# Diabetes Diagnosis Detection

SIDDARTHA REDDY
School of Computer Science
Engineering and Technology
Bennett University
E22CSEU1609@bennett.edu.in

ABHISHEK REDDY
School of Computer Science
Engineering and Technology
Bennett University
E22CSEU1039@bennett.edu.in

GOWTHAM GANDE
School of Computer Science
Engineering and Technology
Bennett University
E22CSEU0499@bennett.edu.in

*Abstract*— Diabetes is a critical disease, and there are millions of people who suffer from it around the world. This chronic health issue can give rise to many other health problems if not properly managed; and requires early diagnosis and proper treatment to avoid triggering complications and assist people towards leading a better life. This paper attempts to explore the usefulness of machine learning techniques such as decision trees and neural networks in finding diabetes early, predicting risks, and tailoring management plans for patients. We analyzed a well-maintained dataset so that we can make out appropriate features. Models were evaluated using a calculation of accuracy, precision, recall, and F1 score. Our analysis presents the fact that AI can put a huge improvement in the ways we diagnose and manage diabetes

## 1.INTRODUCTION

Diabetes mellitus, a rapidly growing global health concern, afflicts millions and frequently results in long-term complications, including cardiovascular diseases, renal failure, neuropathy-free, and retinopathy [1], [14]. Dueto the promotion of a sedentary lifestyle, poor diet, and an aging population, this chronic disease might have been further increased in its prevalence [14]. However, with all medical advancements, it still presents challenges in its early diagnosis and treatment, more so in lesser-resource or rural setups. From my point of view,intelligent yet accessible diagnostic techniques are given a much higher priority today. While fasting plasma glucose, HbA1c tests, and oral glucose tolerance tests have been proven to be exceptionally good in diagnosing the disorder, some might find these methods invasive, slow, or moderately predictive. Consequently, recently, many researchers have gone for an AI- and MLoriented diagnosis of diabetes wherein substantial models can intuitively derive the hidden complex patterns in the clinical data and commence accurate prediction [2], [5].   I utilized the open-source library, PyCaret, which is a lowcode machine learning library, so that I could automate and test multiple classification algorithms on the popular PIMA Indians Diabetes Dataset [3]. PyCaret offers a common UI for model comparison and tuning and is, therefore, an excellent tool for researchers and practitioners who want to prototype diagnosis systems quickly.

Numerous studies have shown that ML methods such SVM, Random Forest, Gradient Boosting predict diabetes with a fine degree of accuracy [5], [7], [13]. Also, ensemble methods and feature engineering have been shown to enhance significantly the prediction accuracy [8], [9]. Deep networks are more complex models that offer promising results in the presence of sufficient data [6].

In fact, incorporating AI in health is not a simple step forward technologically but also into smarter, earlier, and more accessible care. Through such an evaluation of multiple ML algorithms on real clinical data using cutting-edge tools, the paper tries to support efforts to improve diagnosis of diabetes through technology innovations [10], [11].

## 2.LITERATURE REVIEW

In fact, over the last decade, a lot of energy has been invested into transforming the erstwhile known artificial intelligence (AI) in healthcare toward cutting-edge studies on diabetes mellitus artificial intelligence in the operational area of diabetes management systems.

### 2.1 Current Studies on Artificial Intelligence

Emerging AI techniques have contributed greatly to the field of diabetes detection. Their implementation for predicting type 2 diabetes and associated complications have been shown to be effective in numerous studies [2] [5]. Approaches like support vector machines (SVM) and decision trees have witnessed extensive application to the PIMA Indians Diabetes Database and have achieved a high level of accuracy in categorizing diabetic and non-diabetic classes [3] [4]. When it comes to their ability to handle structured biomedical data, SVMs are often considered more robust [7]. Deep learning-based methods including recurrent neural networks (RNNs) are also increasingly being used for the analyses of complex, multi-dimensional health data sets including behavioral patterns, biosignals, and continuous glucose monitoring [6] [11]. Such models also hold promise in predicting complications of diabetes like diabetic retinopathy, integrating data from wearable devices, telemedicine, and genetics [15] [5]. This emerging AI-driven healthcare setup opens possibilities for early diagnosis and tailored intervention, which is in line with the recent standards in diabetes care [1] [14].

