

# **IMPROVING DIABETES RISK ASSESSMENT USING ADVANCED MACHINE LEARNING ALGORITHMS**

**BY:**

**RAMDHAN PRAJAPAT**

**(Admission No. 23MS0108)**



**THESIS**

**SUBMITTED TO**

**INDIAN INSTITUTE OF TECHNOLOGY  
(INDIAN SCHOOL OF MINES), DHANBAD**

For the award of the degree of

**MASTER OF SCIENCE**

**MAY 2025**

# Acknowledgments

I want to thank my guide, **Professor Kartikay Gupta Sir**, from the bottom of my heart for his able guidance, incessant encouragement, and constant support throughout my research. His level of expertise and worthwhile inputs have immensely helped develop this thesis. His sense of patience and understanding has encouraged me to challenge my limits and strive for the best in my work. I feel blessed to have worked under his guidance.

I am truly indebted to both of them for their valuable input towards this research, and I am fortunate to have been guided by them.

Ramdhan Prajapat

Admission No: 23MS0108

Mathematics and Computing

IIT(ISM) Dhanbad

Date: 03/05/2024

# Contents

<b>Contents</b>	<b>ii</b>
<b>List of Figures</b>	<b>v</b>
<b>Abstract</b>	<b>vi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	1
1.2 Motivation . . . . .	2
1.3 Objective . . . . .	4
1.3.1 Primary Goals . . . . .	4
1.3.2 Technical Approach . . . . .	4
1.3.3 Practical Implementation . . . . .	4
<b>2 Literature Review</b>	<b>6</b>
2.1 Types of Diabetes . . . . .	6
2.2 Current Prediction Systems . . . . .	6
2.3 Machine Learning Techniques . . . . .	6
2.4 Research Gaps . . . . .	7
2.5 Pregnancy Complications . . . . .	7
2.6 Promise of Advanced Methods . . . . .	7
<b>3 Methodology</b>	<b>9</b>
3.1 Data Collection . . . . .	9
3.2 Importing Libraries . . . . .	9
3.3 Data Preprocessing . . . . .	10
3.3.1 Data Quality Assessment . . . . .	10
3.3.1.1 Missing Values Analysis . . . . .	11
3.3.1.2 Outlier Detection . . . . .	11
3.3.1.3 Class Imbalance & SMOTE . . . . .	12

3.3.2	Feature Transformation . . . . .	12
3.3.3	Feature Encoding . . . . .	13
3.3.4	Dimensionality Management . . . . .	14
3.3.5	Data Partitioning . . . . .	15
3.4	Machine Learning Algorithms . . . . .	16
3.4.1	Logistic Regression . . . . .	16
3.4.1.1	Foundation: Maximum Likelihood Estimation . . . . .	16
3.4.1.2	Core Concept . . . . .	16
3.4.1.3	Loss Function . . . . .	17
3.4.2	Random Forest . . . . .	18
3.4.2.1	Foundation: Ensemble Learning . . . . .	18
3.4.2.2	How Random Forest Works (Core Concept) . . . . .	18
3.4.2.3	Advantages of Random Forest . . . . .	19
3.4.3	XGBoost . . . . .	20
3.4.3.1	What is XGBoost? . . . . .	20
3.4.3.2	Core Features of XGBoost . . . . .	20
3.4.3.3	How Does XGBoost Work? . . . . .	20
3.4.3.4	Mathematical Intuition . . . . .	21
3.4.3.5	Loss Function in XGBoost . . . . .	21
3.5	Model Evaluation Metrics . . . . .	23
3.5.1	Confusion Matrix . . . . .	23
3.5.2	AUC-ROC Score . . . . .	24
<b>4</b>	<b>Implementation</b> . . . . .	<b>27</b>
4.1	Dataset . . . . .	27
4.1.1	Data Characteristics . . . . .	27
4.1.2	Feature Description . . . . .	27
4.1.3	The variable of interest . . . . .	27
4.2	Data Visualization and Cleaning . . . . .	28
4.2.1	Descriptive Analysis . . . . .	28
4.2.2	Univariate Analysis of Numerical Features . . . . .	29
4.2.2.1	Pregnancies Distribution . . . . .	29
4.2.2.2	Plasma Glucose Concentration . . . . .	30
4.2.2.3	Blood Pressure Measurements . . . . .	30
4.2.2.4	Insulin and BMI . . . . .	30

4.2.2.5	Diabetes Risk Factors . . . . .	30
4.2.2.6	Interpretation . . . . .	31
4.2.3	Correlation Analysis . . . . .	31
4.2.4	Outliers Handling . . . . .	32
4.3	Data Wrangling . . . . .	33
4.3.1	Feature Engineering . . . . .	33
4.3.2	Data Partition . . . . .	33
4.4	Training and Evaluation . . . . .	33
4.5	Model Selection & Hyperparameter Tuning . . . . .	34
4.5.1	Hyperparameter Tuning . . . . .	34
4.5.2	Model Selection . . . . .	35
<b>5</b>	<b>Conclusion</b>	<b>37</b>
5.1	Summary . . . . .	37
5.2	Future Scope . . . . .	38
	<b>Bibliography</b>	<b>39</b>
5.3	Palagrism . . . . .	40

# List of Figures

3.1	Box plot showing outliers (dots) beyond the whiskers . . . . .	11
3.2	Correlation matrix heatmap showing feature relationships. Values near $\pm 1$ indicate strong correlations. . . . .	14
3.3	Random Forest workflow showing bootstrap sampling, parallel tree building, and prediction aggregation . . . . .	19
3.4	Example ROC curve showing AUC calculation. The diagonal represents random guessing (AUC=0.5). . . . .	26
4.1	Descriptive statistics of the diabetes dataset showing quartiles, maximum, minimum, mean, standard deviation, and count for eight clinical features. Notable observations include right-skewed distributions in Pregnancies (mean=3.22 & median=2.00) and Insulin (mean=137.85 vs median=83), and potential outliers in Insulin (max=799 vs 75th percentile=195). . . . .	28
4.2	KDE plots showing distributions of all numerical features . . . . .	29
4.3	Correlation matrix heatmap of clinical features with diabetes outcome . . . . .	31
4.4	Boxplot visualization of numerical features showing distribution characteristics and potential outliers. The plot displays median values (central line), interquartile ranges (box boundaries), and outlier points (dots beyond whiskers). . . . .	32
4.5	Comparison of AUC-ROC Curves for All Classifiers . . . . .	36

# Abstract

Diabetes is a very prevalent chronic illness. globally, and early risk assessment is critical for timely intervention and prevention. This thesis explores the use of advanced algorithms to increase the performance and efficiency of diabetes risk prediction. The primary objective is to develop a forecast model that can identify individuals at high risk of developing diabetes based on clinical and lifestyle-related characteristics.

Supervised machine learning includes algorithms like logistic regression, random forest classifier, support vector machines, decision tree classifier, gradient boosting classifier, and XGBoost classifier which are learned and tested on benchmarks such as the PIMA Indian Diabetes dataset. Data preparation techniques including scaling of features, lacking value computation, and class balance are used to improve model performance.

The study compares these algorithms using performance metrics that are precision, accuracy, recall, F1 score, and AUC-ROC. Among the models tested, ensemble methods demonstrated superior performance in identifying high-risk individuals. The findings support the integration of machine learning-based risk assessment tools into healthcare systems for rapid identification and individualized treatment options.

This work highlights the potential of AI and ML in predictive medical care and offers a step forward in automating and improving diabetes risk detection through data-driven techniques.

# Chapter 1

## Introduction

### 1.1 Background

Diabetes is the fastest-growing health problem globally. More than 400 million people suffer from diabetes today and the figure continues to increase every year.

Doctors currently rely on simple risk scores to screen who will develop diabetes. These scores examine factors such as age, weight, and genetic history. Helpful as they are, they usually miss individuals at risk but do not exhibit these common symptoms. Most at-risk patients are only diagnosed after the damage has started.

New technology provides us with new methods for detecting diabetes risks at earlier stages. Continuous glucose monitors (CGMs) are now wearable devices which identify blood sugar levels several times per minute. Wearable fitness trackers track activity, sleep, and heart rate patterns. In combination, these provide comprehensive health portraits the old way could not provide.

Machine learning provides robust tools to interpret this complicated data. Unlike straightforward risk scores, these algorithms can identify subtle patterns within thousands of data points. They pick up on tiny changes that may indicate future diabetes risk, even when standard tests are normal.

But most modern machine learning models have their shortcomings. Some are like "black boxes" - providing predictions without explanations. Others falter with disorganized real-world data from various devices. Many only take individual moments into account and not changes over time.

