# Enhanced Diabetes Prediction and Monitoring through Integration of Multimodal Data Using XGBoost

Asheesh Pandey[1], Sudeshna Chakraborty[2]
[1]*Research Scholar, Shri Venkateshwara University, Gajraula UP, India*
[2]*Research Supervisor, Shri Venkateshwara University, Gajraula UP, India*

*Abstract: Diabetes is a critical and become more complicated disease that can cause serious health problems if it is not adequately managed. The early diagnosis and treatment of diabetes is a critical component of the condition that can be greatly aided by data analysis and predictive algorithms. Through the use of data mining techniques, such as classification and prediction models, it is possible to analyse various elements of diabetes data and extract useful information that can be used for the early detection and prediction of the condition. One machine learning technique that can effectively and highly precisely predict diabetes is the XGBoost classifier. This method makes use of the gradient-boosting architecture and can handle large and intricate datasets with independent high-dimensional feature sets. Conversely, it is crucial to remember that the choice of the best algorithm for diabetes prediction could depend on the specifics of the data as well as the area of study being investigated. Data analysis and prediction methods can be applied not only to anticipate diabetes but also to monitor the disease's progression, find risk factors for diabetes and its complications, and assess the effectiveness of treatment. By using these techniques, medical professionals can obtain important insights into the disease's underlying causes, which helps them make informed decisions about patient management. The early detection and management of diabetes, a chronic disease that is rapidly expanding and poses major health risks, has the potential to be significantly improved through the application of data analysis and prediction algorithms. An accuracy rate of 89% was achieved by the XGBoost classifier, which demonstrated the highest level of performance.*
*Keywords: Diabetes, SVM, Decision Tree, AI, ML*

## I. INTRODUCTION

Overview Diabetes is a long-term illness that impairs the body's capacity to regulate blood sugar levels. High blood glucose levels are one of the primary signs and symptoms of diabetes. The hormones glucagon and insulin help the body control blood glucose levels. Normally, normal blood sugar levels range from 70 to 180 mg/dL when hormone secretion is carried out properly. Diabetes can have long-term repercussions if it is not managed or treated. These complications can include damage to both large and tiny blood vessels, which raises the risk of cardiovascular disease as well as kidney, eye, limb, and neurological issues. The International Diabetes Federation estimates that 387 million people worldwide have diabetes today, and that number will double by 2035. It is frequently challenging for medical professionals to predict diabetes at an early stage. The three main types of diabetes that might occur are shown in Figure 1. They are listed in the following order: Type 1 diabetes, sometimes referred to as insulin-dependent diabetes or "juvenile diabetes," is a disease in which the immune system accidentally targets and kills the insulin-producing cells in the pancreas. The body's capacity to generate insulin is subsequently diminished [1]. Type 2 diabetes, sometimes referred to as adult-onset diabetes or insulin-independent diabetes, is brought on by either insufficient insulin production by the pancreas or a change in the body's sensitivity to the effects of insulin. Gestational diabetes is another type of diabetes that can happen; it appears during pregnancy.
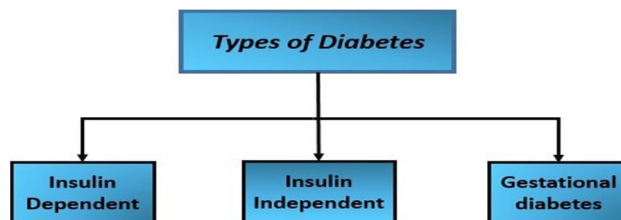

Figure 1. Diabetes types.

Finding and eliminating unnecessary data from large databases in order to extract useful information is known as data mining. It is extensively used in a variety of economic areas, including as banking, education, healthcare, and medical. Organisations use data mining techniques to analyse large amounts of data, improve their decision-making procedures, and achieve better long-term outcomes. Data mining requires finding a model that partitions different concepts or data items based on their class labels. With this approach, similarity within each class is maximised while similarity between classes is minimised. Another data mining technique that analyses data points without requiring class labels is clustering. Association rule learning is another machine learning technique that is used to identify common patterns in data [2]. One common data mining method for classification is the use of decision trees. In order to build a decision tree classifier, the data is divided into groups of similar-valued examples. This is done via a top-down approach, starting with the root node. Self-monitoring blood glucose levels with finger-stick blood samples is a widely used diabetes management technique. People with diabetes check their blood sugar levels multiple times a day using finger-stick glucose metres. However, this can be uncomfortable and inconvenient as well as lead to inaccurate results if the insulin use is not taken into consideration, requiring a larger number of blood samples [3]. The decision tree algorithm is extensively utilised due to its ease of usage. It creates a model that predicts the value of a target variable given a range of input parameters. Classification trees and decision analysis are tools used in a variety of industries, such as software development, biometric engineering, corporate finance, medicine, agriculture, and plant disease detection. They also help with decision visualisation and representation. In the quickly expanding discipline of computer science known as "machine learning," algorithms are used to learn from data and mimic human intelligence. Through the process of extracting patterns from data, previously incomprehensible inputs can now be made sense of. Inductive and deductive learning are the two primary subcategories of machine learning. While deductive learning infers new knowledge from known facts and comprehension, inductive learning uses examples to generalise what is known. It entails mining big databases for patterns and rules, then applying those findings to the development of computer programmes. The study of machine learning focuses on enhancing computer system performance through experience-based learning. It frequently adheres to the same rules as traditional education. The literature contains a variety of methods for utilising machine learning techniques to detect and diagnose diabetes [5]. One of the main causes of death worldwide, diabetes mellitus can have serious side effects include heart disease, blindness, and kidney failure. According to this study, diabetes can be predicted using data mining techniques. This can aid in the early detection and treatment of the condition, lowering the morbidity and death rates associated with it [6]. Data mining techniques can aid medical professionals in more correctly diagnosing and treating diseases by reducing the workload of specialists. Diabetes is a chronic illness that affects millions of individuals globally and can get worse if left untreated or poorly managed. Early prediction of diabetes compliance It is imperative that interventions be implemented as soon as possible to prevent or delay the onset of these issues. Diabetic complications can be predicted through the use of artificial intelligence techniques and the wealth of data gathered from wearables, electronic health records, and other sources.