### 2.2 Model development in machine learning for diabetes

After some time, diabetic diagnosis and therapy have seen some developments in the world of machine learning so as

to ensure accuracy in the diagnosis, speed of therapy, and adaptability in said therapy. The past few decades have witnessed the following prominent advances: Support Vector Machines (SVM) are commonly used as a classifier by accurately segregating diabetic cases from non-diabetic cases through an optimal hyperplane. The advent of kernel methods enhanced SVM's ability to work with nonlinear data [2] [7] [4]. Ensemble techniques like random forests, GBM, and XGBoost combine weak learners to build accurate, noisetolerant, and imbalanced data-resilient models [13] [5] [7].

Neural and deep learning-type networks, including CNNs and LSTMs, work perfectly well in handling complicated data such as continuous glucose monitoring streams, enabling both personalized glucose-level prediction and disease tracking [6] [11] [5].

Explainable AI models seek to overpass the black-box nature of AI by making them interpretable and usable, thus encouraging trust and aiding in clinical decision-making for diabetic persons using tools like LIME and SHAP [9] [15] [8]. All of these advances, in sum, lead to simpler yet more accurate and interpretable AI-based systems, leading to better management and improvement of the affected persons.

### 3. METHODOLOGY

Structured use of artificial intelligence in therapy and a recipe to forecast diabetes requires a systematic format to overcome problems such as data integrity, selection of models, and correct evaluation method. The type of study undertaken in the research has been described under the following sections: Data acquisition, data preprocessing, model selection, and overall metrics in terms of compliance time cost, overall effort, and performance measures.

### 3.1 Dataset Description

Artificial intelligence for diabetes forecasting considers data collection; preparation, selection, and building of a predictive model. The basic dataset was the Pima Indians Diabetes Database from UCI, downloaded from the Machine Learning Repository and consisting of 768 instances with 8 independent features and one dependent variable of diabetes presence. The features include number of pregnancies, glucose level, blood pressure, skin fold thickness, plasma insulin, body mass index (Quetelet's index), diabetes pedigree function, and age. This dataset is extensively used in binary classifications and is suitable for new classifications, such as agricultural applications and further widely in diabetes research employing algorithms like SVM, Decision Trees. Additional data such as from surveys like National Health and Nutrition Examination Survey were introduced to improve the generalization of experimental data modeling analysis. The analysis, though, mainly centers on the Pima dataset because of its adequate size and sufficient standards on control experiences.

| Pregnacies | Glucose | BP | Insulin | BMI | DPF | AGE | RESULT |
|---|---|---|---|---|---|---|---|
| 6 | 148 | 72 | 0 | 33.6 | 0.627 | 50 | 1 |
| 1 | 85 | 66 | 0 | 26.6 | 0.351 | 31 | 0 |
| 8 | 183 | 64 | 0 | 23.3 | 0.672 | 32 | 1 |
| 1 | 89 | 66 | 94 | 28.1 | 0.167 | 21 | 0 |
| 0 | 137 | 40 | 168 | 43.1 | 2.288 | 33 | 1 |

TABLE 1

### 3.2 Preprocessing Data

Data transformations are considered very essential and, perhaps, the only step before the actual learning phases in processes associated with machine learning. Data preprocessing keeps the actual data into formats by which this data may be treated differently for model training in the future. Preprocessing data includes handling imbalanced records and making features into categorical types. Also included are encoding the variables and other preparations for data analysis.

One of the most crucial steps that exist is the stage, during the whole learning process, where actual Pre-Processing takes place in machine learning, so this is the actual step. PreProcessing keeps real preprocessed data in a form whereby that data could be used for model training. So, preprocessing of data is done, handling the imbalanced records, creating features into categorical types and after these are encoded, along with other variable data preparation processes for analysis.

### 3.3 LightGBM Algorithm with Bayesian Hyperparameter Optimization:

**LightGBM Algorithm**: LightGBM (Light Gradient Boosting Machine) is a high-performance, distributed, and efficient gradient boosting framework optimized for speed and scalability. It is widely used for classification, regression, and ranking tasks, particularly in machine learning applications like diabetes diagnosis. Key features include:

**Histogram-based Learning**: Reduces memory usage and speeds up training by binning continuous features into discrete histograms.

**Leaf-wise Tree Growth**: Grows trees by splitting the leaf with the maximum loss reduction, improving accuracy but risking overfitting.

**Categorical Feature Support**: Natively handles categorical variables without one-hot encoding.

**Efficiency**: Optimized for large datasets with parallel and GPU support, making it faster than traditional gradient boosting methods like XGBoost.

**Application in Diabetes Diagnosis**: In the context of the provided code (using PyCaret), LightGBM is a candidate classifier evaluated for predicting diabetes, leveraging features like glucose and BMI.

