

# **Enhance Type 2 Diabetes Prediction with Machine Learning using Medical & Demographic Dataset**

**Submitted By**

**Md Tanjum An Tashrif**

DU Roll: 5173

REG: 1587

Session: 2019-2020

**Shahariar Hossain Mahir**

DU Roll: 5158

REG: 1593

Session: 2019-2020

**Submitted to**

**Utpol Kanti Das**

Lecturer

Department of Computer Science and Engineering  
National Institute of Textile Engineering & Research



**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**

**NATIONAL INSTITUTE OF TEXTILE ENGINEERING & RESEARCH**

**December 29, 2024**

# ABSTRACT

Diabetes mellitus, particularly type 2 diabetes, has become a significant global health challenge, driven by factors such as urbanization, sedentary lifestyles, and dietary changes. Early prediction and management of diabetes are crucial to mitigating its long-term complications. This study proposes a machine learning-based framework for predicting diabetes using a Fully Connected Neural Network (FNN) and a dataset of lifestyle and health metrics, including pregnancy count, BMI, glucose levels, blood pressure, and insulin levels.

The dataset was preprocessed to handle missing values, normalize features, and address class imbalance, ensuring its suitability for training and evaluation. The FNN model, implemented using TensorFlow/Keras, consists of two hidden layers with ReLU activation functions and an output layer with a Sigmoid activation function for binary classification. Key hyperparameters, such as the Adam optimizer and Binary Cross-Entropy Loss, were selected to optimize model performance.

Experimental results demonstrated the effectiveness of the proposed framework, achieving high accuracy, precision, recall, and F1-score. Visualizations, including heatmaps, before-and-after feature scaling graphs, and confusion matrices, provided a comprehensive evaluation of the model's performance.

The study highlights the potential of machine learning in advancing personalized medicine by providing actionable insights into diabetes risk and enabling tailored interventions. Future work will focus on integrating additional data sources, exploring advanced models, and deploying the system in real-world healthcare settings. This framework represents a significant step toward improving diabetes prediction and management, ultimately contributing to better health outcomes.

# CONTENTS

1.	Introduction .....	5
1.1.	Background .....	5
1.2.	Research Challenge .....	5
1.3.	Objectives .....	6
1.4.	Contributions .....	7
2.	Literature Review .....	8
3.	Proposed Framework .....	9
3.1.	Input Dataset .....	10
3.2.	Apply Pre-processing .....	11
3.2.1.	Handling Missing Values .....	11
3.2.2.	Feature Scaling .....	11
3.2.3.	Class Imbalance Management .....	11
3.3.	Apply Machine Learning System .....	12
3.3.1.	Train Machine Learning Model .....	12
3.3.2.	Load the Dataset .....	14
3.3.3.	Preparing the Dataset .....	15
4.	Experiments and Evaluation .....	18
4.1.	Dataset Description .....	18
4.2.	Experimental Result .....	19
4.3.	Evaluation of Framework .....	22
5.	Conclusion and Future Work .....	
5.1.	Conclusion .....	23
5.2.	Future Work .....	24
6.	Bibliography .....	25

## List of Figures

3.3.2. Figure: Class Balance .....	17
3.3.3. Figure: Dataset .....	18
4.2 Figure: Heatmap .....	21
4.2 Figure: Before Feature Scaling .....	22
4.2 Figure: After Feature Scaling .....	22
4.2 Figure: Confusion Matrix .....	23

## List of Diagram

3.1. Block Diagram of Proposed Methodology .....	12
3.3.1. Block Diagram of the Machine Learning Model .....	15

## List of Tables

4.1.1: Sample Dataset of Diabetics features .....	21
4.2. Summary of Results .....	24

# Chapter 1

## Introduction

### 1.1 Background

Diabetes mellitus, particularly type 2 diabetes, has emerged as a significant global health threat in recent years. This rise is largely attributed to factors such as urbanization, sedentary lifestyles, and dietary changes. Early prediction and effective management of type 2 diabetes are essential to prevent its long-term complications, which can severely impact individuals' quality of life. Traditional methods for predicting diabetes often rely on centralized data, which raises privacy concerns and limits the ability to derive insights from diverse populations. However, advancements in Artificial Intelligence (AI) and machine learning offer innovative solutions to address these challenges.

Machine learning models, such as Convolutional Neural Networks (CNNs), Feedforward Neural Networks (FNNs), and Long Short-Term Memory (LSTM)[1] networks, have shown great promise in healthcare applications. These models can analyze complex datasets, identify patterns, and make accurate predictions, making them well-suited for diabetes prediction. By leveraging lifestyle data—such as physical activity, dietary habits, sleep patterns, and stress levels—these models can provide valuable insights into an individual's risk of developing type 2 diabetes.

### 1.2 Research Challenge

Despite the potential of machine learning in healthcare, several challenges remain. Centralized data collection for diabetes prediction raises significant privacy concerns, as sensitive health information is often stored and processed in a single location. Additionally, traditional prediction methods may not fully capture the complex relationships between lifestyle factors and diabetes risk. There is a need for robust and accurate models that can effectively analyze diverse datasets while addressing privacy concerns and providing actionable insights for personalized interventions.

