# DIABETES DETECTION AND DIET PLAN GENERATION USING MACHINE LEARNING

## A PROJECT REPORT

*Submitted by*

**ARCHANA. R**

**HARIPRIYA. R**

**KAVITHA. P**

*In partial fulfillment for the award of the degree*

*of*

**BACHELOR OF ENGINEERING**

*in*

**COMPUTER SCIENCE AND ENGINEERING**

# P. A. COLLEGE OF ENGINEERING AND TECHNOLOGY

(An Autonomous Institution)

Pollachi, Coimbatore Dt. -  642 002

**NOVEMBER 2024**

# P. A. COLLEGE OF ENGINEERING AND TECHNOLOGY

## BONAFIDE CERTIFICATE

Certified that this project report **"DIABETES DETECTION AND DIET PLAN GENERATION USING MACHINE LEARNING"** is the bonafide work of " **ARCHANA R (721721104009), HARIPRIYA R (721721104030), and KAVITHA P (721721104050)"** who carried out the project work under my supervision.


--------------------------------        ---------------------------------

**SIGNATURE**                                 **SIGNATURE**

**Dr. D. CHITRA M.E, Ph.D.,**          **Mrs. K. TAMILSELVI, M.E.,**

Professor                                  Assistant Professor

**HEAD OF THE DEPARTMENT**        **SUPERVISOR**

Department of Computer Science       Department of Computer Science

and Engineering                         and Engineering

P. A. College of Engineering and      P. A. College of Engineering and

Technology                              Technology

Pollachi.                                   Pollachi.


Submitted to the Viva-Voce Examination held on --------------------


-------------------------         -------------------------

**INTERNAL EXAMINER**        **EXTERNAL EXAMINER**

# ACKNOWLEDGEMENT

First and foremost, we thank the GOD ALMIGHTY for blessing us with the mental strength that was needed to carry out the project work. We thank our Chairman **Dr. P. APPUKUTTY, ME, FIE, FIV.,** for his extensive support to successfully carry out this project.

We take privilege in expressing our sincere and heartfelt thanks and gratitude to our beloved Principal **Dr. T. MANIGANDAN, M.E., Ph.D.,** for providing us an opportunity to carry out this project work.

We express our heartfelt thanks to **Dr. D. CHITRA, M.E., Ph.D.,** Professor and Head, Department of Computer Science and Engineering, for her technical guidance and constructive suggestions provided throughout the project work.

We take conceit in express our sincere and deepest thanks to our project guide **Mrs. K. TAMILSELVI, M.E.,** Assistant Professor, Department of Computer Science and Engineering, for her technical guidance, constructive criticism and many valuable suggestions provided throughout the project work.

We take this opportunity to thank and pay gratitude to our Project Coordinator **Mr. T. DINESH KUMAR, M.Tech.,** Assistant Professor, Department of Computer Science and Engineering and all teaching and non-teaching staff members of our department, for their encouragement and valuable suggestion.

We take this opportunity to express our gratitude to our parents, friends, family and other members whose belongings and love have always been with us to carry out this project work successfully.

## ABSTRACT

Diabetes is a global health concern that requires timely intervention to minimize associated complications. The integration of advanced technologies, such as machine learning (ML), has shown significant potential in enhancing the detection and management of diabetes. This framework emphasizes leveraging classification algorithms for accurate diabetes risk assessment, complemented by robust data preprocessing techniques like feature engineering, outlier handling, and hyperparameter optimization. Furthermore, the system promotes personalized healthcare by integrating nutritional insights tailored to individual health metrics. By employing predictive analytics, the methodology supports informed decision-making for healthcare providers and patients alike, ultimately aiming to improve outcomes through early detection and customized dietary guidance.

The approach not only underlines the role of data-driven solutions in modern healthcare but also highlights the adaptability of machine learning systems to provide accessible, efficient, and personalized healthcare tools, fostering advancements in diabetes care and prevention.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF ABBREVATION

| | |
|---|---|
| **AI** | Artificial intelligence |
| **API** | Application program interface |
| **ML** | Machine learning |
| **SVM** | Support vector machine |
| **KNN** | K - nearest neighbors |
| **DT** | Decision tree |
| **RF** | Random forest |
| **GB** | Gradient boosting |
| **BMI** | Body mass index |
| **DPF** | Diabetes pedigree function |
| **CSV** | Comma – separate values |
| **EDA** | Exploratory data analysis |

# CHAPTER 1
# INTRODUCTION

Diabetes is a chronic condition that affects millions worldwide. Early detection is crucial for effective management and prevention of complications. Machine learning (ML) offers a promising solution for diabetes detection by analyzing patient data such as age, blood glucose levels, blood pressure, BMI, and family history. The integration of various ML models, such as Logistic Regression, KNN, and SVM, allows for accurate predictions of diabetes risk based on these factors.

The system generates personalized diet plans for patients. Once diabetes is detected, the diet plan is tailored to the individual's health needs, promoting better management of the condition through dietary adjustments. By combining machine learning for detection and diet generation, the system aims to offer a comprehensive approach to managing diabetes, improving both early diagnosis and patient care.
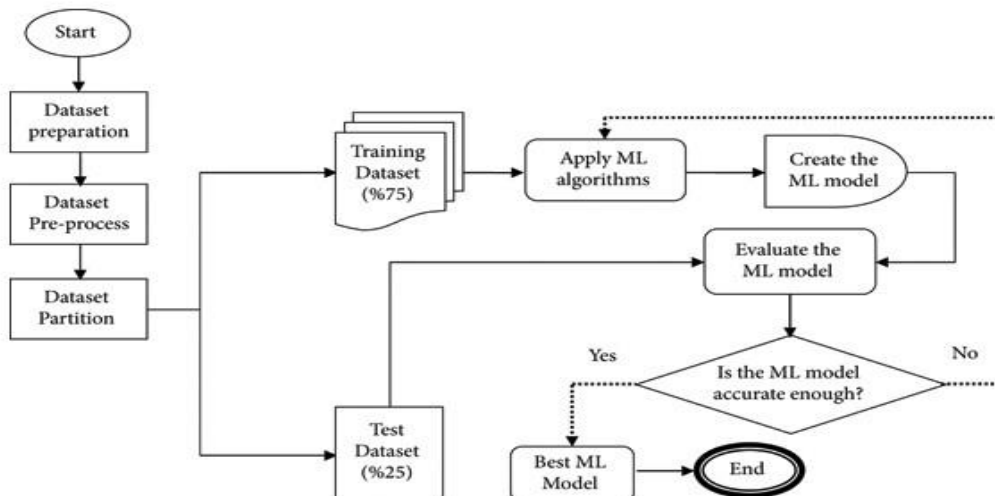


Figure 1.1 A Proposed Technique Using ML for Diabetes Prediction

Figure 1.1 illustrates the proposed technique using Machine Learning for diabetes prediction. This model leverages various machine learning algorithms to analyze input features and predict the likelihood of diabetes, offering a reliable and efficient approach for early diagnosis.

## 1.1 NEED FOR DIABETES DETECTION

Diabetes detection and diet plan generation using machine learning is crucial due to the increasing prevalence of diabetes worldwide. The need for such a system stems from various key factors. Machine learning models can analyze patient data to detect diabetes at an early stage, enabling timely intervention. Early detection is crucial for preventing complications such as cardiovascular diseases, kidney failure, and blindness.

The system generates personalized diet plans for patients based on their specific medical data. This personalized approach enhances the effectiveness of diabetes management, as it takes into account individual health conditions, preferences, and lifestyle factors. Machine learning algorithms such as Logistic Regression, KNN, SVM, and Random Forest, the system ensures high accuracy in predicting diabetes risk.

The automated nature of the system makes it an accessible tool for users, providing diabetes predictions and diet recommendations without requiring frequent consultations with healthcare providers. This increases accessibility to quality healthcare. The diabetes detection, optimizing machine learning models is essential for achieving accurate predictions and efficient performance. This optimization ensures that the system can handle large datasets with minimal computational resources, improving response time and scalability.

The development of the diabetes prediction system involves close collaboration among team members. Code reviews and discussions on the use of machine learning algorithms help foster teamwork, ensuring that all contributors are aligned with the project goals and coding best practices.For team members with less experience in machine learning or healthcare-related projects, code reviews offer valuable learning opportunities. Feedback from experienced developers and data scientists helps junior members enhance their understanding of algorithms, model evaluation, and data preprocessing techniques.

As new data and techniques emerge, the diabetes detection system can be continuously improved. Regular model reviews and updates ensure the system remains accurate, adaptable to new conditions, and capable of providing more personalized health recommendations over time.Ensuring the privacy and security of patient data is a key consideration in healthcare projects. Code reviews help ensure that the diabetes detection system complies with data protection regulations and follows security best practices, safeguarding sensitive health information.

The potential issues during the model training and testing phases, the time spent debugging and troubleshooting is significantly reduced. Early identification of data preprocessing issues, model inconsistencies, or performance bottlenecks helps ensure smoother development cycles, enabling quicker deployment of the diabetes prediction system and its diet planning features.

## 1.2 MACHINE LEARNING

Machine Learning (ML) is a subset of artificial intelligence (AI) that focuses on developing algorithms and models that enable computers to learn

from data. The primary goal is to allow computers to make predictions, decisions, or identifications without being explicitly programmed for a particular task.



Figure 1.2: Concept of ML

Figure 1.2 illustrates the concept of Machine Learning. The process begins with data collection, followed by training the machine using this data. Next, a model is built based on the training, and finally, the model is used for predicting outcomes.

**TYPES OF MACHINE LEARNING**

**i) SUPERVISED LEARNING**

The algorithm is trained on a labeled dataset, where the input data is paired with corresponding output labels. It learns to map inputs to outputs, making predictions on new, unseen data.

**ii) UNSUPERVISED LEARNING**

The algorithm explores patterns and structures within unlabeled data. Clustering and dimensionality reduction are common tasks in unsupervised learning.

### iii) REINFORCEMENT LEARNING

The algorithm learns through interaction with an environment. It receives feedback in the form of rewards or penalties based on its actions, enabling it to optimize its behavior.



