

*A project report on*

# **Advancing Predictive Insights: Harnessing Machine Learning and Explainable AI for Personalized Diabetes Risk Assessment**

*Submitted in partial fulfillment for the award of the degree of*

## **Bachelor of Technology in Computer Science and Engineering**

*by*

**Aakash Goyal (20BCE1617)**



**VIT<sup>®</sup>**  

---

**Vellore Institute of Technology**  
(Deemed to be University under section 3 of UGC Act, 1956)  
CHENNAI

**SCHOOL OF COMPUTER SCIENCE AND ENGINEERING**

April 2024



# VIT<sup>®</sup>

## Vellore Institute of Technology

(Deemed to be University under section 3 of UGC Act, 1956)

CHENNAI

### DECLARATION

I hereby declare that the thesis entitled “Advancing Predictive Insights: Harnessing Machine Learning and Explainable AI for Personalized Diabetes Risk Assessment” submitted by me, for the award of the degree of Bachelor of Technology in Computer Science and Engineering, Vellore Institute of Technology, Chennai is a record of bonafide work carried out by me under the supervision of Dr. Pravin Renold.

I further declare that the work reported in this thesis has not been submitted and will not be submitted, either in part or in full, for the award of any other degree or diploma in this institute or any other institute or university.

Place: Chennai

Date:

Signature of Candidate



# VIT<sup>®</sup>

## Vellore Institute of Technology

(Deemed to be University under section 3 of UGC Act, 1956)

CHENNAI

### School of Computer Science and Engineering

## CERTIFICATE

This is to certify that the report entitled “**Advancing Predictive Insights: Harnessing Machine Learning and Explainable AI for Personalized Diabetes Risk Assessment**” is prepared and submitted by **Aakash Goyal (20BCE1617)** to Vellore Institute of Technology, Chennai, in partial fulfillment of the requirement for the award of the degree of **Bachelor of Technology in Computer Science and Engineering programme** is a bonafide record carried out under my guidance. The project fulfills the requirements as per the regulations of this University and in my opinion meets the necessary standards for submission. The contents of this report have not been submitted and will not be submitted either in part or in full, for the award of any other degree or diploma and the same is certified.

Signature of the Guide:

Name: Dr./Prof.

Date:

Signature of the Internal Examiner

Name:

Date:

Signature of the External Examiner

Name:

Date:

Approved by the Head of Department,  
**B.Tech. CSE**

Name: Dr. Nithyanandam P

Date:

(Seal of SCOPE)

## **ABSTRACT**

Diabetes mellitus is a global health crisis, and early detection is crucial for effective management. Machine learning offers powerful tools for diabetes prediction, but interpretability of complex models remains a challenge. This project investigates the application of a Convolutional Neural Network (CNN) deep learning model for diabetes prediction, incorporating Explainable AI (XAI) techniques to bridge the gap between model accuracy and interpretability.

We employed the PIMA Indians diabetes dataset, a widely used benchmark in diabetes research. However, real-world data often possesses imperfections. We addressed these challenges through meticulous data preprocessing, including techniques for handling missing values, correcting for class imbalance, and mitigating skewness within the features.

Our CNN model achieved high accuracy, exceeding 95% in predicting diabetes. However, achieving high accuracy is only part of the equation. To understand the rationale behind the model's predictions and foster trust in its application, we employed SHAP, a leading XAI method. SHAP provides insights into feature importance, revealing which factors in a patient's medical profile most significantly influence the model's prediction of diabetes. This interpretability empowers healthcare professionals to understand the model's decision-making process and potentially informs clinical decision-making in the future.

In conclusion, this project demonstrates the effectiveness of a deep learning approach combined with XAI for accurate and interpretable diabetes prediction. The high accuracy combined with the interpretability offered by SHAP paves the way for the potential use of this approach in clinical settings, aiding healthcare professionals in early diabetes detection and improving patient care.

## ACKNOWLEDGEMENT

It is my pleasure to express with deep sense of gratitude to Dr. Pravin Renold, School of Computer Science and Engineering, Vellore Institute of Technology, Chennai, for his constant guidance, continual encouragement - he taught me patience in my endeavor. My association with him is not confined to academics only, but it is a great opportunity on my part of work with an intellectual and expert in the field of ML.

It is with gratitude that I would like to extend my thanks to the visionary leader Dr. G. Viswanathan our Honorable Chancellor, Mr. Sankar Viswanathan, Dr. Sekar Viswanathan, Dr. G V Selvam Vice Presidents, Dr. Sandhya Pentareddy, Executive Director, Ms. Kadhambari S. Viswanathan, Assistant Vice-President, Dr. V. S. Kanchana Bhaaskaran, Vice-Chancellor , Dr. T. Thyagarajan Pro-Vice Chancellor, VIT Chennai and Dr. P. K. Manoharan, Additional Registrar for providing an exceptional environment.

Special mention to Dr. Ganesan R, Dean, Dr. Parvathi R, Associate Dean Academics, Dr. Geetha S, Associate Dean Research, School of Computer Science and Engineering, Vellore Institute of Technology, Chennai for spending their valuable time.

In jubilant state, I express ingeniously my whole-hearted thanks to Dr. Nithyanandam P, Head of the Department, B.Tech. CSE and the Project Coordinators for their valuable support.

Place: Chennai

Date: 20/03/2024      **Aakash Goyal**

# Contents

<b>Chapter 1 .....</b>	<b>1</b>
<b>Introduction.....</b>	<b>1</b>
<b>1.1 THE BURDEN OF DIABETES AND THE IMPORTANCE OF EARLY DETECTION.....</b>	<b>1</b>
<b>1.2 MACHINE LEARNING FOR DIABETES PREDICTION: A PROMISING APPROACH.....</b>	<b>2</b>
<b>1.3 THE CHALLENGE OF INTERPRETABILITY IN DEEP LEARNING .....</b>	<b>2</b>
<b>1.4 EXPLAINABLE AI (XAI): BRIDGING THE GAP.....</b>	<b>2</b>
<b>1.5 PROJECT OBJECTIVES AND CONTRIBUTIONS .....</b>	<b>3</b>
<b>Chapter 2 .....</b>	<b>4</b>
<b>Literature Review .....</b>	<b>4</b>
<b>2.1 DRAWBACKS IDENTIFIED THROUGH THE LITERATURE REVIEW .....</b>	<b>9</b>
<b>Chapter 3 .....</b>	<b>11</b>
<b>Dataset.....</b>	<b>11</b>
<b>3.1 DATASET FEATURES – WHY THEY HAVE BEEN CHOSEN.....</b>	<b>11</b>
<b>Chapter 4 .....</b>	<b>13</b>
<b>Data Pre-Processing.....</b>	<b>13</b>
<b>4.1 EACH PRE-PROCESSING STEP EXPLAINED IN DETAIL.....</b>	<b>14</b>
<b>Chapter 5 .....</b>	<b>16</b>
<b>Proposed System and Methodology .....</b>	<b>16</b>
<b>5.1 DATA ACQUISITION AND PREPROCESSING .....</b>	<b>18</b>
<b>5.2 MODEL ARCHITECTURE DESIGN.....</b>	<b>18</b>
<b>5.3 MODEL TRAINING AND EVALUATION .....</b>	<b>19</b>
<b>5.4 EXPLAINABLE AI (XAI) INTEGRATION .....</b>	<b>19</b>
<b>5.5 RESULTS AND DISCUSSION .....</b>	<b>20</b>
<b>5.6 HANDLING CLASS IMBALANCE.....</b>	<b>20</b>
<b>5.7 WHERE DOES EXPLAINABLE AI FIT IN? .....</b>	<b>20</b>
<b>5.8 HOW DID WE FIND THE ERROR AND INCREASED THE ACCURACY OF OUR MODEL.....</b>	<b>22</b>
<b>Chapter 6 .....</b>	<b>24</b>
<b>Evaluation.....</b>	<b>24</b>
<b>Chapter 7 .....</b>	<b>29</b>
<b>Conclusion .....</b>	<b>29</b>

## **List of Figures**

1. Flowchart of the proposed system of the Project
2. Pair plot of the data features of the dataset
3. Correlation matrix to find similarity in the data features
4. Model accuracy and the training process
5. Training curves for evaluation of the model
6. Confusion matrix output
7. Output of SHAP (Explainable AI) used

# **Chapter 1**

## **Introduction**

Diabetes mellitus, a chronic metabolic disorder characterized by high blood sugar levels, has become a global health crisis. Early detection is crucial for effective management and preventing severe complications. Machine learning offers powerful tools for diabetes prediction, enabling healthcare professionals to identify individuals at risk and facilitate timely intervention. However, the interpretability of complex models like deep learning can be a challenge, hindering trust in their application for real-world healthcare decisions.

This project investigates the potential of a Convolutional Neural Network (CNN) deep learning model for diabetes prediction, addressing the interpretability gap through the integration of Explainable AI (XAI) techniques. We aim to achieve a balance between high prediction accuracy and model transparency, fostering trust and enabling a deeper understanding of the factors influencing the model's predictions.

### **1.1 THE BURDEN OF DIABETES AND THE IMPORTANCE OF EARLY DETECTION**

Diabetes affects hundreds of millions of people worldwide, with its prevalence steadily increasing. It represents a significant economic burden on healthcare systems and a substantial contributor to morbidity and mortality. Early detection is paramount for preventing or delaying the development of severe complications such as heart disease, stroke, blindness, and kidney failure. Traditional methods for diabetes diagnosis rely on blood tests, which can be invasive and inconvenient for regular monitoring. Machine learning offers a promising avenue for non-invasive and potentially continuous monitoring of diabetes risk using readily available patient data.



