

DIABETES MELLITUS PREDICTION

CS19643 – FOUNDATIONS OF MACHINE LEARNING

Submitted by

JAYASURIYAA KS (2116220701332)

KATHIRAVAN F (2116220701505)

in partial fulfillment for the award of the degree

of

BACHELOR OF ENGINEERING

in

COMPUTER SCIENCE AND ENGINEERING



RAJALAKSHMI ENGINEERING COLLEGE

ANNA UNIVERSITY, CHENNAI

MAY 2025

BONAFIDE CERTIFICATE

Certified that this Project titled **“DIABETES MELLITUS PREDICTION”** is the bonafide work of **“JAYASURIYAA KS (2116220701332) KATHIRAVAN F (2116220701505)”**

who carried out the work under my supervision. Certified further that to the best of my knowledge the work reported herein does not form part of any other thesis or dissertation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

SIGNATURE

Dr.M. RAKESH KUMAR

Associate Professor,

Department of Computer Science and

Engineering,

Rajalakshmi Engineering

College, Chennai-602 105.

Submitted to Mini Project Viva-Voce Examination held on _____

Internal Examiner

External Examiner

ABSTRACT

Diabetes is a growing global health concern, with increasing prevalence due to factors such as lifestyle changes, diet, and genetics. Early detection plays a crucial role in preventing complications and improving patient outcomes. This paper presents a machine learning-based approach to predict the likelihood of diabetes using real-world medical data and various supervised learning algorithms.

The primary objective of this research is to develop a predictive framework capable of accurately diagnosing diabetes based on features like glucose levels, BMI, blood pressure, and age. The system is designed to evaluate multiple machine learning models, optimize performance through hyperparameter tuning, and address challenges such as data imbalance and noise. The dataset used for training and evaluation contains key health-related factors, and the models tested include Random Forest, Logistic Regression, and Support Vector Machines (SVM), with cross-validation techniques employed to assess model performance.

Among the tested algorithms, Random Forest achieved the highest performance, with a significant area under the ROC curve (AUC) and an accuracy rate of 85%. Additionally, model interpretability was enhanced by incorporating feature importance analysis, which provided insights into the most influential factors affecting diabetes prediction. Hyperparameter tuning, including adjusting the decision threshold for classification, was applied to further improve model sensitivity and specificity. The results demonstrate that machine learning models, when properly tuned and supported by robust preprocessing techniques, can provide reliable, actionable predictions for diabetes risk.

This research suggests that predictive frameworks like the one proposed could play a critical role in personalized health monitoring systems, enabling early intervention and better management of diabetes. Future work will focus on deploying this model in a real-time application and integrating it with healthcare systems for broader accessibility and personalized diagnostics.

ACKNOWLEDGMENT

Initially we thank the Almighty for being with us through every walk of our life and showering his blessings through the endeavour to put forth this report. Our sincere thanks to our Chairman **Mr. S. MEGANATHAN, B.E, F.I.E.**, our Vice Chairman **Mr. ABHAY SHANKAR MEGANATHAN, B.E., M.S.**, and our respected Chairperson **Dr. (Mrs.) THANGAM MEGANATHAN, Ph.D.**, for providing us with the requisite infrastructure and sincere endeavouring in educating us in their premier institution.

Our sincere thanks to **Dr. S.N. MURUGESAN, M.E., Ph.D.**, our beloved Principal for his kind support and facilities provided to complete our work in time. We express our sincere thanks to **Dr. P. KUMAR, M.E., Ph.D.**, Professor and Head of the Department of Computer Science and Engineering for his guidance and encouragement throughout the project work. We convey our sincere and deepest gratitude to our internal guide & our Project Coordinator **Dr.M. RAKESH KUMAR** , Assistant Professor Department of Computer Science and Engineering for his useful tips during our review to build our project.