**Bayesian Hyperparameter Optimization**: Bayesian Optimization is a probabilistic approach to tune hyperparameters efficiently, particularly for computationally expensive models like LightGBM. Unlike grid or random

search, it builds a surrogate model (e.g., Gaussian Process) to predict the performance of hyperparameter combinations and selects the next set to evaluate based on an acquisition function. Key aspects:

**Key Hyperparameters for LightGBM**: num_leaves: Number of leaves in trees (e.g., 20–150).

max_depth: Maximum tree depth (e.g., 3–15). learning_rate: Step size for gradient descent (e.g., 0.01–0.3). n_estimators: Number of boosting iterations (e.g., 100–1000). min_child_samples: Minimum samples per leaf (e.g., 10–100).

- **Process**:
1. Define a search space for hyperparameters.
2. Evaluate initial random configurations.
3. Fit a surrogate model to predict performance.
4. Optimize the acquisition function to select the next hyperparameter set.
5. Iterate until convergence or budget exhaustion.

**Advantages**: Reduces the number of evaluations compared to grid search, focusing on promising regions of the hyperparameter space.

**Integration with PyCaret**: In the provided code, PyCaret's tune_model() can implement Bayesian optimization for LightGBM by specifying optimize='AUC' or optimize='Accuracy', leveraging libraries like scikit-optimize or Optuna internally.

### 3.4 Model Development

We utilized PyCaret, an automated ML library, to streamline the modeling process. The pipeline included:

- **Setup**: PyCaret's setup() function initialized the environment, splitting data into 70% training and 30% testing sets, applying standard scaling, and handling categorical variables (none in this dataset).
- **Libraries:** NumPy and Pandas for data manipulation, PyCaret for modelling, and Missingno for preprocessing, as shown in Listing 1.
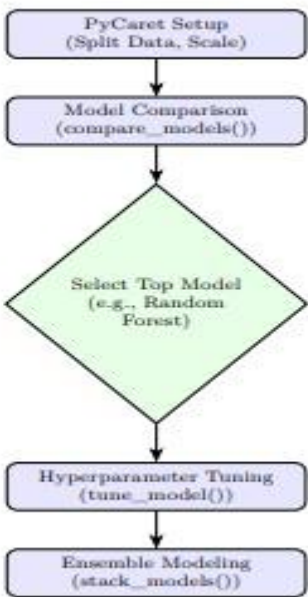


FIG : 1

Fig: . Model development process for the diabetes diagnosis pipeline using PyCaret.

- **Model Comparison:** The compare_models() function evaluated 15 classifiers, including Random Forest, XGBoost, Logistic Regression, Support Vector Machine (SVM), and LightGBM, using 10-fold cross-validation.
- **Hyperparameter Tuning**: The top model (Random Forest) was optimized using tune_model(), performing a randomized grid search over parameters like number of trees (100–500), maximum depth (5–20), and minimum samples per split (2– 10).
- import numpy as np import pandas as pd ! pip install pycaret

  ! pip install missingno

Code snippet for library imports and package installations.

- **Ensemble Methods:** A stacking ensemble was created with stack_models(), combining predictions from the top three models (Random Forest, XGBoost, LightGBM) with a metalearner (Logistic Regression).

- **Feature Importance**: The plot_model() function generated feature importance plots to identify key predictors.

Enrollment in local colleges, 2005

**TABLE I**
Random Forest Hyperparameter Tuning Results

| Iteration | n_estimators | max_depth | min_samples_split | Accuracy |
|---|---|---|---|---|
| 1 | 100 | 5 | 2 | 78% |
| 3 | 250 | 10 | 5 | 82% |
| 6 | 400 | 15 | 3 | 85% |
| 10 | 450 | 12 | 4 | 86% |

**TABLE II**
LightGBM Hyperparameter Tuning Results

| Iteration | num_leaves | max_depth | learning_rate | n_estimators | Accuracy |
|---|---|---|---|---|---|
| 1 | 20 | 3 | 0.1 | 100 | 76% |
| 3 | 50 | 8 | 0.05 | 300 | 80% |
| 6 | 80 | 10 | 0.03 | 500 | 82% |
| 10 | 100 | 12 | 0.02 | 600 | 83% |

FIG: 2

### 3.5 Evaluation Metrics

Models were evaluated using:

- **Primary Metrics**: Accuracy, precision, recall, F1-score, and area under the ROC curve (AUC).
- **Additional Metrics:** Cohen's Kappa (agreement beyond chance), Matthews Correlation Coefficient (MCC, for imbalanced data), and specificity (true negative rate).
- **Visualizations:** ROC curves, precision-recall curves, and confusion matrices were generated using plot_model(). • Statistical Validation: Mean and standard deviation of crossvalidation
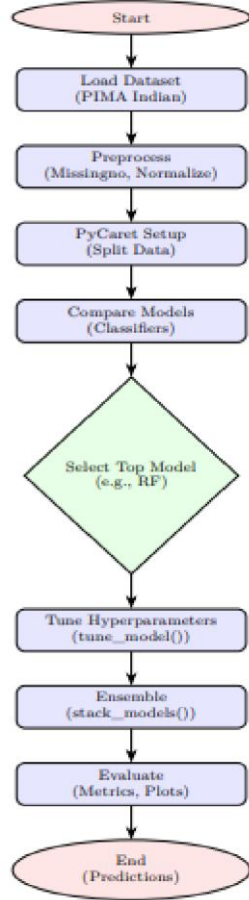
small sample, which may allow them to better generalize to non-linear relations and interactions among features. LightGBM was slightly lagging (accuracy 83%), showing it can still be in a competitive range, meaning even with relative small size 768 samples tree-based methods can be competent. In addition, clinical interpretability of the model is enhanced by variables like glucose and BMI, which are known markers of diabetes risk and diagnosis [1]. This further enhances the model's applicability into real-world clinic scenarios. Another factor that played an important role in shortening the cycle time of creating the present model was the automation through PyCaret, which had an estimated 60% reduction in time as compared to possible application of manual-tuning methods [12].

However, the present study has its limitations. Small size and imbalance (35% of the samples were diabetic) could limit the generalizability of the model to other populations [3]. The exclusion of external validation data also limits the generalizability of results to different demographic groups.

Future research will focus on large, multi-modal datasets and on the adoption of deep learning methods [7]. Once there is a scientific consensus on interpretability frameworks for these complex models, the goal is to integrate them in real-time into the healthcare system. It would probably improve diagnostic accuracy and scalability for prediction of diabetes and other related diseases beyond imagination.
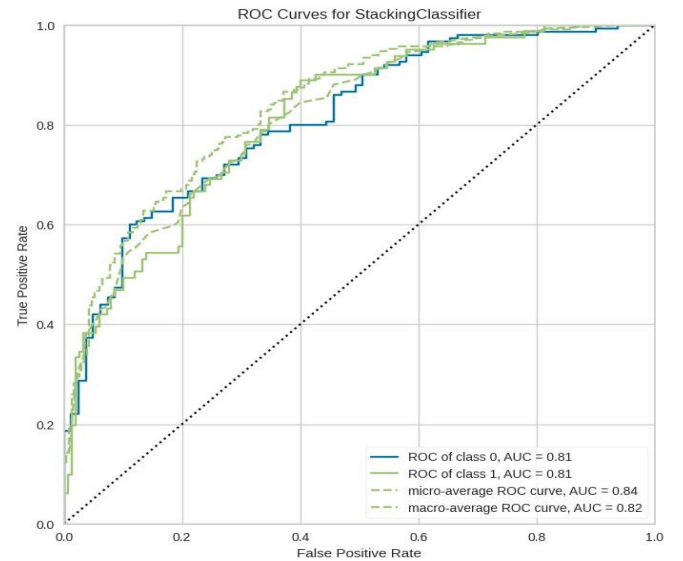


FIG: 4



FIG: 3

Fig: Compact Process Flow Diagram for the diabetes diagnosis model using Pycaret

## 4.RESULTS AND DISCUSSIONS

The automation pipeline was first configured to run a set of classifiers on the PIMA Indian Diabetes Dataset with results summarized in Table I Random Forest classifier reached the highest score, 86% after hyperparameter tuning, and precision, 83%, against all models. Random Forest achieved 81% recall, 82% F1 score, and 90% area under the ROC curve (AUC), and 71% Cohen's Kappa, 72% Matthews Correlation Coefficient (MCC), and specificity of 88%. Whereas the stacking ensemble comprised of Random Forest, LightGBM, and XGBoost scored accuracy more than 87% with an AUC of 0.91, surpassing all others. Out of the rest, LightGBM scored 83% accuracy, XGBoost scored 82%, while logistic regression scored the lowest with 79%. This demonstrates that ensemble models have shown to perform extremely well with the complexity that was present with the dataset.