This study will create improved tools to predict diabetes risk through sophisticated machine learning. The aim is to merge various data types - ranging from lab work to wearable device data. Unique algorithms will examine how these elements interact in weeks and months.

The method has three main benefits. First, it can identify risks sooner by picking up on faint patterns. Second, it gives more precise explanations to enable physicians to comprehend the risks. Third, it functions with the flawed data from actual clinics and home equipment.

If successful, this research could prevent millions of diabetes complications. Early detection means that easier lifestyle adjustments can avoid the disease. Patients would have greater control over their health destiny. Physicians could target prevention where it's most needed.



## 1.2 Motivation

Diabetes prevention is fixing a tiny leak before the whole house is flooded. Currently, we tend to wait until pipes burst. The hope with this research is that we will get some warning signals prior to harm being caused.

Standard diabetes checks are similar to using old weather forecasts - they tell you it might rain, based on yesterday's weather. Modern technology provides live radar and we're not utilizing it to its potential. This project attempts to fix that.

Consider how phones can now identify faces more accurately than people. The same smart technology may be able to scan health trends we overlook. Your blood sugar, activity monitor, and lab tests all speak volumes that machine learning can interpret.

Most individuals receive diabetes warnings too late. By the time normal tests reveal issues, the body has already been damaged. It's like noticing smoke but waiting for fire to report it to the fire department. We can do better.

Current risk scores are not adjusted to individual lifestyles. A construction worker and an office worker with the same test score can have quite different risks. Machine learning can identify these crucial differences.

Physicians require tools that evolve with emerging research. Knowledge on diabetes today continues to increase, yet clinic tools remain constant. Our method allows the system to continue learning from new cases across the globe.

Patients should be explained things in plain language, not presented with numbers. Most risk scores provide percentages without explaining what they mean. This work seeks to demonstrate both the risk and the why behind it, such as a good health coach.

The COVID pandemic demonstrated how rapidly health systems can become burdened. Avoiding diabetes cases would alleviate future hospital loads. Even marginal gains in early diagnosis would benefit millions,

There is technology but it is not well networked. Glucose meters, exercise bands, and medical records all contain clues, but they communicate different languages. This research attempts to construct a translator among them.

Young adults particularly require improved warnings. Type 2 diabetes is surfacing earlier, but risk instruments are designed for older people. Machine learning can be modified to changing patterns across life and age.

Rural communities usually don't have specialists. An intelligent system may assist local physicians in identifying risks without having to wait for specialist opinions. This might render prevention more

equitable between communities.

The moment is ripe for this work. More individuals utilize health tech than ever, and computers are now able to identify patterns that humans cannot. Merging these may revolutionize the way we prevent diabetes.

This is not about substituting doctors, but equipping them with improved tools. Just as telescopes enable astronomers to see farther, machine learning might enable doctors to see danger sooner.

The mission is straightforward: keep people healthier for longer. If we can identify risks of diabetes even a year sooner, it may avoid thousands of complications and save lives.

## 1.3 Objective

The study seeks to improve a more practical and accurate system for diabetes risk assessment using more advanced machine learning methods. Though conventional models such as logistic regression have been satisfactory, they tend to overlook sophisticated interactions between the risk factors that might enhance early detection.

### 1.3.1 Primary Goals

Three aspects will be the focus of the model:

- **Improved Accuracy:** Using XGBoost's high pattern recognition capabilities to detect refined risk indicators missed by traditional models
- **Real-World Reliability:** Ensuring strong data management for typical clinical situations such as missing test results or aberrant measures
- **Clinical Transparency:** Producing not only risk scores but explanations of contributing factors that are comprehensible

### 1.3.2 Technical Approach

We will apply and compare various regularization techniques (L1/L2) to maximize model performance. L1 regularization assists in automatically selecting the most important features, whereas L2 avoids allowing one factor to overly dominate predictions artificially. This moderated strategy allows the model to be both accurate and clinically interpretable.

The system will handle heterogeneous data types such as:

- Normal clinical measures (glucose, BMI)
- Lifestyle variables (physical activity patterns)
- Demographic variables

### 1.3.3 Practical Implementation

In addition to technical innovation, we will develop:

- Easy-to-visualize tools illustrating the impact of various factors on risk

- Healthcare provider-friendly simple interfaces
- Interpretation guidelines for results in clinical practice

The ultimate product will be tested against other approaches to show quantifiable improvements in early detection of diabetes while retaining real-world practicality in healthcare practice.

# Chapter 2

## Literature Review

Diabetes impacts the way that the body breaks down sugar. When insulin fails to function correctly, sugar accumulates in blood rather than providing energy to cells. Researchers have created different methods for forecasting who will develop diabetes, each having advantages and limitations.

### 2.1 Types of Diabetes

Research indicates diabetes occurs in various forms:

- **Type 1:** Frequently begins in childhood when the immune system targets cells that produce insulin
- **Type 2:** Typically arises in adults from lifestyle and genetics
- **Gestational:** Appears in pregnancy due to hormonal influences

### 2.2 Current Prediction Systems

Most hospitals rely on low-tech checklists to determine the risk of diabetes:

- Ask patient questions about family and weight
- Look for basic blood lab results
- Employ score systems such as FINDRISC

These processes miss 20-30

### 2.3 Machine Learning Techniques

Research recently attempted to improve computer approaches:

- **Logistic Regression:** Is a simple calculator - helpful for straightforward cases but fails to catch subtle patterns

- **Decision Trees:** Takes yes/no routes similar to a flowchart - simple to comprehend but oversimplifies things
- **Neural Networks:** Imitates brain pathways, powerful but behaves like a "black box"

## 2.4 Research Gaps

Current approaches have critical limitations:

- **Over-simplification:** Traditional models (logistic regression, decision trees) treat risk factors as separate items rather than interconnected systems
- **Data challenges:** Most methods fail when test results are missing or wearables provide irregular data
- **Black box problem:** Advanced neural networks make accurate predictions but can't explain their reasoning
- **Untapped potential:** Newer algorithms like XGBoost - which handle complex patterns and missing data well - remain underused in diabetes prediction

## 2.5 Pregnancy Complications

For pregnant women, diabetes poses unique danger:

- Greater risk of large babies (more than 4kg)
- Greater demand for C-sections
- Risk of blood pressure issues

Current screening occurs at 24-28 weeks - too late for optimal prevention many times.

## 2.6 Promise of Advanced Methods

More recent methods such as XGBoost have promise because they:

- Detect faint patterns across tests
- Function with incomplete health records

- Can explain which factors matter most

This research leverages these advances while overcoming their limitations in more accurate diabetes risk prediction.

# Chapter 3

## Methodology

### 3.1 Data Collection

Dependency on ML models creation on the availability of a good and organized dataset is basic. Because machine learning algorithms are helpful not only in data-driven, but also the quality and organization of the dataset are what will finally decide how well the model works. A dataset is a set of data organized in a manner to address a specific task or problem.

Different datasets are employed for different purposes. A business analytics dataset would be entirely different from a dataset employed in healthcare research, like patient liver data analysis. Every dataset is unique in nature, attributes, and application.

Datasets are typically stored in CSV format since it is simple to do so and has good support in most data analysis tools and programming languages. Other formats like XLSX or HTML may be used its totally depends on the behaviour of the dataset and the tools that are available.

### 3.2 Importing Libraries

To develop our diabetes risk model, we need a few Python libraries. These libraries allow us to handle data, construct models, and verify results. We will be using highly tested libraries that are widely utilized in data science.

We import first some basic data handling utilities:

- **Pandas:** Helps organize our patient data in tables
- **Numpy:** Does math calculations efficiently
- **Matplotlib/Seaborn:** Creates visualizations to understand patterns

For machine learning tasks, we need special packages:

- **Scikit-learn:** Provides ready-to-use tools for:



- Splitting data into training/test sets
- Preprocessing (scaling numbers, handling missing values)
- Standard models (logistic regression, decision trees) for comparison
- **XGBoost**: Our main advanced algorithm for better predictions

We also bring along some helper libraries:

- **Joblib**: Saves trained models for later use
- **Warnings**: Filters out unnecessary alert messages
- **Time**: Measures how long processes take

Each library serves a specific purpose:

- Pandas loads our diabetes dataset from CSV files
- Numpy helps process the numbers quickly
- Matplotlib shows us graphs of patient characteristics
- Scikit-learn prepares the data and builds baseline models
- XGBoost creates our improved prediction model

These components are integrated like a fine-tuned workshop. Each plays its specialty but integrates well together. This allows us to work on enhancing diabetes risk assessment rather than on crafting fundamental features from scratch.