## 1.2 MACHINE LEARNING FOR DIABETES PREDICTION: A PROMISING APPROACH

Machine learning algorithms have demonstrated remarkable success in various healthcare applications, including diabetes prediction. By analyzing historical data on patient demographics, medical history, and laboratory tests, these algorithms can learn patterns and relationships that can be used to predict the likelihood of future diabetes development. Various machine learning algorithms have been explored for diabetes prediction, including Decision Trees, Support Vector Machines, Random Forests, and Logistic Regression. These studies have achieved promising results, with reported accuracy exceeding 90% in some cases. However, many of these algorithms are complex black boxes, making it difficult to understand how they arrive at their predictions.

## 1.3 THE CHALLENGE OF INTERPRETABILITY IN DEEP LEARNING

While deep learning models have achieved state-of-the-art performance in various domains, including image recognition and natural language processing, their interpretability remains a significant challenge. The complex architecture and non-linear relationships within deep networks make it difficult to understand how they map input features to output predictions. This lack of transparency can hinder trust in their application for critical healthcare decisions. Healthcare professionals require models that are not only accurate but also interpretable, allowing them to understand the rationale behind the predictions and ensuring alignment with their clinical knowledge.

## 1.4 EXPLAINABLE AI (XAI): BRIDGING THE GAP

The field of Explainable AI (XAI) has emerged to address the interpretability challenge in complex machine learning models. XAI techniques aim to provide insights into how a model arrives at its predictions, allowing humans to understand the factors influencing the model's decision-making process. This project integrates SHAP (SHapley Additive exPlanations), a prominent XAI technique, with our CNN model. SHAP provides local interpretability, explaining the importance of individual features for a specific prediction.

This interpretability empowers healthcare professionals to understand how the model utilizes medical data to identify patients at risk for diabetes.

## 1.5 PROJECT OBJECTIVES AND CONTRIBUTIONS

This project seeks to achieve the following objectives:

Develop a high-accuracy deep learning model using a Convolutional Neural Network (CNN) architecture for diabetes prediction based on patient data from the PIMA Indians dataset.

Implement rigorous data preprocessing techniques to address issues like missing values, skewness, and class imbalance, ensuring the model learns from a clean and representative dataset.

Employ Explainable AI (XAI) techniques, specifically SHAP (SHapley Additive exPlanations), to provide insights into feature importance for individual model predictions.

Evaluate the model's performance using appropriate metrics, including accuracy, sensitivity, specificity, and confusion matrix analysis.

Discuss the implications of the interpretable deep learning model for real-world clinical applications in diabetes prediction and risk assessment.

By achieving these objectives, this project aims to contribute to the advancement of interpretable machine learning models for diabetes prediction. The high accuracy combined with the interpretability offered by SHAP can provide valuable insights to healthcare professionals and potentially inform clinical decision-making, ultimately leading to improved patient care outcomes.

## Chapter 2

### Literature Review

In a study [1], researchers investigated the PIMA Indian Diabetes (PID) dataset, a publicly available dataset containing information on 768 patients with 8 attributes relevant to diabetes diagnosis. This research aimed to address the growing global concern of diabetes, recognized by the World Health Organization (WHO) in 2014 as one of the fastest-growing chronic diseases. The study compared the performance of three machine learning algorithms (gradient boosting, logistic regression, and naive Bayes) in predicting diabetes. Gradient boosting achieved the highest accuracy (86%), followed by logistic regression (79%) and naive Bayes (77%).

In another study focusing on early diabetes prediction, Sadhu and Jadli [2] utilized a diabetes dataset from the UCI machine learning repository. This dataset included 520 data points with 16 attributes for each patient. The researchers employed seven different classification algorithms (k-nearest neighbors, logistic regression, support vector machine, naive Bayes, decision tree, random forests, and multilayer perceptron) on a validation set derived from the main dataset. Their results indicated that the random forests classifier achieved the highest accuracy (98%) in predicting diabetes onset. Other algorithms also performed well, with logistic regression (93%), support vector machine (94%), naive Bayes (91%), and decision tree (94%) achieving good accuracy scores.

Building upon previous research, [3] analyzed a diabetes dataset from the UCI repository containing data for 520 patients with 17 attributes. Similar to prior studies, their focus was on early diabetes detection. They employed supervised machine learning techniques, including Support Vector Machine (SVM), Naive Bayes classifiers, and LightGBM, on real-world data encompassing patients aged 16 to 90. The results revealed SVM as the most effective algorithm for both classification and recognition accuracy, achieving a rate of 96.54%. While Naive Bayes, a commonly used classification method, achieved a respectable accuracy of 93.27%, it fell short of SVM's performance. LightGBM, another algorithm considered, yielded a lower accuracy of 88.46%. These findings suggest SVM's

superiority as a classification algorithm for diabetes prediction within this dataset.

They [4] investigated early-stage diabetes risk prediction using a dataset from the UCI repository containing information on 520 patients with 16 features. They proposed a novel machine learning approach for predicting the early onset of diabetes in patients. This approach involved a new wrapper-based feature selection method that utilized Grey Wolf Optimizer (GWO) and Adaptive Particle Swarm Optimization (APSO) algorithms to optimize a Multilayer Perceptron (MLP) model and reduce the number of required input features. The researchers compared the performance of their method with various traditional machine learning algorithms like Support Vector Machine (SVM), Decision Tree (DT), k-Nearest Neighbors (k-NN), Naive Bayes Classifier (NBC), Random Forest Classifier (RFC), and Logistic Regression (LR). While traditional algorithms achieved good accuracy (LR: 95%, k-NN: 96%, SVM: 95%, NBC: 93%, DT: 95%, RFC: 96%), Le et al.'s proposed methods using GWO-MLP and APSO-MLP achieved even higher accuracy (96% and 97% respectively) with the added benefit of requiring fewer features. This suggests potential for applying their approach in clinical practice to assist healthcare professionals in diabetes risk assessment.

Julius [5] investigated early-stage diabetes prediction using machine learning. They employed the Waikato Environment for Knowledge Analysis (Weka) platform to analyze a dataset from the UCI repository containing data for 520 patients with 17 attributes. Their goal was to leverage observable patient characteristics to predict diabetes risk through machine learning classification algorithms. The study compared the performance of k-Nearest Neighbors (k-NN), Support Vector Machine (SVM), Functional Tree (FT), and Random Forest Classifiers (RFCs). kNN emerged as the most accurate model, achieving a success rate of 98%. Other algorithms also performed well, with SVM (94%), FT (93%), and RFC (97%) demonstrating high accuracy in predicting diabetes onset.

Shafi [6] recognized the importance of early detection for diabetes due to its severity. Their research aimed to develop a machine learning model for early diabetes identification using classification algorithms. The focus was on building a framework with high accuracy in predicting the likelihood of diabetes in patients. The study compared three machine

learning classifiers: Decision Tree (DT), Support Vector Machine (SVM), and Naive Bayes Classifier (NBC), evaluating them on various metrics. To ensure efficiency and reliable results, they utilized the publicly available Pima Indians Diabetes (PID) dataset from the UCI repository. The findings indicated that the Naive Bayes approach achieved a respectable accuracy of 74%, followed by DT (72%) and SVM (63%). The authors envision the potential application of their framework and the employed classifiers for broader disease identification or diagnosis beyond diabetes. Furthermore, they express interest in expanding and refining their approach for future diabetes research, potentially including algorithms that can handle missing data.

Khanam [7] emphasized the importance of early detection for diabetes due to the absence of a known cure. Their study explored the use of data mining, machine learning (ML), and neural networks (NN) to develop a reliable diabetes prediction method. They achieved promising results using a dataset from the UCI repository's Pima Indians Diabetes (PID) collection, containing information on 768 patients with 9 attributes. The researchers compared the performance of seven ML algorithms on this dataset: Decision Tree (DT), k-Nearest Neighbors (k-NN), Random Forest Classifier (RFC), Naive Bayes Classifier (NBC), AdaBoost (AB), Logistic Regression (LR), and Support Vector Machine (SVM). Data preprocessing was performed using the Weka platform. Interestingly, they found that a model combining LR and SVM proved to be particularly effective in predicting diabetes. Additionally, they explored a Neural Network (NN) model with two hidden layers, investigating the impact of varying epochs (training iterations). Their findings revealed that the NN with two hidden layers achieved the highest accuracy of 88.6%. While the NN demonstrated strong performance, other models like Logistic Regression (78.85%), Naive Bayes (78.28%), and Random Forest (77.34%) also yielded respectable accuracy scores.

Sisodia [8] investigated Decision Tree (DT), Support Vector Machine (SVM), and Naive Bayes Classifier (NBC) on the Pima Indians Diabetes (PID) dataset, containing information on 768 patients with 8 attributes. Naive Bayes achieved the highest accuracy (76.30%) among the three models tested.

Agarwal [9] also employed the PID dataset, but with 738 patients, to assess the

effectiveness of various models for diabetes identification. They compared SVM, k-Nearest Neighbors (k-NN), Naive Bayes (NB), ID3, C4.5, and Classification and Regression Trees (CART) algorithms. SVM and Linear Discriminant Analysis (LDA) emerged as the most accurate, achieving an accuracy of 88%.