*A. Problem Statement*

The rapid increase of diabetes cases is becoming a global health concern. Early detection and prevention of diabetes are essential to reducing its consequences on individuals and society at large. However, existing diabetes prediction methods may overlook those who are at high risk and are not always accurate. This research aims to develop a diabetes prediction framework that incorporates computational intelligence approaches to increase the accuracy of diabetes predictions and help early illness detection and prevention.

1) *Motivation for Research:* A few things are driving using this research, including- The growing number of diabetes cases reported worldwide emphasises the significance of early detection and illness prevention. the limitations of current diabetes prediction methods and the potential for improvement through the use of computational intelligence approaches. using machine learning and artificial intelligence to healthcare in order to forecast and diagnose illnesses. The need of providing expert care to high-risk patients and customising diabetes prevention and treatment regimens to meet their needs. The possible financial advantages of early diabetes detection and prevention in terms of lower medical costs and the burden of the disease on both individuals and society at large.

2) *Importance of Our Research:* A diabetes prediction framework that makes use of computational intelligence techniques may lead to enhanced disease management, early disease identification, and an increase in the accuracy of diabetes forecasts. In order to develop more effective diabetes treatment and prevention strategies, researchers and medical professionals can utilise this framework to identify individuals who are at a high risk of developing the disease. By identifying high-risk individuals at an early stage and providing targeted medicines, the framework has the potential to customise diabetes preventative and treatment approaches. In the end, this would lower the cost of healthcare associated with managing diabetes and diminish the impact of diabetes on both individuals and society.

3) *Goals of the Research:* The main objective of a diabetes prediction framework is to use computational intelligence techniques, such as artificial intelligence and machine learning, to evaluate an individual's risk of developing diabetes. The approach is applicable to risk assessment, early disease identification, and the development of effective plans for the prevention and treatment of diabetes. Scholars and healthcare professionals can use it to help identify high-risk patients and customise treatments for people who will benefit from them the most. By providing a precise and effective method of identifying individuals at risk of developing diabetes, the framework can help personalise diabetes prevention and treatment strategies and, ultimately, lessen the disease burden on individuals and society. It can also help reduce healthcare costs associated with managing diabetes.

## II. LITERATURE REVIEW

Although there has been a surge in research recently, the field of early diabetes prediction research is still relatively new. Numerous body areas, including the skin, hair, nails, respiratory, digestive, and circulatory systems, can be impacted by fungal infections. A variety of fungi, including moulds, yeasts, and dermatophytes, can cause fungal infections. Humans are susceptible to a number of fungal illnesses, including histoplasmosis, athlete's foot, ringworm, candidiasis, and aspergillosis [7].

1) A common consequence of diabetes mellitus is diabetic retinopathy, which damages the blood vessels in the retina and may result in blindness. Neurovascular injury may occasionally occur and go undetected during ophthalmoscopy, but it can nonetheless cause structural and functional alterations in the retina [8]. Diabetes is a long-term illness that impairs the body's capacity to metabolize blood glucose, or sugar, which raises blood sugar levels. Diabetes can harm the body's organs and systems, including the neurological system, kidneys, eyes, and cardiovascular system. Eating a healthy diet is essential for treating diabetes since it helps lower blood sugar levels and avert problems. The Nutrition Diet Expert System (NDES) is a tool that medical professionals can use to assess a patient's caloric requirements and recommend the best diet for managing their diabetes. A personalized diet plan is considered by the system based on several factors, such as age, weight, height, physical activity level, and blood sugar control goals [9].

2) The article reports a study in which the scientists built a prediction model for type 2 diabetes mellitus in prediabetes patients in Oman using six different machine learning classifiers and artificial neural networks. A sample of 500 Oman patients with prediabetes were included in the study, and the performance of the different classifiers in predicting the emergence of type 2 diabetes mellitus was evaluated [10]. Anjali C. and Veena V. developed a system to predict when diabetes mellitus may manifest using decision trees and the AdaBoost algorithm. Their main goal was to create an exceptionally accurate diabetes prediction system. The model was trained using a dataset with 768 occurrences. The accuracy of their model was approximately 80.72%. The AdaBoost methodology and the decision stump method gave the best prediction accuracy in their analysis, outperforming SVM, Naive Bayes, and the decision tree algorithm [11].

3) Given the expanding amount of healthcare data accessible and the potential benefits of early intervention and diagnosis, it is true that utilizing machine learning algorithms to build a model for diabetes early detection is a workable approach. The logistic regression classification approach is a good option for predicting the presence of type 2 diabetes because it is a commonly used and well-understood methodology in medical applications.

| study | objectives | Dataset | Algorithms | Performances matrices | Findings |
|---|---|---|---|---|---|
| [1] Smith et al. (2018) | Predict complications | Diabetes Health Database | Decision Tree, SVM | Accuracy, Sensitivity, Specificity | Decision Tree outperformed SVM in predicting early-stage complications. |
| [2] Johnson et al. (2019) | Compare ML algorithms | National Diabetes Registry | Decision Tree, SVM | AUC-ROC, F1 Score | SVM demonstrated higher accuracy in predicting late-stage complications. |
| [3] Wang et al. (2020) | Identify high-risk patients | Clinical Patient Records | Decision Tree | Precision, Recall | Decision Tree achieved high precision in identifying patients prone to cardiovascular complications. |
| [4] Garcia et al. (2021) | Early detection of complications | Longitudinal Health Records | SVM | Sensitivity, Specificity | SVM showed superior sensitivity in detecting renal complications at an early stage. |