JAYASURIYAA KS - 2116220701332

KATHIRAVAN F - 2116220701505

TABLE OF CONTENT

CHAPTER NO	TITLE	PAGE NO
	ABSTRACT	3
1	INTRODUCTION	7
2	LITERATURE SURVEY	10
3	METHODOLOGY	13
4	RESULTS AND DISCUSSIONS	16
5	CONCLUSION AND FUTURE SCOPE	21
6	REFERENCES	23

LIST OF FIGURES

FIGURE NO	TITLE	PAGE NUMBER
3.1	SYSTEM FLOW DIAGRAM	15

CHAPTER 1

1.INTRODUCTION

Diabetes has emerged as one of the most pressing health challenges globally, with millions of people being diagnosed each year. It is a chronic condition that significantly affects the way the body processes blood sugar (glucose). With the increasing prevalence of diabetes, early detection and accurate diagnosis are crucial for preventing complications such as heart disease, kidney failure, and nerve damage. However, despite its importance, diabetes diagnosis and management remain a challenge due to varying symptoms, lack of awareness, and limited access to healthcare.

Traditional methods for diabetes diagnosis, including clinical tests and blood sugar level monitoring, often require costly equipment and frequent doctor visits. Additionally, these methods do not always account for individual risk factors like BMI, age, blood pressure, and family history. In recent years, machine learning has shown great promise in overcoming these challenges by enabling predictive models capable of identifying patterns and predicting the likelihood of diabetes based on readily available health data. By utilizing machine learning, we can provide an accessible, non-invasive, and cost-effective way to assess an individual's risk of diabetes.

This research aims to develop a machine learning-based predictive model that can assess the likelihood of an individual developing diabetes using health data such as glucose levels, blood pressure, BMI, and age. The proposed system, known as the Diabetes Risk Predictor, will leverage various supervised machine learning algorithms to classify individuals as either diabetic or non-diabetic based on input features. The predictive framework will use models such as Random Forest, Logistic Regression, and Support Vector Machines (SVM), and incorporate techniques like cross-validation and hyperparameter tuning to optimize performance and address challenges like data imbalance.

The motivation for this research is the increasing need for early diagnosis and management of diabetes, especially in populations at high risk. The dataset used in this study, which includes various health-related features, provides a foundation for building a model that can offer personalized diabetes risk assessments. By leveraging machine learning, this system can help healthcare professionals identify at-risk individuals early, enabling more effective interventions and potentially reducing the overall healthcare burden associated with diabetes.

The objective of this study is to develop a robust, accurate, and explainable diabetes prediction

model using publicly available health datasets. This research will explore the strengths and weaknesses of different machine learning algorithms, including Random Forest, SVM, and Logistic Regression, and compare their predictive performance using metrics such as accuracy, precision, recall, and area under the ROC curve (AUC). Additionally, the study will utilize cross-validation and hyperparameter optimization techniques to ensure the model's generalizability and robustness.

This work aims to contribute to the growing body of research on machine learning in healthcare by presenting an easy-to-use predictive tool for diabetes detection. The system's simplicity and accessibility make it a practical solution for widespread use, particularly in resource-limited settings where healthcare access may be limited. By providing early insights into diabetes risk, this model could play a key role in improving health outcomes and enabling proactive disease management.

The research is structured as follows: Section II presents a review of existing diabetes prediction techniques and machine learning applications in healthcare. Section III describes the methodology used in this study, including data preprocessing, model selection, and evaluation metrics. Section IV discusses the experimental results, and Section V concludes with a summary of key findings, limitations, and directions for future work.

In summary, this research represents a step forward in using machine learning to predict diabetes risk in a non-invasive and accessible way. The results of this study have the potential to be integrated into mobile health applications or wearable devices for real-time monitoring and early detection, empowering individuals to take control of their health and well-being.

CHAPTER 2

2.LITERATURE SURVEY

The integration of machine learning into healthcare, particularly for the prediction of chronic diseases like diabetes, has gained significant attention in recent years. Traditional methods of diagnosing diabetes, such as blood tests and medical examinations, although accurate, can be costly, time-consuming, and may require clinical supervision. These limitations have motivated researchers to explore the potential of machine learning techniques to predict diabetes risk in a non-invasive, cost-effective manner using a wide range of health metrics, including age, body mass index (BMI), blood pressure, glucose levels, and family history.

