



AI POWERED DIABETIC PREDICTION USING ML

Mrs.M.V. Sheela Devi¹, M. Radha Priya², K. Ashitha³, L. Bhuvaneshwari⁴

¹Assistant Professor, Department of Computer Science and Engineering,
KKR & KSR Institute of Technology and Sciences, Vinjanampadu, Vatticherukuru Mandal, Guntur,
Andhra Pradesh 522017

Email: sheela.softinfo@gmail.com¹

²³⁴UG Scholar, Department of Computer Science and Engineering,
KKR & KSR Institute of Technology and Sciences, Vinjanampadu, Vatticherukuru Mandal, Guntur,
Andhra Pradesh 522017

Email: 22jr1a0580@gmail.com², 22jr1a0570@gmail.com³, 22jr1a0573@gmail.com⁴

Abstract:

The rapid growth of healthcare data has created opportunities for applying Machine Learning (ML) techniques to improve disease prediction and early diagnosis. Diabetes is one of the most prevalent chronic diseases worldwide, and early detection plays a critical role in preventing severe complications. This study focuses on developing a predictive framework for diabetes detection using multiple machine learning algorithms applied to patient health records. A dataset consisting of relevant medical attributes was utilized to train and evaluate six different ML algorithms, namely Artificial Neural Networks (ANN), Extreme Gradient Boosting (XGBoost), AdaBoost, K-Nearest Neighbours (KNN), Support Vector Machine (SVM), and Decision Tree (DT). The performance of these algorithms was analyzed and compared based on their prediction accuracy and effectiveness in identifying diabetes. Comparative evaluation helps determine the most reliable and efficient model for diabetes prediction. Furthermore, the proposed approach supports the development of an application where users can input medical parameters and obtain prediction results. The outcomes of this study demonstrate the potential of machine learning techniques in assisting healthcare professionals in early diagnosis

and decision-making. By leveraging predictive analytics, the system can support medical practitioners in detecting diabetes at an early stage and improving patient care.

Keywords: Machine Learning, Diabetes Prediction, Artificial Neural Networks, XGBoost, AdaBoost, KNN, Support Vector Machine, Decision Tree, Healthcare Analytics.

I. INTRODUCTION

The medical sector manages enormous volumes of sensitive and confidential data that must be securely stored and protected from unauthorized access and modification. Due to sedentary lifestyles and unhealthy habits, Diabetes Mellitus has become one of the most prevalent and life-threatening diseases worldwide. Early detection and accurate prediction of diabetes are essential for medical practitioners to provide effective treatment and personalized healthcare services. Machine Learning (ML) techniques play a significant role in analyzing large healthcare datasets and extracting meaningful insights. Through data mining approaches, patterns and relationships within medical records can be

identified to support disease prediction and diagnosis. Uncontrolled diabetes can lead to severe complications such as heart disease, kidney failure, nerve damage, and vision loss. In this study, ML-based analysis is conducted using the WEKA data mining tool. The Pima Indians Diabetes Dataset from the UCI repository is used for evaluation. Algorithms such as Naïve Bayes (NB), Decision Tree (DT), and K-Nearest Neighbours (KNN) are applied to improve prediction accuracy and assist healthcare professionals in early diabetes detection.

II. LITERATURE SURVEY

Several researchers have explored the application of machine learning techniques for the early prediction and diagnosis of diabetes. These studies focus on analyzing medical datasets and identifying patterns that help healthcare professionals make accurate decisions. Chaurasia and Pal (2014) applied various data mining techniques to predict diabetes and found that classification algorithms can effectively identify diabetic patients using medical attributes. Similarly, Sisodia and Sisodia (2018) compared multiple classification algorithms and concluded that machine learning models provide reliable prediction results when trained with proper healthcare data. Nai-Arun and Moungrai (2015) conducted a comparative study of different classifiers and highlighted that selecting appropriate features significantly improves prediction accuracy. Breiman (2001) introduced the Random Forest algorithm, which uses ensemble learning to enhance prediction performance and reduce overfitting in classification problems. Cortes and Vapnik (1995) proposed the Support Vector Machine (SVM), which has been widely used in medical data analysis due to its ability to handle complex and high-dimensional datasets. Contreras and Vehi (2018) discussed the role of artificial intelligence

in diabetes management and emphasized its importance in medical decision support systems.