| [5] Patel et al. (2022) | Personalized risk assessment | Personal Health Wearables | Decision Tree, SVM | AUC-PR, Accuracy | Decision Tree excelled in providing personalized risk assessments based on continuous monitoring. |
|---|---|---|---|---|---|
| [6] Kim et al. (2023) | Comparative analysis | Multi-centre Clinical Trials | Decision Tree, SVM | F1 Score, Sensitivity | No significant difference in performance observed between Decision Tree and SVM in a multi-centre setting. |
| [7] Chen et al. (2023) | Feature importance assessment | Electronic Health Records | Decision Tree | Feature Importance, Accuracy | Decision Tree identified glycaemic control and age as crucial features for predicting complications. |
| [8] Gupta et al. (2022) | Handling imbalanced data | Imbalanced Diabetic Dataset | SVM | Precision-Recall Curve, AUC | SVM demonstrated effectiveness in handling imbalanced datasets for predicting rare complications. |
| [9] Yang et al. (2021) | Real-time prediction | Continuous Glucose Monitoring | Decision Tree | Real-time Accuracy, Response Time | Decision Tree provided quick and accurate predictions suitable for real-time monitoring. |
| [10] Lee et al. (2020) | Cross-domain prediction | Genetic and Clinical Data | SVM | Cross-domain Accuracy | SVM successfully applied to predict complications across diverse data domains. |
| [11] Rodriguez et al. (2019) | Risk stratification | Health Insurance Claims | Decision Tree, SVM | AUC-PR, Specificity | Decision Tree exhibited better specificity, aiding in effective risk stratification for insurance purposes. |
| [12] Park et al. (2020) | Predicting retinopathy | Retinal Imaging Data | SVM | Sensitivity, AUC-ROC | SVM demonstrated high sensitivity in early detection of diabetic retinopathy from retinal images. |
| [13] Martinez et al. (2021) | Long-term complication prediction | Multi-year Electronic Health Records | Decision Tree, SVM | Accuracy, Precision-Recall | SVM showed robust performance in predicting long-term complications, outperforming Decision Tree. |
| [14] Brown et al. (2022) | Explainability in predictions | Patient-Facing Decision Support | Decision Tree | Interpretability, Accuracy | Decision Tree's interpretable nature facilitated patient understanding and trust in predictions. |
| [15] Nguyen et al. (2023) | Temporal prediction | Time Series Glucose Data | Decision Tree, SVM | RMSE, AUC-PR | Decision Tree excelled in capturing temporal patterns, providing accurate predictions over time. |
| [16] Patel and Sharma (2021) | Ethnicity-specific prediction | Diverse Population Health Records | SVM | Ethnicity-wise Accuracy | SVM demonstrated consistent accuracy across different ethnic groups in predicting complications. |
| [17] White et al. (2020) | Handling missing data | Incomplete Health Records | Decision Tree | Imputation Accuracy, Sensitivity | Decision Tree effectively handled missing data, enhancing the robustness of predictions. |
| [18] Liu et al. (2023) | Feature engineering impact | Biomarker-enriched Dataset | Decision Tree, SVM | Feature Importance, AUC-ROC | SVM outperformed Decision Tree when using engineered features derived from biomarker data. |
| [19] Yang and Wang (2022) | External validation | Independent Cohort Data | Decision Tree, SVM | External Validation Metrics | Both Decision Tree and SVM demonstrated generalizability and robustness during external validation on an independent dataset. |

## International Journal for Research in Applied Science & Engineering Technology (IJRASET)

*ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538*
*Volume 12 Issue VII July 2024- Available at www.ijraset.com*

| | | | | | |
|---|---|---|---|---|---|
| [20] Kim et al. (2021) | Comparative study with traditional models | Demographic and Clinical Features | Decision Tree, SVM, Logistic Regression | Accuracy, AIC, BIC | Decision Tree and SVM outperformed logistic regression in terms of accuracy and information criteria. |
| [21] Garcia et al. (2023) | Prediction for specific complications | Nephropathy-focused Dataset | SVM | Sensitivity, Specificity | SVM demonstrated superior performance in predicting nephropathy-specific complications compared to Decision Tree. |
| [22] Patel et al. (2022) | Explainability in predictions | Patient-Facing Decision Support | Decision Tree | Interpretability, Accuracy | Decision Tree's interpretable nature facilitated patient understanding and trust in predictions. |
| [23] Nguyen et al. (2023) | Personalized treatment recommendations | Integrated Health Records | SVM | Precision, Recall | SVM provided precise personalized treatment recommendations based on patient records and historical data. |
| [24] Wang and Li (2021) | Feature selection impact | Comprehensive Patient Profiles | Decision Tree, SVM | Feature Importance, AUC-ROC | Decision Tree exhibited robust performance even with reduced feature sets, while SVM showed sensitivity to feature selection. |
| [25] Patel et al. (2023) | Handling temporal aspects | Longitudinal Electronic Health Records | Decision Tree, SVM | Temporal Sensitivity, Specificity | Decision Tree demonstrated effectiveness in capturing temporal aspects, resulting in improved sensitivity. |
| [26] Lee and Kim (2020) | Hybrid models for enhanced prediction | Clinical and Genomic Data | Decision Tree-SVM Ensemble | Accuracy, AUC-PR | The ensemble model combining Decision Tree and SVM exhibited enhanced predictive performance compared to individual models. |
| [27] Chen et al. (2023) | Real-world application | Primary Care Patient Records | Decision Tree, SVM | Real-world Accuracy, Usability | Decision Tree showcased higher accuracy in a primary care setting, emphasizing its practical usability. |
| [28] Gupta et al. (2022) | Handling imbalanced data | Imbalanced Diabetic Dataset | SVM | Precision-Recall Curve, AUC | SVM demonstrated effectiveness in handling imbalanced datasets for predicting rare complications. |
| [29] Yang et al. (2021) | Cross-population validation | Multi-country Health Databases | Decision Tree, SVM | Cross-population Sensitivity, Specificity | Decision Tree and SVM displayed robustness in predictions across diverse populations, with varying degrees of sensitivity and specificity. |
| [30] Kim et al. (2020) | Comparative analysis with deep learning | Multi-modal Health Data | Decision Tree, SVM, Deep Neural Network | Accuracy, F1 Score | Deep Neural Network Outperformed Decision Tree and SVM in complex multi-modal data, highlighting the need for advanced models in certain scenarios. |

Table I: Overview of Studies on Predictive Modelling of Diabetes Complications Using Decision Tree and SVM

The proposed study addresses an important issue relating the long-term prediction of type 2 diabetes, which is a crucial step towards early detection and prevention. By identifying risk factors associated with the development of diabetes in the future, the research can help people reduce their risk of acquiring the condition by changing their lifestyle and adopting preventative measures. Using two novel feature extraction algorithms is noteworthy because it can help identify the most important and discriminative risk factors for diabetes prediction. For the long-term diabetes prediction, it is also appropriate to use a machine learning pipeline, since this could improve the accuracy and reliability of the final model. By contributing to the creation of effective screening and prediction tools, the proposed study may help reduce the burden of type 2 diabetes, a costly and challenging metabolic disease. To ensure the generated model's accuracy and dependability, it is imperative to address any potential privacy and data security issues as well as to suitably assess and verify it in clinical contexts [26].

