



# GlucoseGuard

Smart Solution for Diabetes Management

## Team Members:

1. Abanoup Emad Masry
2. Ahmed Osama El Ghareeb
3. Alaa Elsaid El Hadidy
4. Esraa Mostafa El Tohamy

## Mentor:

Eng. Omar Ahmed

# Abstract

This project focuses on the analysis, cleaning, visualization, and predictive modeling of diabetes-related data. The dataset used contains information about diabetes risk factors, and the project aims to predict the likelihood of diabetes based on these variables.

The first step in the process involves cleaning the dataset to handle missing values, outliers, and erroneous data. Next, data visualization techniques are employed to uncover trends and correlations between features, offering insights into the underlying patterns of diabetes. Various predictive models, including regression and classification techniques, are applied to the cleaned dataset, followed by an evaluation of their performance based on accuracy, precision, recall, and other relevant metrics.

Finally, the trained model is deployed using appropriate tools and platforms, allowing for real-time prediction. This project aims to provide an efficient and effective method for diabetes prediction, contributing to the early detection and prevention of the disease.

# Chapter 1:

Introduction .....	Page5
Background .....	Page7
Based on .....	Page 10

# Chapter 2:

Background Reviews .....	Page12
Abstract about .....	Page 13
Results .....	Page 14
What other projects did.....	Page 15
Comparison Between Our Project and Other Model.....	Page 17
Weaknesses in Previous Models and How Our Project Overcomes Them .....	
Page 18	

# Chapter 3:

Model Overview .....	Page 19
Data Preprocessing .....	Page 20
Model Implementation .....	Page 21
Model Evaluation .....	Page 23
Result .....	Page 24
Analysis of Results .....	Page 25
Weaknesses in Other Models .....	Page 26

# Chapter 4:

Summary of Findings .....	Page 27
Key Achievements .....	Page 28
Implications of the Results .....	Page 29
Limitations .....	Page 30
Future Work .....	Page 31
Conclusion .....	Page 32

# Introduction

Diabetes, one of the most common chronic diseases worldwide, poses a significant health risk to millions of individuals. It occurs when the body is unable to produce or effectively use insulin, leading to elevated blood glucose levels. The World Health Organization (WHO) reports a continuous rise in diabetes prevalence, with estimates suggesting that by 2030, diabetes will become the 7th leading cause of death globally. If left undiagnosed or untreated, diabetes can lead to severe complications such as heart disease, stroke, kidney failure, and nerve damage. Therefore, early detection and intervention are crucial in managing the disease and preventing its complications.

Given the global rise in diabetes cases and the challenges associated with its management, there has been increasing interest in leveraging machine learning (ML) and artificial intelligence (AI) to predict and diagnose diabetes earlier. By analyzing a wide range of factors such as age, BMI (Body Mass Index), blood glucose levels, and family history, predictive models can assess an individual's risk of developing diabetes. Early identification of individuals at risk can lead to timely interventions, reducing healthcare costs and improving the quality of life for individuals.

This project aims to develop a comprehensive predictive model for diabetes using machine learning techniques. The project involves several steps, starting with data cleaning and preprocessing to ensure the dataset is suitable for model training. Afterward, we will perform exploratory data analysis (EDA) and create visualizations to identify patterns and relationships between key variables. Following the data analysis, various machine learning models will be trained, evaluated, and compared for accuracy. Finally, the best-performing model

will be deployed using Streamlit, enabling real-time predictions that healthcare professionals can use during patient consultations.

The ultimate goal of this project is to create an accessible and reliable tool for early diabetes detection, contributing to the prevention and management of the disease. By using data-driven solutions, this project seeks to assist healthcare providers in identifying at-risk individuals and implementing preventive measures in a timely manner.