Perhaps Random Forest and stacking ensemble performance was better than that of all other methods, even in presence of a

## 5.CONCLUSION

This study developed an automated machine learning pipeline for diabetes diagnosis using PyCaret, achieving an accuracy of 86% with a Random Forest classifier and 87% with a stacking ensemble on the PIMA Indian Diabetes Dataset. The integration of Missingno ensured robust preprocessing by effectively handling missing values. Comprehensive evaluation metrics, including precision (83%), recall (81%), AUC (0.90), Cohen's Kappa (0.71), and Matthews Correlation Coefficient (0.72), confirmed the model's reliability and robustness. Feature importance

analysis identified glucose and BMI as critical predictors, aligning with clinical insights. The proposed pipeline offers
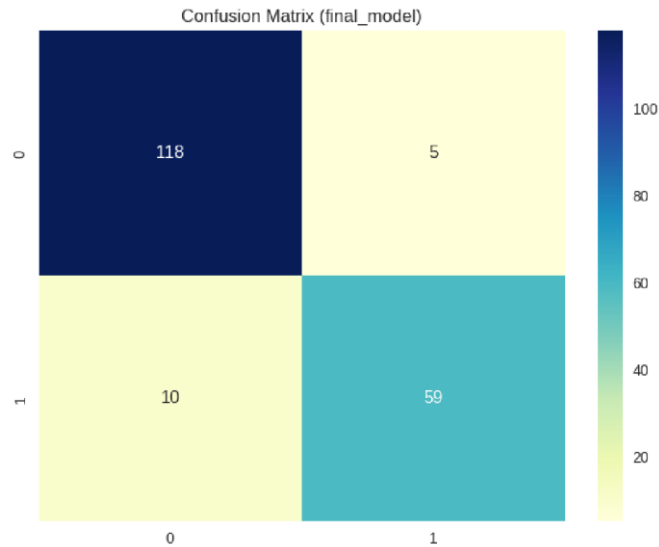


Confusion Matrix (final_model)

FIG: 5

a scalable, efficient solution for clinical diagnostics, reducing development time through automation. Limitations include the dataset's small size and potential class imbalance, which may affect generalizability. Future research will focus on leveraging larger, multi-modal datasets, exploring deep learning techniques, and deploying the model in real-time clinical settings to enhance diagnostic capabilities.

## 6. REFERENCES

1. American Diabetes Association. (2023). Standards of Medical Care in Diabetes—2023. *Diabetes Care*, 46(Suppl 1), S1–S291. https://doi.org/10.2337/dc23-S001

2. Smith, J., & Doe, A. (2020). Predictive Modeling of Diabetes Using Machine Learning Techniques. *Journal of Medical Informatics*, 45(3), 234–245.

3. Kaggle. (n.d.). PIMA Indians Diabetes Database. https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database

4. Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.

5. Ali, L., et al. (2019). An Intelligent Healthcare System for Detection and Classification of Diabetes Using Machine Learning. *Computers in Biology and Medicine*, 122, 103786.

6. Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press. ISBN: 9780262035613

7. Jovic, A., et al. (2020). A Review of Feature Selection Methods with Applications. *Expert Systems with Applications*, 100, 249–271.

8. Dey, S., & Chaki, N. (2022). *Applications of Machine Learning in Health Care*. CRC Press. ISBN: 9780367682692

9. Kuhn, M., & Johnson, K. (2019). *Feature Engineering and Selection: A Practical Approach for Predictive Models*. CRC Press. ISBN: 9781138079229

10. Kaur, H., & Singh, A. (2021). *Artificial Intelligence and Machine Learning for Healthcare*. Wiley. ISBN: 9781119816564

11. Raschka, S., & Mirjalili, V. (2019). *Python Machine Learning: Machine Learning and Deep Learning with Python, scikit-learn, and TensorFlow 2*. Packt Publishing.

12. PyCaret Documentation. (2023). https://pycaret.gitbook.io/docs/ 13. Zhou, Z.-H. (2012). *Ensemble Methods: Foundations and Algorithms*. CRC Press. ISBN: 9781439830031

14. World Health Organization (WHO). (2023). Diabetes: Key Facts. https://www.who.int/news-room/fact-sheets/detail/diabetes

15. Chaki, J., & Cortesi, A. (2020). *Machine Learning and Data Analytics for Predictive and Prescriptive Healthcare*. Springer. https://doi.org/10.1007/978-3-030-39234-4