Figure 1.3: Types of Machine Learning

Figure 1.3 illustrates the types of Machine Learning: Supervised Learning (using labeled data), Unsupervised Learning (finding patterns in unlabeled data), and Reinforcement Learning (learning through interaction to maximize rewards).

### WORKING OF ML

Machine learning works by enabling computers to learn patterns from data and make predictions or decisions without explicit programming. It begins with collecting relevant data, which is then cleaned and preprocessed to handle missing values, remove noise, and format it suitably. Features, which are key characteristics of the data, are selected or engineered to

enhance the model's performance. An appropriate algorithm is chosen based on the task, such as classification, regression, or clustering. The model is trained using a portion of the data, during which it learns by adjusting parameters to minimize errors. Once trained, the model is evaluated using a separate dataset to measure its performance through metrics like accuracy and precision. Adjustments are made through hyperparameter tuning to improve results further. After achieving satisfactory performance, the model is deployed to make predictions on new data. Over time, a feedback loop ensures the model adapts to changes and improves by retraining with updated data, making it more accurate and reliable. This cycle allows machines to solve problems effectively and evolve continuously.
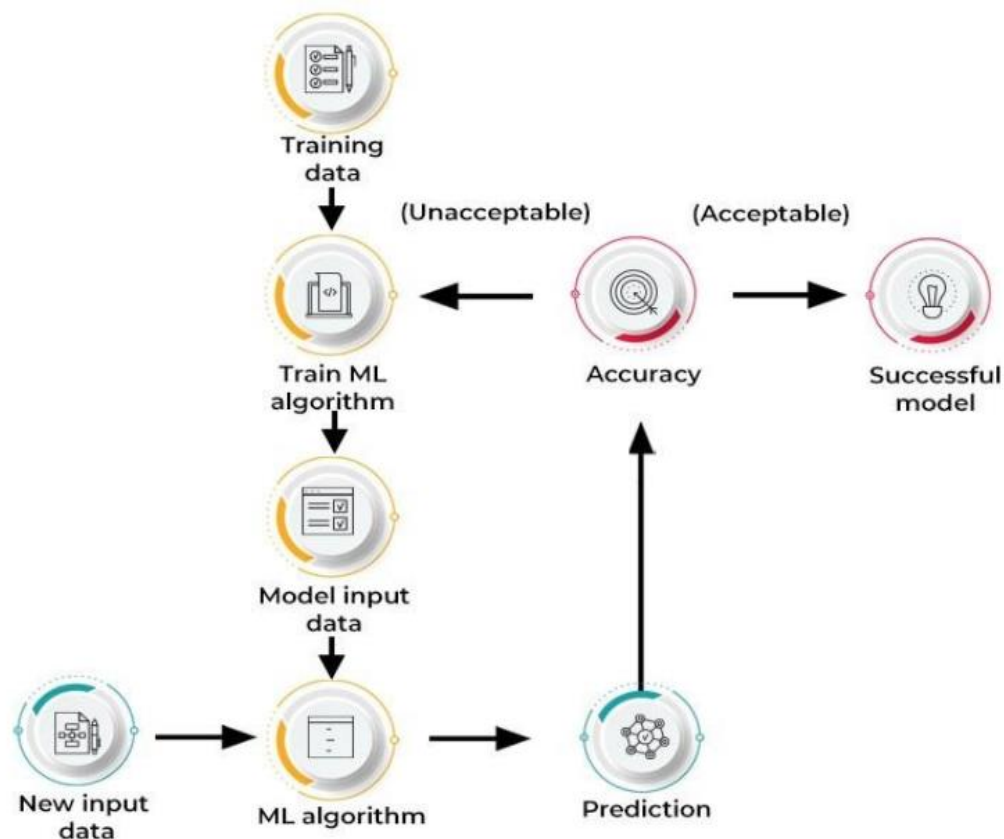


Figure 1.4: Working of ML

Figure 1.4 illustrates the working of Machine Learning. It starts with collecting data, followed by preprocessing and feature extraction.

## 1.3 ADVANTAGE OF ML IN DIABETES DETECTION

Machine learning algorithms streamline the process of diabetes detection by automating data analysis, allowing healthcare professionals to make quicker, data-driven decisions. This reduces the time spent on manual diagnosis, enhancing efficiency in diabetes care. ML models consistently apply the same criteria to evaluate patient data, ensuring that each prediction is based on the same high standards and reducing the chance of inconsistencies that might arise with human judgment. Machine learning models can detect early signs of diabetes from patient data, such as glucose levels and other health parameters, that might go unnoticed through traditional methods. Early identification allows for timely intervention and better management of the condition.

ML models can analyze large datasets quickly and accurately, making them scalable for a wide range of patients. This is particularly useful for handling big data in healthcare systems, ensuring that each individual receives personalized attention without overburdening healthcare professionals. Machine learning models continuously improve as they process more patient data. This adaptive nature allows the system to evolve, refining its predictions over time, and ensuring that it stays relevant and accurate as new patterns emerge in diabetes management. ML-based diabetes detection minimizes the risk of human error in diagnosing the condition. By relying on data-driven predictions, the chances of overlooking critical health indicators are reduced, improving diagnostic accuracy.

ML algorithms can be tailored to suit specific patient populations or healthcare systems. This flexibility allows the system to cater to different demographic groups, considering factors like age, lifestyle, and medical history, ensuring that the diabetes prediction model is highly personalized. Machine learning can automate the diagnostic process, enabling faster and

more accurate diabetes identification. This reduces the workload for healthcare providers and allows for more efficient allocation of resources. By automating data analysis, machine learning reduces the time spent on manual assessment, allowing healthcare professionals to focus on making informed decisions and crafting personalized treatment plans for patients. Machine learning models provide an unbiased, objective analysis of patient data, eliminating personal biases in diagnosis. This leads to more accurate and reliable outcomes in diabetes detection and management.

## 1.4 LIMITATIONS AND CHALLENGES OF ML IN DIABETES DETECTION

Machine learning (ML) models in healthcare, particularly for diabetes detection, face several limitations and challenges. These models rely heavily on high-quality data, and inaccurate, incomplete, or biased data can lead to unreliable predictions. Collecting large datasets with accurate labels is both time-consuming and expensive. Healthcare data is sensitive and governed by strict privacy regulations, such as HIPAA, making it challenging to ensure the security of patient information during model training. Breaches in data security can undermine trust in these technologies. Additionally, many ML models, especially deep learning, operate as "black boxes," making it difficult for healthcare professionals to interpret predictions, which can hinder adoption in medical practice where interpretability is essential.

Another challenge is overfitting, where models perform well on training data but fail to generalize to new, unseen data. This issue necessitates robust validation using diverse datasets to improve generalization. The lack of labeled data is another major hurdle, as annotating medical datasets with patient outcomes is expensive and often inaccessible, limiting the training of reliable models. Diabetes detection models require careful selection of

features like genetics, lifestyle, and medical history. Identifying the most relevant features from a large set is complex and demands domain expertise, adding to the complexity of model development.

ML models must also evolve alongside healthcare practices to remain effective. Continuous monitoring and updating are crucial, as outdated models can lead to inaccurate predictions. Integrating these models with existing medical systems poses additional challenges, as adapting workflows and training staff can incur significant costs. The initial setup and ongoing maintenance of ML systems for diabetes detection can also be expensive, making it difficult for resource-limited healthcare institutions to implement them. Furthermore, biases in ML models, stemming from non-representative training data, can result in unfair and inaccurate predictions, affecting the model's reliability and equity across diverse patient demographics.

## 1.5 TYPES OF MACHINE LEARNING MODELS USED

In the diabetes detection system, machine learning models are used to analyze patient data and predict the likelihood of diabetes. Several machine learning algorithms can be utilized depending on the complexity and type of the data.

### LOGISTIC REGRESSION

Logistic Regression is a widely used algorithm for binary classification tasks, like predicting whether a patient has diabetes or not. It works by applying a logistic function to the input features, producing a probability score between 0 and 1.

Advantages:

(i) Simple to implement and interpret.

(ii) Works well for linearly separable data.

Disadvantages:

(i) Limited to linear relationships.

(ii) May not perform well with highly imbalanced datasets.

## SUPPORT VECTOR MACHINE (SVM)

Support Vector Machine (SVM) is a powerful algorithm that can classify data by finding the hyperplane that best separates the classes. It is effective in both linear and non-linear classifications, using kernel functions to map data into higher dimensions.

Advantages:

(i) Effective for high-dimensional data.

(ii) Robust to overfitting, especially in high-dimensional spaces.

Disadvantages:

(i) Computationally expensive for large datasets.

(ii) Sensitive to the choice of kernel and hyperparameters.

## DECISION TREE

Decision Tree is a non-linear model that splits the data into subsets based on feature values. It recursively partitions the data and forms a tree structure where each node represents a decision based on a feature.

Advantages:

(i) Easy to understand and interpret.

(ii) Handles both numerical and categorical data.

Disadvantages:

(i) Prone to overfitting, especially with deep trees.

(ii) Sensitive to noise in the data.

## RANDOM FOREST

Random Forest is an ensemble learning method that constructs multiple

decision trees and combines their outputs. It reduces overfitting by averaging the results from many trees, making it more accurate and robust.

Advantages:

(i) Reduces overfitting compared to individual decision trees.

(ii) Handles large datasets with high accuracy.

Disadvantages:

(i) Less interpretable than a single decision tree.

(ii) Training time can increase with the number of trees.

## K-NEAREST NEIGHBORS (KNN)

K-Nearest Neighbors (KNN) is a simple, instance-based learning algorithm that classifies data points based on the majority class of their k-nearest neighbors in the feature space. It doesn't require a training phase, making it traight forward.

Advantages:

(i) Simple to implement and understand.

(ii) No training required, only prediction.

Disadvantages:

(i) Computationally expensive at prediction time.

(ii) Sensitive to the choice of k and irrelevant features.

## GRADIENT BOOSTING MACHINES (GBM)

Gradient Boosting is an ensemble learning technique where models are built sequentially. Each new model corrects the errors of the previous one, and predictions are made based on the weighted sum of all models.

