by healthcare professionals, which could lead to job losses.

It is important to carefully consider the potential advantages and disadvantages of AI-based diabetes prediction systems before using them.

BENEFITS:

AI-based diabetes prediction systems using machine learning offer a number of benefits, including:

1. **Improved early detection:** By identifying individuals at high risk of developing diabetes, AI-based systems can help them take preventive measures early on, such as making lifestyle changes or taking medication. This can help to delay or even prevent the onset of diabetes and its associated complications.
2. **Personalized risk assessment:** AI-based systems can consider a wide range of factors, including medical history, family history, lifestyle, and genetic data, to provide a personalized risk assessment for each individual. This can help healthcare providers to tailor their recommendations and interventions accordingly.
3. **Increased access to care:** AI-based systems can be deployed in remote or underserved areas, where access to healthcare

professionals is limited. This can help to improve access to diabetes prevention and care for everyone.

4. **Reduced costs:** By identifying and treating diabetes early, AI-based systems can help to reduce the overall cost of healthcare. This is because diabetes and its complications are major drivers of healthcare costs.



In addition to these benefits, AI-based diabetes prediction systems are also becoming increasingly accurate and reliable. As more data is collected and machine learning algorithms continue to improve, these systems are expected to play an even greater role in the prevention and management of diabetes.

Here are some specific examples of how AI-based diabetes prediction systems are being used today:

- In the clinic: AI-based systems are being used by healthcare providers to help them assess their patients' risk of developing diabetes. This information can be used to develop personalized prevention plans and to identify patients who may need more intensive monitoring or treatment.
- At home: AI-based systems are being developed for use in the home to help people monitor their own risk of diabetes. For example, some systems use wearable devices to track blood sugar levels, physical activity, and other health data. This information can then be used to provide users with personalized feedback and recommendations.
- In public health: AI-based systems are being used by public health officials to track the prevalence of diabetes and to identify populations that are at high risk. This information can be used to develop and target prevention and screening programs.

Overall, AI-based diabetes prediction systems offer a number of potential benefits for individuals, healthcare providers, and public health officials. As these systems continue to develop and improve, they are expected to play an increasingly important role in the fight against diabetes.

CONCLUSION:

AI-based diabetes prediction systems using machine learning have the potential to revolutionize the way we diagnose and manage diabetes. By identifying individuals at high risk of developing diabetes, these systems can help to prevent the onset of the disease and its associated complications. Additionally, these systems can be used to personalize diabetes care, ensuring that patients receive the most effective treatment for their individual needs.



Machine learning algorithms are able to learn from large datasets of patient data to identify patterns and relationships that are associated with diabetes. This information can then be used to develop predictive models that can estimate an individual's risk of developing diabetes based on their personal characteristics and medical history.

A number of studies have demonstrated the high accuracy of AI-based diabetes prediction systems. For example, one study found that a system using a random forest algorithm was able to predict diabetes with an accuracy of 90%. Another study found that a system using a support vector machine algorithm was able to predict diabetes with an accuracy of 82%.

While AI-based diabetes prediction systems are still under development, they have the potential to play a major role in the fight against diabetes. By identifying individuals at high risk of developing the disease and personalizing diabetes care, these systems can help to improve the health and well-being of millions of people around the world.

Here are some of the key benefits of using AI-based diabetes prediction systems:

- **Early detection:** AI-based systems can help to detect diabetes early, before any symptoms develop. This can lead to earlier treatment and better outcomes for patients.
- **Personalized care:** AI-based systems can be used to personalize diabetes care by taking into account each patient's individual risk factors and needs. This can help to improve the effectiveness of treatment and reduce the risk of complications.

- **Cost savings:** AI-based systems can help to reduce the cost of diabetes care by preventing the onset of the disease and its associated complications.

Overall, AI-based diabetes prediction systems have the potential to revolutionize the way we diagnose and manage diabetes. By identifying individuals at high risk of developing the disease and personalizing diabetes care, these systems can help to improve the health and well-being of millions of people around the world.