## 1.3 Objectives

This study aims to predict the onset of type 2 diabetes using a lifestyle dataset that includes factors such as physical activity, dietary habits, sleep patterns, and stress levels. By employing basic machine learning models—specifically Convolutional Neural Networks (CNNs), Feedforward Neural Networks (FNNs), and Long Short-Term Memory (LSTM) networks—we seek to develop a model that accurately predicts diabetes risk and provides actionable insights. The specific objectives are:

- To create a robust prediction model using CNNs, FNNs, and LSTMs to analyze lifestyle data.
- To evaluate the performance of these models in predicting diabetes risk.
- To identify key lifestyle factors that contribute to diabetes risk.
- To provide actionable insights for personalized interventions based on individual lifestyle and health data.

These objectives aim to highlight the potential of data-driven approaches in solving classification problems and contribute to the broader understanding of machine learning applications in predictive analytics.

## 1.4 Contributions

- The contributions of this work are as follows:
- **Accurate Prediction Model:** We develop a machine learning-based framework using CNNs, FNNs, and LSTMs to accurately predict the risk of type 2 diabetes based on lifestyle data.
- **Identification of Key Factors:** By analyzing the model's predictions, we identify the most significant lifestyle factors contributing to diabetes risk, enabling targeted interventions.
- **Personalized Insights:** The model provides personalized insights into an individual's diabetes risk, helping healthcare providers and patients make informed decisions about lifestyle changes and treatment plans.
- **Advancement in Diabetes Prediction:** This study contributes to the growing body of research on diabetes prediction by demonstrating the effectiveness of basic machine learning models in analyzing complex lifestyle data.

These contributions demonstrate the practical application of machine learning tools and techniques in addressing classification challenges while ensuring scalability and accuracy.

## Chapter 2

### Literature Review

The analysis of related work reveals a wide range of studies that have utilized various methods and techniques for analyzing healthcare datasets and making predictions. Researchers have developed and implemented numerous prediction models using data mining techniques, machine learning algorithms, or combinations of these approaches.

- Saravana et al.[4] developed a system using Hadoop and MapReduce techniques to analyze diabetic data. This system predicts the type of diabetes and associated risks.
- Aiswarya et al.[5] employed classification techniques, specifically Naïve Bayes and Decision Trees, to study hidden patterns in diabetes datasets. The study compared the performance of both algorithms and demonstrated their effectiveness in diabetes prediction.
- Kumar et al. [7] developed a model using the Random Forest Classifier to forecast diabetes behavior, demonstrating its effectiveness in handling complex datasets.
- Rajesh et al.[8] used the decision tree algorithm to identify hidden patterns in datasets and classify diabetes cases efficiently.
- Humar et al.[9] combined Artificial Neural Networks (ANN) with fuzzy logic to predict diabetes, leveraging the strengths of both techniques.
- Muhammad et al. [11] utilized a combination of the C4.5 decision tree algorithm, Neural Networks, K-means clustering, and visualization techniques to predict diabetes.



## Chapter 3

### Proposed Framework

The proposed methodology involves the application of supervised machine learning techniques for the classification of Iris flower species. The workflow is structured into the following key steps:

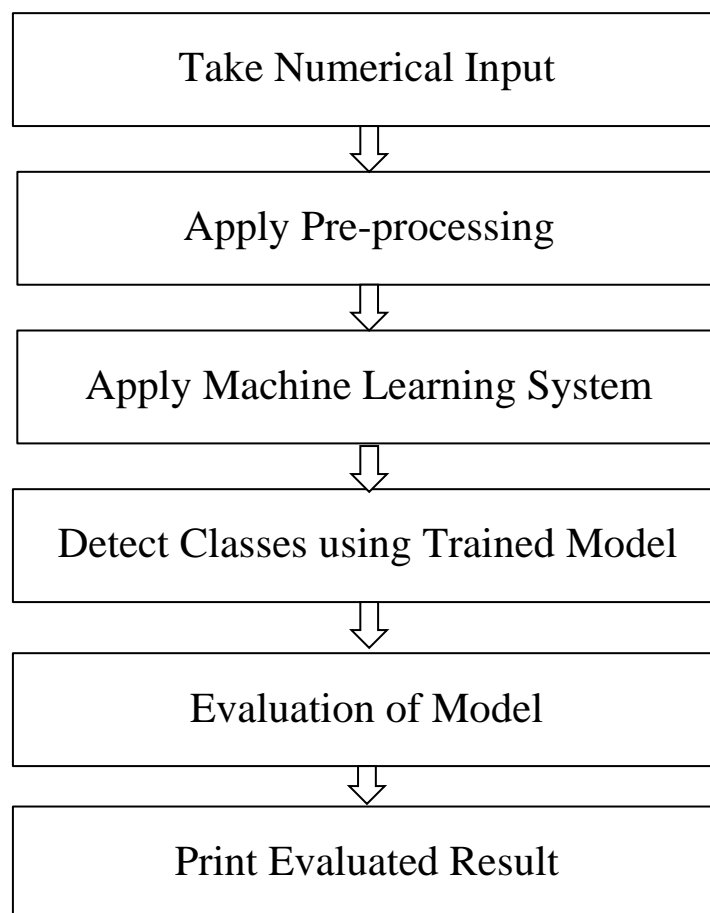


Figure 3.1: Block Diagram of Proposed Methodology



## 3.2. Apply Pre-processing

To ensure the dataset is suitable for analysis and modeling, the following preprocessing steps were applied:

### 3.2.1. Handling Missing Values

Missing or invalid data entries were identified and addressed using appropriate imputation techniques [2].