Advantages:

(i) High predictive accuracy, especially with imbalanced data.

(ii) Can handle a mix of numerical and categorical features.

Disadvantages:

 (i) Prone to overfitting if not tuned correctly.

 (ii) Can be computationally expensive with large datasets.

## NEURAL NETWORKS

Neural networks are modeled after the human brain and consist of layers of interconnected nodes (neurons). Each neuron processes information and passes it to the next layer, allowing the network to learn complex patterns in the data.

 Advantages:

  (i) Capable of learning complex, non-linear relationships.

  (ii) Performs well with large datasets.

 Disadvantages:

  (i) Requires significant computational power.

  (ii) Difficult to interpret, leading to the "black-box" problem.

## NAIVE BAYES

Naive Bayes is a probabilistic model that applies Bayes' theorem, assuming that features are independent given the class. It is particularly useful when there is little data or when the relationships between features are simple.

 Advantages:

  (i) Fast and simple to implement.

  (ii) Performs well with small datasets and when features are independent.

 Disadvantages:

  (i) Assumes feature independence, which is often unrealistic.

  (ii) May not perform well if features are highly correlated.

## 1.6 DIET PLAN GENERATION

The system generates a personalized diet plan, which is a vital component of effective diabetes management. A well-structured diet helps regulate blood glucose levels, improves overall health, and reduces the risk of diabetes-related complications. The diet plan is tailored to each individual's health metrics, such as age, BMI, activity level, and blood glucose levels, ensuring that the nutritional recommendations align with the user's specific needs. It emphasizes maintaining a balance of macronutrients and micronutrients, including carbohydrates, proteins, and fats, to provide adequate energy while controlling carbohydrate intake to prevent blood sugar spikes.

Meal timing and portion control are also key aspects, as eating smaller, regular meals helps stabilize blood sugar levels. Recommendations for meal timing are designed to optimize glucose control throughout the day. The diet plan includes specific food recommendations, such as high-fiber foods, lean proteins, and whole grains, while advising against certain foods, like refined sugars and high-glycemic-index items, to help users make informed dietary choices that support blood sugar regulation.

As the user's health data evolves over time, the diet plan can be updated to reflect new goals or health needs, adopting an adaptive approach that ensures ongoing support and long-term health improvements. Additionally, the diet plan is integrated with the diabetes detection system, providing users with a comprehensive health management solution. Based on the model's predictions and the user's lifestyle information, the system suggests diet plans that complement diabetes treatment and prevention strategies, enhancing overall care.

# CHAPTER 2
# SCHEMES FOR DIABETES DETECTION AND DIET PLAN GENERATION USING MACHINE LEARNING

Sakkayaphop Pravesjit et al. (2020) focused on developing a prediction model for diabetes complications, including eye disease, kidney disease, coronary heart disease, and hyperlipidemia. The model was based on modifying the Iterative Dichotomiser 3 (ID3) decision tree into a binary number vector format. The model's efficiency was evaluated using the 10-fold cross-validation method, which resulted in an accuracy of 92.35%. Furthermore, a mobile application was developed to test the model with general patients, yielding a perfect accuracy rate of 100%. This research demonstrates the potential of machine learning in accurately predicting diabetes complications, offering valuable tools for early intervention and management.

JiMin Liu et al. (2021) highlighted the importance of early detection and treatment in controlling the progression of diabetes, a chronic disease that poses a significant threat to human health. With advancements in artificial intelligence, machine learning models have become increasingly prevalent in disease prediction. However, a single machine learning model often has limitations in diagnostic applications due to its weak generalization ability. To address this, Liu proposed an early diabetes prediction model based on stacking ensemble learning. This model combined Gradient Boosting Decision Tree, Adaboost, and Random Forest as primary learners, with Logistic Regression as a

secondary learner. The results demonstrated that the proposed ensemble model outperformed individual models in predicting early onset diabetes. Despite its high accuracy, the study noted that the timeliness of the model still requires improvement for practical application in real-world scenarios.

Anuj Mangal et al. (2022) proposed the use of Machine Learning (ML) models to predict diabetes, a major global health concern due to rising cases linked to lifestyle and dietary habits. The study emphasizes the challenges of early detection of diabetes using traditional medical methods, which often fail to provide accurate predictions. To address this, Mangal applied two widely-used ML algorithms to train a model based on medical data and symptoms, aiming to enhance early detection of diabetes. The results demonstrated that the Random Forest algorithm achieved an accuracy of 99%, showcasing its potential in predicting diabetes and offering a robust solution for early diagnosis.

Costas Papaloukas et al. (2022) focused on short-term personalized glucose prediction for patients with Type 1 diabetes mellitus (T1DM) as part of diabetes self-management. Accurate glucose prediction is essential for regulating glycemic levels and taking appropriate actions to manage the disease. The study compared two predictive models: an Autoregressive Moving Average (ARMA) model and a Long Short-Term Memory (LSTM) model, evaluating their performance across different prediction horizons. The models were trained and tested on data from 29 real patients, using evaluation metrics such as Root Mean Square Error (RMSE) and Mean Absolute Error (MAE). The results showed that the LSTM model outperformed the ARMA model with RMSE values of 3.13, 6.41, and 8.81 mg/dL, and MAE values of 1.98,

5.06, and 6.47 mg/dL for 5-, 15-, and 30-minute prediction horizons, respectively. This work demonstrates the potential of LSTM models in providing more accurate and reliable glucose predictions for T1DM patients, aiding in better disease management.

Sadia Afrin Shampa et al. (2023) examined the challenges in predicting diabetes due to high blood sugar levels and the impact of missing values and outliers in diabetes datasets. The study utilized diabetes data from Bangladesh, India, and Germany and applied various Machine Learning (ML) models to predict the onset of the disease. The experimental results revealed that boosting ML algorithms like AdaBoost, CatBoost, Gradient Boost, and XGBoost performed particularly well with the Bangladesh dataset, effectively predicting diabetes occurrence. Basic models such as Random Forests and Decision Trees also showed satisfactory performance, as evaluated by performance metrics. The study emphasizes the critical role of early diabetes detection in mitigating associated risks and severity. The findings highlight the effectiveness of boosting algorithms in diabetes prediction and the potential of basic models, contributing to a better understanding of how ML can be applied for early diagnosis and management of diabetes. This research serves as a valuable resource for healthcare professionals and policymakers working toward improved diabetes prediction and public health outcomes.

Srishti Mahajan et al. (2023) explored diabetes as a chronic condition that results from the body's inability to effectively use insulin or produce sufficient insulin to regulate blood glucose levels. The study focused on type 1 and type 2 diabetes, highlighting their prevalence and the increasing global incidence, with the International Diabetes

Federation reporting around 382 million people affected in 2022, and an expected rise to 592 million by 2035. Given diabetes' serious complications, including organ damage and dysfunction, early detection is critical. The research implemented two Machine Learning algorithms, Logistic Regression and Random Forest, to predict diabetes. The performance of both models was evaluated, with the Random Forest algorithm achieving the highest accuracy of 99.03%, showcasing its potential in early diabetes prediction. The study underscores the importance of machine learning in providing reliable tools for diagnosing diabetes and preventing its severe health outcomes.

Jiali Gao et al. (2023) addressed the challenge of traditional, time-consuming methods like blood and urine tests for early diabetes diagnosis. The study applied machine learning algorithms, including KNN, Decision Tree, and Random Forest, to model and predict the Pima Indian diabetes dataset. Performance evaluation using cross-validation and confusion matrix revealed that the Random Forest model outperformed the others, achieving an accuracy of 0.84 and an F1 value of 0.77. The results highlight the significant improvement in prediction accuracy with machine learning, offering valuable insights for early diabetes diagnosis and prevention.

P. Meenakshidevi et al. (2024) aimed to identify and predict diabetes using a combination of four machine learning algorithms: Random Forest (RF), Naive Bayes (NB), Multi-Layer Perceptron (MLP), and Support Vector Machine (SVM). The project integrated SVM with the other algorithms to enhance prediction accuracy. This multi-algorithm framework utilized SVM's robust classification, RF's ensemble learning, MLP's pattern capture, and NB's probabilistic approach to

improve efficiency and address the limitations of traditional SVM-based models. The results showed that the combined models yielded notable accuracy, with RF improving accuracy by 76%, NB by 78%, and MLP by 82%. The study concluded that combining SVM with RF, MLP, and NB offers a comprehensive solution for diabetes prediction, with MLP proving to be the most effective for early diagnosis and treatment, surpassing conventional SVM models in accuracy.

Ifra Shaheen et al. (2024) focused on predicting diabetes using Deep Learning (DL) techniques, particularly ensemble learning approaches, to address the challenges posed by the highly imbalanced Diabetes Prediction Dataset (DPD). The dataset contained 8,500 diabetic patients compared to 91,500 non-diabetic individuals. To mitigate the class imbalance, the Proximity-Weighted Synthetic Oversampling (ProWSyn) technique was applied. Shaheen proposed two ensemble models: a hybrid model called Hi-Le, combining the Highway and LeNet models, and a blending model called HiTCLe, which integrated Highway, LeNet, and Temporal Convolutional Network. The Hi-Le model achieved an accuracy of 94%, with an F1-Score of 96%, precision of 94%, and recall of 95%, surpassing the performance of individual models. The HiTCLe model also outperformed its components, achieving an accuracy of 94% and an F1-Score of 94%, while the individual models achieved accuracies between 89% and 91%. K-Fold Cross Validation and Shapley Additive eXplanations were used to validate the models and analyze feature contributions. The study demonstrated that both ensemble models significantly outperformed their individual counterparts, providing a highly accurate method for early diabetes detection.

# CHAPTER 3

# EARLY DIABETES DETECTION USING DEEP LEARNING

## 3.1 INTRODUCTION

Diabetes, a chronic metabolic disorder, is becoming increasingly prevalent due to poor lifestyles and environmental factors. Early detection and precise prediction are critical to mitigating its effects and preventing severe complications such as kidney failure, heart diseases, and blindness. The existing system focuses on leveraging Deep Learning (DL) and ensemble learning methodologies to predict diabetes accurately and reliably. Two ensemble models, Hi-Le (hybrid) and HiTCLe (blending), are proposed for early diabetes detection, using Explainable Artificial Intelligence (XAI) techniques to enhance interpretability.