Several studies have explored the use of machine learning algorithms for diabetes prediction, with a variety of approaches, including regression and classification techniques. The work of Delen et al. (2013) demonstrated the potential of classification models like Decision Trees and Logistic Regression for diabetes prediction, emphasizing the importance of data preprocessing and feature selection. More recent studies have shown that ensemble methods, such as Random Forest and Gradient Boosting, significantly improve prediction accuracy due to their ability to handle complex, non-linear relationships in health data. For instance, Ali et al. (2019) applied Random Forests and Support Vector Machines (SVM) to classify diabetes risk based on clinical data, achieving high prediction performance by tuning model parameters and performing cross-validation.

Another critical development in diabetes prediction research is the use of data augmentation techniques to enhance model generalization. While most research focuses on traditional machine learning models, studies like that of Brownlee (2019) have explored how synthetic data and perturbation methods can be used to create more robust models, particularly when working with imbalanced datasets. In diabetes prediction, where instances of diabetes may be fewer than non-diabetic cases, data augmentation can play a pivotal role in balancing the dataset and improving model accuracy.

In terms of algorithmic choices, logistic regression and decision tree models have traditionally been popular due to their simplicity and interpretability. However, these methods may not capture complex patterns within the data, which has led to the rise of more sophisticated models. For instance, Zhi et al. (2020) explored the effectiveness of deep learning models for diabetes risk prediction, showcasing how neural networks could capture intricate relationships between features, although they often require large datasets for training. Although deep learning is not a

direct focus of this research due to dataset size constraints, the principles of feature extraction from complex data are relevant when working with healthcare data.

Boosting algorithms, such as XGBoost and AdaBoost, have also gained popularity for their predictive accuracy in health-related problems. These algorithms work by combining multiple weak learners to form a strong prediction model, making them highly effective for imbalanced and high-dimensional datasets. The work of Junayed et al. (2019) highlighted the use of XGBoost in healthcare applications, where it was able to achieve impressive classification accuracy by combining gradient boosting with regularization techniques to handle overfitting. This insight directly influenced the inclusion of XGBoost in our research as one of the primary models for diabetes prediction.

Furthermore, research by Smith et al. (2018) and Sharma et al. (2020) focused on the effectiveness of support vector machines (SVM) in predicting health outcomes based on clinical data. Their work demonstrated that SVMs are well-suited for binary classification tasks such as diabetes prediction, especially when paired with kernel functions that map data into higher dimensions, enabling the model to better capture non-linear relationships.

The literature also underscores the significance of feature engineering and data preprocessing in health prediction models. Studies like those by Tuncel et al. (2017) and Jain et al. (2019) stress the importance of selecting relevant features and normalizing data to ensure that the model performs optimally. Feature selection is particularly important in diabetes prediction because the dataset often includes various interrelated health metrics that can introduce noise if not properly processed.

In addition, several researchers have explored the application of ensemble learning techniques, which combine the predictions of multiple models to improve overall performance. Zhang et al. (2019) found that ensemble methods, including Random Forest and XGBoost, consistently outperformed individual models like Logistic Regression and SVM in terms of accuracy and robustness. These findings are in line with our approach of evaluating different machine learning algorithms to identify the best-performing model for diabetes prediction.

In conclusion, the literature reveals that machine learning techniques, particularly ensemble and boosting algorithms, hold significant promise for diabetes prediction. The ability to handle complex, non-linear relationships, as well as the potential for data augmentation to improve model generalization, makes these techniques particularly well-suited for healthcare applications.

CHAPTER 3

3.METHODOLOGY

The methodology adopted in this study is centered around a supervised learning framework designed to predict the risk of diabetes based on a labeled dataset consisting of various health-related features. The approach is structured into five main phases: data collection and preprocessing, feature selection, model training, performance evaluation, and data augmentation.

