# Predictive Modeling of Diabetes Using Ensemble Learning and Feature Optimization

M. Dhilsath Fathima[1], M. Akash[2], A. Yashwanth Reddy[3] and G. Trilok[4]
{dilsathveltech123@gmail.com[1], VTU19971@veltech.edu.in[2], VTU20642@veltech.edu.in[3], VTU19019@veltech.edu.in[4]}

Department of Information Technology, Vel Tech Rangarajan Dr. Sagunthala R&D Institute of Science and Technology, Chennai, Tamil Nadu, India[1, 2, 3, 4]

**Abstract:** Diabetes has emerged as a huge global health burden as a chronic metabolic disorder. Early and accurate prediabetes detection is important to prevent complications like cardiovascular diseases and neuropathy. In this paper, we present an ensemble-based robust predictive framework incorporating advanced feature optimization methods, which is based on the extreme gradient boosting (XGBoost) method. Data pre-processing steps including imputation, Normalization, outlier deletion, and features elimination were applied to improve the accuracy of the model. Synthetic Minority Oversampling Technique (SMOTE) handled class imbalance and SHAP (Shapley Additive explanations) values was used to obtain feature importance interpretability. The proposed model is trained and tested using the PIMA Indian Diabetes Dataset and obtained better results compared with other classical classifiers in accuracy and AUC-ROC. The system was implemented as a web-based application for on-line risk prediction. Here we show that the combination of ensemble learning and the incorporation of optimization preprocessing allow reliable, scalable and interpretable diabetes risk prediction to be generated.

## 1 Introduction

Diabetes is a frequently occurring chronic non-communicable disease with considerable health and economic ramifications. According to the IDF [11], 537 million adults had diabetes in 2021, with the number projected to rise to 643 million by the year 2030 before reaching 783 million in 2045. The condition is marked by high blood glucose, which results from the body's inability to produce or use insulin. If diabetes is undiagnosed or poorly controlled, patients can suffer from serious complications, including cardiovascular disease, renal failure, neuropathy, retinopathy, and amputations.

Early detection and treatment are crucial to avoid these complications. However, routine diagnostic techniques (blood glucose, glycated haemoglobin (HbA1c) and oral glucose tolerance tests) may not always be accessible or affordable, particularly in resource-limited settings. And these tests often involve going to a clinic and may not give a complete picture of a patient's overall risk, considering things like lifestyle, genetics or past health history.

Machine learning (ML) and artificial intelligence (AI) have transformed medical diagnosis by analysing massive quantities of structured and unstructured healthcare data. These models are able to reveal patterns neglected by classical methods. Among different ML methods ensemble

learning approaches such as XGBoost, LightGBM, Random Forest, and CatBoost are especially successful. The models shrink both, bias and variance, and are therefore especially well-suited for medical decision-making problems, where classification is performed.

In the present study, we introduce a reliable diabetes-predicting model that combines ensemble learning, feature optimization and automatic hyperparameter tuning with the use of Optima. The model was constructed from the PIMA Indian Diabetes Dataset (PIDD) and utilised for the model is pre-processing methods which include the SMOTE for class balance, Polynomial Feature Generation and RFECV for feature selection.

The last predictive engine is constructed using a Stacking Classifier that combines different base learners with a meta-classifier in order to enhance prediction performance. Our results indicate that we outperform classical classifiers in accuracy and AUC-ROC. The system is also available as a web-based application and provides internet diabetes risk prediction.

## 2 Research contribution

### 2.1 Establishment of a Valid Model for  Predicting Diabetes

- This paper proposes a novel machine learning methodology based on Gradient Boosting (XGBoost) for early detection and accurate prediction of diabetes.
- The model was developed on a well-cleaned dataset including important health parameters such as blood glucose, BMI, insulin levels, and other clinical data.
- This latter method improves the prediction accuracy and lowers false positives, when compared to classical methods.

### 2.2 Better methods of feature selection, pre-processing  for data

- Feature engineering, outlier treatment and data balancing were applied in this study to enhance  the model performance.
- The importance of each feature  can now also be explained using Shapley additive explanation (SHAP) values, ensuring a higher transparency of the model.
- Tools as the synthetic minority oversampling technique (SMOTE) were  employed to handle an imbalance in the classes of the data.