## 3.2 EARLY DIABETES DETECTION USING DEEP LEARNING

Early diabetes detection using deep learning utilizes the Diabetes Prediction Dataset (DPD), a highly imbalanced dataset, containing 91,500 non-diabetic and 8,500 diabetic cases. To address this imbalance, the Proximity-Weighted Synthetic Oversampling (ProWSyn) technique was employed to generate synthetic samples for minority classes. The proposed ensemble models, Hi-Le and HiTCLe, demonstrate significant improvements in prediction accuracy.

**Hi-Le Model:** A hybrid model combining Highway and LeNet architectures for capturing spatial and hierarchical relationships in data. It

achieves an accuracy of 94%, F1-Score of 96%, and recall of 95%.

**HiTCLe Model:** A blending model combining Highway, Temporal Convolutional Network (TCN), and LeNet, which integrates temporal and spatial features. It achieves an accuracy of 94% and F1-Score of 94%.

## 3.3 DATA PREPROCESSING AND BALANCING

The Diabetes Prediction Dataset (DPD) forms the foundation of this research. This dataset exhibits a severe class imbalance, with 91,500 non-diabetic cases and only 8,500 diabetic cases. This imbalance could lead to biased predictions if not handled effectively.

To mitigate this, the Proximity-Weighted Synthetic Oversampling (ProWSyn) technique was employed. ProWSyn is a data augmentation approach that generates synthetic samples by focusing on data points near the decision boundary. This technique ensures that the minority class receives sufficient representation, enhancing the model's ability to learn critical patterns related to diabetes.

## 3.4 MODEL

**(i) Hi-Le:** Integrates Highway and LeNet models. Highway networks manage information flow and adaptively process data, while LeNet captures spatial patterns.

**(ii) HiTCLe:** Blends TCN, LeNet, and Highway. TCN captures temporal dependencies, LeNet identifies spatial hierarchies, and Highway fine-tunes predictions.

## 3.5 ALGORITHMS

**(i) Proximity-Weighted Synthetic Oversampling (ProWSyn)**

ProWSyn is a technique that addresses the class imbalance issue by generating synthetic samples for the minority diabetic class. It focuses on data points near decision boundaries, ensuring realistic and meaningful synthetic samples that help the model learn critical patterns effectively. By balancing the dataset, ProWSyn enhances the representation of the minority class, reducing bias and improving prediction accuracy.

**(ii) Highway Network**

The Highway Network is a deep learning architecture that facilitates efficient information flow across layers using gating mechanisms. These mechanisms determine which information to retain or transform, ensuring the preservation of essential features while mitigating issues like vanishing gradients. Its adaptability makes it crucial for processing complex data in the Hi-Le and HiTCLe models.

**(iii) LeNet**

LeNet is a convolutional neural network designed to extract spatial hierarchies and patterns from input data. Using convolutional and pooling layers, it captures low-level and high-level spatial features, making it integral to the proposed models for understanding spatial relationships within the dataset.

**(iv) Temporal Convolutional Network (TCN)**

TCN is a deep learning architecture that captures temporal dependencies in sequential data. By employing dilated convolutions, it effectively models long-range dependencies, making it suitable for analyzing temporal patterns. It is a key component of the HiTCLe model, which

leverages both spatial and temporal features.

**(v) Shapley Additive Explanations (SHAP)**

SHAP enhances model interpretability by quantifying the contribution of individual features to predictions. It uses visualization techniques like beeswarm and dependence plots to provide insights into feature importance and interactions. This makes the model's predictions transparent, aiding stakeholders in understanding the significance of features like glucose and HbA1c levels.

## 3.6 MODEL ARCHITECTURE

To enhance the interpretability of predictions, Shapley Additive Explanations (SHAP) were employed. SHAP provided insights into the contribution of individual features, such as HbA1c levels and glucose levels, to the final predictions. Visualizations such as summary plots, beeswarm plots, and dependence plots offered a deeper understanding of model behavior.
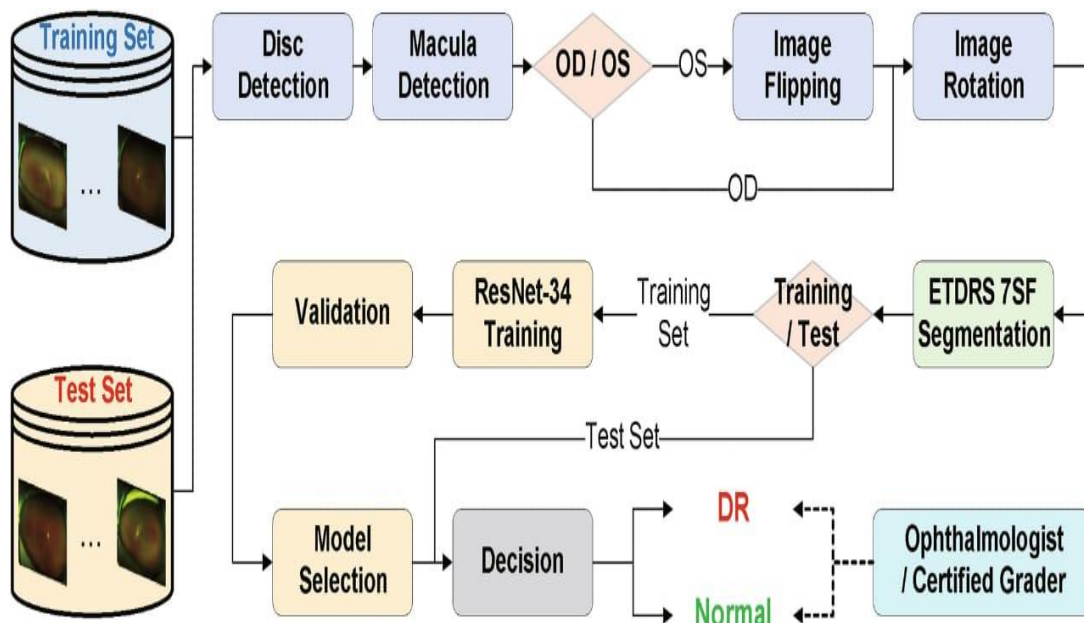


Figure 3.1 Early diabetes detection using ML architecture

Figure 3.1 illustrates the architecture for early diabetes detection using Machine Learning, showcasing data preprocessing, model training, prediction, and evaluation processes.

## 3.7 LIMITATIONS

### (i) Dataset Dependence

The accuracy and generalizability of models like Hi-Le and HiTCLe are significantly influenced by the quality and diversity of the Diabetes Prediction Dataset. Limited or biased data may hinder the model's performance in real-world applications.

### (ii) Computational Complexity

The integration of multiple architectures, such as Highway, LeNet, and TCN, increases the computational requirements, making the training process resource-intensive and potentially time-consuming.

### (iii) Synthetic Data Challenges

While ProWSyn effectively balances the dataset, it may inadvertently introduce noise or overfit the model, especially if the minority class lacks sufficient diversity or represents rare patterns.

### (iv) Limited Temporal Context

Despite the inclusion of TCN for capturing temporal dependencies, the absence of longitudinal data in the dataset restricts the model's ability to predict long-term patterns or trends in diabetes progression.

### (v) Explainability vs. Complexity Trade-Off

Although SHAP enhances the interpretability of model predictions, its complex visualizations may be challenging for non-technical users to fully grasp, limiting its utility for broader audiences.

## (vi) Real-World Deployment

Transitioning the model from a research setting to real-world healthcare applications involves significant challenges, including regulatory compliance, integration with existing systems, and ensuring the privacy and security of patient data.

# CHAPTER 4

# DIABETES DETECTION AND DIET PLAN GENERATION USING MACHINE LEARNING

## 4.1 INTRODUCTION

The diabetes detection and diet plan generation using machine learning system aims to provide an efficient and reliable solution for diabetes detection and management. Early detection is critical for preventing complications associated with diabetes, and the system leverages machine learning algorithms to predict the likelihood of diabetes based on user-provided health metrics. Another key objective is to generate personalized diet plans that align with individual health conditions and dietary preferences, ensuring sustainable diabetes management. By integrating these functionalities into a single platform, the system empowers users to make informed health decisions, reduces dependency on manual healthcare consultations, and promotes proactive health management.

## 4.2 SYSTEM ARCHITECTURE

The architecture of the proposed system is designed for modularity and ease of use. At its core, the system consists of several interconnected components that handle data collection, preprocessing, prediction, and dietary recommendation. Users provide key health metrics such as glucose levels, BMI, blood pressure, insulin levels, and age through a user-friendly web interface. This data is then passed through the preprocessing module, where inconsistencies like missing values or outliers