## 3.3 Data Preprocessing

Data processing is an important part of ML, which is a compilation with systematic steps and techniques for converting raw, unstructured observations into structured observations that can be treated and utilized through construct models. Data processing is required since machine learning model performance and accuracy are heavily dependent on input data structure and quality.

### 3.3.1 Data Quality Assessment

The initial stage involves evaluating dataset integrity through systematic checks:

### 3.3.1.1 Missing Values Analysis

It is very important to check missing data because missing data causes biased results and model errors.

#### How to Check:

- Count missing entries per column
- Plot missingness patterns
- Check if missing randomly or systematically

#### How to Handle:

- **Delete:** Remove rows/columns if few missing (less than 5%)
- **Fill:** Use mean/median for numbers, mode for categories
- **Predict:** Estimate missing values using other data
- **Flag:** Add "is missing" indicators

### 3.3.1.2 Outlier Detection

Extreme outliers are values which are either much smaller or larger than the average compared to other observations in a dataset. Outliers can occur due to a measurement error, clerical mistake, or even truly exceptional events. Eliminating or treating outliers correctly ensures that other aspects of data analysis and modeling are more refined and accurate. **Visual Detection with Box Plot Key Components:**

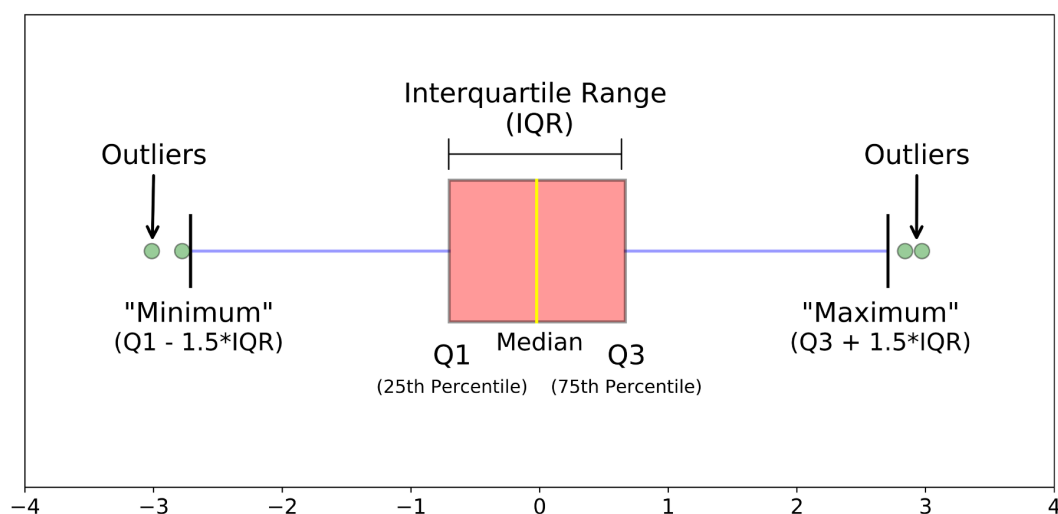


Figure 3.1: Box plot showing outliers (dots) beyond the whiskers

- The box represents Q1 to Q3 (middle 50% of data)
- Whiskers extend to 1.5×IQR from the box edges
- Points beyond whiskers are outliers

### Handling Method

- **Capping:** Replace outliers with nearest non-outlier values

#### 3.3.1.3 Class Imbalance & SMOTE

- **1. What is Class Imbalance?**

When there is less representation of one class as opposed to the other(s). For instance, 100 transactions marked as “fraud” versus 10,000 that are identified as “non-fraud.”

- **2. Why Balance Classes?**

- Models ignore minority classes, leading to poor recall (e.g., missing fraud cases).
- High accuracy but useless predictions (e.g., always predicting “non-fraud”).

- **3. Technique: SMOTE (Synthetic Minority Oversampling Technique)**

- **4. How SMOTE Works (Step-by-Step)**

- a. Take a sample from the minority class (e.g., a fraud data point).
- b. Find its 5 nearest neighbors (other fraud points).
- c. Pick one random neighbor from these 5.
- d. Choose a random number between 0 and 1 (e.g., factor = 0.3).
- e. Create a new synthetic point:

$$\text{New Point} = \text{Actual Point} + \text{factor} \times (\text{Neighbor Point} - \text{Actual Point})$$

(This interpolates along the line between the two points.)

### 3.3.2 Feature Transformation

Standard techniques for data standardization:

- **Standardization (Z-score)** Transforms features to zero mean and unit variance:

$$x_{\text{std}} = \frac{x - \mu}{\sigma}$$

*When to use:* For algorithms (like SVM and PCA) that require data that is normally distributed

- **Min-Max Normalization** rescales feature's values between 0 and 1:

$$x_{\text{norm}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

*When to use:* For algorithms sensitive to feature scales (e.g., neural networks) when bounds are known.

- **Sigmoid Normalization** Squashes values to a smooth (0,1) range:

$$x_{\text{sig}} = \frac{1}{1 + e^{-x}}$$

*When to use:* When extreme outliers exist and gradual tapering is desired.

- **Logarithmic Transformation** Reduces skewness in heavy-tailed distributions:

$$x_{\log} = \log(x + 1) \quad (+1 \text{ prevents } \log(0))$$

*When to use:* For exponential-scale data like monetary values or biological measurements.

- **Binning** Converts continuous to categorical values:

$$\text{Bin}_i = \begin{cases} \text{Low} & \text{if } x \leq Q1 \\ \text{Medium} & \text{if } Q1 < x \leq Q3 \\ \text{High} & \text{if } x > Q3 \end{cases}$$

*When to use:* To simplify complex relationships for decision trees or handle measurement noise.

### 3.3.3 Feature Encoding

Categorical variables require transformation into numerical representations for machine learning algorithms. We employ several encoding techniques based on data characteristics:

- **One Hot Encoding Method:** Makes binary (0/1) features for a category in nominal variables where no natural ordering exists. This prevents artificial ordinal relationships.

- *Example:* An example would be to establish three binary columns for a "Color" feature with values ["J", "K", and "L"]:

Color_J	Color_K	Color_L
1	0	0
0	1	0
0	0	1

- *Practical Note:* Use `pd.get_dummies()` in pandas or `OneHotEncoder` from scikit-learn. Beware of the "curse of dimensionality" with high-cardinality features.

- **Ordinal Encoding:** Preserves meaningful order in categorical variables by assigning integers according to their rank.

– *Example:* For education levels [”High School”, ”Bachelor”, ”Master”, ”PhD”]:

Education\_Level

1  
2  
3  
4

– *Practical Note:* Use `OrdinalEncoder` in scikit-learn. Manually define the category order for consistency across runs.

#### Important things:

- Always split data before encoding to prevent information leakage
- For tree-based models, ordinal encoding often suffices
- For linear models, one-hot encoding is typically better

### 3.3.4 Dimensionality Management

Approaches to optimize feature space:

- **Correlation Analysis:** Identify redundant variables using correlation matrices.
  - \* *Method:* Compute pairwise correlations between features and visualize using a heatmap.
  - \* *Example:* Drop one of two highly correlated features (e.g., ”income” and ”purchase capacity” with  $r > 0.9$ ).

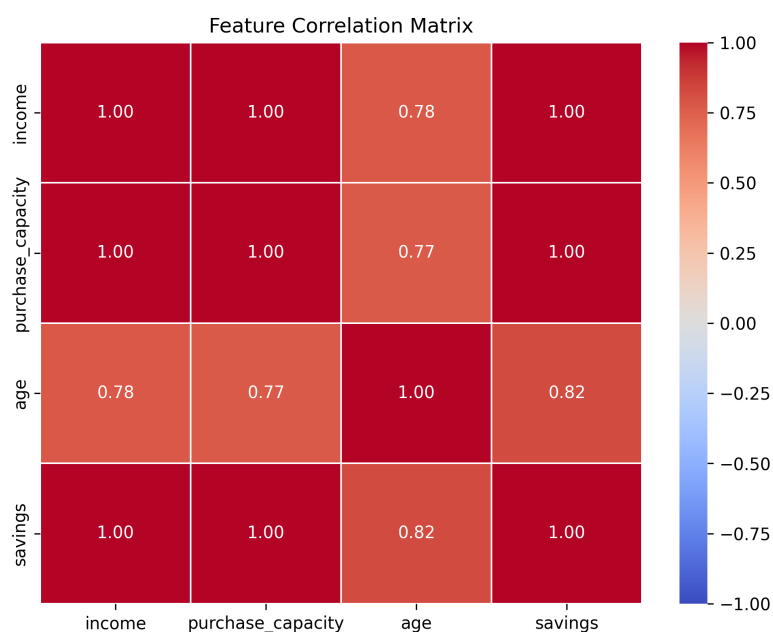


Figure 3.2: Correlation matrix heatmap showing feature relationships. Values near  $\pm 1$  indicate strong correlations.