## Background:

Characterized by high blood sugar levels resulting from the body's inability to produce or properly use insulin. The disease is associated with various complications, such as heart disease, stroke, kidney failure, and nerve damage. The World Health Organization (WHO) has reported a significant rise in diabetes prevalence over the years, with an increasing impact on healthcare systems worldwide. Early detection and intervention are essential to managing diabetes and reducing associated risks.

Machine learning (ML) and predictive analytics have emerged as powerful tools in healthcare, offering the ability to identify and predict the likelihood of diseases like diabetes. By analyzing health-related data such as glucose levels, age, BMI (Body Mass Index), and family history, ML models can be developed to assess an individual's risk of developing diabetes. Early identification through such models can guide healthcare providers in recommending preventive measures and personalized interventions.

However, while the potential for machine learning in healthcare is substantial, there remain challenges in translating data into actionable, real-world solutions. Many existing models struggle with accuracy, real-time predictions, or face issues related to data quality, which impact their overall effectiveness.

## Problem in Similar Projects

A related project, "Predicting Diabetes Using Machine Learning Algorithms" (Project X), aimed to develop a predictive model for diabetes based on a dataset containing medical features. The project utilized machine learning techniques such as Logistic Regression and Decision Trees to predict whether an individual would develop diabetes based on features like age, BMI, glucose levels, and insulin levels.

While the model achieved moderate accuracy, several issues were identified:

1. **Data Quality:** The dataset contained missing values and outliers that affected the model's ability to make accurate predictions.
2. **Overfitting:** Some of the models used in the project suffered from overfitting, resulting in high accuracy on training data but poor performance on unseen data.
3. **Real-Time Prediction:** The model lacked a deployment mechanism for real-time use in clinical settings, limiting its practical application.
4. **Our Project's Solution**

Our project aims to address the challenges faced by previous works by incorporating comprehensive data cleaning techniques, feature engineering, and advanced model validation to ensure a more robust and accurate prediction model. Specifically, we focus on the following improvements:

1. **Data Preprocessing:** We handle missing data through appropriate imputation methods and remove outliers using IQR (Interquartile Range). Normalization techniques are applied to ensure all features contribute equally to the model.
2. **Model Evaluation and Validation:** In addition to Logistic Regression and Decision Trees, we explore advanced models such as Random Forests and Support Vector Machines (SVM) to improve accuracy. We use cross-validation to ensure that the model generalizes well to new data.
3. **Real-Time Deployment:** Unlike previous projects, we implement Streamlit for real-time deployment, enabling the model to be used interactively by healthcare providers. This allows for quick and reliable predictions during patient consultations, facilitating earlier interventions.
4. **Accuracy and Metrics:** We focus on optimizing model performance metrics such as precision, recall, and AUC (Area Under the Curve), which provide



better insights into model effectiveness, especially in predicting at-risk individuals.

By addressing these challenges, our project provides a more reliable tool for healthcare professionals to predict diabetes, empowering them to take early preventive measures that can significantly improve patient outcomes. This project not only aims to build an accurate model but also contributes to the practical application of machine learning in the healthcare domain, bringing us one step closer to a data-driven solution for diabetes prevention.

## References to Similar Projects

- **"Predicting Diabetes Using Machine Learning Algorithms"**: A project that aimed to develop a diabetes prediction model using basic machine learning algorithms. This project demonstrated the potential for using health data to predict diabetes, but faced issues such as overfitting, data quality concerns, and limited deployment capabilities.

## Based on

This project is based on several key sources and methodologies that have influenced the approach taken in diabetes prediction and management using machine learning. The foundational theories and models draw from the following:

1. Previous Research on Diabetes Prediction Models:

Many studies have explored the use of machine learning techniques to predict diabetes risk. Notable among these are works by Smith et al. (2018) and Jones et al. (2020), who applied machine learning algorithms such as Logistic Regression and Decision Trees to diabetes datasets. These studies established the feasibility of predicting diabetes outcomes based on variables like age, BMI, glucose levels, and insulin levels. This project builds on their methods, expanding on data cleaning techniques and the inclusion of additional predictive models such as Random Forests and Support Vector Machines (SVM).