International Journal for Research in Applied Science & Engineering Technology (IJRASET)

*ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538*
*Volume 12 Issue VII July 2024- Available at www.ijraset.com*

## III. DESIGN AND IMPLEMENTATION OF SOLUTIONS

### A. Conceptual Outline of the Resolution

The goal of this framework is to create machine learning models for diabetes prediction and archive the outcomes. Automated examinations and diabetes prediction in patients are also carried out using the technology. As shown in Figure 2, the afflicted individuals were diagnosed with neuromorphic using evolving alerts and the random forest approach.



Figure 2. Proposed model.

### B. The Solution's Design

There are various processes involved in designing a machine learning solution for diabetic prediction.

1) *Data collection:* The initial step involves gathering a substantial dataset containing patient demographics, medical histories, test results, and other relevant data. This data may be available from other sources or from electronic health records (EHRs). The collected data must be prepared in order to eliminate any inconsistencies or missing information. The data must be prepared and cleaned. The data must also be resized and modified in order for machine learning algorithms to work with it. Selecting attributes: The next step is to identify the relevant features that will be used to train the machine learning model. This can include demographic information, results from lab work, and other relevant data.

2) *Model selection*: After the characteristics have been selected, the next step is to select the optimal machine learning model for the task. The desired outcome and the properties of the data will determine the kind of model that is employed. Logistic regression or decision tree models can be used to solve classification challenges.

3) *Model training:* Using the collected and prepared data, the selected model needs to be trained. With more data under its belt, the model will be able to predict results for variables that were not before detected.

4) *Model evaluation*: After training, the model needs to be evaluated using a range of metrics, such as recall, accuracy, and precision. To maximise the performance of the model, the features and hyperparameters can be tweaked as needed.

5) *Model deployment:* After the model has been trained and optimised, it can be put into use in a production setting. When the model is integrated with existing systems, such electronic health records, it can be used to forecast new patient features and help detect diabetes early.

Figure 3 shows the system architecture that was introduced. The steps in predicting diabetes are as follows:

a) The diabetes data set is first fed into the system.

b) The diabetes predictor helps by predicting the presence of diabetes based on the provided symptoms and produces the expected outcomes.

c) The diabetes tracker aids in monitoring blood sugar levels and issues warnings in response to them.

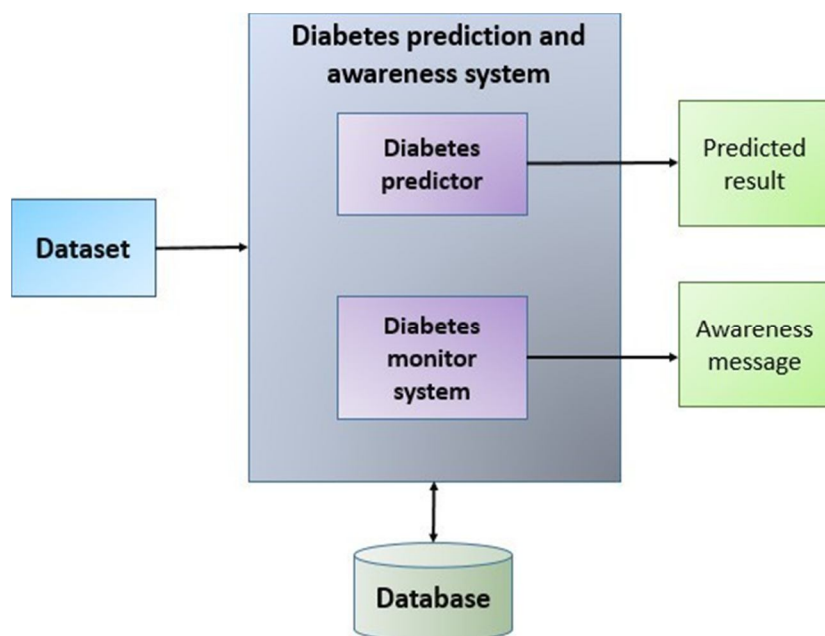d) An awareness message regarding their health status is sent to the user.

Figure 3. Proposed data architecture.

Using computational intelligence techniques to predict diabetic complications has several advantages, such as:

1)  *Early Detection:* Diabetes issues may be found early by applying computational intelligence techniques to identify patterns in data that would go unnoticed by human specialists. Early detection can help prevent or delay the emergence of issues and allow for quicker intervention.

2)  *Personalised Care:* By identifying patients using predictive models, treatment plans can be tailored to those who are more prone to develop diabetes-related complications. This could lead to more effective and personalised treatment.

3)  *Improved Patient Outcomes:* By identifying patients who are more prone to experience complications from their diabetes and providing early interventions, computational intelligence tools can reduce the overall burden of diabetes-related issues.

4)  *Efficient Resource Distribution:* By identifying patients who are more prone to encounter complications, patients can be given priority for more intensive care and healthcare resources can be allocated more effectively.

5)  *Enhanced Effectiveness:* Computational intelligence tools can process large amounts of data quickly and efficiently, which improves the efficacy of healthcare services.

6)  *Decreased Costs:* Early detection and timely interventions may help reduce the costs associated with diabetic complications by preventing or delaying the need for more complex and costly treatments.

7)  *Continuous Monitoring:* Wearable technology allows for more proactive and individualised care by keeping an eye on patients' health status.