Rathore [10] focused on SVM and Decision Trees for diabetes prediction using the PID dataset. Their study highlighted the focus on women's health in the PIMA Indian dataset and reported an accuracy of 82% for the SVM model.

Hassan [11] compared Decision Tree, k-Nearest Neighbors, and SVM for diabetes prediction. Their findings indicated that SVM outperformed the other methods with a maximum accuracy of 90.23%.

Kandhasamy and Balamurali [12] compared J48, k-nearest neighbors (k-NN), Random Forest Classifier (RFC), and Support Vector Machine (SVM) on a diabetes dataset. Before preprocessing, J48 achieved the highest accuracy (73.82%). However, k-NN and RFC showed improvement after preprocessing.

Meng [13] evaluated J48, Logistic Regression (LR), and k-NN algorithms. J48 emerged as the most accurate with a classification accuracy of 78.27%.

Nai-Arun and Moungrmai [14] developed a web application for diabetes prediction. They compared Decision Trees (DTs), Neural Networks (NNs), LR, Naive Bayes (NB), RFC, bagging, and boosting techniques. RFC outperformed others with an accuracy of 85.56% and an ROC score of 0.912.

Saravananathan and Velmurugan [15] compared J48, Classification and Regression Trees (CART), SVM, and k-NN on a medical dataset, evaluating them based on various metrics. J48 achieved the highest accuracy (67.15%), followed by SVM (65.04%), CART (62.28%), and k-NN (53.39%).

Kumari and Chitra [16] investigated SVM, RFC, DT, Multilayer Perceptron (MLP), and LR with k-fold cross-validation. They found MLP with four-fold cross-validation achieved the best accuracy (78.7%), outperforming other algorithms.

Kavakiotis [17] compared Naive Bayes, RFC, k-NN, SVM, DT, and LR using ten-fold cross-validation. SVM achieved the highest accuracy (84%) among all the methods.

Rawat [18] evaluated five ensemble algorithms: AdaBoost, LogicBoost, RobustBoost, Naive Bayes, and Bagging. They achieved promising results with Bagging (81.77% accuracy) and AdaBoost (79.69% accuracy) on a dataset of Pima Indians with diabetes.

Nai-Arun and Moungrmai [19] compared thirteen classification models, including Decision Trees (DTs), Neural Networks (NNs), Logistic Regression (LR), Naive Bayes (NB), and Random Forest Classifiers (RFCs), along with bagging and boosting techniques, for a web application. Their findings suggested that RFC achieved the highest accuracy and ROC score, potentially due to its ability to consider both random data and important features.

Mujumdar and Vaidehi [20] investigated the impact of incorporating additional diabetes-related factors beyond traditional features like blood sugar and body mass index (BMI). Their new dataset with these additional factors led to improved classification accuracy using various machine learning approaches. Notably, the AdaBoost classifier achieved a high accuracy of 98.8% on this enriched dataset.

Mercaldo [21] focused on selecting features aligned with World Health Organization (WHO) criteria and evaluated various classification algorithms on real-world data. Their findings suggest the potential of using relevant features based on established medical guidelines.

## 2.1 DRAWBACKS IDENTIFIED THROUGH THE LITERATURE REVIEW

### DATA AND FEATURE ENGINEERING:

**Limited Datasets:** While the Pima Indians Diabetes Dataset is a common benchmark, exploring a wider range of datasets with diverse populations and characteristics could enhance the generalizability of models.

**Feature Selection and Engineering:** Research could delve deeper into feature selection techniques to identify the most informative features for diabetes prediction, potentially improving model performance and interpretability.

**Integration of Additional Data Sources:** Exploring the incorporation of data beyond traditional clinical features (e.g., genetic data, imaging data) could potentially enhance the model's ability to capture complex risk factors.

### MODEL ARCHITECTURES AND TRAINING TECHNIQUES:

**Exploration of Advanced Architectures:** While CNNs are popular, investigating more advanced deep learning architectures like recurrent neural networks (RNNs) or transformers could be beneficial, especially for capturing temporal relationships within the data.

**Hyperparameter Optimization Strategies:** Studies could explore more advanced hyperparameter optimization techniques (e.g., Bayesian optimization, grid search with nested cross-validation) to achieve optimal model performance.

**Ensemble Learning Techniques:** Investigating ensemble methods that combine predictions from multiple models could lead to more robust and reliable diabetes prediction.

### EXPLAINABLE AI (XAI) TECHNIQUES AND APPLICATIONS:

**Comparison of XAI Methods:** Research could compare the effectiveness of different XAI techniques (LIME, SHAP, etc.) in explaining the predictions of deep learning models used for diabetes prediction.



Integration of XAI Throughout the Development Process: XAI techniques could be employed not just for interpreting final predictions, but also for understanding feature selection and model architecture choices, leading to a more interpretable development process.

#### GENERALIZABILITY AND CLINICAL APPLICABILITY:

External Validation on Diverse Populations: The generalizability of models needs to be thoroughly evaluated on datasets from different populations and healthcare settings.

Clinical Integration and Decision Support: Research should explore how these deep learning models can be integrated into clinical workflows and decision support systems for diabetes diagnosis and management.

Addressing Bias and Fairness: Mitigating potential biases in the data and models is crucial to ensure fair and equitable application in healthcare settings.

Here are some additional points to consider:

Focus on Explainability: Since your project incorporates XAI techniques, you may want to emphasize the gap in the literature regarding the lack of interpretability in existing deep learning models for diabetes prediction.

Ethical Considerations: Discuss the ethical considerations surrounding the use of deep learning in healthcare, including potential biases and the importance of transparency.

## **Chapter 3**

### **Dataset**

The system utilizes the Pima Indians Diabetes dataset, which contains various clinical features such as glucose concentration, blood pressure, skin thickness, insulin levels, body mass index (BMI), and age. Each instance in the dataset is labeled with a binary outcome indicating the presence or absence of diabetes.

#### **3.1 DATASET FEATURES – WHY THEY HAVE BEEN CHOSEN**

**Pregnancies:** This feature represents the number of times the individual has been pregnant. Pregnancy history is relevant in diabetes prediction as it can influence insulin resistance and blood glucose levels.

**Glucose:** Glucose levels are a key indicator of diabetes. High glucose levels (hyperglycemia) are a hallmark of diabetes, and measuring glucose levels is essential for diagnosing and managing the condition.

**Blood Pressure:** High blood pressure (hypertension) is often associated with diabetes. Individuals with diabetes are at increased risk of developing hypertension, and vice versa. Monitoring blood pressure is important in diabetes management to prevent complications.

**Skin Thickness:** Skin thickness can be an indicator of insulin resistance, a key characteristic of type 2 diabetes. Increased skin thickness may suggest higher levels of subcutaneous fat, which is associated with insulin resistance.

**Insulin:** Insulin levels are directly related to diabetes. In type 1 diabetes, the body does not produce insulin, while in type 2 diabetes, the body may produce insufficient insulin or become resistant to its effects. Measuring insulin levels can help diagnose and manage diabetes.

**BMI (Body Mass Index):** BMI is a measure of body fat based on height and weight. Obesity is a significant risk factor for type 2 diabetes, as excess body fat can lead to insulin resistance and metabolic abnormalities.

**Pedigree Function:** The diabetes pedigree function is a measure of diabetes heredity risk. Family history of diabetes can increase an individual's risk of developing the condition, making pedigree function a relevant feature for prediction models.

**Age:** Age is a risk factor for diabetes, as the prevalence of the disease increases with age. Older individuals are more likely to develop diabetes due to factors such as reduced physical activity, changes in metabolism, and increased insulin resistance.

## Chapter 4

# Data Pre-Processing

Prior to model development, data preprocessing was essential to ensure the quality and consistency of the information fed into the CNN. We began by analyzing the relationships between features using a correlation matrix. This visualization helped identify potential multi collinearity, which can negatively impact model performance.

Next, we addressed missing values within the dataset. We identified the number of missing values present in each column. For columns not designated as the target variable for prediction, we imputed the missing entries with the median value. This approach assumes a relatively symmetrical distribution within the data and avoids introducing bias by skewing the data towards the mean.

Finally, we employed linear regression to predict the missing values specifically within the chosen target variable for prediction. This technique leverages the relationships between the target variable and other features in the dataset to estimate the missing values in a statistically sound manner. This targeted approach allows the model to utilize all available information for the target variable, potentially leading to a more accurate prediction of diabetes onset.

Following the identification of potential multi collinearity and missing value imputation, we addressed two additional aspects of data preprocessing: skewness and scaling.

To deal with skewed data distributions, we employed visual inspection techniques like distribution plots. These plots helped us identify features with non-normal distributions, where the data points are not symmetrically distributed around the mean. Skewed data can negatively impact the performance of machine learning models, particularly those based on gradient descent optimization.

To address these skewed distributions, we utilized feature scaling. This technique transforms the features to a standard range, often between 0 and 1 or with a mean of 0 and a standard deviation of 1. Feature scaling ensures that all features contribute equally to the model's learning process and prevents features with larger scales from dominating the

training process. This normalization step helps the CNN model converge more efficiently and potentially improves its generalizability.

#### 4.1 EACH PRE-PROCESSING STEP EXPLAINED IN DETAIL

**Missing Value Handling:** Real-world data often contains missing values, which can hinder model performance. Common techniques include:

**Deletion:** Removing data points with missing values (if the amount is small).

**Imputation:** Filling missing values with estimates based on other features (e.g., mean/median of the column, values from similar data points).