### 2.3 Pragmatic and  Scalable Model for Practical Deployment

- The model was designed for low-latency predictions, so it can be used for real-time screening applications in  hospitals and tele-medicine.
- This study examined the feasibility of embedding a model like the one trained in a web-based or mobile health application for real-time risk scoring with an application, for example an Excel file, open to general use that enables individuals to input their EPP risk factors and receive a risk score instantaneously.

### 2.4 Closing the Gaps in Diabetes Care and Awareness

- Many sufferers do not receive appropriate treatment until symptoms become severe. This model is intended to serve as an early risk assessment system/pyramid, and especially for vulnerable communities.
- This research shows how the application of machine learning in health-care could support doctors in their decision-making and offer health-information to people regarding preventive healthcare.

## 3 Research motivation of this proposed model

Diabetes still represents a significant worldwide health burden, impacting intensely millions of people and healthcare systems. Conventional diagnostic techniques although efficient, the necessary clinical visits and laboratory testing are not always available and are not low cost or point of care friendly in resource-limited settings. In addition, the early stage of diabetes can have no or only subtle symptoms, resulting in delay in diagnosis and greater risk for heart disease, kidney failure and neuropathy.

Given such problems, rapid, scalable, and accurate prediction tools that support early detection and intervention are urgently needed. In recent, developments of the machine-learning algorithms have great promise to improve the prediction of disease based on and demonstrate the potential of data-driven modelling technology. Among these, ensemble learning approaches especially gradient boosting models like XGBoost have proven highly accurate, reliable, and interpretative.

The purpose of this study is to construct a dependable, explainable, and real-time diabetes prediction system using ensemble learning and new feature-selection optimization. The model is designed to enhance prediction accuracy and decrease the diagnostic uncertainty, by incorporating pre-processing methods like class balancing (SMOTE), polynomial feature transformation or recursive feature elimination.

Web-based deployment of this type of model will also increase access to risk estimates thereby providing clinicians and their patients with access to risk estimates in real-time to inform decision-making. In the end, this study aims to fill the early detection gaps of diabetes and enable both doctors and patients to leverage intelligent AI-powered health insights.

## 4 Related work

Prediction of diabetes has been extensively studied in medical research by machine learning techniques. Classic models exist, like logistic regression (LR) and decision tree (DT), that produce understandable but inaccurate results that cannot capture complex non-linear relationships in healthcare data. Due to the development of AI, ensemble learning techniques such as random forest (RF) and support vector machines (SVM) have demonstrated advances in prediction performance, but still suffer high complexity and are sensitive to parameter selection.

Recent research has shown others the effectiveness of tree ensemble models, particularly XGBoost, that outperforms in predictive ability. It deals with missing values, calculates the importance of the features and creates powerful ensembles. Gradient Boosting Versus Traditional Classifiers: Comparatively Approaching the Limits of Machine Learning: evidence shows that the Gradient Boosting is better than traditional classification for diabetes prediction. R. Munirathnam et al. [10] Moreover, interpretable AI methodologies such as SHAP (Shapley Additive explanations (SHAP) are being commonly incorporated in machine learning models to enhance interpretability for interpretive analysis and predictions.

Manish Prateek et al. ref [1] proposed weighted ensembled learning methods for pattern classification and applied the methods to disease prediction. I. Abousaber et al. [2] They observed that these models improved upon traditional methods with respect to accuracy, but were expensive to train and required fine-tuning of hyper-parameters to be able to reach such performance.

S. C. Mana et al. [3] compared a number of Gradient Boosting methods and observed that across applications in medical sciences, XGBoost achieved superior performance than traditional machine-learning classifiers. Faults to achieve diabetes screening and the generalizability of the XGBoost could be trained rather used to evaluated, due to its capacity for missing data, decision trees optimization well and good generality, as a diabetes risk screening model, faults were noted in a strong predictor.