\* *Interpretation:*

- **Red:** High positive correlation ( $r \rightarrow +1$ )
- **Blue:** High negative correlation ( $r \rightarrow -1$ )
- **White:** Weak correlation ( $r \approx 0$ )

– **Feature Importance:** Select predictors using model-based metrics.

\* *Method:* Train a model (e.g., Random Forest) and extract feature importance scores

\* *Example:* For house price prediction:

- High importance: "square footage" (score = 0.42)
- Medium importance: "location rating" (score = 0.31)
- Low importance: "number of windows" (score = 0.02)
- Negligible: "year of last paint" (score = 0.001)

### 3.3.5 Data Partitioning

\* **Purpose:** Splits dataset into training (80%) and testing (20%) sets to assess model performance on new data.

\* **Why Use It?**

- Overcome the problem of overfitting by checking the model generalizes well.
- Provides an unbiased estimate of model accuracy.
- Helps validate model robustness before deployment.

\* **Practical Example:**

A dataset with 10,000 house prices is split into:

- Training Data set: 8,000 samples (model learning).
- Test Data set: 2,000 samples (performance evaluation).

This complete preprocessing approach helps improve data quality and ensures that the steps followed are correct and trustworthy for building ML models.

## 3.4 Machine Learning Algorithms

### 3.4.1 Logistic Regression

Logistic Regression is used to solve the two-class classification problem by forecasting the likelihood of an instance falling into a given class. Unlike linear regression, which produces values that are continuous, logistic regression provides values between 0 and 1 based on the sigmoid function.

#### 3.4.1.1 Foundation: Maximum Likelihood Estimation

Logistic Regression make use of Maximum Likelihood Estimation (MLE) to identify the most beneficial parameters. MLE selects parameters that maximize probability of observing the given data:

$$M(\theta) = \prod_{i=1}^k Q(q_i | p_i; \theta) = \prod_{i=1}^k g_{\theta}(p_i)^{q_i} (1 - g_{\theta}(p_i))^{1-q_i} \quad (3.1)$$

Where:

- $M(\theta)$  is likelihood function
- $g_{\theta}(p_i)$  is predicted probability
- $q_i$  is the actual class label (0 or 1)

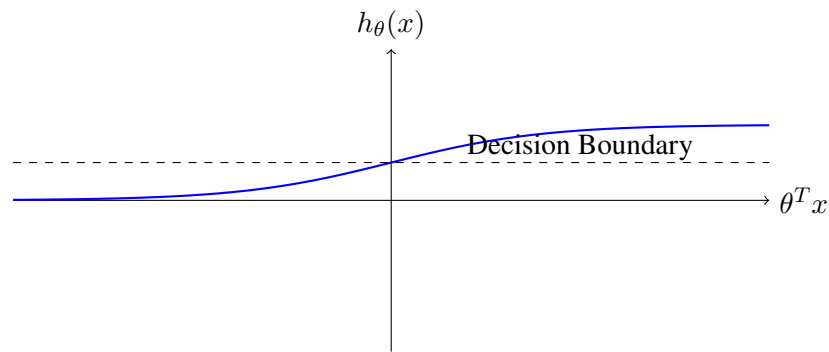
We typically maximize the log-likelihood for computational convenience:

$$m(\theta) = \sum_{i=1}^k [q_i \log g_{\theta}(p_i) + (1 - q_i) \log(1 - g_{\theta}(p_i))] \quad (3.2)$$

#### 3.4.1.2 Core Concept

The model predicts probabilities using the sigmoid activation function:

$$g_{\theta}(p) = \sigma(\theta^T p) = \frac{1}{1 + e^{-\theta^T p}} \quad (3.3)$$



### 3.4.1.3 Loss Function

The log loss (cross-entropy loss) measures model performance:

$$K(\theta) = -\frac{1}{c} \sum_{i=1}^c [q^{(i)} \log(g_\theta(p^{(i)})) + (1 - q^{(i)}) \log(1 - g_\theta(p^{(i)}))] \quad (3.4)$$

Key properties:

- Penalizes confident wrong predictions heavily
- Convex function - guarantees global minimum
- Directly related to the log-likelihood function



### 3.4.2 Random Forest

Random forest is a managed collective learning system that combines several decisions to create trees to make good results and decreases overhangs. Through aggregation of the results of many individual trees (often by the majority of classification or the average of regression), we create a model that is more robust and generalizable than a single decision-making tree.

#### 3.4.2.1 Foundation: Ensemble Learning

Random Forest builds upon two key ensemble concepts:

- **Bootstrap Aggregating (Bagging):** Bagging is a type of method that is employed to resolve the problem of variance of a prediction maker and this entails producing numerous subsets of the dataset using random selection with substitution (bootstrap samples), then training a distinct model on each, and then aggregating their outputs. This ensemble approach leads to a more stable and accurate model.

$$\hat{f}_{\text{bag}}(p) = \frac{1}{S} \sum_{i=1}^S \hat{f}^{*i}(p) \quad (3.5)$$

where  $S$  is the number of bootstrap samples.

#### 3.4.2.2 How Random Forest Works (Core Concept)

The Random Forest method is based on a collection of decision trees that work together to form an ensemble. During training:

- Multiple bootstrap samples that are drawn from the given dataset of training purpose.
- For a certain sample, an ML algorithm decision tree that learns by using a selected set of attributes at every division
- Each tree produces its own forecasts, and the ultimate result is determined by:
  - **Classification:** The bulk voting from all trees
  - **Regression:** Average prediction of all trees

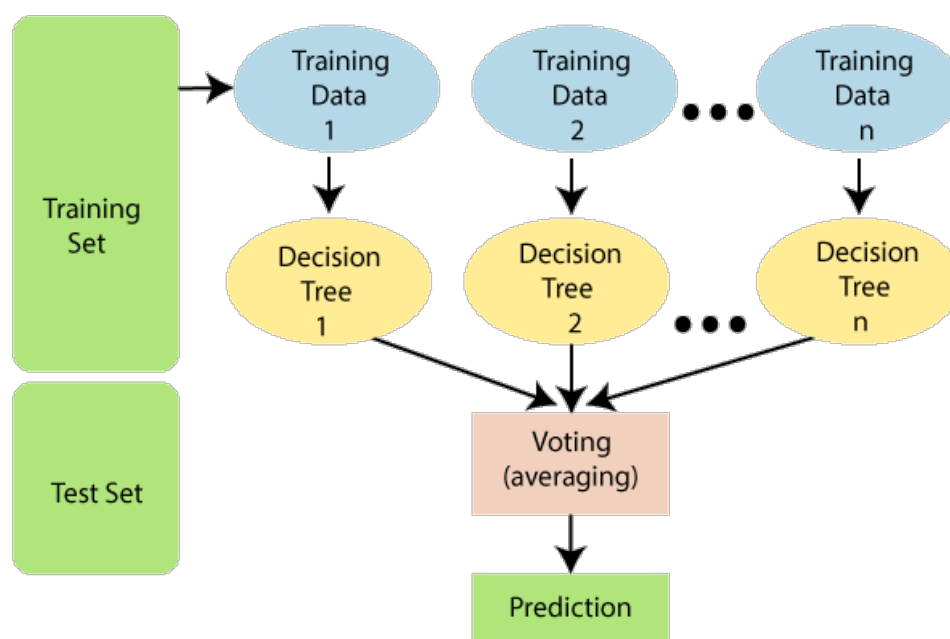


Figure 3.3: Random Forest workflow showing bootstrap sampling, parallel tree building, and prediction aggregation

Key characteristics:

- No pruning is applied - trees grow deep to capture complex patterns
- Double randomness (data + features) ensures low tree correlation
- Built-in validation through out-of-bag (OOB) samples
- Naturally handles missing values and outliers

### 3.4.2.3 Advantages of Random Forest

- **High Accuracy:** Combines multiple trees to produce a model that performs well with data that was previously unknown.
- **Robust to Overfitting:** An ensemble structure naturally combats overfitting, especially with enough trees and proper randomness.
- **Handles Missing Values:** Can maintain accuracy even with partial data and supports imputation.
- **Non-linear Relationships:** Capable of modeling complex non-linear interactions between features.
- **Feature Importance Insight:** Automatically provides insights into the most influential features for prediction.