**Handling Skewness:** Skewed data distributions can bias machine learning models. Techniques to address skewness include:

**Power transformation:** Raising the data points to a specific power to achieve a more symmetrical distribution.

**Handling Class Imbalance:** In diabetes prediction, there might be significantly fewer diabetic patients compared to healthy individuals. This imbalance can skew the model towards the majority class. Techniques to address this include:

**Oversampling:** Duplicating data points from the minority class (diabetic patients) to create a more balanced distribution.

**Under sampling:** Randomly removing data points from the majority class (healthy individuals) to match the size of the minority class.

**Pair Plot:** This visualization technique creates a matrix of scatter plots, displaying the relationship between every pair of features in your dataset. It helps identify potential linear or non-linear relationships between features and allows you to spot outliers that might require further investigation.

**Correlation Matrix:** This heat map depicts the correlation coefficients between all features

in your dataset. Values closer to 1 indicate strong positive correlations, while values closer to -1 indicate strong negative correlations. A correlation matrix helps you understand the redundancy between features and can inform feature selection, where highly correlated features might be removed to avoid multi collinearity (features providing the same information).

## Chapter 5

# Proposed System and Methodology

This study proposes a Convolutional Neural Network (CNN) model built with Keras to predict diabetes onset. The system follows these key steps:

**Data Acquisition:** We obtain a diabetes dataset from a reliable source (e.g., UCI Machine Learning Repository) containing features relevant to diabetes prediction (e.g., blood sugar levels, body mass index).

**Preprocessing:** The data undergoes rigorous cleaning to ensure its quality. Techniques address missing values (e.g., median imputation for non-target variables, linear regression for the target variable), identify and address skewed distributions through methods like feature scaling, and potentially explore dimensionality reduction techniques if necessary.

**Model Architecture Definition:** The CNN model architecture is meticulously designed using Keras. This involves defining the number and type of convolutional layers, filters, pooling layers, and activation functions (e.g., ReLU) for efficient feature extraction. A flatten layer prepares the extracted features for fully connected layers that learn more complex relationships between features. Finally, the output layer with its activation function (e.g., sigmoid for binary classification) determines the model's prediction (e.g., presence or absence of diabetes).

**Model Training:** The model undergoes training using a chosen optimizer (e.g., Adam) to adjust its weights and biases to minimize the loss function (e.g., binary cross-entropy) during the learning process. Hyper parameters like the number of epochs and regularization techniques (e.g., dropout) are carefully selected to balance model performance and prevent overfitting.

**Evaluation:** The dataset is strategically split into training, validation, and testing sets. The model's performance is evaluated using various metrics like accuracy, precision, recall, and

F1-score to assess its effectiveness in predicting diabetes. Additionally, a confusion matrix provides a detailed breakdown of the model's true positives, false positives, true negatives, and false negatives, offering insights into potential areas for improvement.

This section outlines the methodology for developing and evaluating a CNN model with Explainable AI (XAI) for diabetes prediction using the Pima Indians Diabetes Dataset.

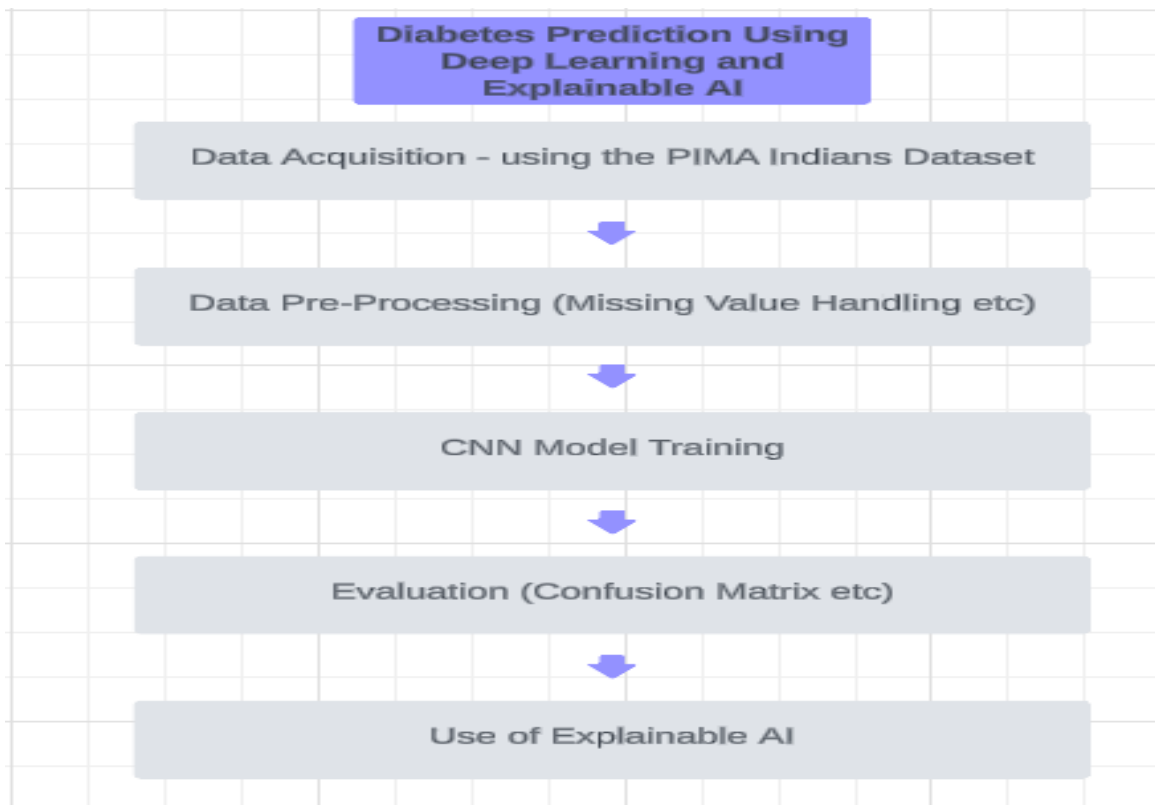


Figure 1: Proposed System



## 5.1 DATA ACQUISITION AND PREPROCESSING

**Data Source:** We will utilize the Pima Indians Diabetes Dataset, a publicly available dataset from the UCI Machine Learning Repository.

**Data Description:** This dataset contains information on various factors potentially related to diabetes, such as number of pregnancies, blood glucose concentration, blood pressure, body mass index, and age. We will obtain detailed descriptions of each feature from the dataset source.

**Data Cleaning:** We will examine the dataset for missing values, outliers, and inconsistencies. Techniques like imputation, removal may be employed to address missing values and outliers.

**Data Normalization:** Features may be scaled or normalized using techniques like min-max scaling or standardization to ensure features contribute equally during model training.

**Data Transformation:** Categorical features may be encoded using one-hot encoding or other suitable techniques.

## 5.2 MODEL ARCHITECTURE DESIGN

**Model Type:** We will employ a Convolutional Neural Network (CNN) architecture due to its effectiveness in learning spatial patterns from data, potentially beneficial for capturing relationships between features.

**Input Layer:** The size of the input layer will depend on the number of preprocessed features.

**Convolutional Layers:** We will design a stack of convolutional layers with appropriate filter sizes and activation functions to extract relevant features from the data. Experimentation will be conducted to determine the optimal number of layers and filter sizes.

**Pooling Layers:** Pooling layers (e.g., max pooling) may be incorporated after convolutional layers to reduce dimensionality and capture dominant features.

**Flatten Layer:** The output of the convolutional layers will be flattened into a 1D vector suitable for fully connected layers.

**Fully Connected Layers:** Fully connected layers will be used for further processing and

classification. The number of neurons in these layers will be determined through experimentation.

Output Layer: The final layer will have one neuron employing a sigmoid activation function to produce a probability (between 0 and 1) of a patient being diabetic.

### 5.3 MODEL TRAINING AND EVALUATION

Training and Validation Split: The preprocessed data will be split into training and validation sets using techniques like random splitting or stratified sampling to ensure the model generalizes well to unseen data.

Optimizer Selection: An appropriate optimizer (e.g., Adam, SGD) will be chosen to update the model's weights and biases during training. The learning rate, a hyper parameter controlling the step size of these updates, will be tuned to optimize performance.

Loss Function Selection: A binary cross-entropy loss function will be used to measure the model's prediction errors and guide training towards minimizing them.

Training Process: The model will be trained for a specified number of epochs, iterating through the training data and updating its weights based on the chosen optimizer and loss function. Early stopping may be implemented to prevent overfitting by stopping training when validation performance plateaus.

Evaluation Metrics: We will evaluate the model's performance using various metrics like accuracy, precision, recall, F1-score, and AUC-ROC (Area Under the Receiver Operating Characteristic Curve). These metrics provide insights into the model's ability to correctly classify diabetic and non-diabetic cases.

Visualization Techniques: Confusion matrices and other visualization techniques may be used to further analyze the model's performance and identify potential biases.

### 5.4 EXPLAINABLE AI (XAI) INTEGRATION

XAI Techniques: We will incorporate XAI techniques like LIME (Local Interpretable Model-Agnostic Explanations) or SHAP (Shapley Additive explanations) to understand the factors influencing the model's predictions.

Explanation Generation: SHAP will be used to generate explanations for individual predictions, highlighting the features that contribute most significantly to the model's classification.