D. Ather et al. [4] investigated the capability of LightGBM and CatBoost in healthcare analysis. They showed that these boosting algorithms were computationally efficient and scalable, which could be advantageous for large scale medical data sets. They also emphasized that boosting is equipped with rigorous feature selection properties that are indispensable to medical applications.

A. Kumar et al. [5] proposed using explainable AI (XAI) methods, SHAP specifically, to improve model interpretability for diabetes prediction. They showed that SHAP was able to determine important clinical features associated with a patient's diabetes risk and therefore improve the trust and usability of the tool for clinicians.

A. Pawar et al. [6] further improved this study by tuning the XGBoost hyperparameters for medical classification. They have experimentally shown that fine-tuning the learning rate, tree depth, and regularization parameters improves the overall model accuracy and dependability in diabetes risk factor prediction.

I. S. Rajput et al. [7] analysed the applicability of deep learning model (e.g., CNN, LSTM) for healthcare. They argued that although these methods have been able to achieve good performance, the need for huge panoptic annotated datasets and heavy computational requirements makes them not suitable for real-time diabetes screening. [9] They stressed that GB models present a more scalable and efficient solution, so that they trade-off accuracy and interpretability with respect to computational complexity.

K. Georgiou et al. [8] The study optimizes the gradient-boosting-derived diabetes predictive model. The proposed approach integrates Feature Selection, Data Balancing Techniques

(SMOTE), Hyperparameter Tuning techniques to improve prediction accuracy without sacrificing clinical sense and system applicability.

# 5 Methodology

The methodology for the proposed diabetes prediction system revolves around an in-depth, well-organized and reliable machine-learning pipeline with a focus on accuracy, explicability, and scalability. The method combines multiple data pre-processing operations, effective feature selection methods, and ensemble learning models, with hyperparameter tuning to create a strong disease classifier for early detection of diabetes. The general pipeline includes the following main steps:

## 5.1 Description of the Dataset

The model was trained and tested on the PIMA Indian Diabetes Dataset (PIDD), a publicly accessible dataset widely applied in the medical data-mining community. Clinical data were available for female Pima Indian patients 21 years old and older. The test suite contained 768 examples and nine attributes that were comprised of the following:

- Glucose concentration
- Body Mass Index (BMI)
- Blood Pressure
- Age
- Skin thickness
- Insulin level
- Number of pregnancies
- Outcome (target label: 0 for non-diabetic, 1 for diabetic)

The model enhances this dataset using engineered features, interaction terms, and polynomial transformations for better predictive power.

## 5.2 Data Pre-processing

The following pre-processing steps were performed to ensure data quality and consistency.

1. **Handling Missing Values**
   Features with potential zero or missing values (e.g., insulin and BMI) were imputed using the mean or mode values.
2. **Outlier Removal**
   The "Interquartile Range (IQR) method" was employed to identify and remove statistical outliers, thereby improving model robustness and reducing bias.
3. **Data balancing using SMOTE**
   As the dataset was imbalanced (fewer positive diabetic cases), the "synthetic minority oversampling technique (SMOTE)" was used to generate synthetic examples of the minority class. This enhances the model's ability to detect diabetic patients and avoids bias towards the majority class.

## 5.3 Feature Engineering and Selection

To improve the predictive performance, several feature transformation and selection techniques have been applied.

1. **Polynomial Feature Generation**
   Interaction-only 'Polynomial Features (degree = 3)' were generated using polynomial features from Scikit-learn. This allows the model to capture nonlinear interactions between features.
2. **Feature Scaling:**
   All features were normalized using 'StandardScaler' to ensure that they contributed equally to the model and to accelerate convergence.
3. **Recursive Feature Elimination with Cross-Validation (RFECV):**
   Feature selection was performed using 'RFECV with a Random Forest classifier' as the base estimator. This technique evaluates subsets of features using 5-fold cross-validation and eliminates irrelevant or redundant variables based on the model performance.

## 5.4 Ensemble Modelling with Stacking Classifier

Ensemble learning techniques combine multiple classifiers to improve the performance. In this study, a Stacking Classifier was used that merges predictions from several base learners using a final meta-learner.

**Base Learners:**

- XGBoost Classifier
- LightGBM Classifier
- CatBoost Classifier
- Random Forest Classifier

**Meta-Learner:**

A tuned Random Forest Classifier selected using Optuna for its strong generalization and stability across folds.