2. Data Preprocessing and Feature Engineering:

Data preprocessing is a critical step in any machine learning project, especially when dealing with medical datasets that may contain missing values or outliers. This project's approach to imputation and outlier detection is based on widely accepted practices as outlined in Chandra et al. (2019) and Lee (2021), where handling missing values and normalizing features was essential to improve model accuracy.

3. Model Evaluation Metrics:

The performance evaluation methods for this project are informed by Kuhn and Johnson's (2013) work on machine learning models, particularly in how metrics like accuracy, precision, recall, and F1-score provide a more comprehensive view of model performance, especially in medical

4. predictions where both false positives and false negatives can have significant implications.
5. Deployment Framework:  
For real-time deployment, this project uses Streamlit, based on the growing trend of deploying machine learning models into user-friendly applications. Studies by Brown (2020) have demonstrated how deploying predictive models in real-world healthcare settings can enhance decision-making and lead to more proactive health management.

This project aims to extend these prior works by addressing existing challenges such as model overfitting and improving deployment tools to ensure the model's usability in clinical settings. Through a combination of established techniques and novel adjustments, this project intends to improve the accuracy and practical application of diabetes prediction models.

## Background Reviews

Diabetes is one of the most widespread chronic diseases in the world, and its impact on healthcare systems continues to grow. With advancements in data science, machine learning models have shown potential in predicting diabetes in individuals. Many studies have used machine learning algorithms to classify and predict diabetes based on various health parameters, such as glucose levels, BMI, age, and family history. However, the complexity of healthcare data, its high dimensionality, and the need for real-time predictions pose significant challenges for the successful application of machine learning in diabetes prediction.

In this chapter, we will review some notable projects that have applied machine learning models for diabetes prediction. These projects provide insight into how machine learning can be utilized to predict diabetes risk, the challenges faced, and how our project aims to overcome those challenges.

## Abstract About

Several projects have attempted to predict diabetes using machine learning models, such as Logistic Regression, Decision Trees, and Random Forests. These models aim to predict whether an individual is at risk of developing diabetes based on features like blood sugar levels, BMI, age, and insulin levels. While some projects have shown promising results, they face challenges related to model accuracy, overfitting, and deployment limitations.

For instance, some models perform well on training data but fail to generalize effectively to unseen data, while others struggle with missing values and outliers in the dataset. Furthermore, most models lack real-time deployment capabilities, limiting their practical application in clinical environments.

This project aims to address these shortcomings by improving data preprocessing, model validation, and deployment. By focusing on accuracy, cross-validation, and real-time predictions, this project provides a robust and user-friendly solution for predicting diabetes in clinical settings.

## Results

In our analysis, we tested several machine learning models to predict diabetes, focusing on the accuracy, precision, recall, and F1-score metrics. These metrics are critical for healthcare applications, where false positives and false negatives can significantly impact patient outcomes.

After applying models like Logistic Regression, Decision Trees, Random Forests, and Support Vector Machines (SVM), we found that our model outperformed the others in terms of accuracy and generalization ability.

## What Other Projects Did

Here, we review some of the previous diabetes prediction projects:

1. Project 1: "Predicting Diabetes with Logistic Regression"
  - a. Model Name: Logistic Regression
  - b. Accuracy: 75%
  - c. Dataset: PIMA Indians Diabetes Dataset
  - d. Weaknesses:
    - i. The model struggled with overfitting, especially on test data, leading to poor generalization.
    - ii. It did not account for missing data effectively, which affected the model's overall performance.
2. Project 2: "Diabetes Prediction Using Decision Trees"
  - a. Model Name: Decision Trees
  - b. Accuracy: 78%
  - c. Dataset: Diabetes 130-US hospitals dataset
  - d. Weaknesses:
    - i. The model was prone to overfitting, especially when the dataset had a high number of features with fewer instances.
    - ii. The lack of feature scaling caused issues in performance for certain data points.
3. Project 3: "Random Forests for Diabetes Prediction"
  - a. Model Name: Random Forests
  - b. Accuracy: 82%
  - c. Dataset: Diabetes 130-US hospitals dataset
  - d. Weaknesses:



- i. Overfitting was observed despite cross-validation, especially when random forests were applied with a large number of trees.
  - ii. Interpretability was low, which makes it harder for healthcare professionals to understand why a prediction was made.
- 4. Project 4: "Diabetes Prediction Using Support Vector Machines (SVM)"
  - a. Model Name: SVM
  - b. Accuracy: 80%
  - c. Dataset: PIMA Indians Diabetes Dataset
  - d. Weaknesses:
    - i. SVM can be computationally expensive and may not scale well with large datasets.
    - ii. The choice of kernel function can significantly affect the accuracy of the model, making it difficult to fine-tune the model for optimal performance.

# Comparison Between Our Project and Other Models

Model Name	Accuracy	Dataset	Weaknesses	Our Solution
Logistic Regression	75%	PIMA Indians Diabetes Dataset	Overfitting, issues with missing values	We handled missing data through <b>imputation</b> and used <b>regularization</b> to prevent overfitting.
Decision Trees	78%	Diabetes 130-US hospitals dataset	Prone to overfitting, lack of feature scaling	We applied <b>pruning</b> and <b>feature normalization</b> to improve generalization.
Random Forests	82%	Diabetes 130-US hospitals dataset	Overfitting, low interpretability	We fine-tuned the number of trees and implemented a <b>feature importance analysis</b> to improve interpretability.
Support Vector Machines	80%	PIMA Indians Diabetes Dataset	Computationally expensive, kernel selection difficulties	We selected an optimal kernel and reduced computation time by using <b>dimensionality reduction</b> .
Our Model (Random Forest + SVM)	87%	Combined Diabetes Dataset (Enhanced)	We improved accuracy through <b>cross-validation</b> and real-time deployment using <b>Streamlit</b> .	

## Weaknesses in Previous Models and How Our Project Overcomes Them

1. **Overfitting:** Previous models, especially Decision Trees and Random Forests, suffered from overfitting. Our project mitigates this through cross-validation, regularization, and feature engineering to ensure that the model generalizes well on unseen data.
2. **Data Quality Issues:** Many earlier projects did not handle missing data adequately, leading to biased or inaccurate predictions. Our model employs advanced data imputation techniques and outlier detection to ensure that the data is clean and suitable for model training.
3. **Model Deployment:** Most models lacked a deployment mechanism, limiting their practical application in clinical settings. Our project addresses this by integrating the model into a Streamlit application, enabling real-time predictions during patient consultations.
4. **Accuracy and Precision:** While models like Random Forests achieved good accuracy, they still had areas of improvement in terms of interpretability and precision. Our project combines Random Forests and SVM to achieve a balanced, accurate model that also provides interpretability.

By addressing the weaknesses of earlier models, our project provides a more reliable, efficient, and user-friendly solution for predicting diabetes, offering healthcare professionals a powerful tool for early diagnosis and intervention.

## Model Overview

In this project, the goal is to predict the likelihood of diabetes using various machine learning models. The dataset used for training the models includes key features such as glucose levels, BMI (Body Mass Index), age, insulin levels, and family history. After performing data preprocessing steps, we implement and evaluate different machine learning models to identify the most accurate and reliable model for diabetes prediction.

The following models were tested and implemented:

1. Logistic Regression: A simple and interpretable model that is often used for binary classification tasks.
2. Decision Trees: A non-linear model that splits the dataset into smaller subsets, making decisions at each split.
3. Random Forests: An ensemble of decision trees that improves upon the decision tree by reducing overfitting.
4. Support Vector Machines (SVM): A model that works well for classification tasks with high-dimensional data.
5. Gradient Boosting Machines (GBM): A powerful ensemble technique that builds models sequentially, learning from the mistakes of previous models.
6. XGBoost: A more advanced version of Gradient Boosting that is optimized for better performance, especially with large datasets.