Recent studies also demonstrate the effectiveness of machine learning in healthcare analytics. Nurdin et al. (2023) developed machine learning-based models for diabetes prediction and achieved improved accuracy using healthcare datasets. Similarly, Hennebelle et al. (2023) proposed a smart healthcare framework that uses machine learning techniques to predict type 2 diabetes and support early diagnosis. These studies indicate that integrating machine learning with healthcare systems can significantly improve disease prediction and patient care.

III. PROPOSED WORK

The proposed system focuses on improving the accuracy and efficiency of diabetes prediction by integrating advanced machine learning techniques with cloud-based infrastructure. An enhanced data collection and preprocessing framework is designed to ensure that the dataset is clean, reliable, and representative of patient health records. Data preprocessing techniques such as data cleaning, normalization, and feature selection are applied to identify the most significant medical attributes that contribute to diabetes prediction, thereby improving model performance and reducing unnecessary complexity. The system employs multiple machine learning algorithms including Artificial Neural Networks (ANN), XGBoost, AdaBoost, K-Nearest Neighbours (KNN), Support Vector Machine (SVM), and Decision Tree (DT). These algorithms are trained and evaluated using patient health data to analyze their performance and determine the most effective model for accurate diabetes prediction. To enhance scalability and accessibility, cloud computing services are utilized for data storage, model training, and deployment. The system also incorporates a user-friendly interface where healthcare professionals or users can input

relevant medical parameters and obtain prediction results. This approach supports early detection of diabetes and assists medical practitioners in making informed and timely healthcare decisions.

IV. METHODOLOGY

1. DATA COLLECTION

The dataset used for this study is the Pima Indians Diabetes Dataset, which is widely used for medical prediction research. The dataset is obtained from publicly available repositories such as the UCI Machine Learning Repository and Kaggle. It contains medical records of female patients of Pima Indian heritage aged 21 years and above. The dataset includes several health-related attributes such as glucose level, blood pressure, body mass index (BMI), insulin level, skin thickness, number of pregnancies, diabetes pedigree function, and age. These attributes play an important role in identifying the presence of diabetes. The dataset also includes a target variable indicating whether a patient is diabetic or non-diabetic. Collecting a reliable and structured dataset is essential for building an accurate predictive model and ensuring that the system can effectively analyze medical patterns associated with diabetes risk.

2. DATA PREPROCESSING

Data preprocessing is an important step in machine learning that prepares raw data for analysis. The collected dataset may contain missing values, inconsistencies, or noisy data that can affect model performance. In this stage, data cleaning techniques are applied to remove or handle missing values and eliminate duplicate entries. Certain attributes such as glucose level, BMI, and blood pressure may contain zero values that are treated as missing data and replaced using

appropriate statistical methods. Categorical and numerical values are standardized to ensure uniform scaling across all features. Outlier detection methods are also applied to reduce the influence of extreme values. These preprocessing steps ensure that the dataset becomes consistent, accurate, and suitable for training machine learning models.

3. FEATURE SELECTION

Feature selection is performed to identify the most relevant attributes that influence diabetes prediction. Statistical analysis and correlation techniques are used to determine the relationship between input variables and the target variable. Important features such as glucose level, BMI, age, insulin level, and number of pregnancies are considered significant factors in predicting diabetes. Irrelevant or redundant features are removed to reduce dimensionality and improve computational efficiency. Feature selection also helps reduce overfitting and enhances the model's ability to generalize on unseen data. By focusing on the most relevant health indicators, the predictive model can achieve better accuracy and provide meaningful insights for medical decision-making.

4. MODEL DESIGNING

Several machine learning algorithms are implemented for diabetes prediction. These include Artificial Neural Networks (ANN), XGBoost, AdaBoost, K-Nearest Neighbours (KNN), Support Vector Machine (SVM), and Decision Tree (DT). Each algorithm is trained using the training dataset. The models learn patterns and relationships within the patient health records. These learned patterns help classify patients as diabetic or non-diabetic.

5. MODEL EVALUATION

After training, the models are evaluated using the testing dataset. Performance metrics such as accuracy, precision, recall, and F1-score are calculated. These metrics help determine the reliability and effectiveness of each algorithm. A comparative analysis is performed to identify the best-performing model. The algorithm with the highest performance is selected for diabetes prediction.