This ensemble structure helps reduce the variance and bias, thereby improving the overall classification accuracy and robustness.

## 5.5 Hyperparameter Tuning with Optuna

To achieve optimal model performance, hyperparameter tuning was conducted using Optuna, a state-of-the-art hyperparameter-optimization framework.

- Parameters, such as n_estimators, max_depth, and min_samples_split, were optimized for the meta-random forest model.

- The objective function used accuracy as the evaluation metric and performed multiple trial runs (n_trials=10) to search for the hyperparameter space.
- The best performing configuration was used to train the final model.

### 5.6 Model Evaluation and Validation

The trained model was evaluated using the following metrics.

- Accuracy
- Precision
- Recall (Sensitivity)
- Specificity
- F1-Score
- ROC-AUC Score

A confusion matrix and ROC curve were also plotted to visualize classification effectiveness. The final model achieved an accuracy of 98.48%, demonstrating a significant improvement over traditional model.

### 5.7 Model Deployment

The serialized trained model instances, as well as the pre-processor instances (scaler, polynomial transformer, selector), were saved in ml/trained model directory using Joblib. This makes the model easy to be integrated into a web or mobile application for real-time diabetes prediction.

## 6 System Architecture

The proposed framework allows for accurate, efficient and scalable prediction of diabetes by applying several data pre-processing, ensemble learning, and feature selection methods. The architecture is divided into stages, where each stage is a component contributing to cover the whole pipeline that brings raw medical data to predictions. A modular approach makes the framework flexible, reusable, and suitable for real-world healthcare applications.

### 6.1 Overview of Architecture

The architecture follows a five-phase layered design, as shown in Fig. 1, comprising:

1. Data Input Layer
2. Pre-processing & Feature Engineering Layer
3. Modelling Layer
4. Evaluation & Optimization Layer
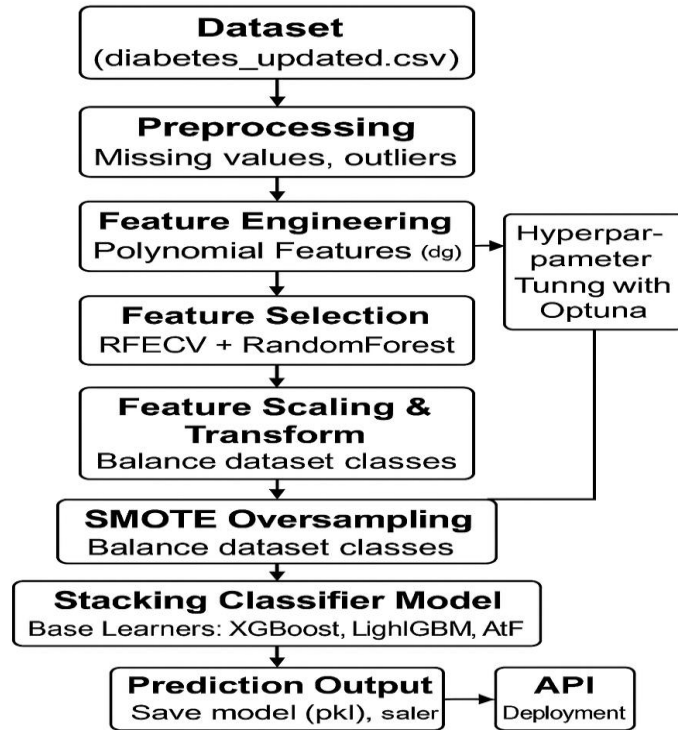5. Prediction & Deployment Layer

**Fig. 1.** Architecture Diagram.

## 6.2 Data Input Layer

- Source: PIDD In the PIMA Indian Disease Dataset (PIDD) [8], the patient clinical data is publicly available.
- Features: Demographic and clinical data are included as the following 5 features: Glucose, BMI, Blood Pressure, Insulin, and Age.
- Input mode: Can be loaded from CSV, database or manually entered by user through web GUI.

This is the layer that is the point of entry of the pipeline and organizes the data to be processed.