are addressed. The processed data is fed into machine learning models that classify users into diabetic or non-diabetic categories.
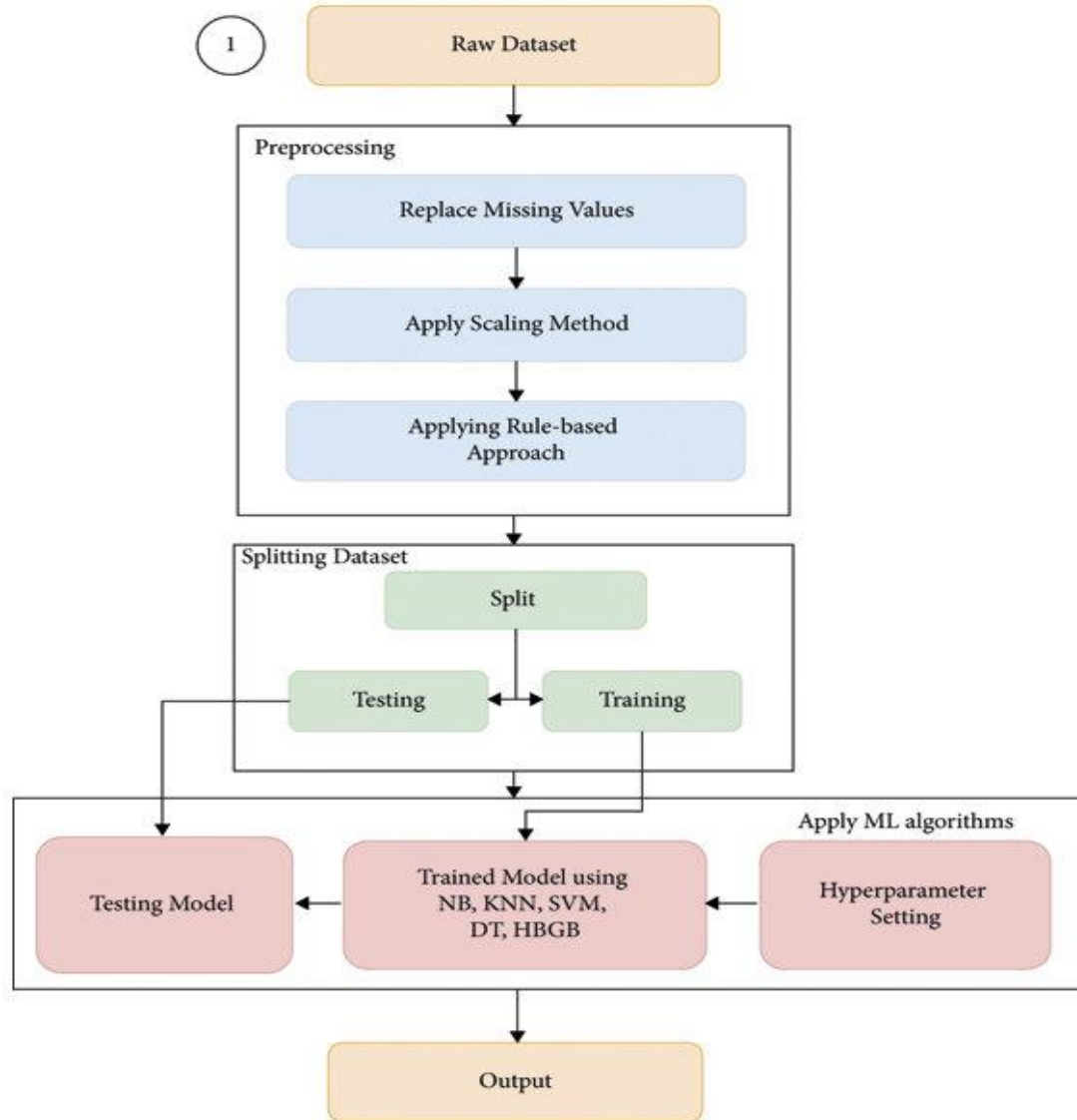


Figure 4.1 System Architecture

Fig 4.1 includes a diet plan generation engine that uses the predictions and other health metrics to create tailored dietary recommendations. The output module displays the results, including predictions and diet plans, along with insights from explainable AI techniques. This ensures transparency in the predictions and builds user trust in the system.

## 4.3 METHODOLOGY

The methodology for this diabetes detection and diet plan generation system follows a structured approach to ensure accuracy, reliability, and user-centric functionality.

### (i) Data Collection and Preprocessing:

Data is initially sourced from publicly available datasets, such as the PIMA Indian Diabetes Dataset, which provides a robust foundation for training the machine learning models. To improve generalizability, data from other relevant health studies may also be integrated. Preprocessing steps include handling missing values. Outlier detection is applied to remove abnormal values that may skew results. Scaling and normalization are used to ensure uniformity across features such as glucose levels, blood pressure and BMI.

### (ii) Feature Engineering:

Key health metrics like glucose levels, BMI, insulin, and age are evaluated, and potentially new features are generated based on domain knowledge, such as the ratio of BMI to glucose levels or age-adjusted risk factors. Feature selection techniques, such as recursive feature elimination and correlation analysis, help refine the dataset to only include highly predictive features, enhancing model performance.

### (iii) Model Training and Evaluation:

The Various machine learning models, including Logistic Regression, Random Forest, and Gradient Boosting, are trained using the preprocessed dataset. The Hyperparameter tuning is conducted using Grid Search or Random Search to optimize each model for performance metrics such as

accuracy, precision and recall. The Cross-validation, specifically K-Fold Cross-Validation, ensures the models generalize well to the unseen data, that reducing the risk of overfitting.

**(iv) Deployment Pipeline:**

The selected model is deployed on the web platform using frameworks like Flask, enabling real-time predictions based on user input. The deployment pipeline also supports automated data updates, allowing for continuous retraining and improvement of the models based on user feedback. A feedback loop enables users to provide insights on the accuracy of both predictions and dietary recommendations, which can further inform retraining cycles and improve system responsiveness over time.

**(v) Diet Plan Personalization:**

Post-prediction, user-specific dietary recommendations are generated. These recommendations are dynamically created based on health metrics, diabetes risk, and user preferences, ensuring the diet plan is not only healthy but also adheres to cultural and personal tastes. A recommendation algorithm, based on user feedback and dietary success rates, refines the personalization over time.

**4.4 MACHINE LEARNING MODELS**

The diabetes detection system employs various machine learning algorithms to identify the model best suited for predicting diabetes risk. Each was evaluated using a structured approach involving data preprocessing, model training, prediction, and performance assessment. The selection of models included Logistic Regression, Naive Bayes, K-Nearest Neighbors (KNN), Decision Tree Classifier, and Support Vector Machine (SVM). The evaluation focused on accuracy and computational efficiency.

**(i) LOGISTIC REGRESSION :**

Logistic Regression was initially implemented as a baseline model due to its simplicity and effectiveness in binary classification tasks. The model demonstrated strong performance, particularly in terms of interpretability, and provided a probabilistic output that aided in understanding predictions. To enhance its performance, hyperparameter optimization was conducted using GridSearchCV, focusing on parameters like the solver method (`'liblinear'`) and regularization strength (`C=0.1`). The optimized Logistic Regression model exhibited high accuracy and served as the primary candidate for deployment.



Figure 4.2 Logistic Regression

The Fig 4.2 shows the respective range of values for the linear regression model. These values are entered as approximate values to represent the relationship between the input features and the predicted output.

**(ii) NAIVE BAYES :**

Naive Bayes, specifically the Gaussian Naive Bayes variant, was chosen for its efficiency and suitability for continuous data such as glucose

levels and BMI. This probabilistic model assumes feature independence, which, although not entirely true for this dataset, allowed for relatively quick and accurate predictions. While the accuracy was competitive, its assumptions limited its ability to capture complex feature interactions.
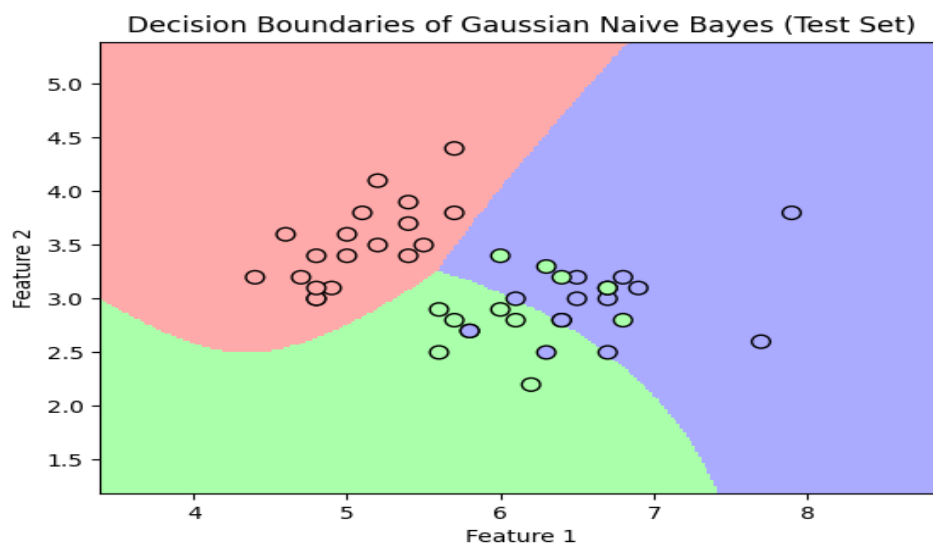


Figure 4.3 Naive bayes

Figure 4.3 illustrates the Naive Bayes algorithm, a probabilistic classifier that applies Bayes' theorem with strong independence assumptions between features. It is used for classification tasks by calculating the probability of each class given the input data.

**(iii) K-NEAREST NEIGHBORS (KNN) :**

The K-Nearest Neighbors (KNN) algorithm was explored to leverage its simplicity and effectiveness in capturing local patterns in the dataset. Using a neighborhood size of `k=3`, the KNN model achieved reasonable accuracy, especially for well-separated classes. However, the algorithm's computational complexity increased with larger datasets, presenting scalability challenges.
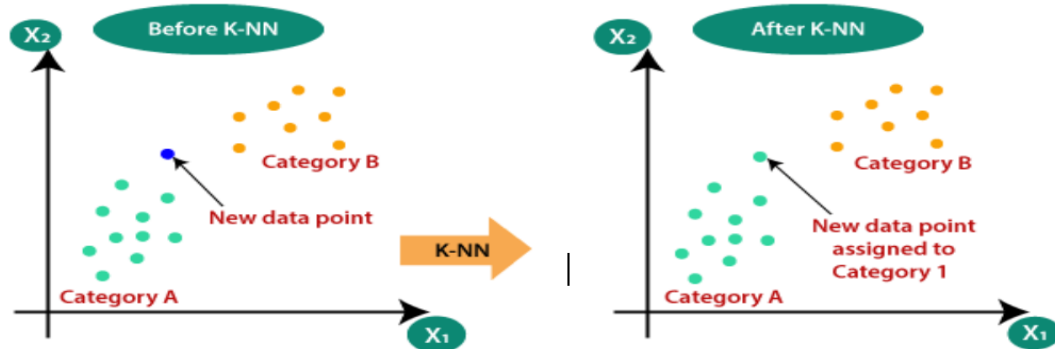
Figure 4.4  K-Nearest Neighbors (KNN)

Figure 4.4 illustrates the K-Nearest Neighbors (KNN) algorithm, which classifies data based on the majority class of its nearest neighbors.

