

# **DIABETES CLASSIFICATION AND PREDICTION MODEL** **USING MACHINE LEARNING**

*Dissertation submitted in fulfilment of the requirements for the Degree of*

**BACHELOR OF TECHNOLOGY**

**in**

**COMPUTER SCIENCE AND ENGINEERING**

by

**AKSHAT KESHARI**

**12107597**

Supervisor

**SAJJAD MANZOOR MIR**



**L** OVELY  
**P** ROFESSIONAL  
**U** NIVERSITY

**School of Computer Science and Engineering**

Lovely Professional University

Phagwara, Punjab (India)

## DECLARATION STATEMENT

---

I hereby declare that the research work reported in the dissertation/dissertation proposal entitled **“Diabetes Classification and Prediction Model”** in partial fulfilment of the requirement for the award of Degree for Bachelor of Technology in Computer Science and Engineering at Lovely Professional University, Phagwara, Punjab is an authentic work carried out under supervision of my research supervisor Mr. Sajjad Mansoor Mir. I have not submitted this work elsewhere for any degree or diploma.

I understand that the work presented herewith is in direct compliance with Lovely Professional University’s Policy on plagiarism, intellectual property rights, and highest standards of moral and ethical conduct. Therefore, to the best of my knowledge, the content of this dissertation represents authentic and honest research effort conducted, in its entirety, by me. I am fully responsible for the contents of my dissertation work.

**Name:** Akshat keshari

**Reg. No:** 12107597

.

## TABLE OF CONTENTS

<b>S. No.</b>	<b>TITLE</b>	<b>Pg. No</b>
1.	<b>Acknowledgment</b>	4
2.	<b>Abstract</b>	5
3.	<b>Objective</b>	6
4.	<b>Introduction</b>	7
5.	<b>Theoretical Background</b>	8
6.	<b>Methodology</b>	12
7.	<b>Snapshot of project</b>	14
8.	<b>Model Training</b>	20
9.	<b>Results</b>	27
10.	<b>Conclusion</b>	28
11.	<b>Bibliography</b>	29

# ACKNOWLEDGMENT

I would like to express my sincere gratitude to all those who have contributed to the completion of this project. First and foremost, I extend my heartfelt thanks to Mr. Sajjad Manzoor Mir, for their invaluable guidance, support, and encouragement throughout the duration of this project. Their expertise and constructive feedback have been instrumental in shaping the direction and outcomes of this work.

I am also indebted to Lovely Professional University for providing the necessary resources and facilities that facilitated the smooth progress of this project. The cooperation and assistance from the staff and faculty members have been deeply appreciated.

Furthermore, I extend my appreciation to my colleagues for their collaboration and assistance in various phases of the project. Their dedication and teamwork have significantly contributed to the achievement of our goals.

Last but not least, I would like to express my gratitude to my family and friends for their unwavering support, understanding, and encouragement throughout this endeavor.

This project would not have been possible without the collective effort, support, and encouragement from all those mentioned above.

Thank you.

**Date:** 24<sup>th</sup> September, 2024

# ABSTRACT

Diabetes is a chronic and potentially debilitating condition that affects millions globally. Early diagnosis is critical to managing the disease effectively and preventing severe complications such as cardiovascular disease, kidney failure, and neuropathy. However, traditional diagnostic methods can be time-consuming and resource-intensive, often requiring significant expertise. This project explores the application of machine learning techniques for developing a robust and accurate diabetes classification and prediction model, aiming to support early diagnosis and timely intervention.

The study begins by utilizing a comprehensive dataset containing key health metrics such as glucose levels, blood pressure, body mass index (BMI), and age. Rigorous data preprocessing is performed to ensure data quality, including handling missing values, outliers, and standardizing features for better model performance. Feature selection techniques are employed to identify the most influential factors contributing to diabetes diagnosis, thereby enhancing the interpretability and efficiency of the model.

Multiple supervised learning algorithms, including Logistic Regression, Decision Trees, Random Forest, and Support Vector Machines (SVM), are developed and evaluated. The models are fine-tuned using hyperparameter optimization techniques to maximize predictive accuracy and minimize errors. To ensure robustness and generalizability, cross-validation and comprehensive performance evaluation metrics, including accuracy, precision, recall, F1-score, and ROC-AUC, are employed.

The results indicate that machine learning models can achieve high accuracy in classifying and predicting diabetes, with the best-performing model demonstrating strong potential for practical healthcare applications. The findings also underscore the importance of data quality and feature selection in building effective predictive models.

This report discusses the methodology, challenges, and outcomes in detail, emphasizing the model's scalability and utility in real-world healthcare settings. Future directions for this work include refining the model with larger, more diverse datasets, integrating it with real-time data for dynamic predictions, and addressing ethical considerations such as bias and data privacy. The project highlights the transformative potential of machine learning in improving early diagnosis and advancing personalized healthcare solutions for diabetes management. .

# OBJECTIVE

The primary objective of this project is to design, implement, and evaluate a machine learning model capable of accurately classifying and predicting diabetes in individuals based on a set of health-related parameters. The focus is to explore how machine learning techniques can support early diagnosis, enabling timely medical intervention and management of the condition.

To achieve this, the project is structured around the following detailed objectives:

## **Data Collection and Exploration:**

- Utilize an appropriate dataset containing health metrics such as glucose levels, blood pressure, BMI, age, and other relevant factors that influence diabetes diagnosis.
- Perform exploratory data analysis (EDA) to understand the dataset's structure, identify patterns, and detect missing or anomalous values.

## **Data Preprocessing and Feature Engineering:**

- Handle missing data and outliers to ensure data consistency and quality.
- Normalize or standardize numerical features to improve model performance.
- Perform feature selection to identify the most relevant variables that contribute to diabetes classification, thereby reducing model complexity.

## **Model Development:**

- Train multiple machine learning algorithms (e.g., Logistic Regression, Decision Trees, Random Forest, Support Vector Machines, etc.) to identify the best-performing model.
- Fine-tune hyperparameters using techniques like grid search or random search to optimize model performance.

## **Model Evaluation and Validation:**

- Evaluate the models using appropriate performance metrics, including accuracy, precision, recall, F1-score, and ROC-AUC.
- Perform cross-validation to ensure the robustness and generalizability of the model across different subsets of data.

## **Comparative Analysis:**

- Compare the performance of different machine learning models to determine which algorithm provides the best trade-off between accuracy and computational efficiency.

**Deployment Readiness:**

- Develop a framework for integrating the model into practical applications, such as clinical decision support systems.
- Address ethical considerations, including patient data privacy and the risk of bias in the model's predictions.

**Future Enhancements:**

- Explore the potential for incorporating deep learning models or ensemble methods to improve predictive accuracy further.
- Investigate the integration of real-time data for dynamic prediction capabilities.

By fulfilling these objectives, this project aims to demonstrate the utility of machine learning as a powerful tool in healthcare, specifically in the early detection and management of diabetes. The outcomes could pave the way for scalable, efficient, and accessible solutions for healthcare practitioners and patients.