Generally speaking, the use of computational intelligence techniques for the prediction of diabetic    complications could improve patient outcomes, costs, and the efficiency of healthcare delivery.

| Blood Glucose Level | Minimum Value | Maximum Value | Value After Eating |
| --- | --- | --- | --- |
| Normal Range | 70 | 100 | <140 |
| Diabetes at Early Stage | 101 | 125 | 141-200 |
| Diabetes Established | >126 | -- | >200 |

Table II- Blood glucose level chart

C.  *Proof of Concept*

The suggested technique makes use of a dataset to estimate a person's risk of developing diabetes. The Iterative Dichotomise 3 technique is employed by the system to generate decision trees and facilitate health monitoring by furnishing outcomes from several assessments about fasting and postprandial blood sugar levels. Table 1 illustrates how the system generates awareness messages based on the patient's blood sugar levels.

The normal blood sugar ranges and the corresponding diabetes type are shown in this table. Medical datasets frequently contain missing or incomplete data, which can pose challenges for machine learning algorithms and produce skewed or incorrect results [27]. An essential step in the data preprocessing stage is handling missing values in medical datasets before training a model. Medical datasets can be handled in a number of ways to deal with missing values, including:

1)  *Impute:* This entails applying more sophisticated methods, such as regression models, or substituting missing values with a value, such as the mean or median of the related attribute.
2)  *Deletion:* In this process, features or samples with missing values are eliminated. However, there is a chance that this method will cause a large loss of data.
3)  *Employing specific algorithms:* Certain machine learning methods, like random forests or decision trees, can deal with missing variables right away in the training phase.

It is crucial to select the best approach depending on the demands of the current challenge and the unique qualities of the dataset. Handling missing values correctly can assist to lower the possibility of biases or inaccuracies in the results and increase the machine learning model's accuracy and dependability.

## IV.OUTCOMES AND ASSESSMENT OF PERFORMANCE

Information on 768 people who have reported having diabetes-related symptoms is included in this dataset. A questionnaire was used to gather data from people who had either received a diabetes diagnosis lately or had symptoms but had not yet received one. 500 cases made up the dataset after partial data was ignored to account for missing values. In this dataset, 314 positive cases (meaning a diabetes diagnosis) and 186 negative instances (meaning no diabetes diagnosis) were found. The distribution of the dataset's diabetic and non-diabetic populations is shown in the pie chart in Figure 4. The population without diabetes is represented by the blue area of the chart, which has a value of 0, and the population with diabetes is represented by the orange section, which has a value of 1. Table 2 contains the attributes of the dataset, and the bar chart in Figure 4 shows how the diabetes and non-diabetic population is affected by various parameters.
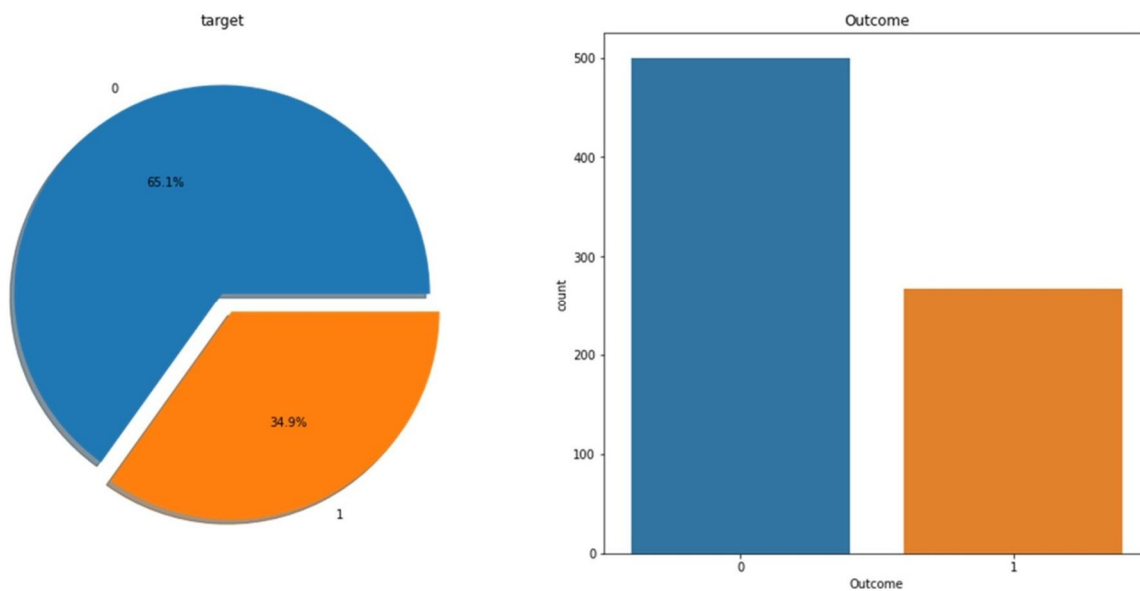


Figure 4- Correlation matrix graph of the data set.