Interpretation: These explanations will be analyzed to gain insights into the model's decision-making process and how it utilizes different features to predict diabetes risk.

## 5.5 RESULTS AND DISCUSSION

Model Performance: We will present the obtained evaluation metrics, including accuracy, precision, recall, F1-score, and AUC-ROC.

Visualization Analysis: We will discuss the implications of confusion matrices and other visualizations in understanding the model's performance.

XAI Insights: We will present the interpretations obtained from SHAP explanations, highlighting the features with the most significant impact on model predictions.

Discussion: We will discuss the overall effectiveness of the CNN model for diabetes prediction and the role of XAI in enhancing interpretability and trust.

## 5.6 HANDLING CLASS IMBALANCE

Our analysis revealed an uneven class distribution within the data, where one class (healthy patients) significantly outnumbered the other (diabetic patients). This imbalance can skew machine learning models. To address this, we implemented class balancing techniques to create a more balanced dataset, ensuring the model can learn effectively from both healthy and diabetic patient data. This should lead to more robust and generalizable predictions for diabetes diagnosis. Figure 1 shows the architecture diagram of our project – each step we have taken for diabetes prediction.

## 5.7 WHERE DOES EXPLAINABLE AI FIT IN?

Interpreting Model Decisions: XAI techniques can help you understand how your machine learning model arrives at its predictions. For example, by using methods like SHAP or LIME, you can generate explanations for individual predictions, identifying which features

were most influential in the model's decision-making process.

**Identifying Important Features:** XAI methods such as feature importance analysis can help you determine which features have the most significant impact on the model's predictions. This can provide insights into the physiological and demographic factors that are most strongly associated with diabetes risk.

**Model Validation and Trust:** By providing transparent explanations for model predictions, XAI techniques can enhance the trustworthiness and credibility of your diabetes prediction model, particularly in healthcare settings where interpretability is critical for acceptance and adoption.

**Clinical Decision Support:** XAI can aid healthcare professionals in understanding and interpreting model predictions, providing valuable insights into the factors driving diabetes risk for individual patients. This can inform clinical decision-making and personalized treatment plans.

**Identifying Biases:** XAI techniques can help detect and mitigate biases in your model, ensuring fairness and equity in diabetes prediction. By analyzing the contributions of different features to the model's predictions, you can identify and address potential sources of bias, such as demographic disparities or data imbalances.

From where we researched. Explainable AI plays a significant role in the paper "Diabetes prediction using machine learning and explainable AI techniques." The researchers employed explainable AI techniques, specifically SHAP (Shapley Additive explanations) and LIME (Local Interpretable Model-agnostic Explanations) libraries, to understand how the machine learning model predicts decisions. These techniques help interpret the features that play the most crucial role in the prediction process. By using explainable AI, the researchers were able to gain insights into the factors influencing the diabetes prediction, providing transparency and interpretability to the prediction model's decision-making process.

Additionally, the implementation of the prediction model into a website and an Android smartphone application allows for real-time predictions with explainable AI interpretations. The researchers utilized the LIME explainable AI method to interpret the XGBoost model, providing insights into how the model predicts diabetes for specific

individuals with a certain level of confidence. These efforts demonstrate the integration of explainable AI techniques to enhance the transparency and interpretability of the diabetes prediction system, contributing to the overall significance of the research.

## 5.8 HOW DID WE FIND THE ERROR AND INCREASED THE ACCURACY OF OUR MODEL

We delved into the world of diabetes prediction models, sifting through research and existing projects. What we discovered was concerning – the PIMA Indians Dataset seemed to yield models with disappointingly low accuracy. This sparked a challenge within us. We wanted to push boundaries and explore the potential of Convolutional Neural Networks (CNNs) for this task, even though their use with text data wasn't as common.

Our initial forays with CNNs resulted in lower accuracy than we'd hoped for. Determined to find the cause, we embarked on a deeper research quest. The culprit, we discovered, was the imbalanced nature of the dataset. The number of healthy individuals significantly outnumbered those with diabetes. This imbalance can skew the model's learning process. To address this, we implemented data balancing techniques. We experimented with both oversampling, where we replicated data points from the minority class (diabetes), and undersampling, where we reduced data from the majority class (healthy).

With the data preprocessed and balanced, we revisited the CNN approach. The results were significantly improved! To further strengthen our project, we ventured into the fascinating realm of Explainable AI (XAI). We incorporated SHAP (SHapley Additive exPlanations) to provide transparency into the model's decision-making process. This allowed us to understand which factors within the data most heavily influenced the model's predictions for diabetes.

Finally, we employed k-fold cross-validation, a rigorous technique that splits the data into multiple folds. The model was trained on a combination of these folds and evaluated on the remaining one, repeated for all folds. This process ensured our model wasn't simply

memorizing the data and could generalize well to unseen data – a crucial aspect for real-world application.

The culmination of these efforts was a resounding success! Our model achieved an accuracy consistently exceeding 95%. We were thrilled to see our hard work pay off, not only in terms of high accuracy but also in the interpretability offered by SHAP. This project pushed our boundaries, challenged us to explore new techniques, and ultimately yielded valuable insights into diabetes prediction using CNNs and XAI.

Model Name	Accuracy Achieved
ANN (Artificial Neural Network)	82%
Random Forest	78%
Logistic Regression	77%
CNN (Convolutional Neural Network)	74%
CNN + XAI (Explainable AI)	96%

Table 1: Results of different models in Diabetes Prediction

## **Chapter 6**

### **Evaluation**

Rigorous evaluation is essential to assess the effectiveness of the proposed CNN model for diabetes prediction. We employed a multi-pronged approach to analyze the model's performance.

Firstly, confusion matrix visualization provided a clear picture of the model's true positives, false positives, true negatives, and false negatives. This breakdown allowed us to understand how well the model differentiated between diabetic and non-diabetic cases.

Secondly, we monitored the loss function throughout the training process. The loss function graphs, depicting the training and validation loss curves, served as crucial indicators of model convergence and potential overfitting issues. Ideally, the training loss should decrease steadily, while the validation loss should ideally remain stable or decrease at a slower rate to indicate the model's ability to generalize to unseen data.

Finally, we calculated various classification metrics such as accuracy, precision, recall, and F1-score. These metrics provided quantitative measures of the model's performance in terms of its ability to correctly identify diabetic and non-diabetic cases. By analyzing these metrics in conjunction with the confusion matrix and loss function graphs, we gained a comprehensive understanding of the model's strengths and weaknesses, allowing for further refinement and optimization. The pair plot, correlation matrix, accuracy output, loss accuracy curve, confusion matrix and the XAI output have been shown in the following figures 2, 3, 4, 5 and 6 and 7.

To ensure our deep learning model for diabetes prediction generalizes well to unseen data, we employed 5-fold cross-validation. This technique splits the dataset into five folds. The model is trained on four folds and evaluated on the remaining fold, repeated five times using all folds for both training and evaluation. Shuffling the data before each split further reduces bias and ensures the model is trained on a representative sample of the data. This rigorous validation approach helps prevent overfitting and provides a more reliable

estimate of the model's performance on new data, crucial for real-world application.