**(iv) DECISION TREE CLASSIFIER:**

Decision Tree Classifier was implemented to explore its capability of modeling non-linear relationships in the data. This algorithm provided a clear and interpretable decision-making process by creating a tree structure based on feature thresholds. While it excelled in capturing intricate patterns, the model was prone to overfitting without additional regularization techniques, such as pruning.
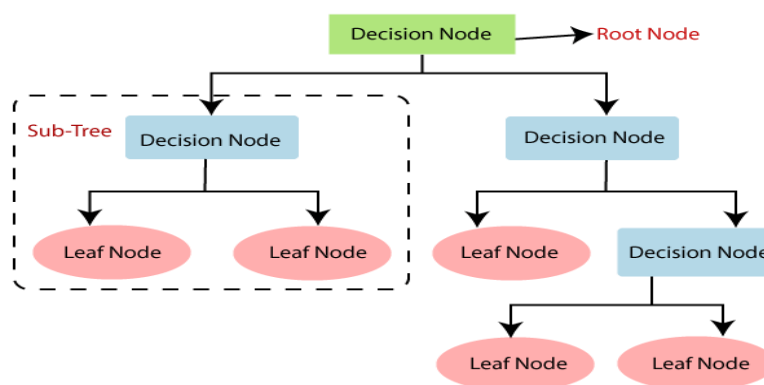


Figure 4.5 Decision Tree Classifier

Figure 4.5 illustrates the Decision Tree Classifier, which splits data into branches based on feature values, making decisions at each node.

**(v) SUPPORT VECTOR MACHINE (SVM):**

Support Vector Machine (SVM) was employed with a polynomial kernel to explore its potential in handling non-linear relationships. SVM offered strong performance by finding an optimal hyperplane to separate classes in a high-dimensional space. However, the computational cost and complexity of hyperparameter tuning posed challenges, particularly for large datasets.
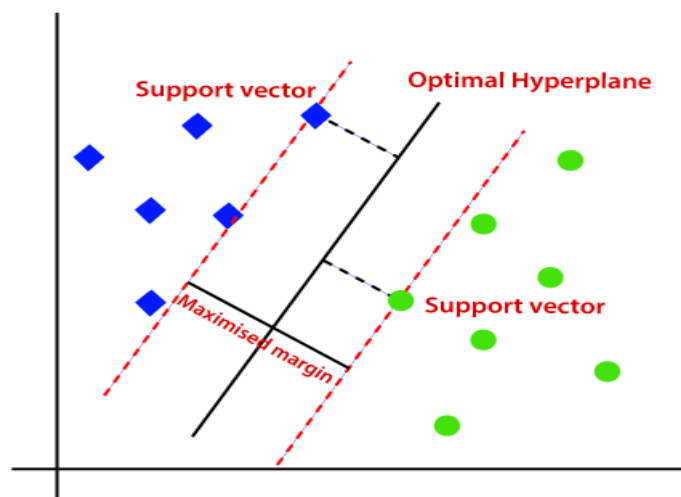


Figure 4.6 Support vector machine

Figure 4.6 illustrates the Support Vector Machine (SVM) algorithm, which finds the optimal hyperplane that best separates data into different classes.

**4.5 DIET PLAN GENERATION**

The diet plan generation module is a vital component of the proposed system. It uses user-specific health metrics, such as BMI, glucose levels, and lifestyle factors, to provide personalized dietary recommendations. The generated diet plans focus on balancing macronutrients like carbohydrates, proteins, and fats to stabilize blood sugar levels and improve overall health. For instance, a user with high glucose levels might receive recommendations

to consume high-fiber foods, lean proteins, and whole grains while avoiding sugary and processed foods.

The system also considers individual preferences and cultural dietary habits, ensuring that the recommendations are practical and sustainable. Feedback mechanisms are incorporated into the diet plan generation process, allowing the system to refine its recommendations based on user responses.

## 4.6 WEB-BASED INTERFACE

The web interface of the proposed system serves as the primary interaction point for users. Developed using Flask, the interface is designed to be intuitive and user-friendly. It allows users to input their health metrics, receive diabetes predictions, and access personalized dietary recommendations. The interface also incorporates visualizations, such as bar graphs and pie charts, to present the predictions and diet plans in an easily understandable format.

A key feature of the interface is the integration of explainable AI insights. Users can view the importance of various health metrics, such as glucose levels and BMI, in influencing their diabetes risk. This transparency builds trust in the system and helps users understand the rationale behind the recommendations. The web interface is also designed to be responsive, ensuring compatibility across devices, including desktops, tablets, and smartphones.

## 4.7 ADVANTAGES

The integration of multiple machine learning models ensures high predictive accuracy, allowing for early and reliable detection of diabetes risk. By automating diabetes risk prediction, the system reduces the need for regular manual screenings, which can be costly and time-consuming. By

incorporating user-specific health metrics and dietary preferences, the system provides personalized diet plans that cater to individual health needs. This level of customization promotes better adherence to dietary recommendations, ultimately supporting sustainable diabetes management and improved health outcomes. The use of SHAP values and other explainable AI techniques helps users understand the importance of their health metrics in the prediction process. This transparency fosters trust in the system, empowering users to take informed actions based on the results.

Designed with modularity in mind, the system can be easily expanded to include additional functionalities, such as real-time monitoring through IoT devices or mobile integration. This scalability ensures that the system remains adaptable to technological advancements and user needs. The system reduces dependency on healthcare providers for regular diabetes management, offering a cost-effective alternative for individuals. By automating the diagnostic and dietary recommendation processes, the system alleviates the burden on healthcare resources, making it an affordable and efficient solution for diabetes management. By incorporating user feedback, the system continuously refines its predictions and dietary recommendations, ensuring that it evolves to meet individual health needs effectively. This dynamic approach enhances user satisfaction and improves the accuracy of predictions over time.

# CHAPTER 5

# SYSTEM SPECIFICATION

## 5.1 HARDWARE REQUIREMENTS

▶ Processor Type      : Intel(R) Core (TM) i3

▶ Speed        : 2.30GHZ

▶ RAM        :8 GB RAM

▶ Hard disk       : 1 TB

## 5.2 SOFTWARE REQUIREMENTS

▶ Operating System    : Windows 10

▶ Tool        : Visual Studio

▶ Coding Language    : Python

## 5.3 TECHNOLOGIES USED

The system integrates key technologies such as Flask, Pandas, scikit-learn, and Matplotlib to deliver efficient and accurate diabetes detection and personalized diet plan generation. Flask enables a seamless web interface, Pandas handles data preprocessing and manipulation, scikit-learn powers machine learning algorithms, and Matplotlib provides detailed visualizations for insights. Together, these technologies ensure the system is robust, user-friendly, and capable of real-time interaction.

### 5.3.1 FLASK

Flask is a lightweight and flexible web development framework that facilitates the creation of the interactive user interface for the proposed system. It allows users to input their health data, such as glucose levels, BMI, and insulin levels, and receive predictions and dietary recommendations in real-time. Flask serves as the backend framework that connects the user interface with machine learning models, ensuring smooth data flow between the frontend and backend.

Flask's scalability enables the system to accommodate future enhancements, such as integrating real-time health monitoring through IoT devices. The framework supports the use of RESTful APIs, allowing data to be sent to and retrieved from the machine learning models efficiently. Flask's modular structure makes it easy to extend the application with additional features, such as user authentication or historical health data tracking.

By providing an intuitive platform for user interaction, Flask ensures a seamless experience for users, allowing them to engage with the system effortlessly. Its compatibility with other Python libraries, such as Pandas and scikit-learn, ensures efficient communication between components, further enhancing the system's performance and reliability.

### 5.3.2 PANDAS

Pandas is a powerful Python library for data manipulation and analysis, playing a crucial role in preprocessing the health data provided by users. In the context of diabetes detection, Pandas is used to clean and prepare the input data, addressing issues such as missing values, outliers, and inconsistent

formats. This ensures that the data is ready for analysis by machine learning algorithms, improving the accuracy and reliability of predictions.

Pandas' Data Frame structure provides an efficient way to store and manipulate large datasets, such as those used for training the machine learning models. Tasks like filtering data, aggregating values, and applying transformations are performed seamlessly, enabling the system to handle complex datasets efficiently. Additionally, Pandas supports the integration of new data sources, making it possible to update the system with the latest health metrics and research findings.

By leveraging Pandas, the system ensures that health data is consistently formatted and processed, laying the foundation for accurate machine learning predictions and personalized diet recommendations.

### 5.3.3 SCIKIT-LEARN

scikit-learn is a versatile library for machine learning, providing tools for implementing and training predictive models. In this project, scikit-learn is used to develop models for diabetes detection, such as Logistic Regression, Random Forest, and Gradient Boosting. These models analyze user health data and classify individuals as diabetic or non-diabetic with high accuracy.

The library offers robust tools for feature selection, enabling the system to identify key predictors like glucose levels, BMI, and insulin. Hyperparameter tuning and cross-validation techniques in scikit-learn ensure that the models are optimized for performance, minimizing overfitting and enhancing generalization. The library also provides metrics for evaluating

model performance, such as accuracy, precision, and recall, which are critical for ensuring reliable predictions.

scikit-learn's seamless integration with data preprocessing libraries like Pandas and visualization tools like Matplotlib allows for a streamlined workflow, from data input to prediction. Its efficiency and scalability make it a core technology for building the machine learning models in the proposed system.

# CHAPTER 6

# IMPLEMENTATION AND RESULT

## 6.1 MODULE DESCRIPTION

## 6.1.1 FLOW DIAGRAM

The model training process involves feeding the extracted features into an LSTM network to learn the temporal relationships between hand landmarks during gestures. The LSTM model is trained on a labeled dataset, where each gesture is associated with a specific sign. During training, the model learns to predict the corresponding gesture from the hand landmark data. The training is done in multiple epochs to optimize the model's performance, using techniques like cross-validation to improve generalization.