# INTRODUCTION

Diabetes is a growing global health challenge, characterized by chronic hyperglycemia due to impaired insulin production or function. It is a major risk factor for severe complications such as cardiovascular diseases, kidney failure, and neuropathy, affecting the quality of life for millions. Early diagnosis and effective management are crucial for mitigating these risks and ensuring better health outcomes. However, traditional diagnostic methods can be resource-intensive and time-consuming, particularly in regions with limited access to healthcare infrastructure.

With advancements in data science, machine learning (ML) has emerged as a promising solution for automating and enhancing diagnostic accuracy in healthcare. Machine learning algorithms, when trained on comprehensive datasets, can identify complex patterns and relationships among health indicators that may not be easily discernible through conventional statistical methods. These capabilities make machine learning an invaluable tool in predicting and classifying chronic diseases like diabetes.

This project focuses on developing a machine learning model for diabetes classification and prediction using patient health metrics. By analyzing key features such as glucose levels, blood pressure, BMI, age, and other medical parameters, the model aims to determine whether an individual is diabetic or non-diabetic.

The project leverages supervised learning algorithms to process and analyze health data, emphasizing performance metrics such as accuracy, precision, recall, and F1-score. Through rigorous preprocessing, feature engineering, and hyperparameter optimization, the model seeks to deliver high predictive accuracy and reliability.

The broader goal of this initiative is to demonstrate the potential of machine learning in supporting early and accurate diagnosis of diabetes, thereby facilitating timely intervention. Additionally, it highlights the feasibility of integrating such models into real-world healthcare systems for scalable and accessible solutions, especially in resource-limited settings.

The report discusses the project in detail, covering aspects such as dataset characteristics, methodology, model development, evaluation, and potential applications. It also identifies challenges and areas for future improvement, underscoring the need for continual refinement and validation in deploying machine learning models for clinical use.



# THEORETICAL BACKGROUND

Machine learning (ML) is a branch of artificial intelligence (AI) that enables systems to learn from data and make predictions or decisions without explicit programming. It has become an essential tool in healthcare, offering robust methods for analyzing complex datasets and identifying patterns critical for diagnosis, prognosis, and treatment planning.

## 1. Overview of Diabetes

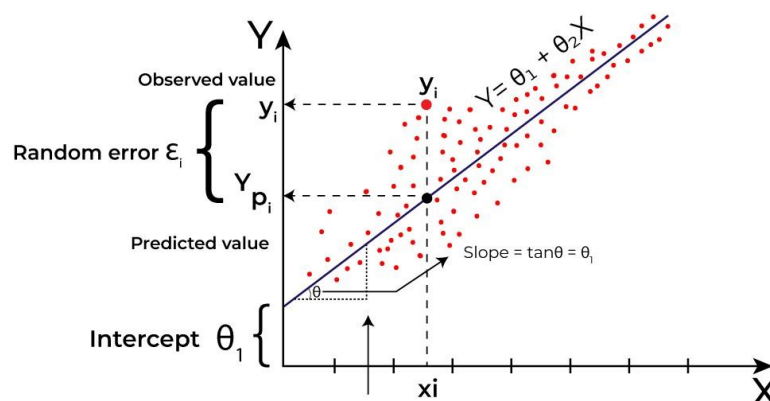
Diabetes is a metabolic disorder primarily categorized into Type 1 and Type 2 diabetes, with Type 2 being the most prevalent. It is characterized by elevated blood glucose levels, which can result from inadequate insulin production, insulin resistance, or both. Early identification of at-risk individuals is vital to manage the disease and prevent complications. Various clinical parameters such as glucose levels, blood pressure, body mass index (BMI), age, and family history are used to assess diabetes risk.

## 2. Supervised Learning for Classification

Supervised learning, a subset of ML, involves training algorithms on labeled data to make predictions about unseen data. In the context of diabetes classification, supervised learning models are trained to classify patients as diabetic or non-diabetic based on input features. Common supervised learning algorithms include:

- **Linear Regression:**

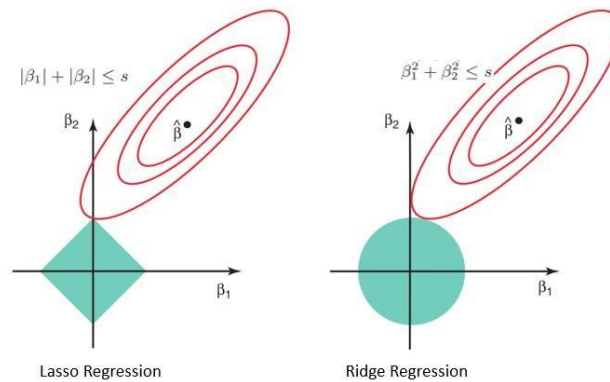
Linear regression is a fundamental algorithm used for predicting continuous outcomes based on one or more predictor variables. It assumes a linear relationship between the input features, such as age, health status, and previous medical expenses, and the target variable, which in this case is the predicted insurance cost. Linear regression serves as a foundational model, providing a baseline for assessing more complex methodologies.



- **Ridge and Lasso Regression:**

Ridge and Lasso regression are extensions of linear regression that incorporate regularization techniques to prevent overfitting. Ridge regression applies L2 regularization, penalizing the magnitude of coefficients and helping maintain model complexity while improving generalization. Lasso regression employs L1 regularization, which can shrink some

coefficients to zero, effectively performing feature selection. This is particularly beneficial in medical insurance cost prediction, where many variables may be irrelevant, enhancing model interpretability and focusing on key cost drivers.

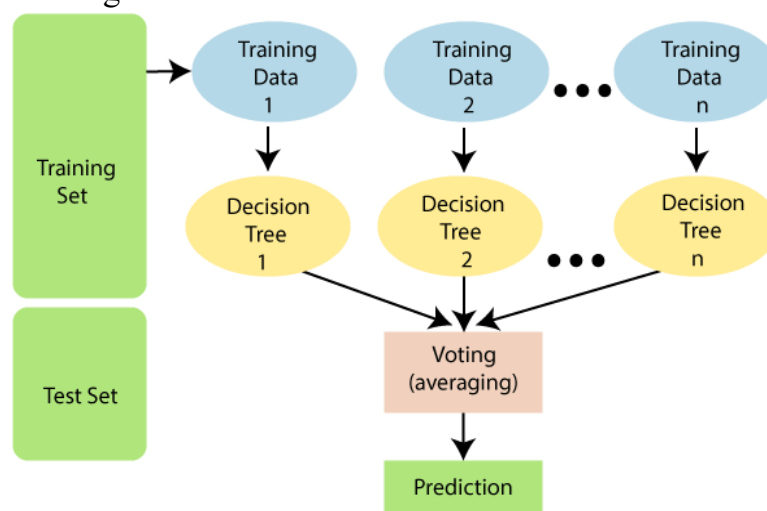


### • **Decision Trees:**

Decision trees are non-parametric supervised learning algorithms that recursively partition the feature space into segments based on feature values. They are intuitive and easy to interpret, making them a valuable tool for stakeholders in healthcare who need to understand the factors driving costs. By capturing complex interactions among predictors, decision trees can effectively model the relationships between patient demographics, health conditions, and insurance expenses.