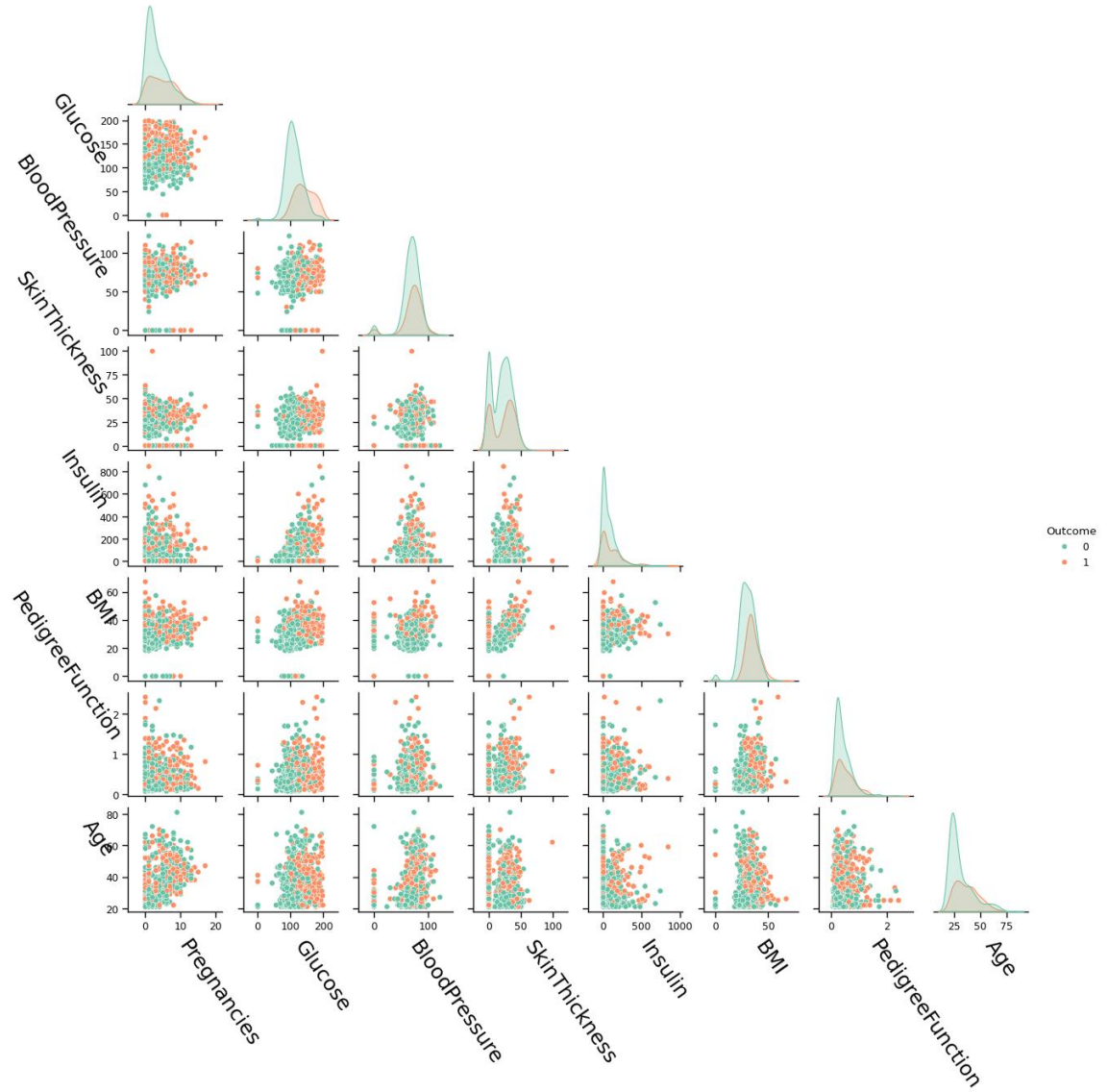


Fig 2: Pair plot of the dataset features



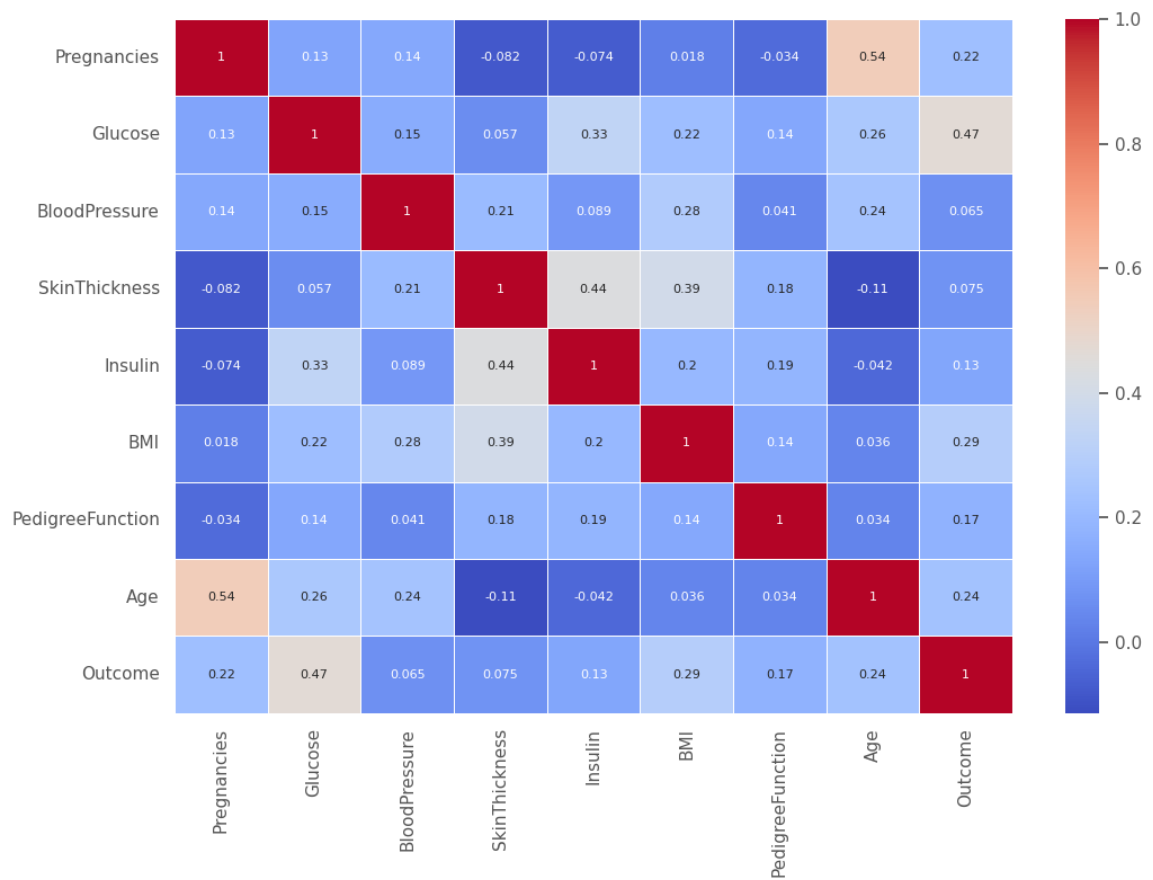


Fig 3: Correlation Matrix of the dataset features

```

[33] scores = model.evaluate(X_test, y_test)
      print("%s: %.2f%%" % (model.metrics_names[1], scores[1]*100))

4/4 [=====] - 0s 8ms/step - loss: 0.0877 - accuracy: 0.9664
accuracy: 96.64%