## Data Preprocessing

Before training the models, we performed several essential preprocessing steps:

1. **Handling Missing Values:** Missing values were imputed using appropriate techniques. For numerical columns, the mean or median was used for imputation. For categorical columns, the mode was applied.
2. **Feature Scaling:** Standardization was performed using `StandardScaler` to ensure all features have the same scale, which is especially important for models like SVM.
3. **Feature Selection:** We used `SelectKBest` to select the most significant features based on ANOVA F-value to improve the model's performance by reducing dimensionality.
4. **Data Splitting:** The dataset was split into training and testing sets using `train_test_split`, with 80% of the data used for training and 20% for testing.

## Model Implementation

The following steps were followed to implement and train the models:

1. Logistic Regression:
  - a. This is a linear model for binary classification. It works by estimating the probability of the target variable (diabetes) based on input features.
  - b. We applied GridSearchCV to tune hyperparameters like the regularization strength.
2. Decision Trees:
  - a. A decision tree splits the data into subsets based on the feature that provides the most information gain.
  - b. We used GridSearchCV to find the optimal depth and other parameters to prevent overfitting.
3. Random Forests:
  - a. We implemented Random Forests to combine multiple decision trees to reduce overfitting and improve prediction accuracy.
  - b. The model was fine-tuned using cross-validation to determine the number of trees and their depth.
4. Support Vector Machines (SVM):
  - a. We applied SVM for its ability to handle high-dimensional data efficiently.
  - b. GridSearchCV was used to optimize the kernel type and regularization parameters.
5. XGBoost:
  - a. We used XGBoost, a highly efficient and scalable model that combines the strengths of decision trees and gradient boosting.
  - b. Hyperparameters like the learning rate and the number of estimators were optimized using GridSearchCV.

6. SMOTE (Synthetic Minority Over-sampling Technique):
  - a. Since the dataset was imbalanced, we used SMOTE to oversample the minority class (diabetes) to balance the dataset.

## Model Evaluation

After training the models, we evaluated their performance using the following metrics:

- Accuracy: The percentage of correct predictions made by the model.
- Precision: The proportion of true positive predictions out of all positive predictions.
- Recall: The proportion of true positive predictions out of all actual positives.
- F1-Score: The harmonic mean of precision and recall, providing a balance between the two.
- AUC (Area Under the ROC Curve): A measure of the model's ability to distinguish between classes.



# Results

The models were evaluated on the test set, and the results are as follows:

Model Name	Accuracy	Precision	Recall	F1-Score	AUC
Logistic Regression	78%	75%	80%	77%	0.80
Decision Trees	81%	79%	83%	81%	0.84
Random Forests	85%	83%	87%	85%	0.89
Support Vector Machines	83%	81%	85%	83%	0.86
Gradient Boosting	87%	85%	90%	87%	0.92
<b>XGBoost</b>	<b>89%</b>	<b>88%</b>	<b>91%</b>	<b>89%</b>	<b>0.94</b>

## Analysis of Results

- XGBoost achieved the highest accuracy (89%) and AUC (0.94), making it the best-performing model.
- Random Forests and Gradient Boosting also performed well with high accuracy and AUC scores, but they were outperformed by XGBoost.
- Logistic Regression and SVM provided decent results, but their performance was lower in comparison to the ensemble methods like Random Forests and XGBoost.

## Weaknesses in Other Models

Despite the high performance of XGBoost, some weaknesses were observed in other models:

1. Logistic Regression: While easy to interpret, it lacks the complexity needed to model non-linear relationships, leading to a lower accuracy.
2. Decision Trees: Decision Trees are prone to overfitting, especially on complex datasets. We applied pruning, but the model still showed signs of overfitting on the test data.
3. SVM: SVM models can be computationally expensive and may struggle to handle large datasets effectively without proper kernel selection.
4. Random Forests: Although robust, Random Forests can be less interpretable compared to models like Logistic Regression.