Figure 6.1: System Flow Diagram

Figure 1.1 illustrates the proposed technique using Machine Learning for diabetes prediction. This model leverages various machine learning algorithms to analyze input features and predict the likelihood of diabetes, offering a reliable and efficient approach for early diagnosis.

## 6.1.2 COLLECTING DATASET

The dataset for training and testing machine learning models is sourced from publicly available datasets like the PIMA Indian Diabetes Dataset and additional health records. These datasets include vital features such as glucose levels, BMI, blood pressure, insulin levels, and family history, which are essential for diabetes prediction.

The dataset undergoes thorough cleaning and annotation to ensure accuracy and reliability. Missing values are handled using statistical imputation techniques, and outliers are treated to prevent skewed model predictions. The cleaned and annotated dataset is then split into training and testing sets to validate the models effectively.

| Pregnancie | Glucose | BloodPres | SkinThickn | Insulin | BMI | DiabetesP | 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 |
| 10 | 168 | 74 | 0 | 0 | 38 | 0.537 | 34 | 1 |
| 10 | 139 | 80 | 0 | 0 | 27.1 | 1.441 | 57 | 0 |
| 1 | 189 | 60 | 23 | 846 | 30.1 | 0.398 | 59 | 1 |
| 5 | 166 | 72 | 19 | 175 | 25.8 | 0.587 | 51 | 1 |
| 7 | 100 | 0 | 0 | 0 | 30 | 0.484 | 32 | 1 |
| 0 | 118 | 84 | 47 | 230 | 45.8 | 0.551 | 31 | 1 |
| 7 | 107 | 74 | 0 | 0 | 29.6 | 0.254 | 31 | 1 |
| 1 | 103 | 30 | 38 | 83 | 43.3 | 0.183 | 33 | 0 |
| 1 | 115 | 70 | 30 | 96 | 34.6 | 0.529 | 32 | 1 |
| 3 | 126 | 88 | 41 | 235 | 39.3 | 0.704 | 27 | 0 |
| 8 | 99 | 84 | 0 | 0 | 35.4 | 0.388 | 50 | 0 |
| 7 | 196 | 90 | 0 | 0 | 39.8 | 0.451 | 41 | 1 |
| 9 | 119 | 80 | 35 | 0 | 29 | 0.263 | 29 | 1 |
| 11 | 143 | 94 | 33 | 146 | 36.6 | 0.254 | 51 | 1 |
| 10 | 125 | 70 | 26 | 115 | 31.1 | 0.205 | 41 | 1 |
| 7 | 147 | 76 | 0 | 0 | 39.4 | 0.257 | 43 | 1 |

Figure 6.2 Collecting dataset

The Figure 6.2 illustrates that the collecting of dataset helps Machine Learning models to find patterns in people and make predictions that they are at risk of having Diabetes.

## 6.1.3 FEATURE EXTRACTION

Feature extraction involves identifying the most relevant health metrics contributing to diabetes prediction. Features like glucose levels, BMI, and age are selected based on their correlation with diabetes risk. Statistical methods and domain knowledge are used to ensure that only meaningful features are included in the training process.

The extracted features are scaled and normalized to ensure consistency across the dataset. This step is crucial for algorithms like Logistic Regression and Random Forest, as it improves the models' ability to interpret the data. Feature importance analysis, such as SHAP (SHapley Additive exPlanations), is also conducted to highlight the contribution of each feature to the prediction process.



Figure 6.3 Distribution of Variables

Figure 6.3 illustrates the Distribution of variables among the dataset for diabetes prediction using machine learning models.



Figure 6.4 Univarient Analysis

Figure 6.4 illustrates that Univarient Analysis looks at one variable at a time to understand its patterns.
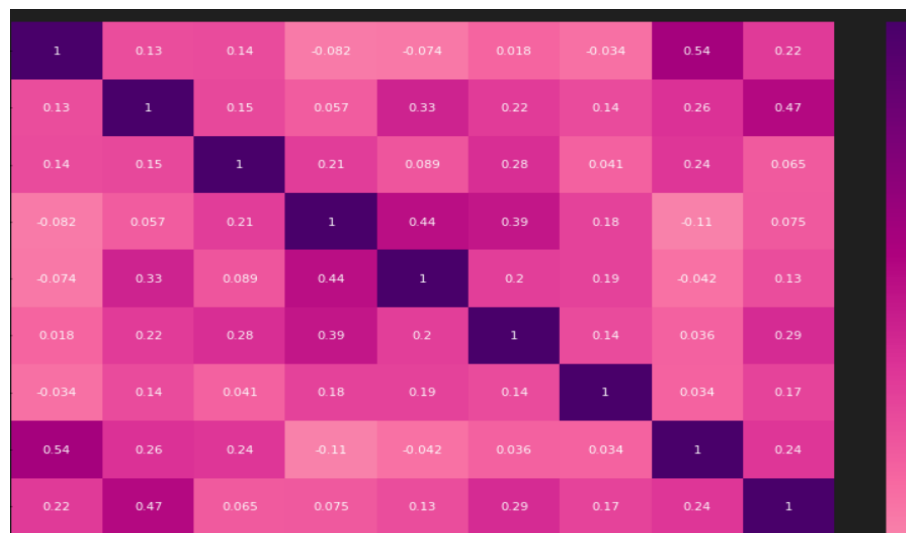


Figure 6.5 Bivarient Analysis

Figure 6.5 illustrates that Bivarient Analysis examines the relationship between two variables to understand its patterns.

### 6.1.4 TRAIN MODEL

The model training process involves feeding the preprocessed and feature-engineered data into various machine learning algorithms. Models such as Logistic Regression, Random Forest, and Gradient Boosting a trained to classify users as diabetic or non-diabetic based on their health metrics.

Each model is trained over multiple iterations (epochs), using techniques like cross-validation to improve generalization and prevent overfitting. Hyperparameter tuning is performed to optimize the models for accuracy and reliability. The training phase also involves evaluating model performance using metrics like accuracy, precision, recall, and F1-score to select the best-performing algorithm.

### 6.2 SOURCE CODE

**App.py**

```
from flask import Flask, render_template, request

import random

app = Flask(_name_)

@app.route('/')

def home():

    return render_template('index.html')

@app.route('/form')
```

```python
def form():

    return render_template('form.html')

@app.route('/result', methods=['POST'])

def result():

    age = float(request.form['Age'])

    glucose = float(request.form['Glucose'])

    blood_pressure = float(request.form['BloodPressure'])

    insulin = float(request.form['Insulin'])

    bmi = float(request.form['BMI'])

    skin_thickness = float(request.form['SkinThickness'])

    # Determine blood sugar range

    blood_sugar_range = calculate_blood_sugar_range(glucose)

    # Simple logic to determine diabetes prediction

    prediction = 'positive' if glucose > 140 or bmi > 30 else 'negative'

    # Get diet plan and quotes based on the result

    data = get_diet_plan(prediction, glucose, bmi)

    data['blood_sugar_range'] = blood_sugar_range   # Add range to the
data dictionary

    # Generate a positive quote for users with a positive result
```

```python
    quote = get_positive_quote() if prediction == 'positive'

else

 None

    return render_template('result.html', data=data, quote=quote)

def calculate_blood_sugar_range(glucose):

    if glucose < 100:

        return "Normal (Under 100 mg/dL)"

    elif 100 <= glucose <= 125:

        return "Prediabetes Range (100-125 mg/dL)

 else:

        return "Diabetes Range (Over 125 mg/dL)"

def get_diet_plan(prediction, glucose, bmi):

    diet_data = {}

    if prediction == 'positive':

        if glucose > 200:

            severity = "Severe"

            diet_data = {

                'severity': severity,

                'food': ["Lean meats, whole grains, legumes"],
```

```
            'fruits': ["Berries, Apples (in moderation)"],

            'vegetables': ["Leafy greens, broccoli"],

            '5210_rule': ["5 servings of fruits/vegetables", "1 hour physical
activity", "2 hours or less screen time", "0 sugary drinks"],

        }

    elif glucose > 140:

        severity = "Moderate"

        diet_data = {

            'severity': severity,

            'food': ["Whole grains, beans, legumes"],

            'fruits': ["Apples, Pears, Citrus"],

            'vegetables': ["Spinach, Kale"],

            '5210_rule': ["5 servings of fruits/vegetables", "1 hour physical
activity", "2 hours or less screen time", "0 sugary drinks"],

        }

    else:

        severity = "Mild"

        diet_data = {

            'severity': severity,

            'food': ["Whole wheat, Oats, Nuts"],
```

```python
            'fruits': ["Berries, Apples, Pears"],

            'vegetables': ["Leafy greens, Tomatoes"],

            '5210_rule': ["5 servings of fruits/vegetables", "1 hour physical
activity", "2 hours or less screen time", "0 sugary drinks"],

            }
        else: diet_data = {

            'severity': 'Negative',

            'food': ["Whole grains, fruits, lean proteins, healthy fats"],

            'fruits': ["Berries, Apples, Pears"],

            'vegetables': ["Carrots, Spinach, Cabbage"],

            '5210_rule': ["5 servings of fruits/vegetables", "1 hour physical
activity", "2 hours or less screen time", "0 sugary drinks"]  }

        return diet_data

    def get_positive_quote():

        quotes = [

            "Your health is an investment, not an expense.",

            "Believe in yourself and your ability to overcome challenges.",

            "Take care of your body; it's the only place you have to live.",

            "Knowledge is the best prescription for a healthier life.",

            "Healthy living is not a goal; it's a journey.",
```

