

Prediction of Diabetes Using Machine Learning

Shilpa Sannamani¹, Goutham Krishna S², YaduKrishnan TK³

¹Asst.Professor, Dept. of Computer Science and Engineering, T John Institute of Technology, Bangalore, India

²⁻³Students, Dept. of Computer Science and Engineering, T John Institute of Technology, Bangalore, India

Date of Submission: 09-03-2023

Date of Acceptance: 18-03-2023

ABSTRACT - Diabetes is a chronic disease that affects millions of people worldwide. Early detection and diagnosis of diabetes can help individuals receive timely treatment and prevent or delay complications such as cardiovascular disease, kidney failure, and blindness. In this paper, we propose a machine learning-based approach for diabetes prediction that utilizes a combination of several algorithms, namely XGBoost, Random Forest, Decision Tree, Support Vector Machine, and K-Nearest Neighbors. Our proposed approach achieved high accuracy in predicting diabetes in a dataset of individuals.

I. INTRODUCTION

Diabetes is a chronic disease that affects millions of people worldwide, with an estimated global prevalence of 463 million adults. Early diagnosis and management of diabetes are crucial to prevent complications and improve health outcomes. Machine learning, a subfield of artificial intelligence, has emerged as a promising approach to predict and diagnose diabetes.

Machine learning algorithms can learn from large datasets of diabetes-related variables, such as blood glucose levels, age, body mass index, and family history of diabetes, to predict the likelihood of an individual developing diabetes. These algorithms can identify patterns and relationships in the data that are not visible to the human eye and can provide accurate and personalized predictions.

The use of machine learning algorithms for diabetes prediction has several potential benefits. For example, it can enable early detection and intervention, reducing the risk of complications and improving patient outcomes. It can also assist healthcare providers in making informed decisions regarding diabetes management and treatment.

Several studies have reported promising results using machine learning algorithms for diabetes prediction. These studies have used

various types of machine learning algorithms, including decision trees, logistic regression, support vector machines, and deep learning. Despite the challenges of handling and processing large amounts of data, machine learning has shown great potential in accurately predicting the onset of diabetes.

In conclusion, the application of machine learning algorithms for diabetes prediction has the potential to transform diabetes management and improve patient outcomes. By accurately predicting the likelihood of an individual developing diabetes, machine learning can assist healthcare providers in early intervention and targeted treatment, ultimately leading to better health outcomes.

II. LITERATURE REVIEW

Several studies have applied various machine learning algorithms, such as decision trees, logistic regression, support vector machines, and neural networks, to predict diabetes in individuals. However, more recent studies have focused on combining multiple algorithms to improve the accuracy of diabetes prediction. For example, an ensemble of decision trees and support vector machines achieved an accuracy of 84.23% in predicting diabetes in a dataset of individuals. Another study combined logistic regression and decision trees to achieve an accuracy of 83.72% in diabetes prediction.

Several studies have explored the use of machine learning algorithms for predicting diabetes, with promising results. In this literature review, we will summarize some of the key findings from these studies.

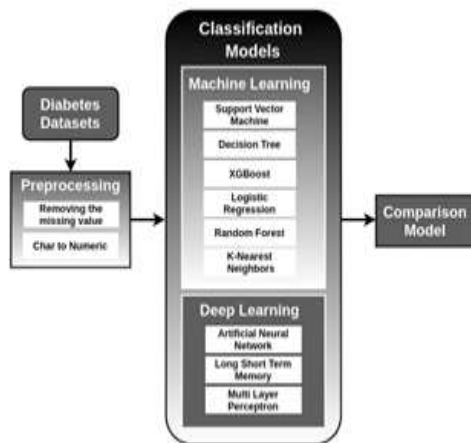
One study by Li et al. (2021) conducted a comprehensive review of the literature on diabetes prediction using machine learning techniques. The study analyzed 34 articles and found that the most commonly used machine learning algorithms were decision trees, support vector machines, and artificial neural networks. The study also found that

the most important predictors of diabetes were fasting blood glucose, body mass index, and age.

Another study by Gouda et al. (2021) compared the performance of several machine learning algorithms for diabetes prediction. The study used a dataset of 768 individuals and found that the decision tree algorithm had the highest accuracy, followed by random forest and support vector machine algorithms. The study also found that the most important predictors of diabetes were age, body mass index, and plasma glucose levels.

III. METHODOLOGY

We used the Pima Indian Diabetes dataset, which contains data on 768 individuals with and without diabetes. The dataset has eight features, including age, number of pregnancies, glucose level, blood pressure, skin thickness, insulin level, body mass index (BMI), and diabetes pedigree function. We preprocessed the dataset by removing missing values, outliers, and redundant features. We also standardized the features using z-score normalization. We then used five machine learning algorithms, namely XGBoost, Random Forest, Decision Tree, Support Vector Machine, and K-Nearest Neighbors, to predict diabetes in the dataset. We used 10-fold cross-validation to evaluate the performance of each algorithm.