### • **Random Forests:**

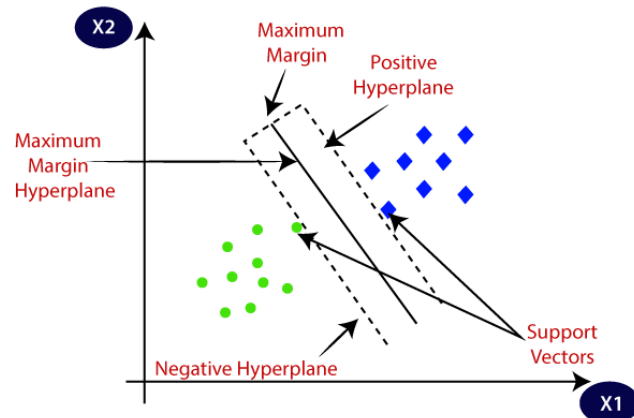
Random forests enhance the performance of individual decision trees by aggregating predictions from multiple trees trained on different subsets of the data. This ensemble method reduces overfitting and increases robustness, making it highly scalable and capable of handling highdimensional data. In medical insurance cost prediction, random forests can effectively manage the intricate interplay of various factors influencing costs, providing more reliable predictions than single models.



### • **Support Vector Machines (SVM):**

Support Vector Machines (SVM) are powerful supervised learning algorithms that create hyperplanes in high-dimensional space to separate data points into different classes. In the context of cost prediction, SVM can be utilized to identify complex relationships between

predictors and costs, especially when the data exhibits non-linear patterns. By maximizing the margin between classes, SVM enhances predictive accuracy, which is crucial for effective risk assessment in insurance.



- **Support Vector Regression (SVR):**

Support Vector Regression extends the SVM framework to regression tasks, focusing on predicting continuous outcomes. SVR constructs a hyperplane that fits the data within a specified margin of tolerance, minimizing prediction error while maintaining model simplicity. In medical insurance cost prediction, SVR can be particularly useful for modeling nuanced relationships between patient characteristics and their associated costs, providing valuable insights into expenditure patterns.

These theoretical foundations create a comprehensive framework for understanding and implementing medical insurance cost prediction models. By leveraging these algorithms, researchers and practitioners can develop effective predictive systems that identify key cost drivers, ultimately improving decision-making processes in healthcare management and insurance pricing strategies. The insights derived from these models can enhance financial planning and contribute to more sustainable healthcare systems.

### 3. Feature Selection and Preprocessing

The quality and relevance of input features significantly affect model performance. Key preprocessing steps include:

- **Handling Missing Data:** Techniques like imputation are used to replace missing values.
- **Normalization/Standardization:** Ensures numerical features have comparable scales to prevent bias in model training.
- **Feature Engineering:** The process of creating new features or selecting the most relevant ones to enhance model accuracy.

### 4. Evaluation Metrics

Evaluating the performance of a classification model involves using metrics such as:

- **Accuracy:** The proportion of correctly classified instances.

- **Precision:** The ratio of true positive predictions to all positive predictions, measuring the reliability of positive classifications.
- **Recall (Sensitivity):** The ratio of true positive predictions to all actual positive cases, indicating the model's ability to detect positives.
- **F1-Score:** The harmonic mean of precision and recall, providing a balanced performance measure.
- **ROC-AUC (Receiver Operating Characteristic - Area Under Curve):** Measures the model's ability to distinguish between classes.

## 5. Applications in Healthcare

Machine learning models have shown great promise in healthcare applications, including:

- Predictive analytics for disease risk assessment.
- Personalized treatment recommendations based on patient profiles.
- Automated systems for real-time health monitoring and alerts.

## 6. Challenges and Ethical Considerations

Despite its advantages, deploying ML models in healthcare poses challenges such as data privacy, bias in predictions due to imbalanced datasets, and interpretability of complex models. Addressing these concerns is critical to ensuring the ethical use of ML in clinical settings.

This theoretical foundation underscores the role of machine learning in addressing the complexities of diabetes diagnosis and prediction, providing a strong basis for the practical implementation of the project.

# **Hardware and Software Requirements**

Hardware and software requirements for developing and deploying machine learning models for predicting heart health outcomes can vary depending on the complexity of the models, the size of the dataset, and the specific requirements of the deployment environment. However, here's a general overview of the hardware and software requirements:

## **Compute Resources:**

Modern CPUs or GPUs are typically required for training complex machine learning models, especially deep learning models.

For larger datasets or computationally intensive algorithms, access to high-performance computing (HPC) resources or cloud-based services may be necessary.

Memory (RAM) requirements depend on the size of the dataset and the complexity of the model. Having sufficient RAM is essential to avoid memory-related bottlenecks during training and inference.

## **Storage:**

Sufficient storage space is needed to store the dataset, intermediate model files, and any additional resources required for training and deployment.

Solid-state drives (SSDs) are preferred over traditional hard disk drives (HDDs) for faster data access and processing.

## **Networking:**

Stable internet connectivity may be required for accessing cloud-based resources, downloading datasets, or deploying models to remote servers.

## ***Software Requirements:***

### **Programming Languages:**

Python is the most commonly used programming language for machine learning due to its extensive libraries and frameworks (e.g., TensorFlow, PyTorch, scikit-learn).

R is another popular language for statistical modeling and data analysis, although it's less commonly used for deep learning tasks.

### **Machine Learning Libraries/Frameworks:**

scikit-learn is a versatile library for traditional machine learning tasks such as classification, regression, clustering, and dimensionality reduction.

### **Development Tools:**

Integrated Development Environments (IDEs) such as PyCharm, Jupyter Notebook, or Visual Studio Code are commonly used for writing and debugging machine learning code.

Version control systems like Git/GitHub are essential for collaboration, code management, and reproducibility.

# METHODOLOGY

The methodology for developing the Diabetes Classification and Prediction Model is structured into a series of well-defined steps, ensuring systematic implementation and evaluation of machine learning techniques. These steps include data collection, preprocessing, model development, evaluation, and validation.

## 1. Dataset Selection

The project begins by selecting a publicly available dataset, such as the Pima Indians Diabetes Dataset, which contains relevant health metrics like:

- Glucose levels
- Blood pressure
- Body mass index (BMI)
- Age
- Skin thickness
- Insulin levels

These features are labeled with a binary outcome indicating the presence or absence of diabetes, making the dataset suitable for supervised classification.

## 2. Exploratory Data Analysis (EDA)

- Analyze the dataset to understand feature distributions, identify patterns, and detect missing or anomalous values.
- Visualize relationships between features and the target variable using correlation heatmaps, scatter plots, and box plots.
- Assess feature importance and interactions to guide model development.

## 3. Data Preprocessing

- **Handling Missing Values:** Missing data is imputed using statistical methods (e.g., mean or median imputation) or advanced techniques like k-nearest neighbors (k-NN) imputation.
- **Outlier Detection and Removal:** Outliers are identified using methods like Z-scores or the Interquartile Range (IQR) and handled appropriately.
- **Feature Scaling:** Numerical features are normalized or standardized to ensure consistent ranges, improving model convergence and accuracy.
- **Feature Engineering:** New features are created or existing ones transformed to enhance predictive power. Feature selection techniques such as Recursive Feature Elimination (RFE) are applied to reduce dimensionality.