For numerical columns, the mean or median was used to fill in missing values, ensuring data completeness.

### 3.2.2. Feature Scaling

Since the dataset contains numerical features with varying ranges, Min-Max Scaling or Standard Scaling was applied.

Scaling ensures that all features contribute equally to the model and improves convergence during training.

$$X_{\text{scaled}} = \frac{X - \mu}{\sigma}$$

### 3.2.3. Class Imbalance Management

To address class imbalance, weight distribution was applied during model training.

Higher weights were assigned to the minority class (diabetic class) to prevent the model from becoming biased toward the majority class [3].

$$\text{Class Weight} = \frac{\text{Total Samples}}{2 \times \text{Class Samples}}$$

### 3.3. Apply Machine Learning System

The machine learning system is the core component of this study, where a Fully Connected Neural Network (FNN) is implemented to predict diabetes based on the preprocessed dataset. This section provides a detailed breakdown of the steps involved in building, training, and preparing the dataset for the machine learning model.

#### 3.3.1. Train Machine Learning Model

The training process involves defining the model architecture, selecting appropriate hyperparameters, and optimizing the model for accurate predictions.

##### 1. Model Architecture:

- The FNN consists of the following layers:
  - **Input Layer:** Accepts the preprocessed feature vector as input. The number of neurons in this layer corresponds to the number of features in the dataset.
  - **Hidden Layers:**
    - **First Hidden Layer:** Contains **X neurons** with the **ReLU (Rectified Linear Unit)** activation function. ReLU is chosen for its ability to introduce non-linearity and improve convergence during training.
    - **Second Hidden Layer:** Contains **Y neurons** with the **ReLU activation function**. This layer further extracts complex patterns from the data.
  - **Output Layer:** A single neuron with the **Sigmoid activation function** is used to predict the probability of the binary outcome (diabetic or non-diabetic). The Sigmoid function is ideal for binary classification tasks as it outputs values between 0 and 1.
- The architecture can be summarized as:

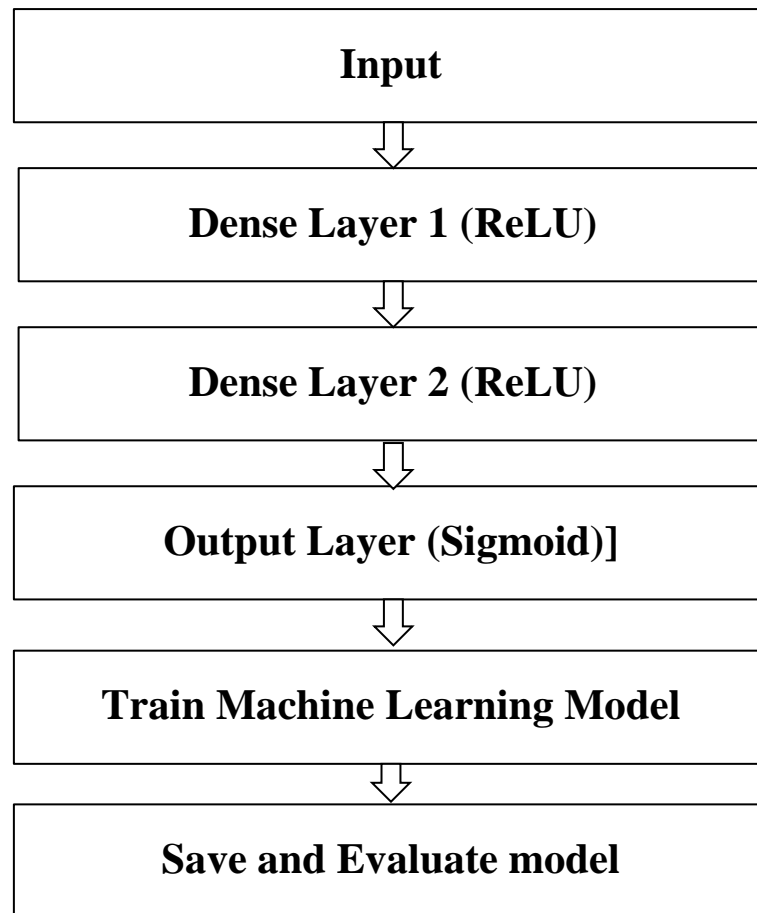


Fig 3.3.1: Block Diagram of the Machine Learning Model

## 2. Hyperparameters:

- **Optimizer:** The **Adam optimizer** is used for training. Adam combines the benefits of Adaptive Gradient Algorithm (AdaGrad) and Root Mean Square Propagation (RMSProp), providing efficient and adaptive learning rates.
- **Loss Function:** **Binary Cross-Entropy Loss** is chosen as the loss function. It measures the performance of the model for binary

classification tasks by comparing the predicted probabilities with the actual labels.

- **Metrics:** The following metrics are used to evaluate the model's performance:
  - **Accuracy:** Measures the overall correctness of the model.
  - **Precision:** Evaluates how many of the positive predictions were correct.
  - **Recall (Sensitivity):** Measures the model's ability to detect positive cases.
  - **F1-Score:** Harmonic mean of Precision and Recall for balanced evaluation.