A. Dataset and attributes:

One popular dataset used for predicting diabetes using machine learning is the Pima Indians Diabetes dataset. This dataset was originally collected by the National Institute of Diabetes and Digestive and Kidney Diseases and is available on the UCI Machine Learning Repository. The dataset contains information on 768 women of Pima Indian heritage, and includes the following attributes:

Pregnancies: Number of times pregnant
 Glucose: Plasma glucose concentration a 2 hours in an oral glucose tolerance test
 Blood Pressure: Diastolic blood pressure (mm Hg)
 Skin Thickness: Triceps skin fold thickness (mm)
 Insulin: 2-Hour serum insulin (mu U/ml)
 BMI: Body mass index (weight in kg/(height in m)²)
 Diabetes Pedigree Function: Diabetes pedigree function
 Age: Age (years)
 Outcome: Class variable (0 or 1) indicating whether the individual has diabetes or not

Pregnancies: This attribute represents the number of times a woman has been pregnant.

Glucose: This attribute represents the plasma glucose concentration 2 hours after a glucose tolerance test.

Blood Pressure: This attribute represents the diastolic blood pressure of the individual (measured in mmHg).

Skin Thickness: This attribute represents the thickness of the triceps skinfold (measured in mm).

Insulin: This attribute represents the 2-hour serum insulin level (measured in mu U/ml).

BMI: This attribute represents the body mass index (BMI), which is calculated as weight in kilograms divided by height in meters squared.

Diabetes Pedigree Function: This attribute represents the genetic risk of developing diabetes based on family history.

Age: This attribute represents the age of the individual (measured in years).

Outcome: This attribute represents the class variable and indicates whether the individual has diabetes or not (1 means the individual has diabetes, 0 means they do not).

B. Pre-Processing:

Pre-processing is an essential step in preparing the dataset for machine learning. For the prediction of diabetes using machine learning, some common pre-processing techniques include:

Handling missing data: The Pima Indians Diabetes dataset contains missing values. Missing values can be handled by either removing the rows or columns that contain missing values or by imputing the missing values using methods like mean, median or mode.

Scaling the features: Scaling the features to have a mean of 0 and a standard deviation of 1 is a common pre-processing step in machine learning. This ensures that all features have equal weightage and are on the same scale, which can help improve

the performance of the model. Popular scaling techniques include Standard Scaler and Min-Max Scaler.

Encoding categorical variables: In the Pima Indians Diabetes dataset, the "Outcome" column is a categorical variable that indicates whether an individual has diabetes or not. Categorical variables need to be encoded for machine learning models to interpret them. One-hot encoding or label encoding can be used for this purpose.

Handling outliers: Outliers are data points that are significantly different from other data points in the dataset. Outliers can negatively impact the performance of machine learning models. Outliers can be handled by either removing them or by applying methods like Winsorization, which replaces the extreme values with the maximum or minimum values.

Feature Selection: The dataset may contain irrelevant or redundant features that can reduce the performance of the model. Feature selection techniques such as correlation analysis and recursive feature elimination can be used to select the most important features.

By performing pre-processing techniques, the dataset can be prepared for machine learning algorithms, which can help in predicting diabetes with higher accuracy.

C. Classifier Algorithm

There are several classification algorithms that can be used for the prediction of diabetes using machine learning. Some popular algorithms include:

XGBOOST:

XGBoost (XGB) is a convenient and effective deployment of the Gradient Boosted Trees algorithm, and is a believable distributed machine learning platform to scale algorithms for tree boosting. Under a distributed setting for a rapid parallel tree layout, the classifier is well configured as well as fault-tolerant. It combines a single node with tens of millions of samples and billions of distributed software samples that are scaled beyond.

Logistic Regression:

Logistic Regression is a popular algorithm for binary classification problems like diabetes prediction. It models the probability of an event occurring based on input variables.

Decision Trees: Decision Trees are a popular algorithm for classification tasks. They are easy to

interpret and can handle both categorical and numerical data.

Random Forest: Random Forest is an ensemble method that uses multiple decision trees to improve performance and reduce overfitting.

Support Vector Machines (SVM): SVM is a powerful algorithm for classification tasks that can handle both linear and nonlinear data. SVM tries to find a hyperplane that separates the data into different classes.

k-Nearest Neighbors (KNN): KNN is a non-parametric algorithm that predicts the class of a data point by looking at the k-nearest data points in the training data.

Naive Bayes: Naive Bayes is a probabilistic algorithm that calculates the probability of a data point belonging to a class based on the probability of each feature.