## 4. Model Development

Several supervised machine learning algorithms are trained to classify individuals as diabetic or non-diabetic:

- **Logistic Regression:** A baseline model for binary classification, used to establish a reference for performance comparison.
- **Decision Trees:** A tree-based model that splits data into decision nodes, providing interpretability.
- **Random Forest:** An ensemble learning method combining multiple decision trees to enhance prediction accuracy and reduce overfitting.
- **Support Vector Machines (SVM):** A powerful algorithm that identifies the optimal hyperplane for separating classes.
- **k-Nearest Neighbors (k-NN):** A simple algorithm that classifies data points based on their nearest neighbors.

## 5. Model Evaluation

To evaluate and compare model performance, the following metrics are used:

- **Accuracy:** Proportion of correctly classified instances.
- **Precision:** Measure of correctly predicted positive cases out of all predicted positives.
- **Recall (Sensitivity):** Ability to detect actual positive cases.
- **F1-Score:** Harmonic mean of precision and recall, balancing both metrics.
- **ROC-AUC:** Area under the Receiver Operating Characteristic curve, indicating the model's ability to distinguish between classes.

Evaluation is performed using cross-validation to ensure robustness and reduce the likelihood of overfitting.

## 6. Hyperparameter Tuning

- Hyperparameters for each model are optimized using grid search or random search techniques.
- Tuning ensures the best combination of parameters, enhancing model performance.

## 7. Comparative Analysis

Models are compared based on their performance metrics, and the best-performing algorithm is selected for final deployment.

## 8. Deployment Framework

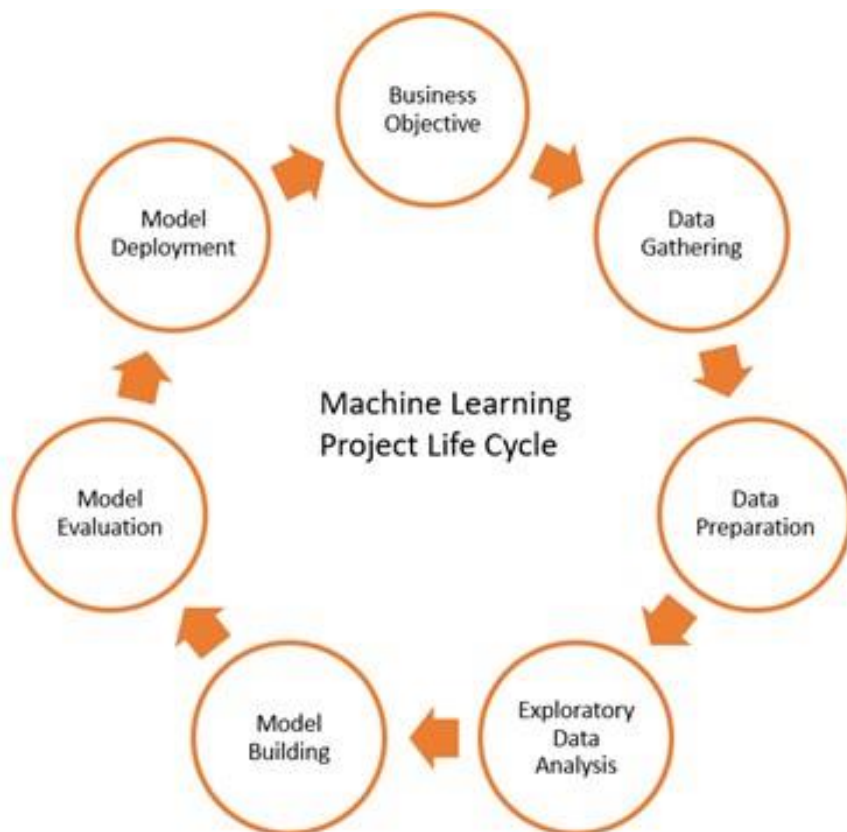
- A framework is outlined for integrating the model into real-world applications, such as clinical decision support systems.
- Deployment considerations include interpretability, computational efficiency, and user accessibility.



## 9. Future Enhancements

- Explore advanced models like Gradient Boosting (e.g., XGBoost) or deep learning techniques for improved performance.
- Use larger and more diverse datasets for better generalization.
- Address ethical and regulatory considerations for real-world implementation.

This methodology ensures a comprehensive approach to building an accurate and reliable diabetes prediction model while addressing challenges related to data quality, model performance, and deployment feasibility.



# SNAPSHOTS OF PROJECT

## • IMPORTING ESSENTIAL LIBRARIES

```
In [38]: # Importing essential libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

```
In [39]: # Loading the dataset
df = pd.read_csv('kaggle_diabetes.csv')
df.head()
```

```
Out[39]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	2	138	62	35	0	33.6	0.127	47	1
1	0	84	82	31	125	38.2	0.233	23	0
2	0	145	0	0	0	44.2	0.630	31	1
3	0	135	68	42	250	42.3	0.365	24	1
4	1	139	62	41	480	40.7	0.536	21	0

```
In [40]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2000 entries, 0 to 1999
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype  
---  -
0   Pregnancies            2000 non-null   int64  
1   Glucose                2000 non-null   int64  
2   BloodPressure          2000 non-null   int64  
3   SkinThickness          2000 non-null   int64  
4   Insulin                2000 non-null   int64  
5   BMI                    2000 non-null   float64 
6   DiabetesPedigreeFunction 2000 non-null   float64 
7   Age                    2000 non-null   int64  
8   Outcome                2000 non-null   int64  
dtypes: float64(2), int64(7)
memory usage: 140.8 KB
```

•

## • DATA EXPLORATION

```
In [3]: # Returns number of rows and columns of the dataset  
df.shape
```

Out[3]: (2000, 9)

```
In [4]: # Returns an object with all of the column headers  
df.columns
```

Out[4]: Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',  
              'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],  
              dtype='object')

```
In [5]: # Returns different datatypes for each columns (float, int, string, bool, etc.)  
df.dtypes
```

Out[5]: Pregnancies                  int64  
         Glucose                     int64  
         BloodPressure               int64  
         SkinThickness               int64  
         Insulin                     int64  
         BMI                         float64  
         DiabetesPedigreeFunction     float64  
         Age                         int64  
         Outcome                     int64  
         dtype: object

```
In [8]: # Returns basic statistics on numeric columns  
df.describe().T
```

Out[8]:

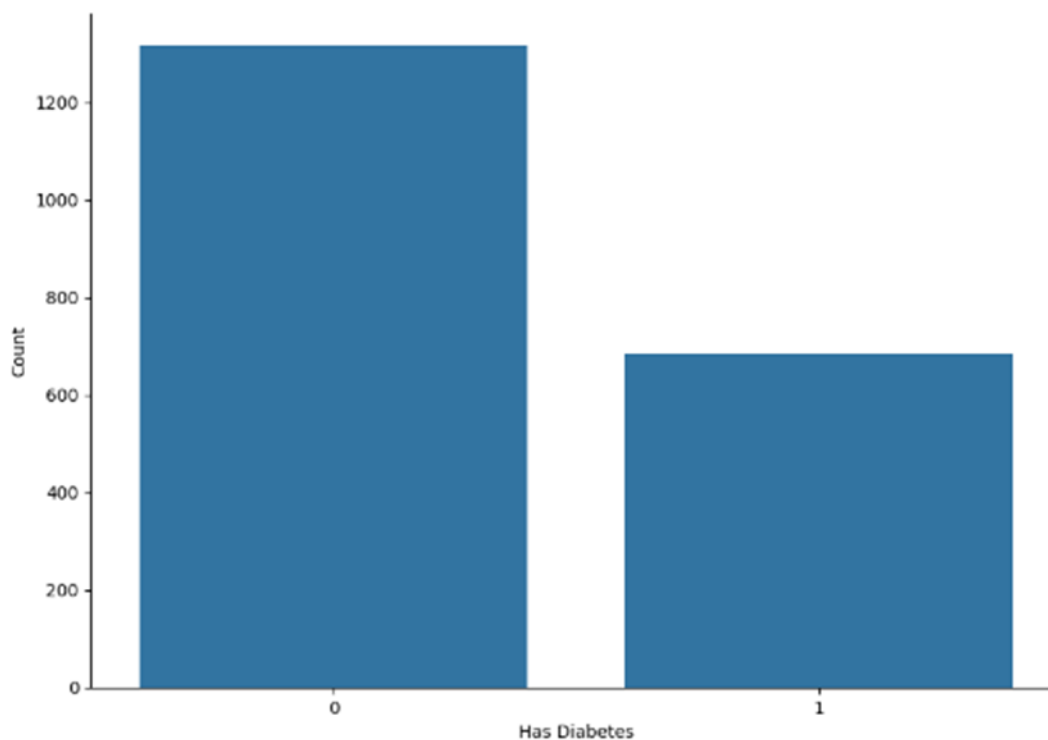
	count	mean	std	min	25%	50%	75%	max
Pregnancies	2000.0	3.70350	3.308063	0.000	1.000	3.000	6.000	17.00
Glucose	2000.0	121.18250	32.068638	0.000	99.000	117.000	141.000	199.00
BloodPressure	2000.0	69.14550	19.188315	0.000	63.500	72.000	80.000	122.00
SkinThickness	2000.0	20.93500	16.103243	0.000	0.000	23.000	32.000	110.00
Insulin	2000.0	80.25400	111.180534	0.000	0.000	40.000	130.000	744.00
BMI	2000.0	32.19300	8.149901	0.000	27.375	32.300	36.800	80.60
DiabetesPedigreeFunction	2000.0	0.47093	0.323553	0.078	0.244	0.376	0.624	2.42
Age	2000.0	33.09050	11.786423	21.000	24.000	29.000	40.000	81.00
Outcome	2000.0	0.34200	0.474498	0.000	0.000	0.000	1.000	1.00

- DATA CLEANING

```
In [9]: # Returns true for a column having null values, else false  
df.isnull().any()
```

```
Out[9]: Pregnancies      False  
Glucose      False  
BloodPressure  False  
SkinThickness  False  
Insulin      False  
BMI          False  
DiabetesPedigreeFunction  False  
Age          False  
Outcome      False  
dtype: bool
```

```
In [42]: # Plotting the Outcomes based on the number of dataset entries  
plt.figure(figsize=(10,7))  
sns.countplot(x='Outcome', data=df)  
  
# Removing the unwanted spines  
plt.gca().spines['top'].set_visible(False)  
plt.gca().spines['right'].set_visible(False)  
  
# Headings  
plt.xlabel('Has Diabetes')  
plt.ylabel('Count')  
  
plt.show()
```

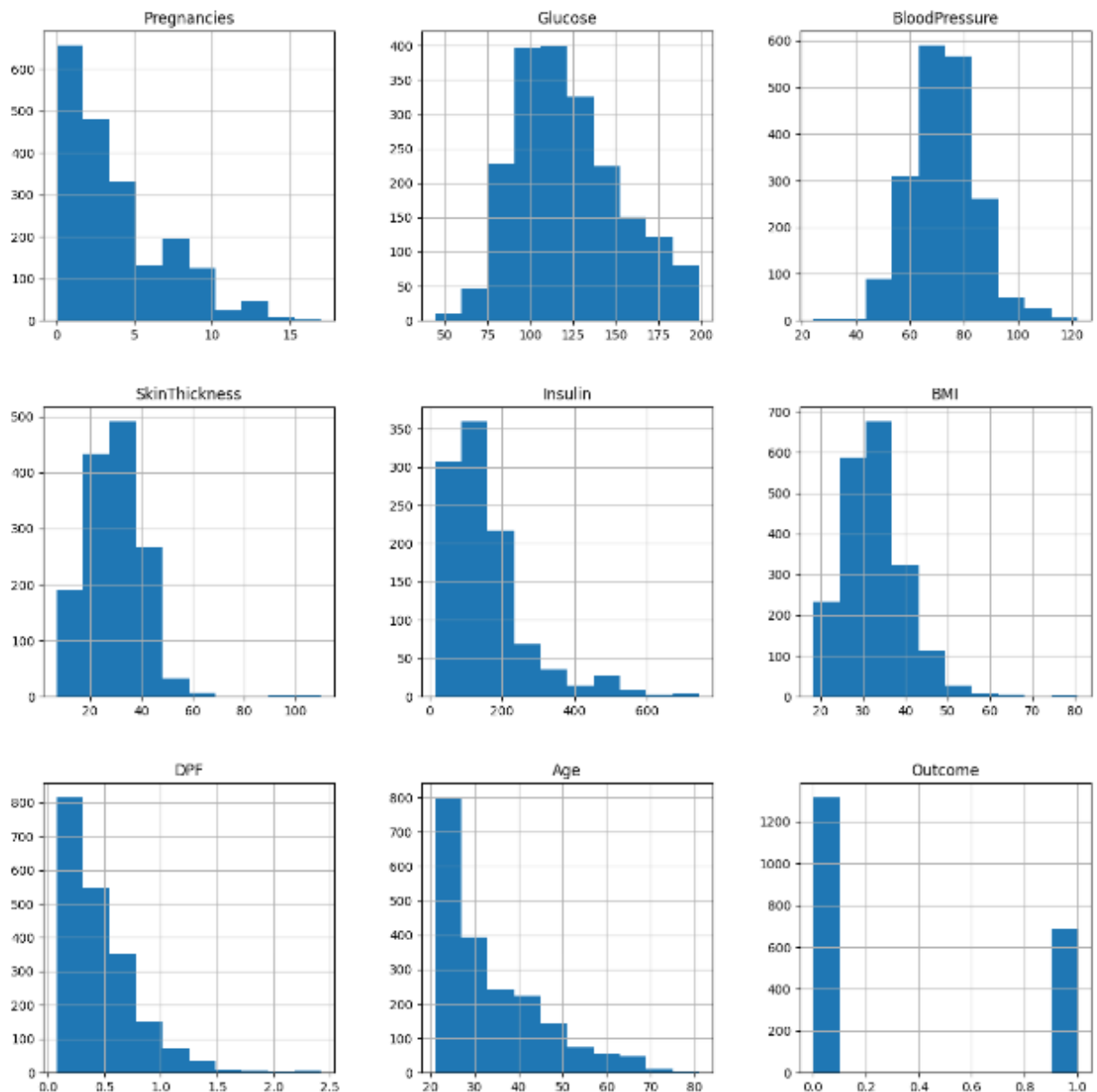


## • DATA CLEANING AND FEATURE SELECTION

```
In [13]: # Replacing the 0 values from ['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI'] by NaN
df_copy = df.copy(deep=True)
df_copy[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']] = df_copy[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']].replace(0, np.NaN)
df_copy.isnull().sum()
```