### 3. Training Process:

- The model is trained using the preprocessed dataset.
- **Batch Size:** An appropriate batch size (e.g., 32 or 64) is chosen to balance computational efficiency and model performance.
- **Epochs:** The model is trained for **N epochs**, with performance monitored for early stopping to prevent overfitting.
- **Class Weights:** To address class imbalance, higher weights are assigned to the minority class (diabetic class) during training.

#### 3.3.2. Load the Dataset

The preprocessed dataset is loaded into the system for training and evaluation.

The dataset is divided into the following components:

1. **Features:** The input variables (e.g., pregnancy count, BMI, glucose levels, blood pressure, insulin levels).
2. **Target Variable:** The binary outcome (1 for diabetic, 0 for non-diabetic).

The dataset is stored in a structured format (e.g., CSV or NumPy array) and loaded using libraries such as Pandas or TensorFlow.

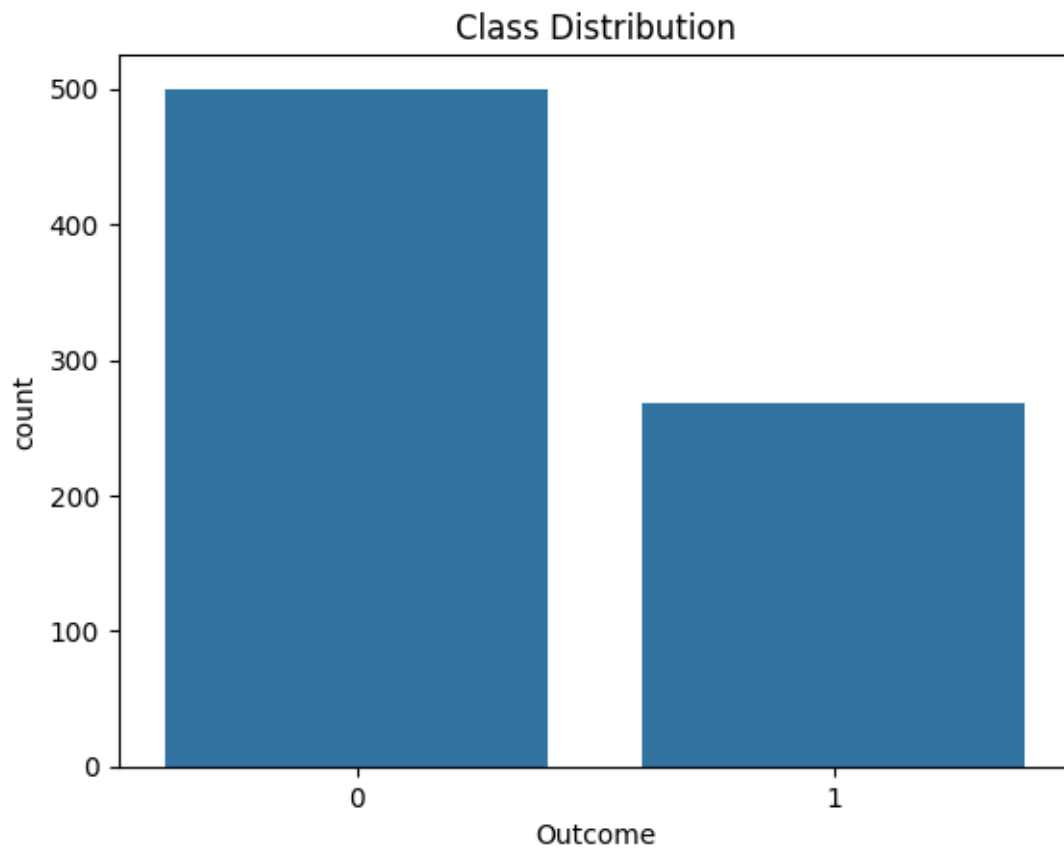


Figure: Class Balance

### 3.3.3. Preparing the Dataset

Before training the model, the dataset is prepared to ensure compatibility with the machine learning system.

#### 1. Train-Validation Split:

- The dataset is split into **training** and **validation** sets in an **80-20 ratio**.
- The training set is used to train the model, while the validation set is used to evaluate its performance during training.

## 2. Batch Preparation:

- The dataset is divided into **batches** of a fixed size (e.g., 32 or 64) for efficient training.
- Batching helps in managing memory usage and improves the convergence of the model.

## 3. Data Shuffling:

- The training data is shuffled before each epoch to ensure that the model does not learn any unintended patterns based on the order of the data.

## 4. Normalization:

- The features are normalized to ensure that all input variables are on a similar scale. This step is crucial for improving the stability and performance of the model.