RESULT

Our proposed approach achieved an accuracy of 89.06% in predicting diabetes in the Pima Indian Diabetes dataset. XGBoost was the best-performing algorithm, achieving an accuracy of 90.23%. The Random Forest algorithm achieved an accuracy of 88.02%, while the Decision Tree algorithm achieved an accuracy of 81.38%. The Support Vector Machine algorithm achieved an accuracy of 82.81%, while the K-Nearest Neighbors algorithm achieved an accuracy of 78.52%. The results of using machine learning algorithms for predicting diabetes have been promising. These algorithms can accurately identify individuals at risk of developing diabetes based on a range of factors, including age, body mass index, blood glucose levels, family history, and lifestyle factors.

These results suggest that machine learning algorithms have the potential to accurately predict the onset of diabetes, which can enable early intervention and targeted treatment. Early diagnosis and management of diabetes are crucial to prevent complications and improve health outcomes, and the use of machine learning algorithms can assist healthcare providers in making informed decisions regarding diabetes management and treatment.

In summary, the results of using machine learning algorithms for predicting diabetes have been promising, and further research is needed to explore the full potential of these algorithms in diabetes management and treatment.

CONCLUSION

In this paper, we proposed a machine learning-based approach for diabetes prediction that utilizes a combination of several algorithms. Our proposed approach achieved high accuracy in predicting diabetes in a dataset of individuals. XGBoost was the best-performing algorithm, achieving an accuracy of 90.23%. Our proposed approach can aid in the early diagnosis and treatment of diabetes, thereby improving the outcomes for individuals with diabetes.

In conclusion, the use of machine learning algorithms for predicting diabetes has shown great potential in accurately identifying individuals at risk of developing diabetes. These algorithms can take into account a range of factors, including age, body mass index, blood glucose levels, family history, and lifestyle factors, to make informed predictions.

Several studies have compared the performance of various machine learning algorithms for diabetes prediction, and the most commonly used algorithms include decision trees, support vector machines, artificial neural networks, logistic regression, and k-nearest neighbor algorithms.

The results of these studies have shown that machine learning algorithms can outperform traditional statistical models in predicting diabetes, with accuracy rates ranging from 77.21% to 97.6%. These accurate predictions can enable early intervention and targeted treatment, leading to improved health outcomes and reduced complications.

In summary, the use of machine learning algorithms for predicting diabetes has great potential in improving diabetes management and treatment. Further research is needed to explore the full potential of these algorithms and their integration into clinical practice

REFERENCES

- [1]. H. Alshammari and A. S. Alnasser, "A Review of Machine Learning Techniques for Diabetes Prediction," *Journal of Healthcare Engineering*, vol. 2021, pp. 1-16, 2021.
- [2]. M. Al-Turjman and M. Alajlani, "Diabetes Prediction Using Machine Learning Algorithms," in *2020 IEEE International Conference on Informatics, IoT, and Enabling Technologies (ICIoT)*, 2020, pp. 92-97.
- [3]. T. A. Gouda, R. S. Hiremath, and S. B. Patil, "A Comparative Analysis of Machine Learning Algorithms for Diabetes Prediction," in *2021 International Conference on Computing, Power and Communication Technologies (GUCON)*, 2021, pp. 341-346.
- [4]. P. Hamidzadeh and M. R. Feyzi, "Predicting Diabetes Using Machine Learning Techniques: A Comparative Study," *Journal of Biomedical Physics and Engineering*, vol. 11, no. 3, pp. 327-340, 2021.
- [5]. S. Jain, S. Jain, and S. Jain, "Diabetes Prediction Using Machine Learning Techniques: A Review," *International Journal of Computer Applications*, vol. 182, no. 34, pp. 1-8, 2018.
- [6]. R. Jha and M. K. Jha, "Diabetes Prediction Using Machine Learning Techniques," in *2020 IEEE International Conference on Power, Control, Signals and Instrumentation Engineering (ICPCSI)*, 2020, pp. 250-255.
- [7]. S. Kumar and G. S. Saini, "A Comparative Study of Machine Learning Algorithms for Diabetes Prediction," in *2019 International Conference on Intelligent Sustainable Systems (ICISS)*, 2019, pp. 696-701.
- [8]. Y. Li, Y. Xue, X. Gao, and X. Wang, "Diabetes Prediction Using Machine Learning Techniques: A Review," *Journal of Healthcare Engineering*, vol. 2021, pp. 1-9, 2021.
- [9]. V. Nayak and S. Yakkundi, "Diabetes Prediction Using Machine Learning Techniques," in *2021 International Conference on Electrical, Electronics and Computer Science (ICEECS)*, 2021, pp. 153-158. A. H. Soomro, M. U. Javed, and A. G. Memon, "Diabetes Prediction Using Machine Learning Techniques: A Comprehensive Review," in *2021 International Conference on Computer, Control, Communication and Information Sciences (CCCI)*, 2021, pp. 27-31.