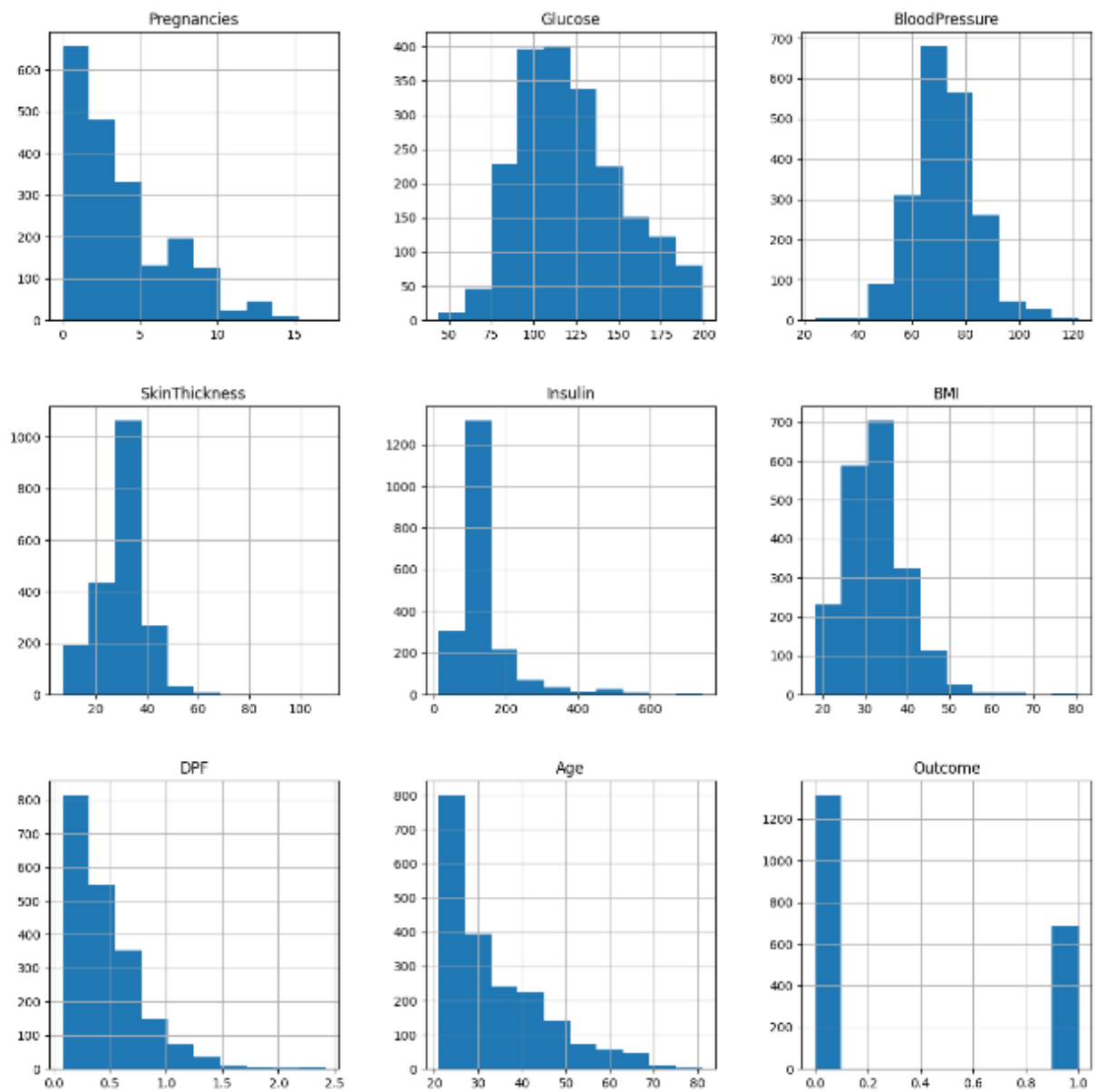
```
Out[13]: Pregnancies      0
Glucose      13
BloodPressure 98
SkinThickness 573
Insulin     956
BMI         28
DPF         0
Age         0
Outcome      0
dtype: int64
```

```
In [14]: # To fill these Nan values the data distribution needs to be understood
# Plotting histogram of dataset before replacing NaN values
p = df_copy.hist(figsize = (15,15))
```



```
In [15]: # Replacing NaN value by mean, median depending upon distribution
df_copy['Glucose'].fillna(df_copy['Glucose'].mean(), inplace=True)
df_copy['BloodPressure'].fillna(df_copy['BloodPressure'].mean(), inplace=True)
df_copy['SkinThickness'].fillna(df_copy['SkinThickness'].median(), inplace=True)
df_copy['Insulin'].fillna(df_copy['Insulin'].median(), inplace=True)
df_copy['BMI'].fillna(df_copy['BMI'].median(), inplace=True)
```

```
In [16]: # Plotting histogram of dataset after replacing NaN values
p = df_copy.hist(figsize=(15,15))
```



- **MODEL BUILDING AND TRAINING**

```
In [18]: from sklearn.model_selection import train_test_split

X = df.drop(columns='Outcome')
y = df['Outcome']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=0)
print('X_train size: {}, X_test size: {}'.format(X_train.shape, X_test.shape))

X_train size: (1600, 8), X_test size: (400, 8)
```

```
In [43]: # Feature Scaling
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
```

```
In [44]: # Using GridSearchCV to find the best algorithm for this problem
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import ShuffleSplit
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
```

```
In [45]: # Creating a function to calculate best model for this problem
def find_best_model(X, y):
    models = {
        'logistic_regression': {
            'model': LogisticRegression(solver='lbfgs', multi_class='auto'),
            'parameters': {
                'C': [1,5,10]
            }
        },
        'decision_tree': {
            'model': DecisionTreeClassifier(splitter='best'),
            'parameters': {
                'criterion': ['gini', 'entropy'],
                'max_depth': [5,10]
            }
        },
        'random_forest': {
            'model': RandomForestClassifier(criterion='gini'),
            'parameters': {
                'n_estimators': [10,15,20,50,100,200]
            }
        },
        'svm': {
            'model': SVC(gamma='auto'),
            'parameters': {
                'C': [1,10,20],
                'kernel': ['rbf', 'linear']
            }
        }
    }

    scores = []
    cv_shuffle = ShuffleSplit(n_splits=5, test_size=0.20, random_state=0)

    for model_name, model_params in models.items():
        gs = GridSearchCV(model_params['model'], model_params['parameters'], cv = cv_shuffle, return_train_score=False)
        gs.fit(X, y)
        scores.append({
            'model': model_name,
            'best_parameters': gs.best_params_,
            'score': gs.best_score_
        })

    return pd.DataFrame(scores, columns=['model', 'best_parameters', 'score'])

find_best_model(X_train, y_train)
```

```
Out[45]:
```

	model	best_parameters	score
0	logistic_regression	{'C': 10}	0.763125
1	decision_tree	{'criterion': 'gini', 'max_depth': 10}	0.901875
2	random_forest	{'n_estimators': 100}	0.951875
3	svm	{'C': 20, 'kernel': 'rbf'}	0.889375