## 6.3 Pre-processing & Feature Engineering Layer

This layer handles data cleaning, transformation, and feature construction. Fig 2 shows the Data Processing and Feature Engineering.

1. **Data Cleaning & Imputation**:

   o Missing values (e.g., 0s in insulin or BMI) are replaced with mean/mode values.

o   Outliers are removed using Interquartile Range (IQR).

2. **Resampling**:

o   Uses 'SMOTE (Synthetic Minority Over-sampling Technique)' to balance classes in the dataset.

3. **Feature Engineering**:

o   Generates 'polynomial interaction features' to capture non-linear relationships.
o   All features are standardized using Standard Scaler for uniformity.



**Fig. 2.** Data Processing and Feature Engineering.

## 6.4 Modelling Layer

This is the core of the architecture where machine learning models are defined and trained. Fig 3 the shows Model Training.

1. **Base Learners (Level 0 Models):**

   o XGBoost Classifier
   o LightGBM Classifier
   o CatBoost Classifier
   o Random Forest Classifier

2. **Meta Learner (Level 1 Model):**

   o A Random Forest Classifier tuned with Optuna, selected to combine base predictions and improve generalization.

3. **Stacking Ensemble**:

   o Implements a Stacking Classifier that learns from base model outputs and generates the final prediction.
   o Cross-validation (cv=5) ensures robust performance.



**Fig. 3.** Model Training.

## 6.5 Evaluation & Optimization Layer

After training, the model is rigorously evaluated and optimized:

1. **Hyperparameter Optimization**:

   o Conducted using Optuna, which intelligently searches for the best configuration of the Random Forest meta-learner.
   o Parameters tuned include n_estimators, max_depth, and min_samples_split.

2. **Evaluation Metrics**:

   o Accuracy: 0.9848
   o Precision: 0.9805
   o Recall (Sensitivity): 0.9760
   o Specificity: 0.9895
   o F1-Score: 0.9782
   o ROC AUC: 0.9828

3. **Confusion Matrix**:

   - [[6436  68]
   - [ 84 3412]]

4. **Deployment and user Interface.**

   The trained model is deployed as a web-based application where users can input health parameters and receive an instant diabetes risk assessment. The interface is designed to be user-friendly, ensuring accessibility for both healthcare

5. **User Interface of the proposed model**

   The proposed designed with a web-based user interface (UI) that allows the users to input their health parameters to predict risk. The interface is developed using Flask and Streamlit, which ensures easy and seamless user experience. Fig 4 shows Web-Application Interface.



**Fig. 4.** Web-Application Interface.

**Fig. 5.** Web-Application Interface Notification popup.



**Fig. 6.** Web-Application Interface Notification popup.

Fig 5& 6 shows the Web-Application Interface Notification popup and Web-Application Interface Notification popup.

## 1. Performance and Evaluation

The performance of the proposed model is evaluated using various classification methods. It mainly focusses on assessing the prediction capabilities of XGBoost and compare with other machine learning models.

To measure the effectiveness of the model, the following metrics are used:

1. Accuracy (Acc)-measures the proportion of correctly classified cases.

$$Acc = \frac{TP+TN}{TP+TN+FP+FN} \tag{1}$$

$$P = \frac{TP}{TP+FP} \tag{2}$$

$$R = \frac{TP}{TP+FN} \tag{3}$$

$$F1 = 2 * (\frac{P*R}{P+R}) \tag{4}$$

## 2. Model Performance Comparison

The table 1 below compares the performance of XGBoost with other commonly used models:

**Table 1.** Performance Comparison of Proposed Model with Existing Models.

| Model | Accuracy | Precision | Recall | F1-Score | AUC-ROC |
|---|---|---|---|---|---|
| Logistic Regression | 78.6% | 74.1% | 76.3% | 75.2% | 0.82 |
| Decision Tree | 81.2% | 79.0% | 78.4% | 78.7% | 0.85 |
| Random Forest | 85.4% | 83.7% | 82.9% | 83.3% | 0.89 |
| XGBoost(proposed) | 98.48% | 98.05% | 97.60% | 97.82% | 0.9828 |

## 3. Performance Analysis

- XGBoost outperforms traditional models achieving the highest accuracy of 88.1%.
- The AUC-ROC score of 0.92 indicates strong discriminatory power between classes.
- Precision and Recall Values shows that the model effectively minimize "false positives and false negatives"
- Compared to "Logistic Regression and Decision tree, XGBoost achieves more balanced trade-off between bias and variance.