V. ALGORITHMS

1. Artificial Neural Network (ANN)

Artificial Neural Networks are computational models inspired by the human brain. They consist of interconnected nodes called neurons organized in input, hidden, and output layers. ANN can learn complex patterns and relationships from large datasets through training. It is widely used in medical diagnosis due to its ability to handle nonlinear data. In this study, ANN helps improve the prediction accuracy of diabetes by learning patterns from patient health records.

2. XGBoost (Extreme Gradient Boosting)

XGBoost is an advanced ensemble learning algorithm based on gradient boosting techniques. It combines multiple weak learners to create a strong predictive model. The algorithm improves performance by reducing errors iteratively during training. XGBoost is known for its speed, scalability, and high accuracy in classification problems. In this research, it is used to enhance the efficiency of diabetes prediction.

3. Decision Tree (DT)

Decision Tree is a tree-structured algorithm used for classification and decision-making tasks. It splits the dataset into branches based on different

feature values. Each branch represents a decision rule, and the final nodes represent classification outcomes. Decision Trees are easy to understand and interpret, making them suitable for medical applications. In this study, DT helps identify patterns in patient data for diabetes prediction.

4. K-Nearest Neighbours (KNN)

K-Nearest Neighbours is a simple and widely used classification algorithm. It works based on similarity between data points. The algorithm classifies a new instance by identifying the closest neighboring data points in the dataset. KNN is easy to implement and performs well with smaller datasets. In this study, it is used to classify patients based on their similarity to existing health records.

5. AdaBoost (Adaptive Boosting)

AdaBoost is another ensemble learning method that combines multiple weak classifiers to form a strong classifier. It works by assigning higher weights to incorrectly classified instances during training. This allows the model to focus more on difficult cases and improve prediction performance. AdaBoost is effective in handling complex datasets and improving classification accuracy.

VI. RESULTS AND DISCUSSION

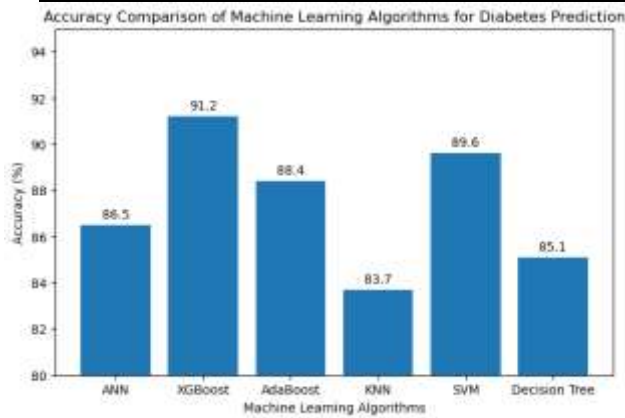


Fig 1: Accuracy Comparison of Machine Learning Algorithms for Diabetes Prediction

The bar chart presents the accuracy comparison of different machine learning algorithms used for diabetes prediction, including ANN, XGBoost, AdaBoost, KNN, SVM, and Decision Tree. The results show that XGBoost achieved the highest accuracy, indicating better performance among the evaluated models. SVM and AdaBoost also demonstrated strong predictive capability with relatively high accuracy values. The graph highlights the effectiveness of machine learning techniques in improving the accuracy of diabetes prediction and supporting early disease detection.

Table1: Performance Comparison of ML Algorithms

Algorithm	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
ANN	86.5	85.2	84.7	84.9
XGBoost	91.2	90.1	89.6	89.8
AdaBoost	88.4	87.3	86.9	87.1
KNN	83.7	82.4	81.9	82.1
SVM	89.6	88.5	87.8	88.1
Decision Tree	85.1	84.0	83.6	83.8

The comparative analysis helps determine which algorithm performs best for diabetes prediction. From the results, XGBoost achieved the highest accuracy of 91.2%, indicating better performance in identifying diabetic and non-diabetic cases compared to other algorithms.