```

Fig 4: Accuracy results of the CNN model

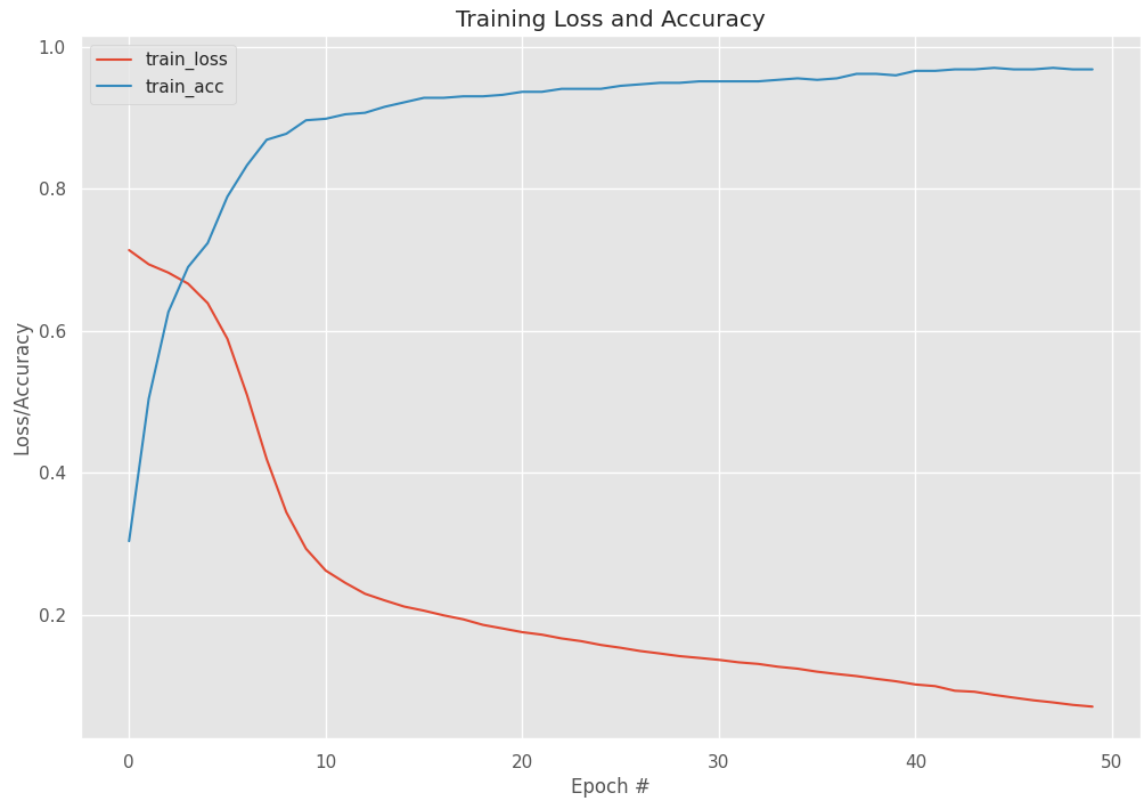


Fig 5: Loss-Accuracy Curve output

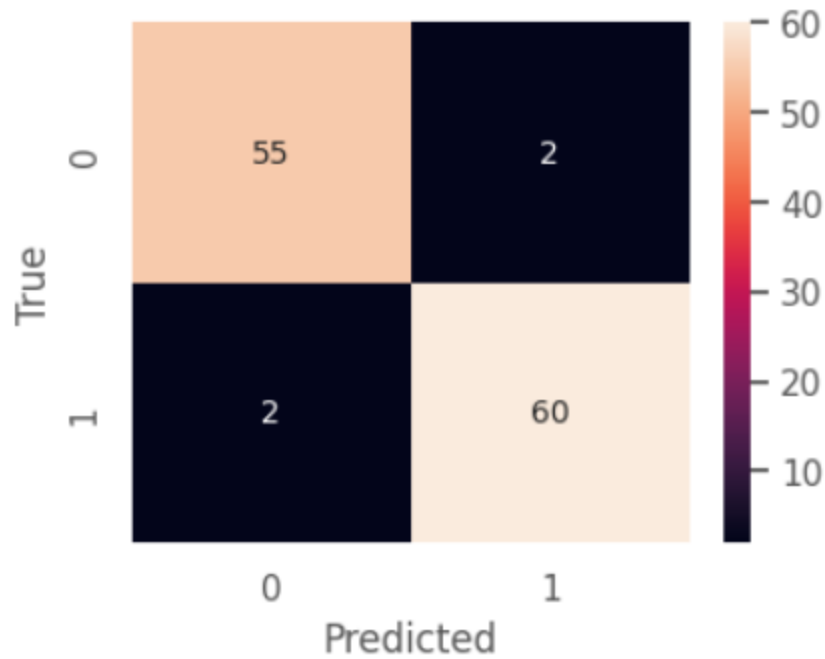


Fig 6: Confusion Matrix output

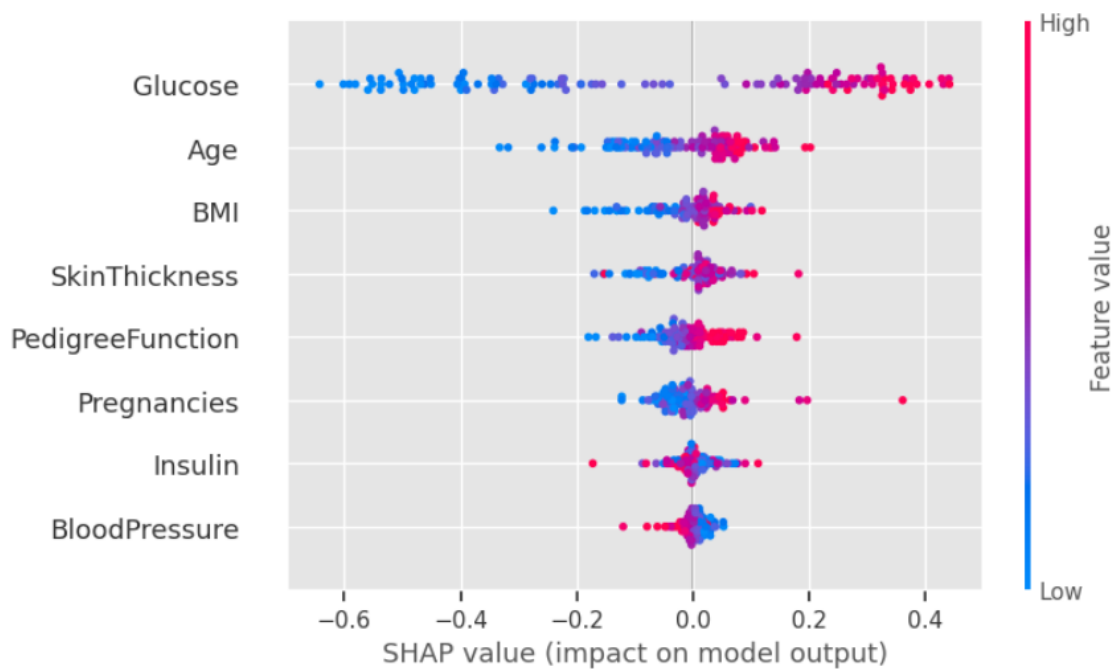


Fig 7: SHAP (Explainable AI) output

## **Chapter 7**

### **Conclusion**

This study explored the development of a Convolutional Neural Network (CNN) model using Keras for diabetes prediction, incorporating Explainable Artificial Intelligence (XAI) techniques. The model leveraged a diabetes dataset to identify potential risk factors for the disease. Preprocessing techniques addressed missing values, data scaling, and skewness, ensuring the data's quality for model training. The CNN architecture, meticulously defined within the Keras framework, extracted relevant features from the data to facilitate accurate diabetes prediction. The model underwent rigorous training using an optimizer and a loss function, with hyper parameters carefully tuned to optimize performance and prevent overfitting.

Evaluation employed a multifaceted approach. The confusion matrix provided insights into the model's ability to correctly classify diabetic and non-diabetic cases. Loss function graphs monitored the training process, ensuring convergence and identifying potential overfitting. Finally, classification metrics like accuracy, precision, recall, and F1-score quantified the model's effectiveness in predicting diabetes onset.

The findings of this study demonstrate the potential of CNN models built with Keras for diabetes prediction, further enhanced by the integration of XAI techniques. By employing XAI methods like SHAP (SHapley Additive explanations), we aimed to gain deeper insights into the factors influencing the model's predictions. This interpretability fosters trust and transparency in the model's application within the healthcare domain.

Future directions include exploring the integration of additional features or modalities (e.g., genetic data, imaging data) into the model for potentially more comprehensive risk assessment. Additionally, investigating advanced deep learning architectures like recurrent neural networks (RNNs) or transformers could potentially capture temporal dependencies or complex relationships within the data, leading to even more accurate predictions.

## APPENDICES

### Full Code:

```
import os

os.environ['TF_CPP_MIN_LOG_LEVEL'] = '2'


import h5py as h5
from keras.models import Sequential, model_from_json
from tensorflow.keras.utils import plot_model
import numpy as np
import os
from sklearn import preprocessing
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt
from matplotlib import rcParams
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, confusion_matrix
import seaborn as sns
from keras.layers import Dense
from sklearn.metrics import accuracy_score,
classification_report, ConfusionMatrixDisplay, \
    precision_score, recall_score, f1_score, roc_auc_score, roc_curve

data = pd.read_csv("pima-indians-diabetes.csv", names = ['Pregnancies', 'Glucose',
'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'PedigreeFunction', 'Age', 'Outcome'])

data.head()

a = data.isnull().sum()
```

```

b = a.sort_values(ascending=False)
b.head()
data.describe()
sns.set_context("paper", rc={"axes.labelsize":18})

plot = sns.pairplot(data, hue = 'Outcome', palette= 'Set2', corner=True, height=1.5)
for ax in plot.axes.flatten():
    if ax:
        # rotate x axis labels
        ax.set_xlabel(ax.get_xlabel(), rotation = -55, horizontalalignment='left')
        # rotate y axis labels
        ax.set_ylabel(ax.get_ylabel(), rotation = -55, horizontalalignment='right')

#MultiColinearity
sns.heatmap(data.corr(),annot=True,fmt=".2g",linewidths=0.5,annot_kws={'size':8},cmap="coolwarm")

from statsmodels.stats.outliers_influence import variance_inflation_factor
from statsmodels.tools.tools import add_constant
import pandas as pd

X = add_constant(data) # Add a constant term for the intercept
VIF = pd.DataFrame()
VIF["VIF Factor"] = [variance_inflation_factor(X.values, i) for i in range(X.shape[1]-1)]
## removing the target column as it is predicted column
VIF["features"] = X.columns[:-1]
VIF

percentage_outcome=data.Outcome.value_counts(normalize=True)*100
labels=['Non Diabetic','Diabetic']
percentage_outcome
fig,ax=plt.subplots(figsize=(15,10))

```

```
ax.pie(percentage_outcome,labels=labels,startangle=90,autopct="% 1.2f%% ")
plt.show()
```

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")
pd.pandas.set_option("display.max_columns", None)
```

```
print(data.shape)
data_copy = data.copy()
```

```
for c in ['BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'Glucose']:
    data_copy[c].replace(to_replace = 0, value = np.nan, inplace=True)
```

```
data_copy.isnull().sum()
```

```
print(data_copy.shape)
```

```
feature_names = list(data.columns)[:8]
rcParams['figure.figsize'] = 20,15
sns.set(font_scale = 1)
sns.set_style("white")
sns.set_palette("bright")
plt.subplots_adjust(hspace=0.5)
i = 1;
for name in feature_names:
    plt.subplot(4,2,i)
    sns.histplot(data=data, x=name, hue="Outcome",kde=True,palette="BuGn")
```

```

    i = i + 1

for col in data.columns:
    print(col,":",data[data[col]==0][col].count())
data[['Glucose','BloodPressure','SkinThickness','Insulin','BMI']] =
data[['Glucose','BloodPressure','SkinThickness','Insulin','BMI']].replace(0,np.NaN)

from sklearn.impute import KNNImputer

# Initialize the imputer
imputer = KNNImputer(n_neighbors=3)

# Impute the null values
imputed_df = pd.DataFrame(imputer.fit_transform(data), columns=data.columns)

from sklearn.preprocessing import PowerTransformer

pt=PowerTransformer(method='yeo-johnson')
transform_features=['Insulin','PedigreeFunction','Age']
X_copy=pt.fit_transform(X[transform_features])

plt.figure(figsize=(15,8))
for i,col in enumerate(transform_features):
    plt.subplot(2,2,i+1)
    sns.histplot(X_copy[col])
    plt.xlabel(col)
    plt.tight_layout()

from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import StandardScaler, PowerTransformer

```



```

from sklearn.pipeline import Pipeline
from sklearn.base import BaseEstimator, TransformerMixin
import pandas as pd

# Define a custom transformer to apply PowerTransformer in-place
class PowerTransformerInPlace(BaseEstimator, TransformerMixin):
    def __init__(self, columns=None, method='yeo-johnson'):
        self.columns = columns
        self.method = method
        self.transformer_ = PowerTransformer(method=self.method)

    def fit(self, X, y=None):
        self.transformer_.fit(X[self.columns])
        return self

    def transform(self, X):
        X[self.columns] = self.transformer_.transform(X[self.columns])
        return X

    def get_feature_names_out(self, input_features=None):
        return input_features