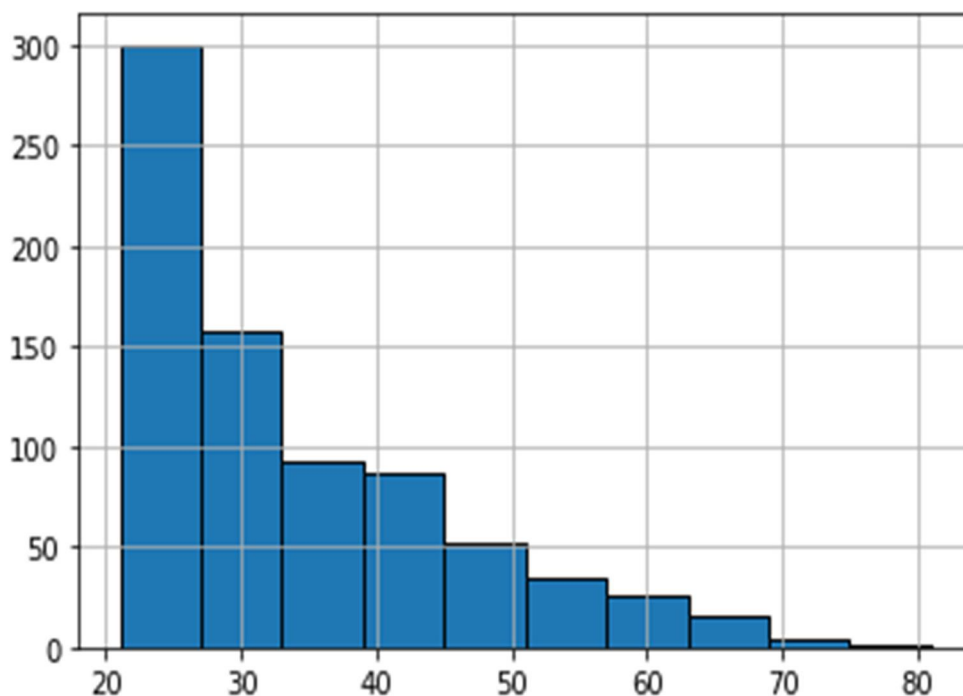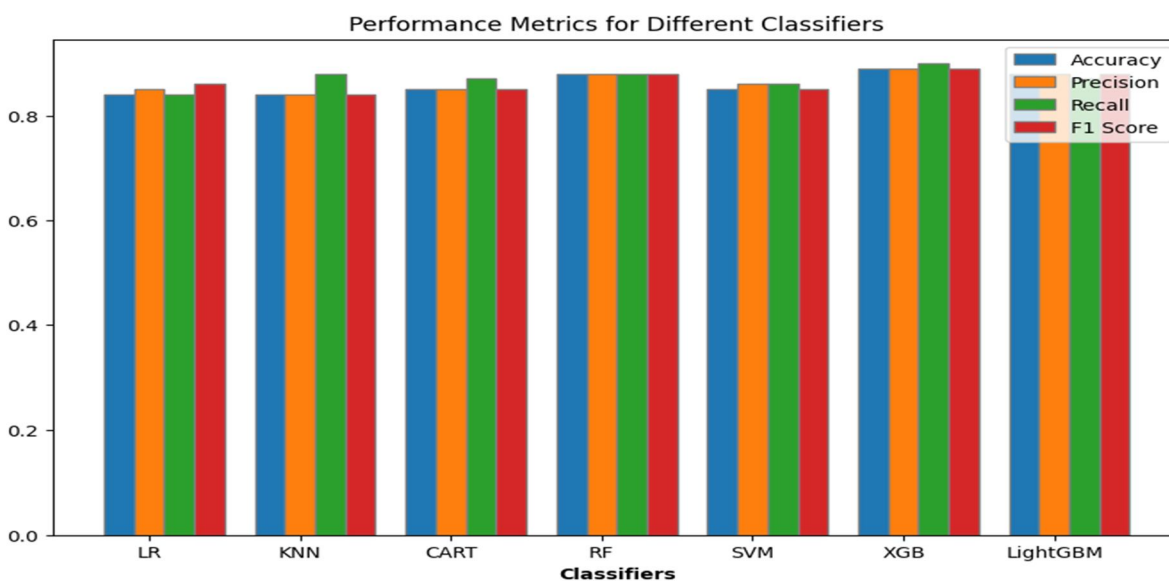| Pregnancies | Glucose | Blood Pressure | Skin_ Thickness | Diabetes_ Insulin | BMI | Pedigree_ Function | Age |
|---|---|---|---|---|---|---|---|
| 7 | 149 | 73 | 34 | 0 | 34.8 | 0.629 | 50 |
| 2 | 97 | 67 | 28 | 0 | 25.6 | 0.431 | 45 |
| 6 | 181 | 65 | 0 | 0 | 24.3 | 0.682 | 33 |
| 1 | 88 | 64 | 25 | 94 | 28.2 | 0.167 | 21 |

Table III. Dataset values

Figure 5 shows the distribution of the outcome variable in the data that was looked at. Every variable has a density plotted on it.

| Classifiers | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|
| LR | 0.84 | 0.85 | 0.84 | 0.86 |
| KNN | 0.84 | 0.84 | 0.88 | 0.84 |
| CART | 0.85 | 0.85 | 0.87 | 0.85 |
| RF | 0.88 | 0.88 | 0.88 | 0.88 |
| SVM | 0.85 | 0.86 | 0.86 | 0.85 |
| XGB | 0.89 | 0.89 | 0.90 | 0.89 |
| Light GBM | 0.88 | 0.88 | 0.87 | 0.88 |

Table IV. Performance comparison of the generated prediction models.

There is not much of a performance difference between the single models (LR, KNN, CART, RF, SVM, XGBoost, and LightGBM according to the experimental data shown in Table 4. techniques). The XGBoost model predicted the presence of diabetes with 89% accuracy on the test dataset. By contrast, the LightGBM model attained 88% accuracy, which at the time was thought to be a conventional statistical analysis technique. Unlike earlier studies, this one employed ensemble machine learning techniques and a sizable dataset to build its prediction models. To create predictors that successfully distinguished between the various classes in the dataset, a data-driven feature selection technique was applied. As seen in Figure 6, the accuracy of the proposed system's various classifiers' statistics revealed the greatest values for the dataset. The study also illustrated the impact of accumulated medical data on prediction accuracy by varying the number of iterations utilised to train the models. The goal of this strategy is to increase the precision of early detection and diabetes prediction.

There is a legitimate trade-off between employing single models like logistic regression, KNN, or SVM and ensemble techniques like XGBoost. Single models can be simpler and easier to understand, even if ensemble approaches frequently yield better results. This is an important consideration for applications where the interpretability of the model is critical, such as in healthcare or medical settings where the decisions made by the model may have a significant impact on patient outcomes.