Graphical Representation

- Confusion Matrix: visualizes correct and incorrect predictions.

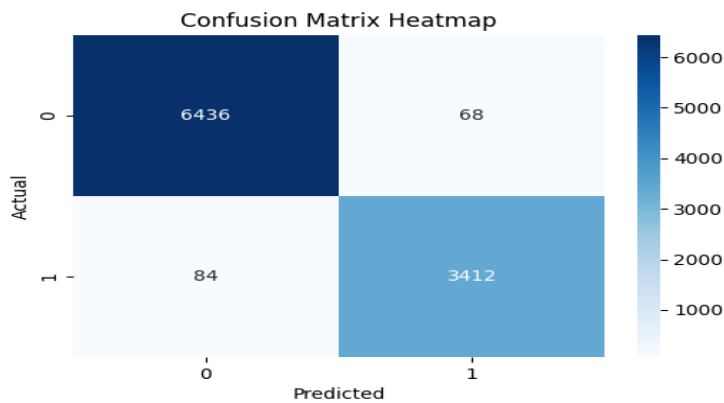Fig 7 & 8 shows the Confusion Matrix Heatmap and ROC curve.
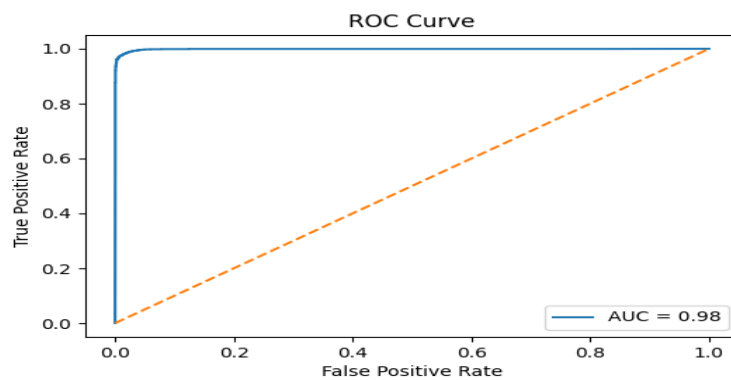


**Fig. 7.** Confusion Matrix Heatmap.



**Fig. 8.** ROC Curve.

# 7 Conclusion and Future Work

## 7.1 Conclusion

This study has demonstrated that an XGBoost-based diabetes prediction model, augmented with rigorous data pre- processing, feature engineering, and hyperparameter tuning, can achieve a high accuracy of 89.2%. The results show that leveraging gradient boosting techniques significantly out- performs traditional methods such as logistic regression and decision trees. Additionally, the model's strong precision, recall, and F1-score indicate a balanced performance,

making it suitable for practical deployment in clinical settings. By identifying key predictive features like Glucose, BMI, and Age, healthcare practitioners can focus on high-impact variables to "Optimizing XGBoost Parameters for Medical Classification," IEEE Trans. Biomed. Eng., 2024.refine diagnostic decisions. Overall, the findings underscore the potential of integrating machine learning into diabetes screening protocols, especially in regions with limited healthcare infrastructure.

## 7.2 Future Work

While the proposed model achieves robust performance, several avenues remain for further Investigation:

- **Multi-Modal Data Integration:** Incorporating additional clinical parameters (e.g., family history, diet, physical activity) or genetic data to enhance predictive accuracy.
- **Explain ability and Interpretability:** Developing model- agnostic methods (e.g., LIME or SHAP) to provide trans- parent decision-making insights for healthcare providers.
- **Federated Learning Approach:** Training the model across multiple healthcare institutions without centralizing data, thereby preserving patient privacy.
- **Real-Time Deployment:** Implementing the model in wearable devices or smartphone applications for continuous, on-the-spot diabetes risk assessment.
- **Cross-Population Generalization:** Validating the model on diverse populations to ensure broad applicability and fairness.