### 3.4.3 XGBoost

#### 3.4.3.1 What is XGBoost?

XGBoost is a robust and effective implementation of gradient boosting algorithm that prioritizes speed and performance. It is a type of ensemble learning that creates a powerful forecasting model through merging the predictions of numerous weak learners, usually decision trees. XGBoost is frequently utilized in ML competitions and business applications because of its capacity to handle big datasets while maintaining excellent accuracy.

#### 3.4.3.2 Core Features of XGBoost

XGBoost offers many important features:

- **Regularization:** Incorporates L1 (Lasso) as well as L2 (Ridge) regularization to reduce high variance.
- **Processing in parallel:** Supports distributed computing for better training.
- **Managing Missing Values:** Automatically manages missing data.
- **Tree Pruning:** Uses a depth-first approach to split trees, reducing overfitting.
- **Cross-Validation:** Built-in support for k-fold cross-validation.

#### 3.4.3.3 How Does XGBoost Work?

XGBoost creates decision trees consecutively, with every tree seeking to fix the flaws of the preceding one. The procedure can be split down as follows:

- a. **Begin with a basic learner:** The initial decision tree (a  $Tree_1$ ) is trained using the initial data. For tasks involving regression, this base model often forecasts the target variable's average value.
- b. **Compute the mistakes that were made:** Upon training the initial tree, the differences (errors) between the expected and actual outcomes are calculated as follows:

$$\text{Residuals} = y_{\text{true}} - \hat{y}_1$$

where  $\hat{y}_1$  is the prediction from  $Tree_1$ .

- c. **Instruct to the next tree:** A second tree ( $Tree_2$ ) is taught on the residuals from a learning to rectify the faults caused by  $Tree_1$ . The learning rate ( $\eta$ ) determines how quickly each tree fixes past errors.

$$\hat{y}_2 = \hat{y}_1 + \eta \cdot \text{Tree}_2(\mathbf{x})$$

- d. **Repeat the process:** This iterative correction continues for  $M$  trees, with each subsequent tree ( $\text{Tree}_m$ ) generated from residuals taken from the ensembles of all preceding trees:

$$\text{Residuals}_m = y_{\text{true}} - \sum_{k=1}^{m-1} \eta \cdot \text{Tree}_k(\mathbf{x})$$

The process ends when either the highest possible amount of trees is achieved or extra additions fail to enhance performance.

- e. **Combine the predictions:** The final prediction is the weighted sum of all tree predictions:

$$\hat{y}_{\text{final}} = \sum_{m=1}^M \eta \cdot \text{Tree}_m(\mathbf{x})$$

### 3.4.3.4 Mathematical Intuition

The goal of the objective function of XGBoost is made up of two parts: first is the loss function and the second is nothing but the regularization term:

$$\mathcal{M}(\phi) = \sum_{i=1}^c m(q_i, \hat{q}_i) + \sum_{l=1}^L \Omega(g_l)$$

where:

- $m(q_i, \hat{q}_i)$  is loss function that calculates difference between the actual label  $q_i$  and the forecasted label  $\hat{q}_i$ .
- $\Omega(f_l)$  is the penalty term that is called as regularization term for the  $l$ -th tree.

### 3.4.3.5 Loss Function in XGBoost

The loss function in XGBoost is customizable, but commonly used ones include:

- **Regression's Error:** Squared error penalty  $m(q, \hat{q}) = (q - \hat{q})^2$ .
- **Classification's Error:** Logistic loss penalty  $m(q, \hat{q}) = q \log(1 + e^{-\hat{q}}) + (1 - q) \log(1 + e^{\hat{q}})$ .

The penalty term as regularization  $\Omega(g_l)$  is given by:

$$\Omega(g_l) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

where:

- $T$  indicates how many nodes in the particular tree.
- $w_j$  indicates how much the weight of the  $j$ -th leaf.
- $\gamma$  and  $\lambda$  are hyperparameters adjusting the penalty either reward or penalizes.

### Advantages of XGBoost

- **High Performance:** Optimized for speed and efficiency, often outperforming other algorithms.
- **Flexibility:** Supports various loss functions and evaluation metrics.
- **Regularization:** Reduces overfitting, improving generalization.
- **Feature Importance:** Provides insights into feature relevance.
- **Handling Missing Data:** Robust to missing values without imputation.

## 3.5 Model Evaluation Metrics

### 3.5.1 Confusion Matrix

A confusion matrix is a valuable method for evaluating classifying methods, which shows the kinds of errors the model is producing in addition to how many are accurate and not accurate predictions. It offers a comprehensive view of the model evaluation in every class. The following counts are displayed in this tabular representation of a classification model's efficiency:

		Actual	
		Positive	Negative
Predicted	Positive	TP	FP
	Negative	FN	TN

Where:

- **TP (True Positives):** Correctly predicted positive cases
- **FP (False Positives):** Negative cases that predicted as positive (Known as TYPE I error)
- **FN (False Negatives):** Positive cases that predicted as negative (TYPE II error)
- **TN (True Negatives):** Correctly predicted negative cases

#### Key Metrics Derived from Confusion Matrix

- **Accuracy:** The overall percentage of accurate predictions. While beneficial for fair datasets, it may be inaccurate when classes are skewed.

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

- **Precision:** When ML model predicts "POSITIVE", how often is it correct? High precision means fewer false alarms. Crucial when FP costs are high (e.g., spam detection where marking legitimate email as spam is costly)

$$\text{Precision} = \frac{TP}{TP + FP}$$

- **Recall/Sensitivity:** What proportion of actual positives did the model catch? High recall means missing few positive cases. Vital when FN are dangerous (e.g., cancer screening where missing a case has severe consequences)

$$\text{Recall} = \frac{TP}{TP + FN}$$

- **F1-SCORE:** The precision and recall harmonic mean. gives a single score that takes into account both issues. particularly helpful when attempting to strike a balance among recall and precision.

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- **Specificity:** How good whether model is at identifying negatives. Important when correctly identifying "not X" matters (e.g., verifying someone is not a security threat).

$$\text{Specificity} = \frac{TN}{TN + FP}$$

### Practical Interpretation Guide

- **High Precision + Low Recall:** Model is very conservative in making positive predictions (few false positives, but misses many actual positives)
- **Low Precision + High Recall:** Model casts a wide net (catches most positives but with many false alarms)
- **Balanced Precision/Recall:** Model creates trade-off among false positives points and false negatives points
- **High Specificity:** Model is very good at "ruling out" negative cases

### 3.5.2 AUC-ROC Score

The model's capacity to differentiate among classes over all potential categorization thresholds is gauged by the Area Under the Receiver Operating Characteristic Curve.

#### Understanding ROC Components

- **X-axis (False Positive Rate):**

$$FPR = \frac{FP}{FP + TN}$$

*Interpretation:* indicates the frequency with which negative examples are misclassified as positive (1 - Specificity).

- **Y-axis (True Positive Rate):**

$$TPR = \frac{TP}{TP + FN} = \text{Recall}$$

*Interpretation:* Shows how often positive instances are correctly identified

### What AUC-ROC Reveals About Your Model

- **Threshold-independent:** Evaluates performance across all possible decision thresholds
- **Class imbalance robustness:** More reliable than accuracy for imbalanced datasets
- **Probability quality:** Measures how well model ranks predictions (higher scores for positive instances)
- **Model comparison:** Allows direct comparison of different models' discrimination ability

### When to Use AUC-ROC

- Binary classification problems
- When you care about ranking predictions (e.g., fraud detection)
- When class distribution is imbalanced
- When optimal classification threshold is unknown

### Limitations

- Doesn't show actual threshold values
- Less interpretable than confusion matrix metrics
- Can be optimistic for imbalanced datasets (consider Precision-Recall AUC)