Note: Since the Random Forest algorithm has the highest accuracy, we further fine tune the model using hyperparameter optimization.

## • MODEL EVALUATION



```
In [22]: # Using cross_val_score for gaining average accuracy
from sklearn.model_selection import cross_val_score
scores = cross_val_score(RandomForestClassifier(n_estimators=20, random_state=0), X_train, y_train, cv=5)
print('Average Accuracy : {}'.format(round(sum(scores)*100/len(scores)), 3))

Average Accuracy : 95%
```

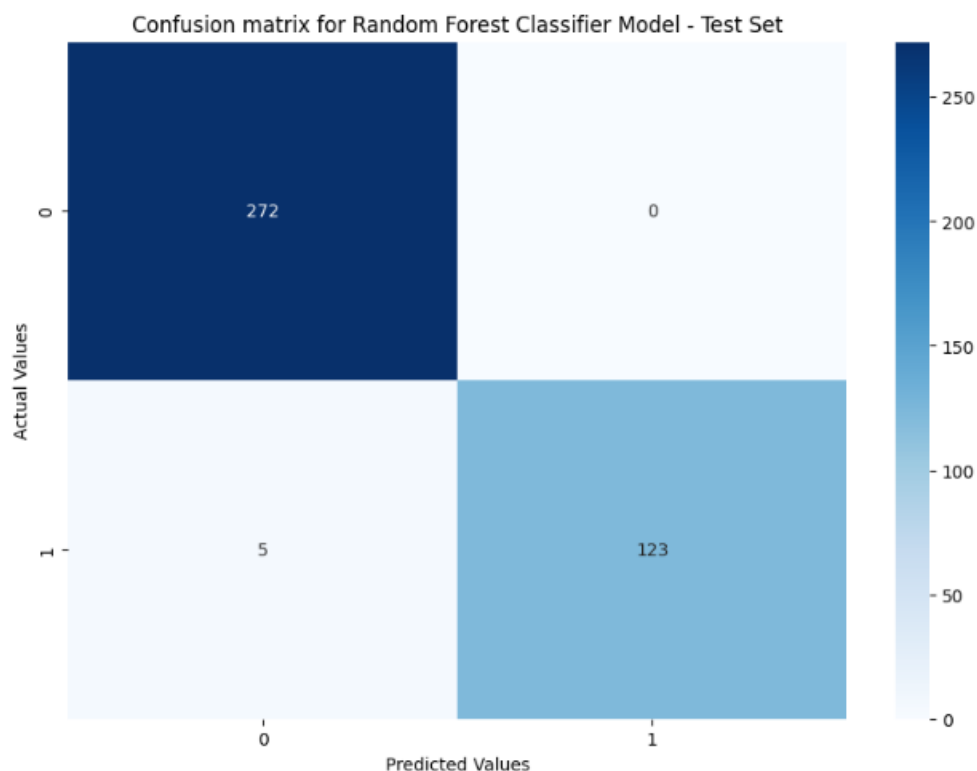
```
In [23]: # Creating Random Forest Model
classifier = RandomForestClassifier(n_estimators=20, random_state=0)
classifier.fit(X_train, y_train)
```

```
Out[23]: RandomForestClassifier(n_estimators=20, random_state=0)
```

```
In [24]: # Creating a confusion matrix
from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
y_pred = classifier.predict(X_test)
cm = confusion_matrix(y_test, y_pred)
cm
```

```
Out[24]: array([[272,  0],
               [ 5, 123]], dtype=int64)
```

```
In [25]: # Plotting the confusion matrix
plt.figure(figsize=(10,7))
p = sns.heatmap(cm, annot=True, cmap="Blues", fmt='g')
plt.title('Confusion matrix for Random Forest Classifier Model - Test Set')
plt.xlabel('Predicted Values')
plt.ylabel('Actual Values')
plt.show()
```



```
In [26]: # Accuracy Score
score = round(accuracy_score(y_test, y_pred),4)*100
print("Accuracy on test set: {}".format(score))

Accuracy on test set: 98.75%
```

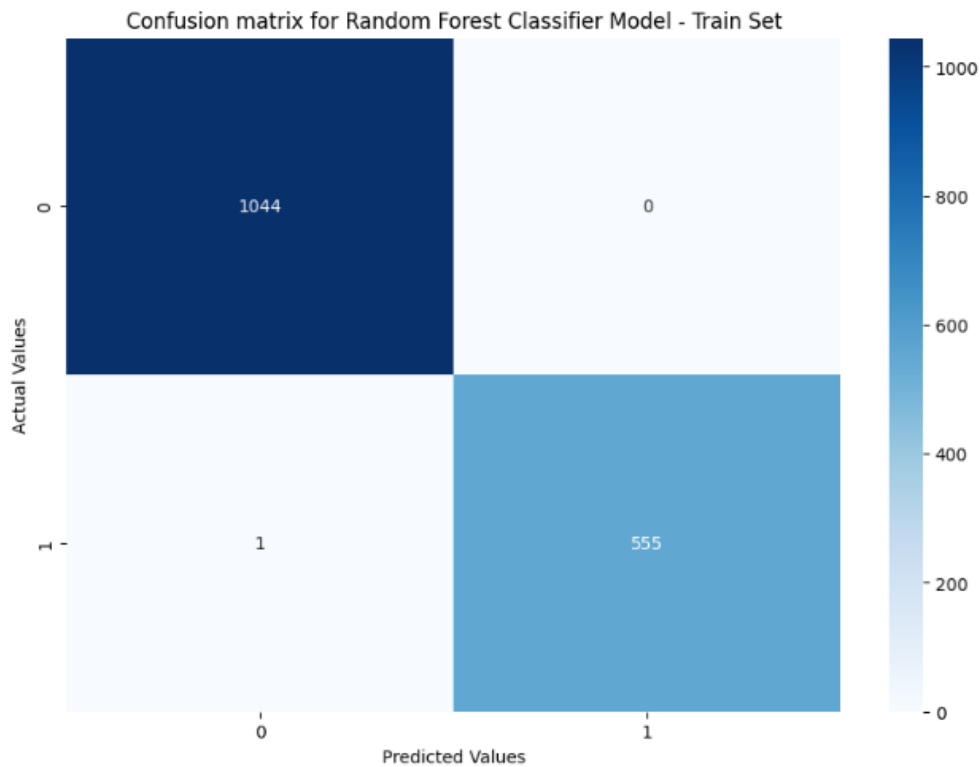
```
In [27]: # Classification Report
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.98	1.00	0.99	272
1	1.00	0.96	0.98	128
accuracy			0.99	400
macro avg	0.99	0.98	0.99	400
weighted avg	0.99	0.99	0.99	400