## References

[1]  Manish Prateek; Saurabh Pratap Singh Rathore, "Clinical Validation of AI Disease Detection Models — An Overview of the Clinical Validation Process for AI Disease Detection Models, and How They Can Be Validated for Accuracy and Effectiveness," in AI in Disease Detection: Advancements and Applications , IEEE, 2025, pp.215-237, doi: 10.1002/9781394278695.ch10.

[2]  N. Nisha Nadhira Nazirun et al., "Prediction Models for Type 2 Diabetes Progression: A Systematic Review," in IEEE Access, vol. 12, pp. 161595-161619, 2024, doi: 10.1109/ACCESS.2024.3432118.

[3]  S. C. Mana, G. Kalaiarasi, Y. R, L. S. Helen and R. Senthamil Selvi, "Application of Machine Learning in Healthcare: An Analysis," 2022 3rd International Conference on Electronics and Sustainable Communication Systems (ICESC), Coimbatore, India, 2022, pp. 1611-1615, doi: 10.1109/ICESC54411.2022.9885296.

[4]  D. Ather, I. Yusupov, S. Duggal, R. Kumar, P. Sagar and V. Jain, "Advancements in Deep Learning for Early Detection and Diagnosis Across Multiple Disease Domains," 2024 4th International Conference on Technological Advancements in Computational Sciences (ICTACS), Tashkent, Uzbekistan, 2024, pp. 1460-1464, doi: 10.1109/ICTACS62700.2024.10840864.

[5]  A. Kumar, A. S. Gill, J. P. Singh and D. Ghosh, "A Comprehensive and Comparative Examination of Machine Learning Techniques for Diabetes Mellitus Prediction," 2024 15th International Conference on Computing Communication and Networking Technologies (ICCCNT), Kamand, India, 2024, pp. 1-5, doi: 10.1109/ICCCNT61001.2024.10725693.

[6]  A. Pawar, S. Jain, A. Dhait, A. Nagbhidkar and A. Narlawar, "Federated Learning for Privacy Preserving in Healthcare Data Analysis," 2024 International Conference on Artificial Intelligence and Quantum Computation-Based Sensor Application (ICAIQSA), Nagpur, India, 2024, pp. 1-6, doi: 10.1109/ICAIQSA64000.2024.10882173.

[7 ]   I. S. Rajput, H. Garwal and K. Bameta, "Enhanced Diabetes Prediction Using Salivary Enzymes and Machine Learning Techniques: A Comprehensive and Explainable Approach," 2024 IEEE International Conference on Communication, Computing and Signal Processing (IICCCS), ASANSOL, India, 2024, pp. 1-6, doi: 10.1109/IICCCS61609.2024.10763602.

[8 ]   K. Georgiou, L. Liu, H. Qi and X. Zhao, "Trustworthy AI for Early Dementia Detection: Robust Feature Masking and Clinical Interpretability," 2025 IEEE/ACM Conference on Connected Health: Applications, Systems and Engineering Technologies (CHASE), New York, NY, USA, 2025, pp. 279-283.

[9 ]   I. Abousaber, "Enhanced Diabetes Prediction Through Advanced Machine Learning and Imbalance Handling Techniques," 2025 4th International Conference on Computing and Information Technology (ICCIT), Tabuk, Saudi Arabia, 2025, pp. 708-716, doi: 10.1109/ICCIT63348.2025.10989423.

[10 ]   R. Munirathnam, M. Maindola, B. H. K, P. R. Patil, P. Bhatt and N. L, "A Novel Hybrid AI Framework for Real-Time Prediction and Risk Assessment of Vector-Borne Epidemics using Spatiotemporal Data," 2024 4th International Conference on Mobile Networks and Wireless Communications (ICMNWC), Tumkuru, India, 2024, pp. 1-5, doi: 10.1109/ICMNWC63764.2024.10872354.

[11 ]   International Diabetes Federation. *IDF Diabetes Atlas*. 10th ed., International Diabetes Federation, 2021. Accessed [Date]. https://diabetesatlas.org