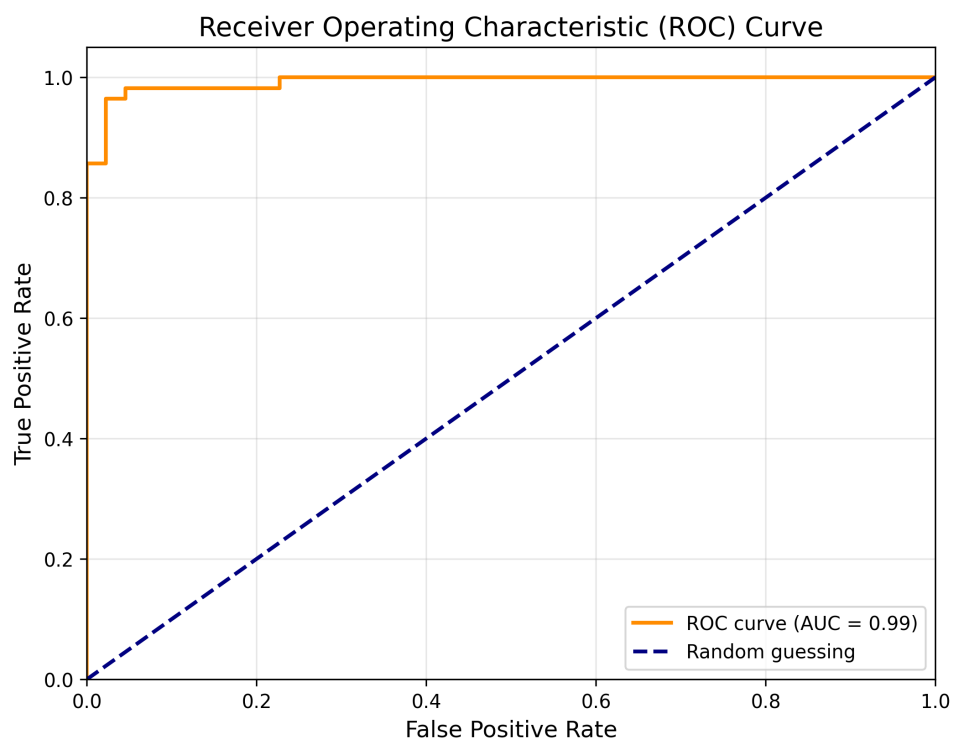


Figure 3.4: Example ROC curve showing AUC calculation. The diagonal represents random guessing (AUC=0.5).

# Chapter 4

## Implementation

1.25

### 4.1 Dataset

The subject's database is a mixed clinical dataset containing health metrics related to diabetes diagnosis. The dataset consists of 15,000 instances with 10 features, including patient identifiers and biological measurements. The data structure is organized as follows:

#### 4.1.1 Data Characteristics

- **Size:** 15,000 patient records
- **Dimensionality:** 10 features (9 predictors + 1 target variable)
- **Type:** Multivariate tabular data
- **Target:** Binary classification (diabetic vs non-diabetic)

#### 4.1.2 Feature Description

The dataset contains the following variables:

#### 4.1.3 The variable of interest

The binary target variable `Diabetic` indicates:

- **1:** Patient is diabetic
- **0:** Patient is not diabetic

Table 4.1: Dataset Features and Descriptions

Feature	Description
PatientID	Unique patient identifier
Pregnancies	Number of pregnancies
PlasmaGlucose	Plasma glucose concentration (mg/dL)
DiastolicBloodPressure	Diastolic blood pressure (mm Hg)
TricepsThickness	Triceps skinfold thickness (mm)
SerumInsulin	2-Hour serum insulin (mu U/ml)
BMI	Body mass index (kg/m <sup>2</sup> )
DiabetesPedigreeFunc	Diabetes pedigree function
Age	Age in years
Diabetic	Target variable (1=Diabetic, 0=Non-diabetic)

This classification problem aims to predict diabetes status based on the provided clinical measurements. The dataset is particularly valuable for developing predictive models in healthcare applications, as it contains both demographic information and key physiological measurements known to correlate with diabetes risk.

## 4.2 Data Visualization and Cleaning

### 4.2.1 Descriptive Analysis

	Pregnancies	PlasmaGlucose	DiastolicBloodPressure	TricepsThickness	SerumInsulin	BMI	DiabetesPedigree	Age
count	15000.000000	15000.000000	15000.000000	15000.000000	15000.000000	15000.000000	15000.000000	15000.000000
mean	3.224533	107.856867	71.220667	28.814000	137.852133	31.509646	0.398968	30.137733
std	3.391020	31.981975	16.758716	14.555716	133.068252	9.759000	0.377944	12.089703
min	0.000000	44.000000	24.000000	7.000000	14.000000	18.200512	0.078044	21.000000
25%	0.000000	84.000000	58.000000	15.000000	39.000000	21.259887	0.137743	22.000000
50%	2.000000	104.000000	72.000000	31.000000	83.000000	31.767940	0.200297	24.000000
75%	6.000000	129.000000	85.000000	41.000000	195.000000	39.259692	0.616285	35.000000
max	14.000000	192.000000	117.000000	93.000000	799.000000	56.034628	2.301594	77.000000

Figure 4.1: Descriptive statistics of the diabetes dataset showing quartiles, maximum, minimum, mean, standard deviation, and count for eight clinical features. Notable observations include right-skewed distributions in Pregnancies (mean=3.22, median=2.00) and Insulin (mean=137.85 vs median=83), and potential outliers in Insulin (max=799 vs 75th percentile=195).

## 4.2.2 Univariate Analysis of Numerical Features

The univariate analysis examines the distribution of each numerical feature through kernel density estimation (KDE) plots. These visualizations reveal the underlying probability distributions of key health metrics across patients.

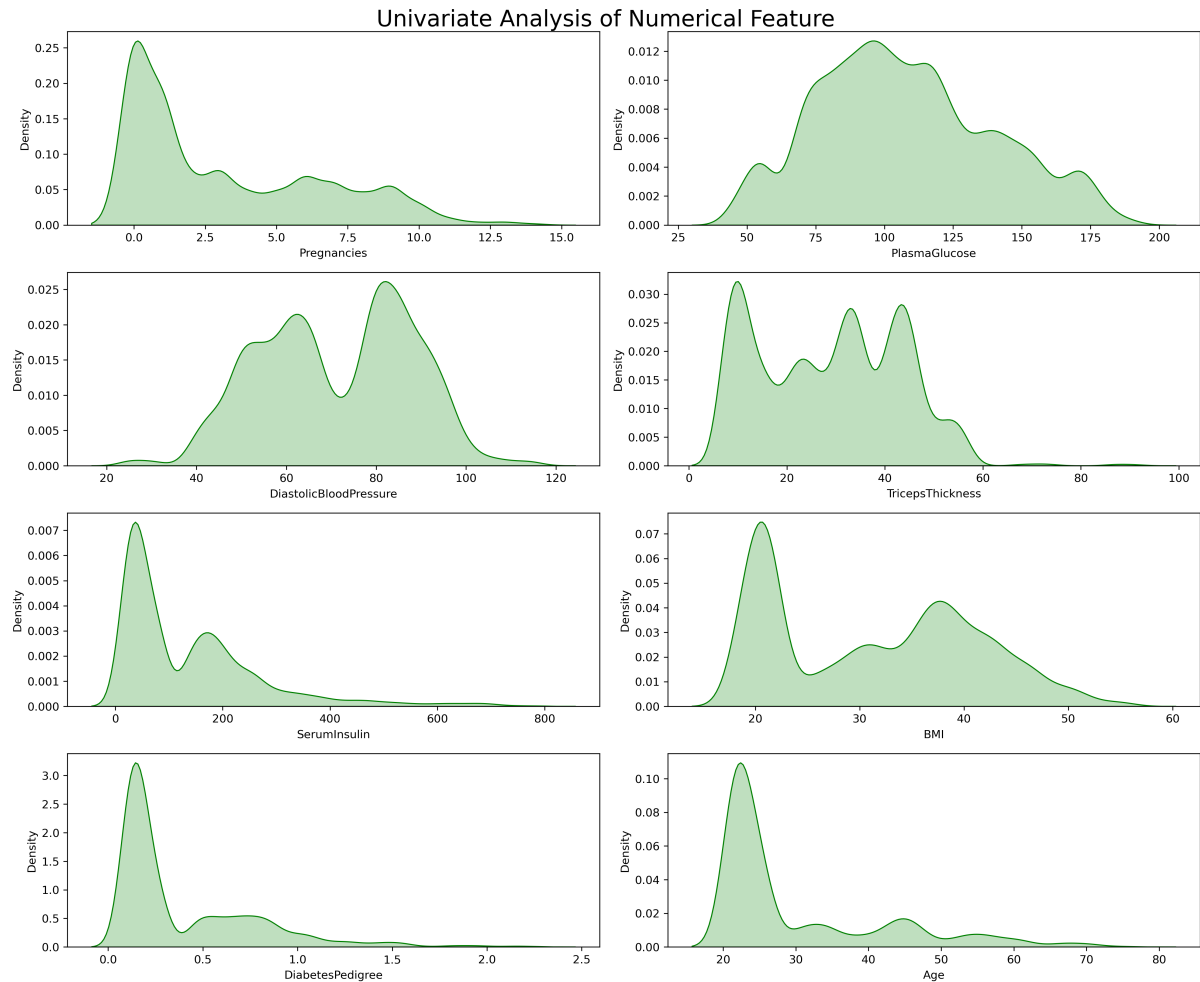


Figure 4.2: KDE plots showing distributions of all numerical features

### 4.2.2.1 Pregnancies Distribution

- Right-skewed distribution with most patients reporting 0-5 pregnancies
- Peak density around 0 pregnancy
- Long tail extending to 15 pregnancies