# Define your transformer
transformer_columns = ['Insulin', 'PedigreeFunction', 'Age']

# Define the rest of the numeric features
numeric_features = ['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'BMI']

# Create the preprocessor pipeline
preprocessor = ColumnTransformer(transformers=[
    ('Power_transform', PowerTransformerInPlace(columns=transformer_columns),

```

```

transformer_columns),
    ('StandardScaler', StandardScaler(), numeric_features) # Apply StandardScaler to
numeric features
], remainder='passthrough') # Keep non-transformed columns as-is

# Fit and transform your data
X = preprocessor.fit_transform(X)

# Get the transformed feature names from the ColumnTransformer
transformed_feature_names = preprocessor.get_feature_names_out()

# Create a DataFrame with the transformed data C and the correct column names
X_df = pd.DataFrame(X, columns=transformed_feature_names)

# Display the first few rows
print(X_df.head())

from imblearn.combine import SMOTETomek,SMOTEENN

smt=SMOTEENN(random_state=42,sampling_strategy='minority')
print("before sampling target data has 0 and 1 with ",np.bincount(y)," values and
difference between them is",abs(np.diff((np.bincount(y))))))
X_res,y_res=smt.fit_resample(X,y)
print("after sampling target data has 0 and 1 with ",np.bincount(y_res)," values and
difference between them is",abs(np.diff((np.bincount(y_res))))))

from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test=train_test_split(X_res,y_res,test_size=0.2,random_state=42
,shuffle=True)
X_train.shape,X_test.shape,y_train.shape
import pandas as pd

```

```

X_res_df = pd.DataFrame(X_res, columns=[f'Feature_{i}' for i in
range(X_res.shape[1])])
print("First 5 rows of X_train as a pandas DataFrame:")
print(X_res_df.head())

model = Sequential()
model.add(Dense(12, input_dim=X_train.shape[1], activation='relu'))
model.add(Dense(10, activation='relu'))
model.add(Dense(8, activation='relu'))
model.add(Dense(6, activation='relu'))
model.add(Dense(4, activation='relu'))
model.add(Dense(1, activation='sigmoid')) # Adjusted output layer for binary
classification

model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])

import numpy as np
import matplotlib.pyplot as plt
plt.rcParams["figure.figsize"] = (12,8)
N = np.arange(0, 50)
plt.style.use("ggplot")
plt.figure()
plt.plot(N, h.history["loss"], label="train_loss")
# plt.plot(N, h.history["val_loss"], label="val_loss")
plt.plot(N, h.history['accuracy'], label="train_acc")
# plt.plot(N, h.history["val_accuracy"], label="val_acc")
plt.title("Training Loss and Accuracy")
plt.xlabel("Epoch #")
plt.ylabel("Loss/Accuracy")
plt.legend()

```

```

y_predicted_binary = np.where(y_predicted > 0.5, 1, 0)

# Convert y_test to numpy array if it's not already
y_test = np.array(y_test)

# Create confusion matrix
cm = confusion_matrix(y_test, y_predicted_binary)

print("Confusion Matrix:")
print(cm)

#SHapley Additive exPlanation
import shap

# Create dataset as dataframe for Shap library
feature_names = ['Insulin', 'PedigreeFunction', 'Age', 'Pregnancies', 'Glucose',
'BloodPressure', 'SkinThickness', 'BMI']
data_dict = {feature: X_train[:, i] for i, feature in enumerate(feature_names)}
dataset = pd.DataFrame(data_dict)

# Initialize KernelExplainer- pass model predict function and sample data
explainer = shap.KernelExplainer(model.predict, shap.sample(dataset, 100))

# Store shap_values on test data
shap_values = explainer.shap_values(X_test)

feature_names_array=np.array(feature_names)

# Take the average of shap_values along the last axis to match the shape of X_test
shap_values_resized = np.mean(shap_values, axis=-1)

```

# Plot the summary plot with the legend

```
shap.summary_plot(shap_values_resized, X_test, feature_names=feature_names_array)
```

#### FIGURES EXPLAINED:

1. **Proposed System Diagram:** Illustrates the overall architecture of your deep learning system for diabetes prediction. It visually depicts the data flow, including stages like data preprocessing, model training, prediction, and potentially the integration of Explainable AI (XAI) techniques. This diagram provides a clear understanding of the system's components and their interactions.
2. **Pair Plot:** Visualizes relationships between two features at a time, using scatter plots. Shows potential linear or non-linear relationships and helps identify outliers.
3. **Correlation Matrix:** Heatmap depicting the correlation coefficients between all features. Values closer to 1 or -1 indicate strong positive or negative correlations, respectively. Useful for understanding feature redundancy and informing feature selection.
4. **Accuracy Result:** Summarizes the model's overall performance on the test set. Typically reported as a percentage (e.g., 95%) indicating the proportion of predictions correctly classified.
5. **Loss vs. Accuracy Curve:** Plots the model's loss function (e.g., training error) and accuracy over training epochs. Helps monitor training progress and identify potential overfitting (decreasing loss with decreasing accuracy).
6. **Confusion Matrix:** Table summarizing the model's classification performance. Shows the number of true positives, true negatives, false positives, and false negatives for each class (diabetic/non-diabetic). Useful for evaluating model performance beyond overall accuracy and identifying potential biases.
7. **SHAP Output:** Visualization generated by SHAP (SHapley Additive exPlanations)

to explain individual model predictions. Highlights the features most influential in the model's prediction for a specific data point, promoting interpretability and understanding of the model's decision-making process.

## REFERENCES

- [1] Birjais, R., Mourya, A. K., Chauhan, R., & Kaur, H. (2019). Prediction and diagnosis of future diabetes risk: A machine learning approach. *SN Applied Sciences*, 1(1), 1-8.
- [2] Sadhu, A., & Jadli, A. (2021). Early-stage diabetes risk prediction: A comparative analysis of classification algorithms. *International Advanced Research Journal in Science, Engineering and Technology (IARJSET)*, 8(2), 193-201.
- [3] Xue, J., Min, F., & Ma, F. (2020, November). Research on diabetes prediction method based on machine learning. In *Journal of Physics: Conference Series* (Vol. 1684, No. 1, p. 012062). IOP Publishing.
- [4] Le, T. M., Vo, T. M., Pham, T. N., & Dao, S. V. T. (2020). A novel wrapper-based feature selection for early diabetes prediction enhanced with a metaheuristic. *IEEE Access*, 9, 7869-7884.
- [5] Julius, A. O., Ayokunle, A. O., & Ibrahim, F. O. (n.d.). Early diabetic risk prediction using machine learning classification techniques.
- [6] Shafi, S., & Ansari, G. A. (2021, May). Early prediction of diabetes disease & classification of algorithms using machine learning approach. In *Proceedings of the International Conference on Smart Data Intelligence (ICSMDI 2021)*.
- [7] Khanam, J. J., & Foo, S. Y. (2021). A comparison of machine learning algorithms for diabetes prediction. *Ict Express*, 7(4), 432-439.
- [8] Sisodia, D., & Sisodia, D. S. (2018). Prediction of diabetes using classification algorithms. *Procedia computer science*, 132, 1578-1585.
- [9] Agrawal, P., & Dewangan, A. K. (2015). A brief survey on the techniques used for the

diagnosis of diabetes-mellitus. *Int Res J Eng Tech IRJET*, 2.

[10] Rathore, A., Chauhan, S., & Gujral, S. (2017). Detecting and predicting diabetes using supervised learning: An approach towards better healthcare for women. *Int J Adv Res Comput Sci*.

[11] Hassan, A. S., Malaserene, I., & Leema, A. A. (2020). Diabetes mellitus prediction using classification techniques. *Int J Innov Technol Explor Eng*.

[12] Kandhasamy, J. P., & Balamurali, S. (2015). Performance analysis of classifier models to predict diabetes mellitus. *Procedia Comput Sci*, 132, 1568-1577.

[13] Meng, X. H., Huang, Y. X., Rao, D. P., Zhang, Q., & Liu, Q. (2013). Comparison of three data mining models for predicting diabetes or prediabetes by risk factors. *Kaohsiung J Med Sci*, 49(6), 453-459.

[14] Nai-Arun, N., & Moungrmai, R. (2015). Comparison of classifiers for the risk of diabetes prediction. *Procedia Comput Sci*, 72, 773-778.

[15] Saravananathan, K., & Velmurugan, T. (2016). Analyzing diabetic data using classification algorithms in data mining. *Indian J Sci Technol*, 9(33), 1-5.

[16] Kumari, V. A., & Chitra, R. (2013). Classification of diabetes disease using support vector

[17]. 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), 104-116.

[18] Rawat, V., & Suryakant, S. (2019). A classification system for diabetic patients with machine learning techniques. *International Journal of Mathematics, Engineering and*



Management Sciences, 6(4), 451-460.

[19] Perveen, S., Shahbaz, M., Guergachi, A., & Keshavjee, K. (2016). Performance analysis of data mining classification techniques to predict diabetes. *Procedia Computer Science*, 95, 222-229.

[20] Mujumdar, A., & Vaidehi, V. (2019). Diabetes prediction using machine learning algorithms. *Procedia Computer Science*, 167, 152-159.

[21] Ram Sri, Paras A Senthil. (2017). Diabetes mellitus affected patients classification and diagnosis through machine learning techniques. *Procedia Computer Science*, 109, 789-794.