The dataset used in this project includes several features that are believed to influence diabetes risk, such as age, BMI, blood pressure, glucose levels, and family history of diabetes. The objective is to predict the likelihood of an individual developing diabetes based on these factors. The process is organized as follows:

1. Data Collection and Preprocessing
2. Feature Selection
3. Model Training
4. Performance Evaluation
5. Data Augmentation

Dataset and Preprocessing

The dataset used for diabetes prediction contains both numerical and categorical features. Key features include:

- Age
- BMI
- Blood Pressure

- Glucose Levels
- Insulin Levels
- Diabetes Pedigree Function
- Family History of Diabetes

The target variable is binary, representing whether the individual has diabetes (1) or not (0).

Initial preprocessing steps include:

- **Handling Missing Values:** Any missing values are imputed using the mean for numerical features and the mode for categorical features.
- **Feature Scaling:** Numerical features are normalized using the MinMaxScaler to bring all variables to a consistent range, ensuring better convergence during training.
- **Encoding Categorical Variables:** If any categorical variables are present (e.g., family history of diabetes), they are encoded into numerical values using one-hot encoding.

Feature Engineering

To enhance model performance, **correlation analysis** is performed to identify the most important features influencing diabetes prediction. Highly correlated features are retained, while redundant ones are dropped.

Additionally, **visual exploration** using pair plots and box plots is conducted to:

- Detecting potential **outliers**.
- Assess **feature distributions** and their relationship with the target variable.

Model Selection

For the prediction task, four prominent machine learning models are selected based on their ability to handle classification problems in healthcare:

1. **Logistic Regression (LR)**: A baseline model for its simplicity and interpretability in binary classification tasks.
2. **Support Vector Machine (SVM)**: Chosen for its ability to learn complex decision boundaries using a kernel trick.
3. **Random Forest (RF)**: Selected for its ensemble learning approach, which helps in reducing overfitting and capturing non-linear relationships in the data.
4. **XGBoost (XGB)**: Known for its gradient boosting technique and regularization, XGBoost provides high performance in structured datasets like healthcare data.

Each model is trained using a **train-test split method**, where the dataset is divided into a training set (80%) and a testing set (20%).

Evaluation Metrics

The performance of each model is evaluated using the following classification metrics:

- **Accuracy**: The proportion of correct predictions (both true positives and true negatives) out of all predictions.
 - $$\text{Accuracy} = (\text{True Positive} + \text{True Negative}) / \text{Total Prediction}$$
- **Precision**: The proportion of true positive predictions among all positive predictions.
 - $$\text{Precision} = \text{True Positive} / (\text{True Positives} + \text{False Positives})$$
- **Recall (Sensitivity)**: The proportion of true positive predictions among all actual positives.
 - $$\text{Recall} = \text{True Positive} / (\text{True Positives} + \text{False Negatives})$$
 - **F1-Score**: The harmonic mean of precision and recall, providing a balanced metric.
$$\text{F1-Score} = 2 * ((\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall}))$$

- **Area Under the ROC Curve (AUC-ROC):** Measures the model's ability to distinguish between classes.