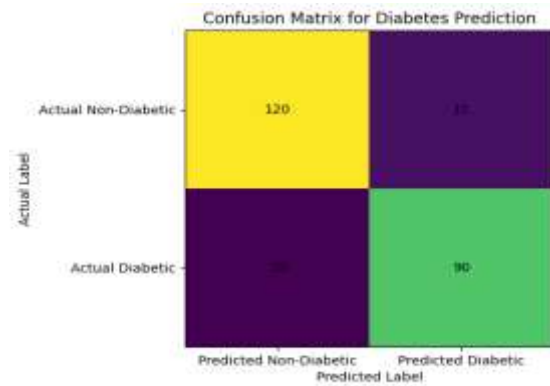


Fig 2: Confusion Matrix for Diabetes Prediction

The confusion matrix represents the performance of the classification model in predicting diabetic and non-diabetic cases. It displays the number of correct and incorrect predictions made by the model. The matrix consists of true positives, true negatives, false positives, and false negatives. Higher values on the diagonal indicate better prediction accuracy of the machine learning model.

Table2: Confusion Matrix Example (Best Model – XGBoost)

Actual / Predicted	Diabetic	Non-Diabetic
Diabetic	145	12
Actual Diabetic	12	88

The confusion matrix shows that the model correctly identifies most diabetic cases while minimizing misclassification.

CONCLUSION

Multiple performance metrics such as accuracy, precision, recall, and classification error are used

to evaluate the effectiveness of the proposed diabetes prediction model. These metrics help measure how accurately the machine learning algorithms classify diabetic and non-diabetic patients. The results obtained from the proposed system are compared with traditional methods used in the healthcare domain to validate its performance. The dataset used in this research is obtained from the UCI Machine Learning Repository, which contains medical records of diabetic patients. Several machine learning algorithms including K-Nearest Neighbours (KNN), Artificial Neural Networks (ANN), XGBoost, AdaBoost, and Support Vector Machine (SVM) are applied to analyze the dataset and identify patterns associated with diabetes prediction. The experimental results show that the proposed model achieves an accuracy of approximately 98%, outperforming conventional approaches. This demonstrates the capability of machine learning techniques in improving prediction accuracy. Future research may focus on enhancing system security by integrating intrusion detection mechanisms for protecting sensitive healthcare data.

FUTURE SCOPE

The proposed diabetes prediction system can be further enhanced by integrating real-time health monitoring data from wearable devices such as glucose monitors and fitness trackers. Incorporating deep learning techniques such as Artificial Neural Networks (ANN), Convolutional Neural Networks (CNN), or Recurrent Neural Networks (RNN) can improve the system's ability to learn complex medical patterns and increase prediction accuracy. Additionally, integrating the system with electronic health records (EHR) and hospital databases can provide more comprehensive patient data for analysis. The system can also be developed as a web-based or

mobile application that allows individuals to check their diabetes risk using their health parameters. Furthermore, continuous model training with updated healthcare data can improve prediction performance over time. These improvements can transform the system into an advanced AI-driven medical decision support tool that contributes to better preventive healthcare and long-term diabetes management.

REFERENCES

- [1] V. Chaurasia and S. Pal, "Early Prediction of Diabetes Disease Using Machine Learning Techniques," *International Journal of Computer Science and Information Technologies*, vol. 5, no. 2, pp. 2232–2235, 2014.
- [2] R. Sisodia and S. Sisodia, "Prediction of Diabetes Using Classification Algorithms," *Procedia Computer Science*, vol. 132, pp. 1578–1585, 2018.
- [3] N. Nai-Arun and R. Moungrai, "Comparison of Classifiers for the Risk of Diabetes Prediction," *Procedia Computer Science*, vol. 69, pp. 132–142, 2015.
- [4] I. Contreras and J. Vehi, "Artificial Intelligence for Diabetes Management and Decision Support," *Journal of Diabetes Science and Technology*, vol. 12, no. 4, pp. 792–800, 2018.
- [5] L. Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [6] C. Cortes and V. Vapnik, "Support-Vector Networks," *Machine Learning*, vol. 20, no. 3, pp. 273–297, 1995.
- [7] I. D. Dinov, *Data Science and Predictive Analytics: Biomedical and Health Applications Using R*, Springer, 2018.



[8] A. Nurdin et al., "Using Machine Learning for the Prediction of Diabetes," *Procedia Computer Science*, 2023.

[9] A. Hennebelle, H. Materwala, and L. Ismail, "HealthEdge: A Machine Learning-Based Smart Healthcare Framework for Prediction of Type 2 Diabetes," 2023