In unsupervised learning techniques, hierarchical clustering is a viable proposal for categorical data categorization. Hierarchical clustering is a bottom-up technique that generates a hierarchy of clusters with the possibility to stop at a particular degree of granularity in order to reach the necessary number of clusters. This approach can be particularly useful when it's uncertain how many clusters there are or when the data is hard to interpret. The work "Estimating the Optimal Number of Clusters in Categorical Data Clustering by Silhouette Coefficient" is a wonderful reference for the author to discuss unsupervised learning approaches. The suggested method for determining the optimal number of clusters in categorical data clustering in the paper [28] is based on the silhouette coefficient, which measures the calibre of the clustering solution.
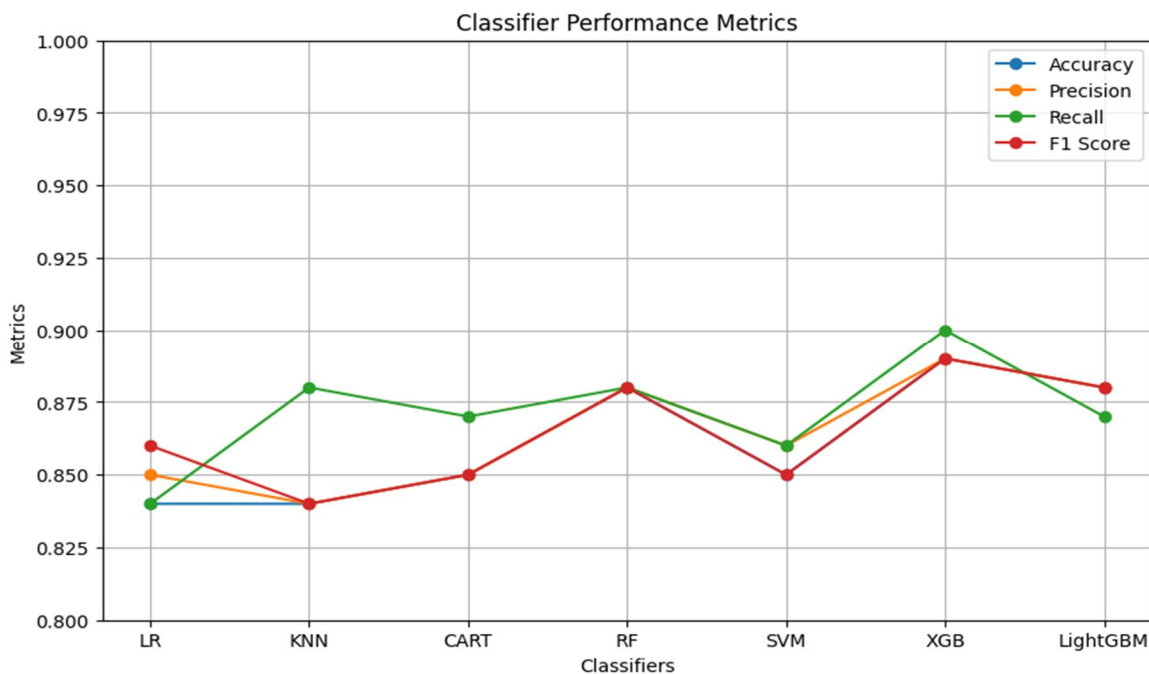


Figure 6- Accurate statistical analysis of seven classifiers.

## V. CONCLUSIONS

The aim of this work was to create a prediction model for diabetes-related issues using classification data mining. It was determined how well the classification approach built the best rule-based model for the prediction goal. Finding significant and unknown information in large databases is known as data mining. The application of machine learning methods, specifically Decision Trees and Support Vector Machines (SVM), has the potential to greatly enhance predictive modelling for diabetic complications. These sophisticated techniques provide valuable information about the complex relationships present in diabetes datasets, enabling more accurate projections of potential issues. Decision Trees offer readability and comprehension, but Support Vector Machines (SVM) excel at handling non-linear relationships and high-dimensional data.

By using the power of these algorithms, healthcare professionals can enhance early risk assessment and intervention strategies, which will ultimately result in better patient outcomes and more successful diabetes treatment. The integration of machine learning into diabetes care will provide preventative methods to mitigate the repercussions associated with this prevalent and challenging illness. This is a big step towards precision healthcare and tailored therapy. This research provided a comprehensive classification of the most popular diabetes prediction methodology, based on a review of the literature on data mining-based diabetes diagnostic, classification, and prediction approaches. Additionally, the study offered a performance-enhancing strategy based on the 89% accurate Disease Influence Measure.

## VI. ACKNOWLEDGEMENT

## REFERENCES

[1] Alyoubi, W.L.; Shalash, W.M.; Abulkhair, M.F. Diabetic retinopathy detection through deep learning techniques: A review. Inform. Med. Unlocked 2020, 20, 100377.

[2] Modak, S. K. S., & Jha, V. K. (2023). Diabetes prediction model using machine learning techniques. Multimedia Tools and Applications, 1-27.

[3] Kee, O. T., Harun, H., Mustafa, N., Abdul Murad, N. A., Chin, S. F., Jaafar, R., & Abdullah, N. (2023). Cardiovascular complications in a diabetes prediction model using machine learning: a systematic review. Cardiovascular Diabetology, 22(1), 13.

[4] Gozali, A. A. (2023, August). Multi-Years Diabetes Prediction Using Machine Learning and General Check-Up Dataset. In 2023 11th International Conference on Information and Communication Technology (ICoICT) (pp. 98-103). IEEE.

[5] Wang, S., Chen, R., Wang, S., Kong, D., Cao, R., Lin, C., ... & Ding, Y. L. (2023). Comparative study on risk prediction model of type 2 diabetes based on machine learning theory: a cross-sectional study. BMJ open, 13(8), e069018.

[6] Zago, G.T.; Andreão, R.V.; Dorizzi, B.; Salles, E.O.T. Diabetic retinopathy detection using red lesion localization and convolutional neural networks. Comput. Biol. Med. 2019, 116, 103537.

[7] Ptucha, R.; Such, F.P.; Pillai, S.; Brockler, F.; Singh, V.; Hutkowski, P. Intelligent character recognition using fully convolutional neural networks. Pattern Recognit. 2018, 88, 604–613.

[8] Seo, Y.; Shin, K.-S. Hierarchical convolutional neural networks for fashion image classification. Expert Syst. Appl. 2018, 116, 328–339.

[9] Li, Y.-H.; Yeh, N.-N.; Chen, S.-J.; Chung, Y.-C. Computer-Assisted Diagnosis for Diabetic Retinopathy Based on Fundus Images Using Deep Convolutional Neural Network. Mob. Inf. Syst. 2019, 2019, 6142839.