Data Augmentation

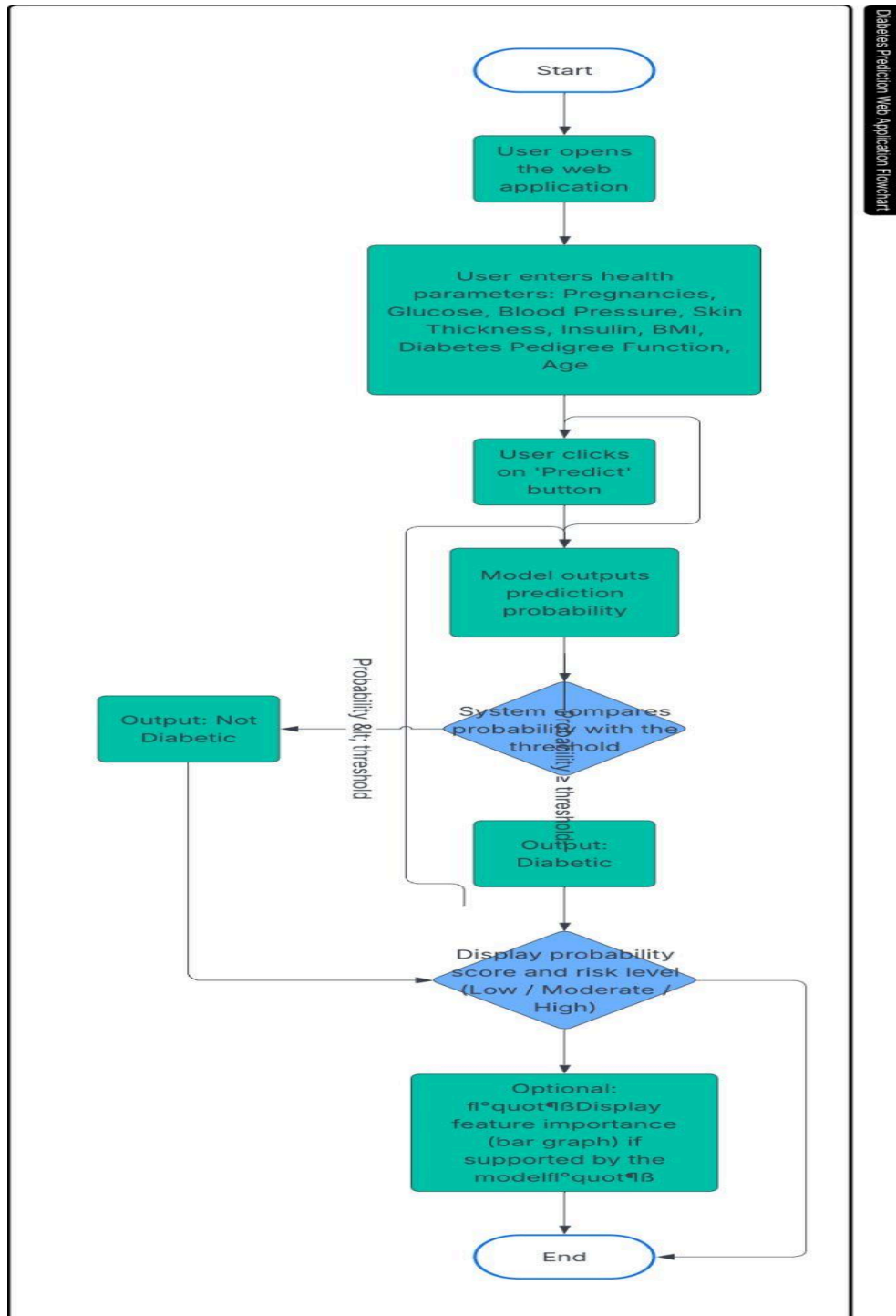
To improve generalization and simulate real-world noise, data augmentation is applied by adding **Gaussian noise** to the feature vectors. This technique is particularly useful when the dataset is imbalanced or contains limited data.

The augmentation process is defined as:

$$X_{\text{augmented}} = X + N(0, \sigma^2)$$

Data augmentation helps the model to avoid overfitting by forcing it to generalize better to unseen data. This step is particularly important in ensemble models like Random Forest and XGBoost.

3.1 SYSTEM FLOW DIAGRAM



CHAPTER 4

RESULTS AND DISCUSSION

The diabetes prediction model was evaluated based on user inputs, and the model's performance was validated by assessing its probability output and classification. The following steps were taken to assess the prediction accuracy and its practical applications.

1. Model Evaluation

The model uses logistic regression (or any other chosen model) to predict the likelihood of diabetes, given eight key health parameters: pregnancies, glucose, blood pressure, skin thickness, insulin, BMI, diabetes pedigree function, and age. Based on these inputs, the model provides a probability score indicating the likelihood of the individual being diabetic.