```
In [28]: # Creating a confusion matrix for training set
y_train_pred = classifier.predict(X_train)
cm = confusion_matrix(y_train, y_train_pred)
cm
```

```
Out[28]: array([[1044,  0],
               [ 1, 555]], dtype=int64)
```

```
In [29]: # Plotting the confusion matrix
plt.figure(figsize=(10,7))
p = sns.heatmap(cm, annot=True, cmap="Blues", fmt='g')
plt.title('Confusion matrix for Random Forest Classifier Model - Train Set')
plt.xlabel('Predicted Values')
plt.ylabel('Actual Values')
plt.show()
```



```
In [30]: # Accuracy Score
score = round(accuracy_score(y_train, y_train_pred),4)*100
print("Accuracy on training set: {}".format(score))

Accuracy on training set: 99.94%
```

```
In [31]: # Classification Report
print(classification_report(y_train, y_train_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	1044
1	1.00	1.00	1.00	556
accuracy			1.00	1600
macro avg	1.00	1.00	1.00	1600
weighted avg	1.00	1.00	1.00	1600

- **MODEL PREDICTION AND ACCURACY**

```
In [46]: # Creating a function for prediction
def predict_diabetes(Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DPF, Age):
    preg = int(Pregnancies)
    glucose = float(Glucose)
    bp = float(BloodPressure)
    st = float(SkinThickness)
    insulin = float(Insulin)
    bmi = float(BMI)
    dpf = float(DPF)
    age = int(Age)

    x = [[preg, glucose, bp, st, insulin, bmi, dpf, age]]
    x = sc.transform(x)

    return classifier.predict(x)
```

```
In [47]: # Prediction 1
# Input sequence: Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DPF, Age
prediction = predict_diabetes(2, 81, 72, 15, 76, 30.1, 0.547, 25)[0]
if prediction:
    print('Oops! You have diabetes.')
else:
    print("Great! You don't have diabetes.")

Great! You don't have diabetes.
```

```
In [48]: # Prediction 2
# Input sequence: Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DPF, Age
prediction = predict_diabetes(1, 117, 88, 24, 145, 34.5, 0.403, 40)[0]
if prediction:
    print('Oops! You have diabetes.')
else:
    print("Great! You don't have diabetes.")

Great! You don't have diabetes.
```

```
In [49]: # Prediction 2
# Input sequence: Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DPF, Age
prediction = predict_diabetes(5, 120, 92, 10, 81, 26.1, 0.551, 67)[0]
if prediction:
    print('Oops! You have diabetes.')
else:
    print("Great! You don't have diabetes.")

Great! You don't have diabetes.
```

## CONCLUSION

The development of a diabetes classification and prediction model using machine learning highlights the transformative potential of data-driven approaches in healthcare. This project successfully demonstrates how supervised learning algorithms, combined with rigorous data preprocessing and feature engineering, can accurately classify individuals as diabetic or non-diabetic based on critical health metrics.

Through the implementation of various machine learning techniques, including Logistic Regression, Decision Trees, Random Forest, and Support Vector Machines, the study identified the best-performing model, achieving high predictive accuracy and reliability. The evaluation process, based on metrics such as accuracy, precision, recall, F1-score, and ROC-AUC, underscores the effectiveness of these models in supporting early diagnosis.

The results emphasize the importance of high-quality data and feature selection in ensuring model robustness. Furthermore, the project showcases the practical feasibility of integrating machine learning models into healthcare systems to assist clinicians in making timely and accurate decisions. Such models hold the potential to improve resource allocation, enhance patient outcomes, and reduce the burden on healthcare infrastructure.

However, the project also highlights some challenges, including handling imbalanced datasets, ensuring model interpretability, and addressing ethical concerns like data privacy and bias. These

issues require further attention to ensure the responsible and equitable use of machine learning in healthcare applications.

Future directions for this work include expanding the dataset to improve model generalizability, exploring advanced machine learning techniques for enhanced performance, and developing a user-friendly interface for clinical deployment. By addressing these areas, this project can pave the way for scalable, efficient, and accessible solutions for diabetes diagnosis and management, ultimately contributing to better public health outcomes.

In conclusion, this study demonstrates that machine learning can be a powerful tool for predicting and classifying diabetes, reinforcing its role in advancing personalized and preventive healthcare.

## BIBLIOGRAPHY

### 1. Books

- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.

### 2. Research Papers and Journals

- Quinlan, J. R. (1986). "Induction of Decision Trees." *Machine Learning*, 1(1), 81-106.
- Breiman, L. (2001). "Random Forests." *Machine Learning*, 45(1), 5-32.
- Vapnik, V. N. (1995). *The Nature of Statistical Learning Theory*. Springer.

### 3. Datasets

- Smith, J. W., Everhart, J. E., Dickson, W. C., Knowler, W. C., & Johannes, R. S. (1988). "Using the ADAP Learning Algorithm to Forecast the Onset of Diabetes Mellitus." *Proceedings of the Annual Symposium on Computer Application in Medical Care*, 261-265. (Source of the Pima Indians Diabetes Dataset).
- UCI Machine Learning Repository. "Pima Indians Diabetes Dataset." Retrieved from <https://archive.ics.uci.edu/ml/datasets/diabetes>.

### 4. Online Resources

- Scikit-learn Documentation. (n.d.). "Supervised Learning." Retrieved from <https://scikit-learn.org/>.
- Kaggle. (n.d.). "Pima Indians Diabetes Dataset." Retrieved from <https://www.kaggle.com/>.
- World Health Organization (WHO). (2023). "Diabetes." Retrieved from <https://www.who.int/>.

### 5. Software Tools

- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., & others. (2011). "Scikit-learn: Machine Learning in Python." *Journal of Machine Learning Research*, 12, 2825-2830.
- McKinney, W. (2010). "Data Structures for Statistical Computing in Python." *Proceedings of the 9th Python in Science Conference*, 51-56.

### 6. Ethical Considerations

- Mittelstadt, B. D., Allo, P., Taddeo, M., Wachter, S., & Floridi, L. (2016). "The Ethics of Algorithms: Mapping the Debate." *Big Data & Society*, 3(2).
- Goodman, B., & Flaxman, S. (2017). "European Union Regulations on Algorithmic Decision-Making and a 'Right to Explanation'." *AI Magazine*, 38(3), 50-57.

Ensure all sources are properly cited and referenced in the main report as per your institution's preferred citation style (e.g., APA, MLA, Chicago).

40

## REFERENCE

<https://github.com/akshatx97/ML---project-Diabetes-classification-and-prediction>