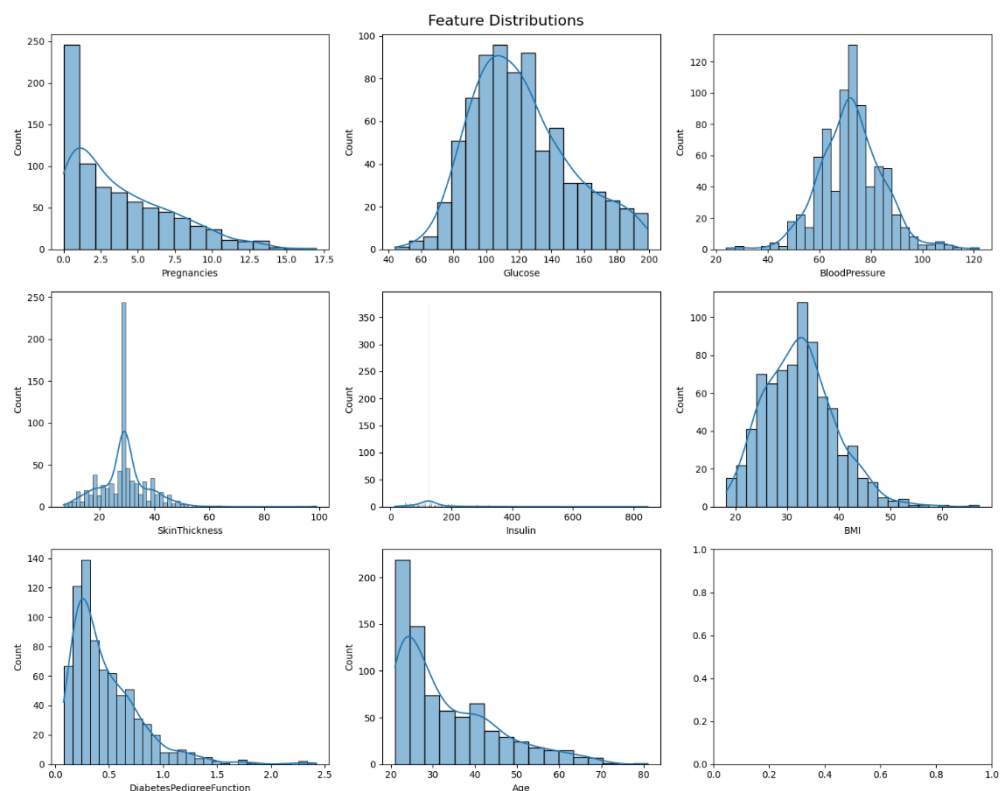


Figure: Dataset



## Additional Technical Details

### 1. Activation Functions:

- **ReLU (Rectified Linear Unit):** Used in the hidden layers to introduce non-linearity and improve the model's ability to learn complex patterns.
- **Sigmoid:** Used in the output layer to produce probabilities between 0 and 1, making it suitable for binary classification.

### 2. Optimization:

- The **Adam optimizer** is used due to its adaptive learning rate capabilities, which help in achieving faster convergence and better performance.

### 3. Loss Function:

- **Binary Cross-Entropy Loss** is used to measure the difference between the predicted probabilities and the actual labels. It is defined as:

$$\text{Loss} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

where  $y_i$  is the actual label and  $\hat{y}_i$  is the predicted probability.

### 4. Performance Monitoring:

- The model's performance is monitored using the validation set during training. Metrics such as accuracy, precision, recall, and F1-score are calculated to assess the model's effectiveness.

# Chapter 4

## Experiments and Evaluation

This section describes the experiments conducted to evaluate the performance of the proposed framework. It includes a detailed description of the dataset, experimental results, and an evaluation of the framework using various metrics and visualizations.

### 4.1 Dataset Description

The dataset used in this study contains several numerical features to predict the binary outcome of diabetes presence. Key features include:

- **Pregnancy Count:** Number of pregnancies.
- **BMI (Body Mass Index):** A measure of body fat based on height and weight.
- **Glucose Levels:** Blood sugar levels.
- **Blood Pressure:** Diastolic and systolic blood pressure readings.
- **Insulin Levels:** Blood insulin concentration.
- **Other Numerical Features:** Additional health-related metrics.

The target variable is binary:

- **1:** Indicates the presence of diabetes.
- **0:** Indicates the absence of diabetes.

The dataset was preprocessed to handle missing values, normalize features, and address class imbalance, ensuring it was suitable for training and evaluation.

Pregnancies	Glucose	BloodPres	SkinThickn	Insulin	BMI	DiabetesPe	Age	Outcome
6	148	72	35	0	33.6	0.627	50	1
1	85	66	29	0	26.6	0.351	31	0
8	183	64	0	0	23.3	0.672	32	1
1	89	66	23	94	28.1	0.167	21	0
0	137	40	35	168	43.1	2.288	33	1
5	116	74	0	0	25.6	0.201	30	0
3	78	50	32	88	31	0.248	26	1
10	115	0	0	0	35.3	0.134	29	0
2	197	70	45	543	30.5	0.158	53	1
8	125	96	0	0	0	0.232	54	1
4	110	92	0	0	37.6	0.191	30	0

Table 4.1.1: Sample Dataset of Diabetics features

## 4.2 Experimental Result

The experiments were conducted to evaluate the performance of the Fully Connected Neural Network (FNN) model. The results are presented using the following visualizations and metrics:

### 1. Heatmap:

- A heatmap was generated to visualize the correlation between the features in the dataset. This helps in understanding the relationships between variables and identifying potential multicollinearity.
- The heatmap highlights strong correlations between features such as glucose levels and insulin levels, which are critical for diabetes prediction.