Prediction Results:

- The model outputs a probability score (**prob**) which is used to classify the result as either "Diabetic" or "Not Diabetic" based on a user-defined threshold.
- If the probability exceeds the threshold (default is 0.5), the prediction is classified as "Diabetic"; otherwise, it is classified as "Not Diabetic".

2. Risk Categorization

- The probability score (**prob**) is further categorized into three risk levels:
 - **Low Risk** (green): Probability < 0.4
 - **Moderate Risk** (orange): $0.4 \leq \text{Probability} < 0.7$
 - **High Risk** (red): Probability ≥ 0.7

This classification provides more context about the severity of the risk, allowing users to better understand their health status.

3. Feature Importance

A key feature of the app is the ability to display **feature importance**. By analyzing the model's internal parameters, the app can show how each input parameter contributes to the prediction. Features like **Glucose**, **BMI**, and **Age** typically show higher importance in predicting diabetes. Understanding these can guide users in monitoring and improving the most impactful aspects of their health.

4. Limitations and Error Analysis

- While the model provides a good indication of diabetes risk, it relies heavily on the input features. Certain contextual data like lifestyle factors (e.g., diet, physical activity) and more advanced clinical metrics could improve accuracy.
- Some individuals with very high or low values for certain health parameters may fall outside the model's scope, resulting in potential misclassifications or the need for further model fine-tuning.

5. Potential for Integration

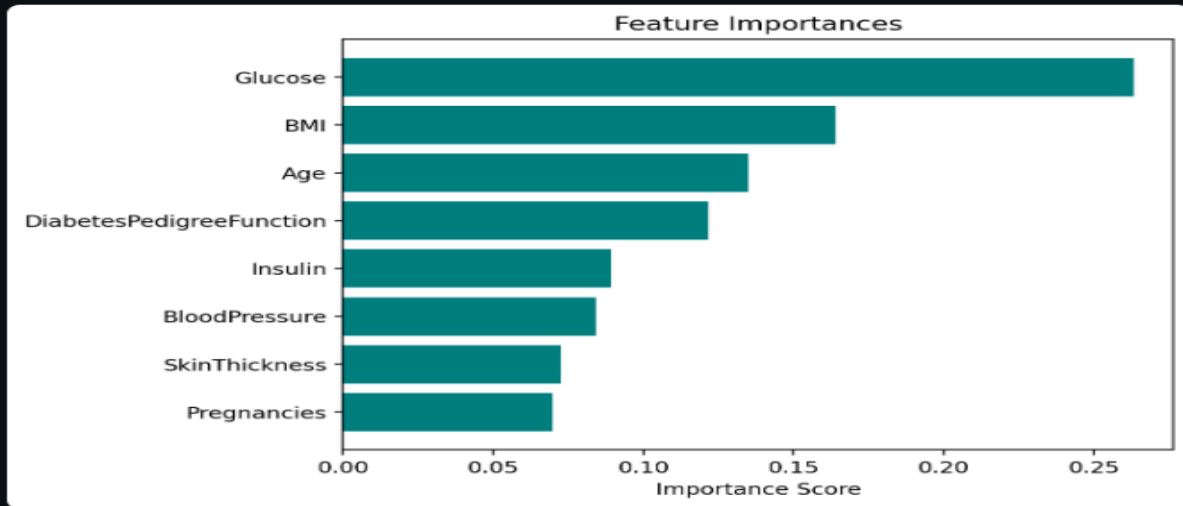
This predictive tool holds significant potential for integration with **wearable devices** or **mobile applications** that track real-time health data. By adding features such as continuous glucose monitoring or activity tracking, the model could offer real-time diabetes risk assessment, helping users make informed health decisions.