"Small changes can lead to a lifetime of wellness",

]

return random.choice(quotes)

if _name_ == "_main_":
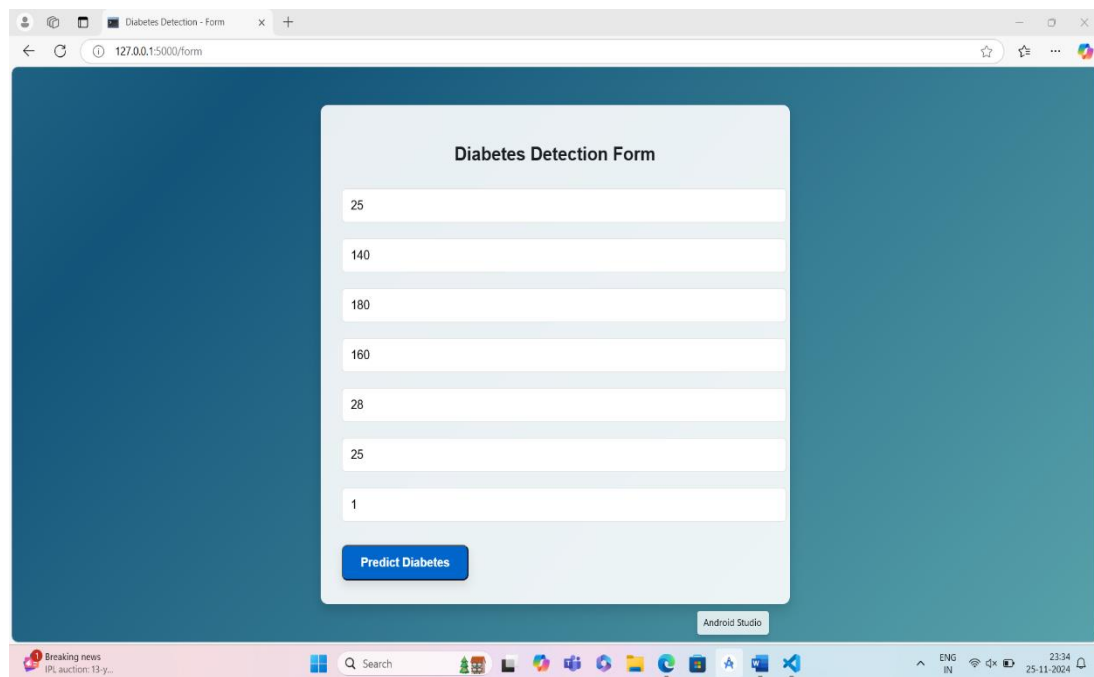
app.run(debug=True)

## 6.3 RESULT



Figure 6.6 Home Page

Figure 6.6 shows the Home page of the result which is used to start prediction by start test.

Figure 6.7 Diabetes Detection Form

Figure 6.7 shows the Diabetes Detection Form is used to predict diabetes in people by giving their details in the form for further process.



Figure 6.8 Form with values

Figure 6.8 shows the Diabetes Detection Form with input data values from user for prediction of diabetes level.
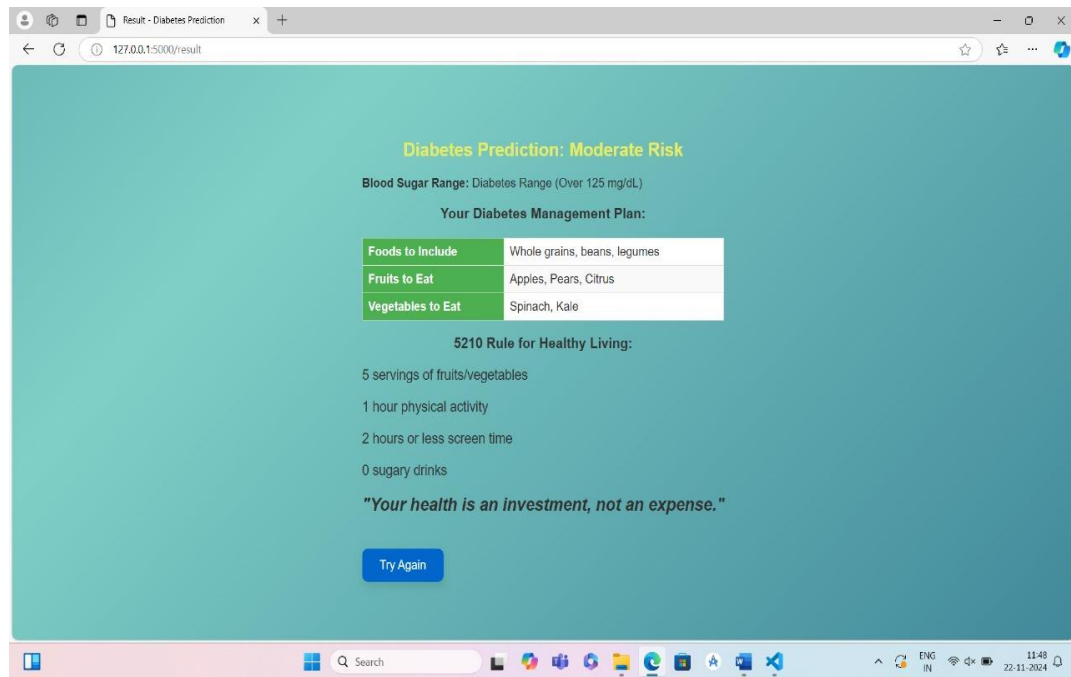


Figure 6.9 Prediction result and Diet plan

Figure 6.9 shows the Diabetes Prediction Result with blood sugar range predicted. This also generates Diet plan according to level of diabetes predicted and also gives rules for healthy living.

# CHAPTER 7

# CONCLUSION AND FUTURE ENHANCEMENT

## 7.1 CONCLUSION

The Diabetes Detection and Diet Plan Generation System effectively combines machine learning algorithms with user-centric design to meet the critical need for early diabetes diagnosis and personalized dietary management. Using models like Logistic Regression and Random Forest, the system analyzes user health metrics, including glucose levels and BMI, to deliver accurate predictions.

Diabetes Detection achieved its objectives of accuracy, scalability, and ease of use. The web-based interface, developed with Flask, ensures a seamless user experience, allowing real-time interaction and insights. Additionally, explainable AI techniques enhance transparency, helping users understand the factors influencing their results and building trust in the system.

The Project highlights the potential of machine learning in healthcare innovation. It provides a practical, scalable solution to diabetes management and establishes a foundation for future advancements in personalized healthcare technologies. Once trained, the model can be used to classify real-time gestures accurately.

**7.2 FUTURE ENHANCEMENT**

Future enhancements for the Diabetes Detection and Diet Plan Generation System include integrating wearable devices like continuous glucose monitors (CGMs) and fitness trackers to enable real-time health monitoring. A mobile application can be developed to make the system more accessible, allowing users to input health data and access predictions and diet recommendations on the go. Additionally, adapting diet plans to regional and cultural preferences will improve user adherence and satisfaction.

Expanding the system's capabilities to analyze additional health metrics, such as cholesterol levels and blood pressure, will provide more comprehensive insights. Incorporating a more diverse dataset with varied demographics and health profiles will enhance the model's robustness and reliability. Integration with healthcare providers can facilitate collaboration, allowing professionals to validate and refine predictions and recommendations.

Lastly, extending the system to predict related health risks, such as hypertension and cardiovascular diseases, will broaden its scope. Advanced visualization techniques can also be introduced to present predictions and dietary recommendations in a more user-friendly manner, improving the overall experience and effectiveness of the system.

**REFERENCES**

1. ANUJ MANGAL, VINOD JAIN (2022) – "Machine Learning-Based Early Detection of Diabetes Using Random Forest Algorithm".

2. AYATUN NESA, SADIA AFRIN SHAMPA, MD. SAIFUL ISLAM (2023) – "Machine Learning-Based Diabetes Prediction: A Cross-Country Perspective''.

3. CHITRADEVI, SUPRIYA, N SUBHASH CHANDRA, CHITRADEVI T N, HAIDER ALABDELI (2024)-"Diabetes Mellitus Prediction and Classification Using Firefly Optimization Based Support Vector Machine"

4. COSTAS PAPALOUKAS, DAPHNE N. KATSAROU, ELENI I. GEORGA, MARIA CHRISTOU, STELIOS TIGAS, DIMITRIOS I. FOTIADIS (2022) – "Short-Term Glucose Prediction in Type 1 Diabetes Mellitus Using ARMA and LSTM Models".

5. IFRA SHAHEEN (2024) – "Deep Learning-Based Diabetes Prediction Using Hi-Le and HiTCLe Ensemble Models with ProWSyn for Class Imbalance Mitigation".

6. JIALI GAO, NA HU (2023) – "Machine Learning-Based Early Diabetes Diagnosis Using Pima Indian Diabetes Dataset".

7. JIMIN LIU, LUHAO FAN, QUANQIU JIA, LONGRI WEN, CHENGFENG SHI (2021) – "Early Diabetes Prediction Using

Stacking Ensemble Learning with Gradient Boosting Decision Tree, Adaboost, Random Forest, and Logistic Regression".

8. MADHUMITA PAL, SMITA PARIJA, GANAPATI PANDA (2021) – "Improved Prediction of Diabetes Mellitus Using Machine Learning-Based Approach".

9. P. MEENAKSHIDEVI, T R. LOGESH, G. NAVAYUGAN, M. SUGESH KANNAN (2024) – "Integration of Random Forest, Naive Bayes, Multi-Layer Perceptron, and SVM for Diabetes Prediction".

10. SADIA AFRIN SHAMPA, MD. SAIFUL ISLAM, AYATUN NESA (2023) – "Boosting Algorithms for Early Diabetes Prediction Using Multi-National Datasets".

11. SAKKAYAPHOP PRAVESJIT, KRITTIKA KANTAWONG, SUPAN TONGPHET, PANU BHROMMALEE, NAPA RACHATA (2020) – "Prediction Model for Diabetes Complications Using Modified ID3 Decision Tree and Mobile Application Testing".

12. SRISHTI MAHAJAN, PRADEEPTA KUMAR SARANGI, ASHOK KUMAR SAHOO, MUKESH ROHRA (2023) – "Diabetes Prediction Using Logistic Regression and Random Forest for Early Detection and Prevention".

13. SRISHTI MAHAJAN, PRADEEPTA KUMAR SARANGI, ASHOK KUMAR SAHOO, MUKESH ROHRA (2023) – "Diabetes Mellitus Prediction Using Supervised Machine Learning Techniques".