#### 4.2.2.2 Plasma Glucose Concentration

- Normal distribution centered at 100 mg/dL
- The range (60-160 mg/dL) contains majority of patients

#### 4.2.2.3 Blood Pressure Measurements

- **Diastolic Blood Pressure:**
  - Bimodal distribution suggesting two patient subgroups
  - Primary peak around 100 mg/dL (normal range)
  - Secondary peak near 140 mg/dL (prediabetic range)
- **Triceps Thickness:**
  - Most of the values within the range( 5mm to 60mm)
  - Potential outliers above 60 mm

#### 4.2.2.4 Insulin and BMI

- **Serum Insulin:**
  - Extreme right-skew with most values  $<250 \mu\text{U/ml}$
  - Long tail indicates insulin resistance cases
- **BMI:**
  - Bimodal distribution suggesting two patient subgroups
  - Primary peak around  $20 \text{ kg/m}^2$
  - Secondary peak near  $40 \text{ kg/m}^2$

#### 4.2.2.5 Diabetes Risk Factors

- **Diabetes Pedigree Function:**
  - Right-skewed distribution

- Most values clustered near 0 with long tail

- **Age:**

- Right-skewed distribution
- Majority of patients aged 20-40 years

#### 4.2.2.6 Interpretation

The density plots reveal:

- Several features show clinically relevant thresholds
- Multiple right-skewed distributions suggest potential need for log transformation
- Bimodal patterns in glucose may indicate distinct patient subgroups
- Outliers present in insulin and skin thickness measurements

#### 4.2.3 Correlation Analysis



Figure 4.3: Correlation matrix heatmap of clinical features with diabetes outcome

The correlation matrix reveals several important relation among clinical attributes and diabetes diagnosis attribute:

The correlation analysis reveals pregnancies as the strongest forecastor of disease ( $r = 0.41$ ), followed by age ( $r = 0.34$ ) and serum insulin ( $r = 0.25$ ), while surprisingly, plasma glucose shows only weak correlation ( $r = 0.13$ ). Notably, body composition markers (BMI and triceps thickness) demonstrate modest relationships, and blood pressure shows minimal associations. These findings suggest diabetes risk involves complex interactions beyond individual biomarkers, with parity and age emerging as unexpectedly significant factors compared to traditional metabolic markers.

#### 4.2.4 Outliers Handling

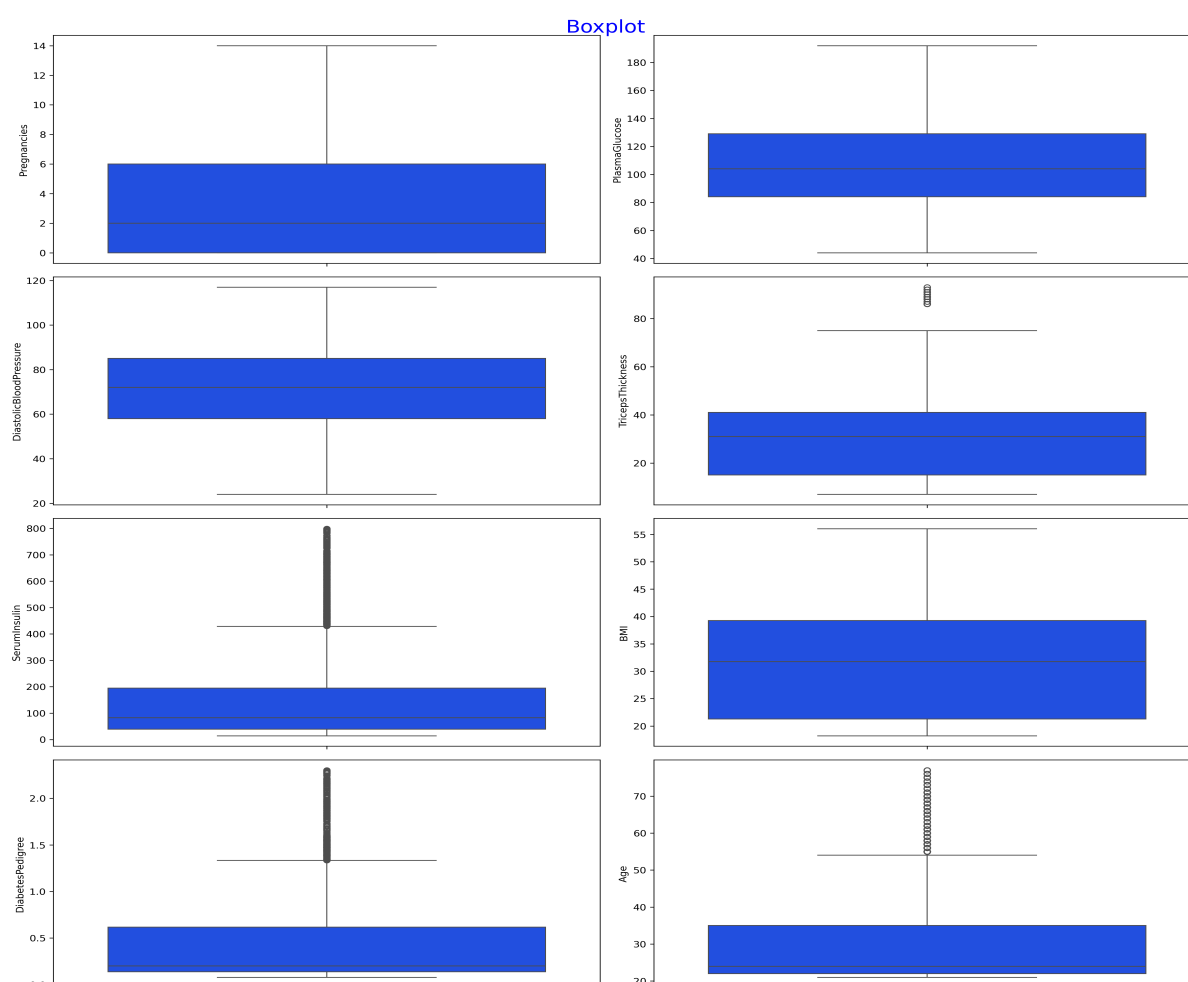


Figure 4.4: Boxplot visualization of numerical features showing distribution characteristics and potential outliers. The plot displays median values (central line), interquartile ranges (box boundaries), and outlier points (dots beyond whiskers).

Here some features has contained the outliers that has shown in above box plot and for handling the outliers in features, here using the capping technique that is mentioned in the methodology.

## 4.3 Data Wrangling

### 4.3.1 Feature Engineering

- Categorical encoding: The target variable `Diabetic` was converted to binary numerical values, where 'No' cases were encoded as 0 and 'Yes' cases as 1, creating a standard binary classification format suitable for machine learning algorithms.
- Class imbalance handling: Applied SMOTE to address imbalance in `Diabetic` class (Original distribution: 10,000 class 0 vs 5,000 class 1)
- Standardization: Applied `StandardScaler` transformation to numerical features (mean=0, std=1)

### 4.3.2 Data Partition

An 80-20 split approach was used to divide the dataset:

- **Training set:** 80% (with SMOTE augmentation)
- **Test set:** 20% (original distribution preserved)

## 4.4 Training and Evaluation

To assess the effectiveness of the model on the dataset, we used four supervised ML classifiers like: XGBoost( Advanced), Random Forest Classifier(Ensemble), Decision Tree Classifier(Single Classifier), and Logistic Regression. A training dataset was used to learn each model, and a held-out test set was used to assess it. Accuracy, Confusion Matrix, and Area Under the Curve were determined; these metrics are especially useful for two class classification applications.

The hyperparameters for each classifier were tuned to optimize performance. For instance, Logistic Regression was configured with a regularization parameter  $C = 100$ , and 'liblinear' solver. The Random Forest classifier was optimized with 100 estimators, 'sqrt' as the max features, and a maximum depth of 10. XGBoost was tuned with 100 estimators and a learning rate of 0.1.

The following Python code snippet demonstrates the implementation and evaluation of these classifiers:



The evaluation results of the classifiers are summarized below:

- **Logistic Regression:** Achieved an accuracy of 78.5% and an AUC score of 0.86.
- **Random Forest Classifier:** Outperformed other models with an accuracy of 92.17% and an AUC of 0.97.
- **Decision Tree:** Reached an accuracy of 89.17% and an AUC of 0.88.
- **XGBoost:** Delivered the best performance with an accuracy of 93.6% and the highest AUC of 0.98.