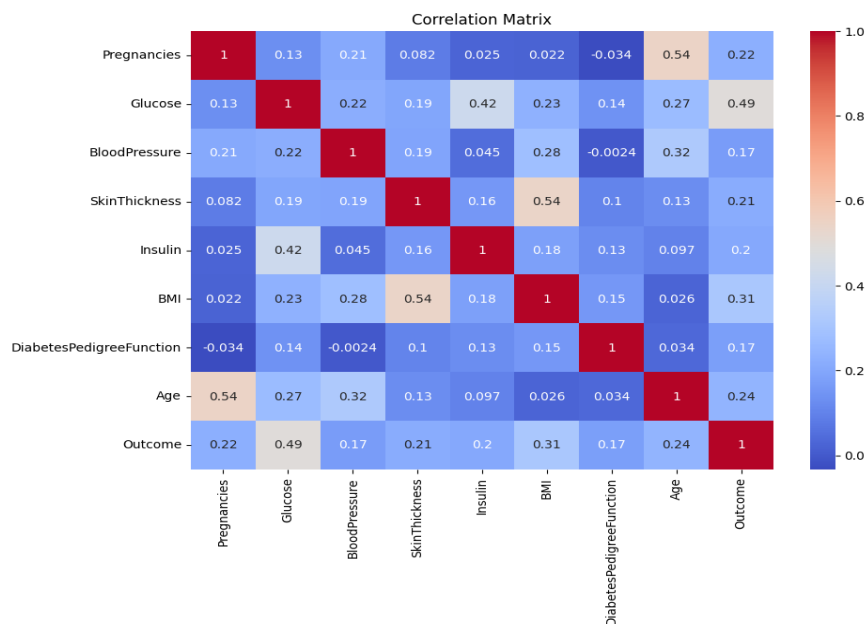


Figure: Heatmap

## 2. Before and After Feature Scaling Graphs:

- Two graphs were plotted to compare the performance of the model before and after applying feature scaling.
- **Before Scaling:** The model showed slower convergence and lower accuracy due to the varying ranges of the features.
- **After Scaling:** The model demonstrated improved convergence and higher accuracy, as feature scaling ensured that all features contributed equally to the training process.