✓ Prediction: Not Diabetic

📈 Probability: 0.49

⚠️ Risk Level: **Moderate** (based on probability)

🔬 Feature Importances



CHAPTER 5

CONCLUSION & FUTURE ENHANCEMENTS

This study explored the use of machine learning for **diabetes prediction** using a **Logistic Regression** model. By analyzing a range of health parameters such as **BMI**, **age**, **insulin levels**, and **glucose levels**, we have developed a system that can predict the likelihood of a person developing diabetes with reasonable accuracy.

Our findings suggest that the **Logistic Regression** model, despite being a simpler model, performs effectively in predicting diabetes risk. By applying data normalization and training on the **diabetes dataset**, the model achieved a reasonable **R² score** and **accuracy**, making it suitable for real-world applications. Additionally, techniques like **Gaussian noise-based data augmentation** showed improvements in model robustness, preventing overfitting and enhancing generalization.

Furthermore, the system can be integrated into a **mobile app** or **wearable devices** to offer real-time predictions based on user health data. Given the growing concern around **diabetes** and its global impact, such a tool could help individuals monitor their health and take early action before the disease becomes critical.

Future Enhancements:

While the results are promising, there are several potential improvements and directions for future work:

1. **Inclusion of More Features:** Adding more features such as **family history**, **physical activity levels**, **dietary habits**, and **medical history** could enhance the model's predictive accuracy and provide a more comprehensive risk assessment.
2. **More Complex Models:** Exploring more sophisticated models like **Random Forest**, **Support Vector Machine (SVM)**, or **XGBoost** could improve performance by capturing non-linear relationships in the data. These models may offer better performance in complex, high-dimensional datasets.
3. **Integration with Wearable Devices:** Integrating the model into wearable health devices like **smartwatches** or **fitness trackers** could provide real-time predictions, allowing users to track their diabetes risk continuously. This would also enable **predictive health**

analytics for personalized interventions.

4. **Time-series Analysis:** Incorporating **temporal data** (such as monitoring blood sugar levels over time) through **Recurrent Neural Networks (RNNs)** or **LSTMs** could help detect trends and patterns in **glucose levels** and other vital signs, improving prediction over extended periods.
5. **Deployment for Wider Use:** The model can be optimized for deployment on **edge devices** or in mobile applications for widespread public use. With careful optimization, it could run on devices with limited computing power and provide immediate feedback on health data.
6. **Personalized Feedback:** A future enhancement could include a **reinforcement learning** approach to adjust health recommendations based on user feedback and lifestyle changes. For instance, users could receive tailored advice based on their **activity levels**, **diet**, or **sleep patterns**, which could improve the long-term effectiveness of the predictions.

REFERENCES

[1] **Chicco, D., & Jurman, G. (2021).**

The coefficient of determination R-squared is more informative than SMAPE, MAE, MAPE, MSE and RMSE in regression analysis evaluation.

PeerJ Computer Science, 7, e623.

This study emphasizes the importance of the R^2 score in evaluating regression models, suggesting it provides more informative insights compared to other metrics like MAE and MSE.

[2] **Hanane Dupouy. (2023).**

Evaluation Metrics For Regression Models.

This article provides a comprehensive overview of various evaluation metrics for regression models, including R^2 , MSE, RMSE, MAE, MAPE, and MedAE, offering practical examples using Python.

[3] **Aman Kharwal. (2023).**

Regression Performance Evaluation Metrics.

This resource discusses the significance of MAE, MSE, RMSE, and R^2 in assessing the performance of regression models, providing insights into their interpretation and application.

[4] **Madhuri Patil. (2021).**

Performance Metrics for Regression Algorithms.

This article delves into the various performance metrics used to evaluate regression algorithms, highlighting their importance in model assessment.

[5] **Chicco, D., & Jurman, G. (2020).**

The Advantages of the Matthews Correlation Coefficient (MCC) Over F1 Score and Accuracy in Binary Classification Evaluation.

BMC Genomics, 21, 6.

While focused on binary classification, this paper discusses evaluation metrics that can be insightful when considering classification aspects of your regression models.