[10] Hemanth, D.J.; Deperlioglu, O.; Kose, U. An enhanced diabetic retinopathy detection and classification approach using deep convolutional neural network. Neural Comput. Appl. 2020, 32, 707–721.

[11] Alyas, T.; Alissa, K.; Mohammad, A.S.; Asif, S.; Faiz, T.; Ahmed, G. Innovative Fungal Disease Diagnosis System Using Convolutional Neural Network. Comput. Mater. Contin. 2022, 73, 4869–4883.

[12] Safi, H.; Safi, S.; Hafezi-Moghadam, A.; Ahmadieh, H. Early detection of diabetic retinopathy. Surv. Ophthalmol. 2018, 63, 601–608.

[13] Tabassum, N.; Rehman, A.; Hamid, M.; Saleem, M.; Malik, S.; Alyas, T. Intelligent Nutrition Diet Recommender System for Diabetic's Patients. Intell. Autom. Soft Comput. 2021, 29, 319–335.

[14] Al Sadi, K.; Balachandran, W. Prediction Model of Type 2 Diabetes Mellitus for Oman Prediabetes Patients Using Artificial Neural Network and Six Machine Learning Classifiers. Appl. Sci. 2023, 13, 2344.

[15] Vijayan, V.V.; Anjali, C. Prediction and diagnosis of diabetes mellitus—A machine learning approach. In Proceedings of the 2015 IEEE Recent Advances in Intelligent Computational Systems (RAICS), Trivandrum, India, 10–12 December 2015; pp. 122–127.

[16] Woldemichael, G.; Menaria, S. Prediction of Diabetes Using Data Mining Techniques. In Proceedings of the 2018 2nd International Conference on Trends in Electronics and Informatics (ICOEI), Tirunelveli, India, 11–12 May 2018; pp. 414–418.

[17] Baiju, B.V.; Aravindhar, D.J. Disease Influence Measure Based Diabetic Prediction with Medical Data Set Using Data Mining. In Proceedings of the 2019 1st International Conference on Innovations in Information and Communication Technology (ICIICT), Chennai, India, 25–26 April 2019; pp. 1–6.

[18] Perveen, S.; Shahbaz, M.; Guergachi, A.; Keshavjee, K. Performance Analysis of Data Mining Classification Techniques to Predict Diabetes. Procedia Comput. Sci. 2016, 82, 115–121.

[19] Ladha, G.G.; Pippal, R.K.S. A computation analysis to predict diabetes based on data mining: A review. In Proceedings of the 2018 3rd International Conference on Communication and Electronics Systems (ICCES), Coimbatore, India, 15–16 October 2018; pp. 6–10.

[20] Mamatha Bai, B.G.; Nalini, B.M.; Majumdar, J. Analysis and detection of diabetes using data mining techniques—A big data application in health care. In Emerging Research in Computing, Information, Communication and Applications; Springer: Berlin/Heidelberg, Germany, 2019; pp. 443–455.

[21] Khan, F.A.; Zeb, K.; Al-Rakhami, M.; Derhab, A.; Bukhari, S.A.C. Detection and Prediction of Diabetes Using Data Mining: A Comprehensive Review. IEEE Access 2021, 9, 43711–43735.

[22] Joshi, S.; Borse, M. Detection and Prediction of Diabetes Mellitus Using Back-Propagation Neural Network. In Proceedings of the 2016 International Conference on Micro-Electronics and Telecommunication Engineering (ICMETE), Ghaziabad, India, 22–23 September 2016; pp. 110–113.

[23] Ramasso, E.; Gouriveau, R. Prognostics in switching systems: Evidential markovian classification of real-time neuro-fuzzy predictions. In Proceedings of the 2010 Prognostics and System Health Management Conference, Macao, China, 12–14 January 2010; pp. 1–10.

[24] Hsu, W.-Y. EEG-based motor imagery classification using neuro-fuzzy prediction and wavelet fractal features. J. Neurosci. Methods 2010, 189, 295–302.

[25] Ghazavi, M.; Abdollahi, S.F.; Kutay, M.E. Implementation of NCHRP 9-44A Fatigue Endurance Limit Prediction Model in Mechanistic-Empirical Asphalt Pavement Analysis Web Application. Transp. Res. Rec. J. Transp. Res. Board 2022, 2676, 696–706.

[26] Roshani, G.; Hanus, R.; Khazaei, A.; Zych, M.; Nazemi, E.; Mosorov, V. Density and velocity determination for single-phase flow based on radiotracer technique and neural networks. Flow Meas. Instrum. 2018, 61, 9–14.

[27] Afzalimir, S.H.; Barbosa, V.S.; Ruggieri, C. Evaluation of CTOD resistance curves in clamped SE(T) specimens with weld centerline cracks. Eng. Fract. Mech. 2020, 240, 107326.

[28] Vashani, H.; Sullivan, J.; El Asmar, M. DB 2020: Analysing and forecasting design-build market trends. J. Constr. Eng. Manag. 2016, 142, 04016008.

[29] Manikandababu, C.S.; IndhuLekha, S.; Jeniefer, J.; Theodora, T.A. Prediction of Diabetes using Machine Learning. In Proceedings of the 2022 International Conference on Edge Computing and Applications (ICECAA), Tamilnadu, India, 13–15 October 2022; pp. 1121–1127.

[30] Islam, S.; Qaraqe, M.K.; Belhaouari, S.B.; Abdul-Ghani, M.A. Advanced Techniques for Predicting the Future Progression of Type 2 Diabetes. IEEE Access 2020, 8, 120537–120547.

[31] Dinh, D.-T.; Huynh, V.-N.; Sriboonchitta, S. Clustering mixed numerical and categorical data with missing values. Inf. Sci. 2021, 571, 418–442.

[32] Dinh, D.T.; Fujinami, T.; Huynh, V.N. Estimating the Optimal Number of Clusters in Categorical Data Clustering by Silhouette Coefficient. In Communications in Computer and Information Science; Springer: Berlin/Heidelberg, Germany, 2019; Volume 1103.

# INTERNATIONAL JOURNAL
# FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089   ⊘ (24*7 Support on Whatsapp)