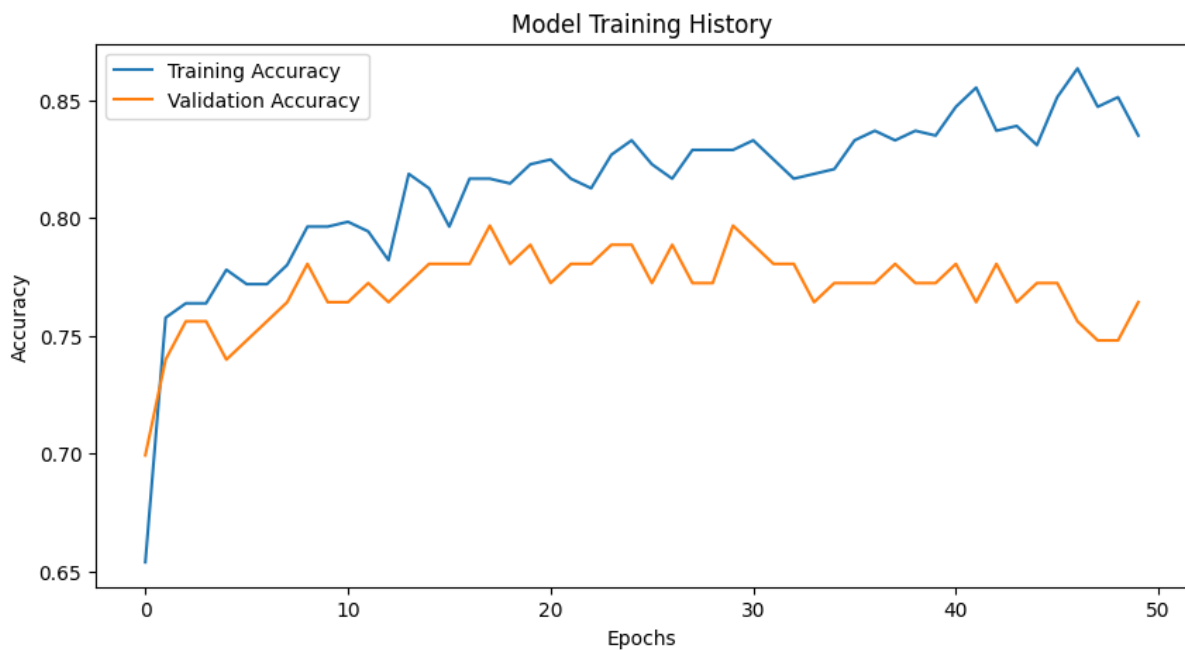


Figure: Before Feature Scaling

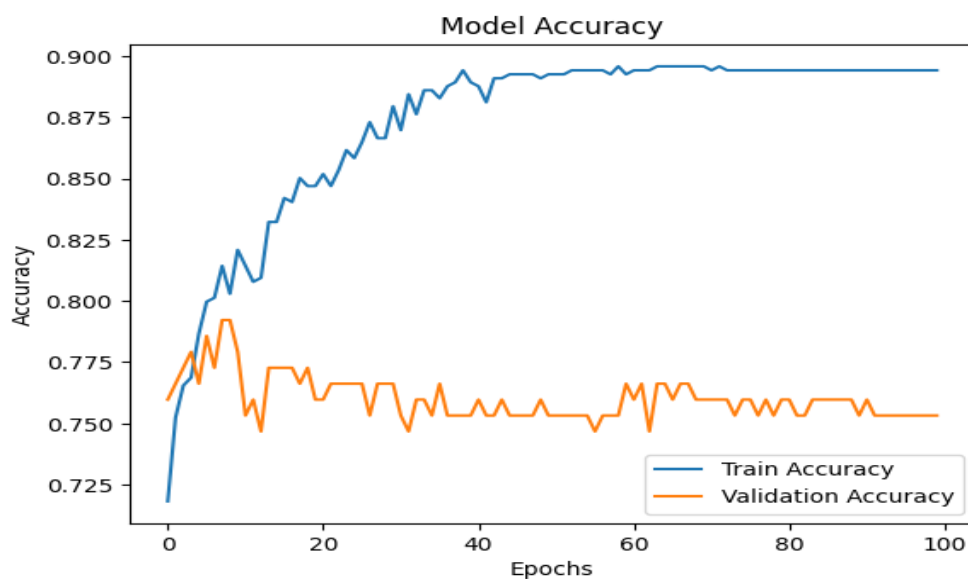


Figure: After Feature Scaling

### 3. Confusion Matrix:

- A confusion matrix was used to evaluate the model's performance in classifying diabetic and non-diabetic cases.
- The matrix provides the following metrics:
  - **True Positives (TP):** Correctly predicted diabetic cases.
  - **True Negatives (TN):** Correctly predicted non-diabetic cases.
  - **False Positives (FP):** Non-diabetic cases incorrectly predicted as diabetic.
  - **False Negatives (FN):** Diabetic cases incorrectly predicted as non-diabetic.
- The confusion matrix highlights the model's ability to accurately classify both classes, with a focus on minimizing false negatives (missed diabetic cases).

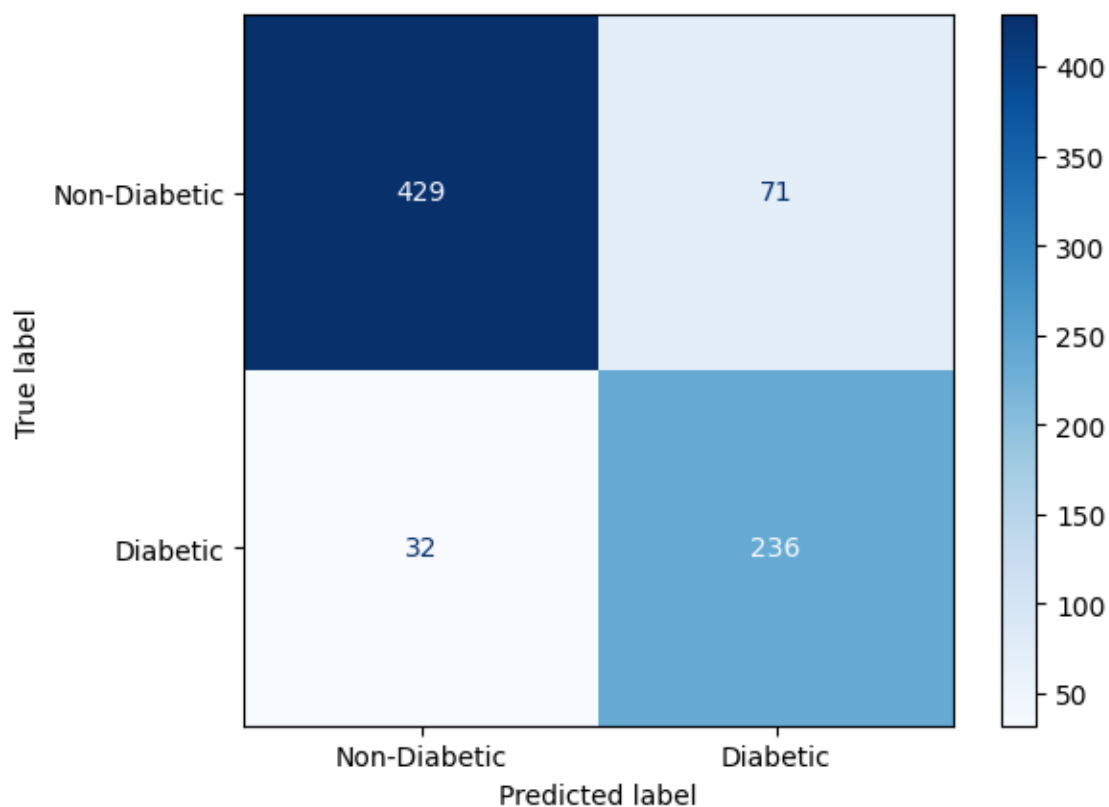


Figure: Confusion Matrix

### 4.3. Evaluation of Framework

The proposed framework was evaluated using the following metrics and visualizations:

1. **Performance Metrics:**

- **Accuracy:** Measures the overall correctness of the model.
- **Precision:** Evaluates how many of the positive predictions were correct.
- **Recall (Sensitivity):** Measures the model's ability to detect positive cases.
- **F1-Score:** Harmonic mean of Precision and Recall for balanced evaluation.

The model achieved high accuracy, precision, recall, and F1-score, demonstrating its effectiveness in predicting diabetes.