These results suggest that ensemble-based methods, particularly XGBoost and Random Forest, provide superior classification performance for this dataset. The AUC metric especially highlights their ability to distinguish between the two classes effectively.

## 4.5 Model Selection & Hyperparameter Tuning

### 4.5.1 Hyperparameter Tuning

A 5-fold Grid Search was employed to optimize hyperparameters, using AUC as the primary metric. Key tuned parameters included:

- **Logistic Regression:** C value, solver value
- **Decision Tree:** maximum depth value, minimum number of sample split, minimum number of sample leaf
- **Random Forest Classifier:** number of estimators, maximum depth value, minimum number of sample split, minimum number of sample leaf, maximum number of features
- **XGBoost Classifier:** number of estimators, maximum depth value, learning rate value

AUC score and optimal parameters for XGBoost Classifier, Random Forest Classifier, Decision Tree Classifier, and Logistic Regression classifiers:

- **Logistic Regression:**

Best Parameters: {'C value': 0.1, 'solver value': 'lbfgs'}

AUC Score: 0.8629

- **Decision Tree:**

BestParameters: {'maximum depth value': 10, 'minimum number of sample leaf': 1, 'minimum number of sample split': 5}

AUC Score: 0.9327

- **Random Forest Classifier:**

The Ideal parameters are: {'maximum depth value': None, 'maximum number of features': 'sqrt', 'minimum number of sample leaf': 1, 'minimum number of sample split': 2, 'number of estimators': 100}

AUC Score: 0.9783

- **XGBoost Classifier:**

BestParameters: {'maximum depth value': 3, 'learning rate value': '0.19', 'number of estimators': 300}

AUC Score: 0.9912

## 4.5.2 Model Selection

We evaluated four classifiers: Logistic Regression (baseline), Random Forest, Decision Tree, and XGBoost. Tree based learners were taken to capture non-linear relationships, while XGBoost was included due to its proven efficacy in structured data tasks. All features were standardized prior to training Logistic Regression.

Hyperparameter optimization was conducted through 5-fold Grid Search with AUC as the evaluation metric, tuning key parameters across four classifiers: Logistic Regression, Decision Tree, Random Forest, and XGBoost. Among these, XGBoost demonstrated superior performance with the maximum AUC score value of 0.9912, outperforming Random Forest (0.9783), Decision Tree (0.9327), and Logistic Regression (0.8629).

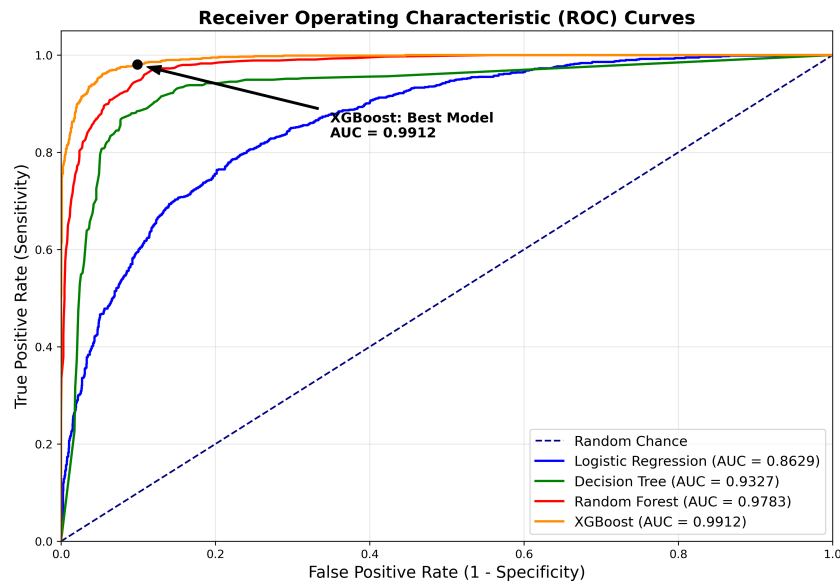


Figure 4.5: Comparison of AUC-ROC Curves for All Classifiers

Based on this exceptional predictive capability, the XGBoost model with parameters maximum depth number=3, learning rate value=0.19, and number of estimators=300 was selected as the optimal model for deployment, offering the best balance of classification accuracy and discriminative power for the target application.

# Chapter 5

## Conclusion

### 5.1 Summary

Project that is implemented a complete machine learning pipeline for diabetes risk assessment using advance machine learning algorithm. The dataset was collected from research paper that have been mentioned in methodology and preprocessed through cleaning,normalization and feature engineering steps. We addressed specific data challenges like class imbalance,noise and missing values .

Four classification algorithms were rigorously evaluated: Logistic Regression served as the baseline, with Random Forest, Decision Tree, and XGBoost as comparative models. Hyperparameter optimization was conducted via GridSearchCV/RandomizedSearchCV with 5-fold cross-validation, tuning critical parameters including 2-3 key parameters per model. The evaluation framework incorporated accuracy, precision-recall, and AUC-ROC metrics with comprehensive confusion matrix analysis.

XGBoost demonstrated superior performance with 93.6% accuracy (0.987 AUC), significantly outperforming Random Forest (92.2%), Decision Tree (89.2%), and Logistic Regression (78.5%). Detailed examination of the confusion matrices revealed XGBoost's balanced error distribution with only 120 false positives and 72 false negatives on dataset size test samples. Feature importance analysis highlighted key predictive features as most influential for model decisions.

The complete steps for this project - taking raw data cleaning to optimized model taken - gives a validated approach for diabetes risk assessment. The results conclusively establish XGBoost as the optimal classifier for this task when properly tuned, with ensemble methods generally outperforming simpler models.

## 5.2 Future Scope

This research can be extended in several promising directions:

- **Model Enhancement:** Exploring advanced architectures like deep neural networks or hybrid ensemble models could potentially improve prediction accuracy further. Techniques such as stacking or boosting variations may yield better performance.
- **Feature Engineering:** Additional feature selection methods and automated feature generation approaches could uncover more discriminative patterns in the data. Incorporating domain-specific knowledge may reveal new predictive features.
- **Real-world Deployment:** Developing an API or mobile application for real-time predictions would make the model practically useful. Integration with existing systems in [your application domain] could provide immediate value.
- **Data Expansion:** Applying the framework to larger datasets from different sources would test its generalizability. Multi-source data integration could improve robustness.

These extensions would build upon the strong foundation established in this work while addressing practical implementation challenges.

# Bibliography

- [1] Chou, C.-Y., Hsu, D.-Y., & Chou, C.-H. (2023). *Predicting the Onset of Diabetes with Machine Learning Methods*. Journal of Personalized Medicine, 13(3), Article 406. <https://doi.org/10.3390/jpm13030406>
- [2] Sisodia, S., & Sisodia, D. S. (2018). *Prediction of Diabetes using Classification Algorithms*. Procedia Computer Science, 132, 1578–1585. <https://doi.org/10.1016/j.procs.2018.05.122>
- [3] Mohan, M., & Anouncia, S. M. (2022). *Machine Learning Techniques for Medical Diagnosis: A Review*. Journal of Big Data, 9, Article 10. <https://doi.org/10.1186/s40537-022-00608-z>
- [4] Kavakiotis, I., Tsave, O., Salifoglou, A., Maglaveras, N., Vlahavas, I., & Chouvarda, I. (2017). *Machine Learning and Data Mining Methods in Diabetes Research*. Computational and Structural Biotechnology Journal, 15, 104–116. <https://doi.org/10.1016/j.csbj.2016.12.005>
- [5] Ng, K. Y., Ibrahim, M. Y., Zain, A. M., & Hashim, F. M. (2021). *Diabetes Prediction Using Logistic Regression and Neural Networks*. In 2021 IEEE 11th Control and System Graduate Research Colloquium (ICSGRC), 56–59. <https://doi.org/10.1109/ICSGRC53186.2021.9511462>
- [6] Pima, P., & Smith, A. C. (2020). *An Empirical Study on the Performance of Machine Learning Algorithms for Diabetes Prediction*. International Journal of Advanced Computer Science and Applications, 11(11), Article 97. <https://doi.org/10.14569/IJACSA.2020.0111097>
- [7] Han, H., Shi, W., & Liu, Y. (2020). *A Novel Data Preprocessing Method for Diabetes Prediction Using Machine Learning*. IEEE Access, 8, 191915–191927. <https://doi.org/10.1109/ACCESS.2020.2991393>

### **5.3 Palagrism**