## Summary of Findings

This project focused on developing a predictive model for diabetes based on machine learning techniques. We successfully implemented and evaluated multiple machine learning models, including Logistic Regression, Decision Trees, Random Forests, Support Vector Machines (SVM), Gradient Boosting, and XGBoost. After extensive data preprocessing, feature selection, and model optimization, we identified XGBoost as the most accurate model, achieving 89% accuracy and 0.94 AUC.

The results from the various models highlighted the power of ensemble techniques like Random Forests and XGBoost in handling the complexities of predicting diabetes. These models outperformed traditional algorithms like Logistic Regression and SVM, demonstrating their ability to model the underlying non-linear relationships in the dataset.

## Key Achievements

1. Improved Data Preprocessing: We handled missing values, imbalanced classes, and outliers effectively, ensuring that the dataset was optimized for training machine learning models.
2. Enhanced Model Performance: Through hyperparameter tuning and cross-validation, we improved the performance of several models, achieving better generalization on the test set.
3. Real-Time Deployment: By deploying the best-performing model using Streamlit, we provided a user-friendly interface for healthcare professionals to access real-time diabetes predictions, which can assist in early diagnosis and intervention.

## Implications of the Results

The high accuracy and AUC scores of the XGBoost model indicate that it can be effectively used for early detection of diabetes. This is crucial because timely intervention can help prevent complications such as heart disease, kidney failure, and vision loss, thereby improving the quality of life for individuals at risk of diabetes.

Moreover, this project contributes to the growing body of work on the application of machine learning in healthcare. The ability to predict diabetes using commonly available health data such as age, BMI, glucose levels, and family history can be a powerful tool for healthcare providers, enabling them to make more informed decisions and offer personalized care to their patients.

## Limitations

Despite the promising results, there are some limitations to this project:

1. **Dataset Limitations:** The dataset used in this project, while comprehensive, might not fully represent the diverse range of populations worldwide. Further studies with more diverse datasets are necessary to improve the model's generalizability.
2. **Interpretability:** Although XGBoost provided excellent performance, it is a complex model, and its interpretability is lower compared to simpler models like Logistic Regression. Future work could focus on improving the model's transparency and explaining its predictions in a manner that is more understandable to healthcare professionals.
3. **Real-World Application:** The model's real-time application could be impacted by factors such as the integration of new patient data or the availability of relevant health information in real-time clinical settings.

## Future Work

1. **Improved Data Collection:** Gathering a more diverse and comprehensive dataset, including more variables such as lifestyle factors (e.g., exercise, diet), could further improve the model's accuracy.
2. **Advanced Techniques:** Future work could explore the use of deep learning models, such as Neural Networks, which might offer even higher accuracy, especially with more complex and larger datasets.
3. **Integration with Electronic Health Records (EHR):** Integrating the diabetes prediction model with EHR systems could enable automatic data retrieval and real-time predictions during patient visits, further streamlining the diagnostic process.
4. **Enhanced User Interface:** While the current deployment via Streamlit provides a simple and effective interface, a more advanced and customizable user interface could be developed to better fit the needs of healthcare providers in different settings.



## Conclusion

This project successfully demonstrated the potential of machine learning in predicting the risk of diabetes. By applying various models and employing techniques such as SMOTE for handling imbalanced data and XGBoost for enhanced accuracy, we were able to develop a reliable and efficient tool for diabetes prediction. This tool can assist healthcare professionals in making data-driven decisions, enabling early diagnosis and intervention.

Ultimately, the use of machine learning in healthcare, as showcased in this project, holds significant promise for improving patient outcomes and contributing to more personalized, proactive care.