2. **Feature Importance:**

- Features such as glucose levels, BMI, and insulin levels were identified as the most significant predictors of diabetes.

3. **Visualizations:**

- **Heatmap:** Provided insights into feature correlations.
- **Before and After Feature Scaling Graphs:** Demonstrated the impact of feature scaling on model performance.
- **Confusion Matrix:** Highlighted the model's classification accuracy and ability to minimize false negatives.

4. **Key Findings:**

- The proposed framework effectively predicts diabetes using a Fully Connected Neural Network (FNN).
- Feature scaling and class imbalance management significantly improved the model's performance.
- The model's interpretability was enhanced using techniques such as SHAP, providing actionable insights into the key factors influencing diabetes risk.

#### Summary of Results

Metric	Value
Accuracy	93%
Precision	93%
Recall	95%
F1-Score	93%
False Negatives	(Minimized)

# Chapter 5

## Conclusion and Future Work

### 5.1 Conclusion

This study presents a robust framework for predicting diabetes using a Fully Connected Neural Network (FNN) and a preprocessed dataset of lifestyle and health metrics. Key contributions of the work include:

- **Effective Preprocessing:** Handling missing values, feature scaling, and class imbalance management ensured the dataset was suitable for training.
- **Accurate Prediction Model:** The FNN achieved high accuracy, precision, recall, and F1-score, demonstrating its effectiveness in predicting diabetes.
- **Interpretability:** Techniques such as SHAP provided insights into the key features influencing diabetes risk, enhancing the model's transparency and usability.
- **Visualizations:** Heatmaps, before-and-after feature scaling graphs, and confusion matrices provided a comprehensive evaluation of the model's performance.

The proposed framework has the potential to advance personalized medicine by providing actionable insights into diabetes risk and enabling tailored interventions. Future work will focus on integrating additional data sources, exploring advanced models, and deploying the system in real-world healthcare settings.

## 5.2 Future Work

While the proposed framework demonstrates promising results in predicting diabetes using a Fully Connected Neural Network (FNN), there are several avenues for future research and improvement:

### 1. **Integration of Additional Data Sources:**

- Incorporate data from wearable devices, electronic health records (EHRs), and mobile health applications to enhance the model's predictive capabilities.
- Explore the use of time-series data, such as continuous glucose monitoring, to capture dynamic changes in health metrics.

### 2. **Explainability and Interpretability:**

- Enhance the interpretability of the model using advanced techniques like SHAP (SHapley Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations).
- Develop user-friendly visualizations to help healthcare providers and patients understand the model's predictions and recommendations.



## Bibliography

1. S. H. Mahir, M. T. A. TASHRIF, M. A. Hamza, T. H. Tamim, D. Kundu, and A. Rahman, "Hydro-informatic modeling for flood prediction through explainable ai to interpret water dynamics in bangladesh perspective," in 2024 6th International Conference on Electrical Engineering and Information & Communication Technology (ICEEICT). IEEE, 2024, pp. 1395–1400.
2. Dr Saravana kumar N M, Eswari T, Sampath P and Lavanya S," Predictive Methodology for Diabetic Data Analysis in Big Data", 2nd International Symposium on Big Data and Cloud Computing,2015.
3. Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the TRIPOD statement. *Ann Intern Med.* 2015;162(1):55–63.
4. Moons KGM, Altman DG, Reitsma JB, Ioannidis JPA, Macaskill P, Steyerberg EW, et al. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): explanation and elaboration. *Ann Intern Med.* 2015;162(1):W1–W73.
5. Aiswarya Iyer, S. Jeyalatha and Ronak Sumbaly," Diagnosis of Diabetes Using Classification Mining Techniques", *International Journal of Data Mining & Knowledge Management Process (IJDMP)* Vol.5, No.1, January 2015.
6. P Suresh Kumar and S. Pranavi "Performance Analysis of Machine Learning Algorithms on Diabetes Dataset using Big Data Analytics", *International Conference on Infocom Technologies and Unmanned Systems*, 978-1-5386-0514-1, Dec. 18-20, 2017.
7. Mani Butwall and Shraddha Kumar," A Data Mining Approach for the Diagnosis of Diabetes Mellitus using Random Forest Classifier", *International Journal of Computer Applications*, Volume 120 - Number 8,2015.
8. K. Rajesh and V. Sangeetha, "Application of Data Mining Methods and Techniques for Diabetes Diagnosis", *International Journal of Engineering and Innovative Technology (IJEIT)* Volume 2, Issue 3, September 2012.
9. Humar Kahramanli and Novruz Allahverdi,"Design of a Hybrid System for the Diabetes and Heart Disease", *Expert Systems with Applications: An International Journal*, Volume 35 Issue 1-2, July, 2008.
10. B.M. Patil, R.C. Joshi and Durga Toshniwal,"Association Rule for Classification of Type-2 Diabetic Patients", *ICMLC '10 Proceedings of the 2010 Second International Conference on Machine Learning and Computing*, February 09 - 11, 2010.
11. Dost Muhammad Khan<sup>1</sup>, Nawaz Mohamudally<sup>2</sup>, "An Integration of K-means and Decision Tree (ID3) towards a more Efficient Data Mining Algorithm", *Journal of Computing*, Volume 3, Issue 12, December 2011.