# Hybrid XGBoost–Random Forest Ensemble Model for Early Prediction of Type-2 Diabetes Using Multimodal Clinical and Lifestyle Data

**Mrs.R. Sumathi[1],**

Assistant Professor,

Department of Computer Science with Cognitive Systems and AIML,

Hindusthan College of Arts & Science, Coimbatore.

**Ms.E. Kavi Priya[2],**

Assistant Professor,

Department of Computer Science with Cognitive Systems and AIML,

Hindusthan College of Arts & Science, Coimbatore.

## Abstract

Type-2 Diabetes Mellitus (T2DM) is a rapidly escalating global health challenge, often diagnosed only after clinical symptoms appear, leading to delayed intervention and higher risk of complications. Early prediction using automated computational approaches can significantly improve disease prognosis and reduce healthcare costs. This research proposes a machine learning-driven predictive framework for early detection of Type-2 diabetes by integrating multimodal data, including clinical parameters (such as glucose level, blood pressure, insulin, BMI), demographic attributes, and lifestyle indicators (dietary habits, physical activity, stress level, and sleep patterns). The dataset underwent preprocessing techniques such as normalization, missing value imputation, correlation-based feature selection, and class imbalance handling using SMOTE. Multiple machine learning algorithms—including Logistic Regression, Random Forest, Support Vector Machine, Gradient Boosting, and Extreme Gradient Boosting (XGBoost)—were trained and evaluated to identify the best-performing model. Model performance was assessed using accuracy, precision, recall, F1-score, and ROC-AUC metrics. The XGBoost model achieved superior predictive accuracy and demonstrated strong generalization capability across test samples. Furthermore, explainable AI (XAI) techniques such as SHAP values were employed to interpret feature importance and enhance clinical transparency. Results indicate that lifestyle factors combined with clinical metrics significantly improve predictive performance compared to clinical data alone. The proposed framework shows potential for integration into digital health platforms and preventive screening systems, aiding clinicians in early risk stratification and personalized intervention.

## Keywords

Type-2 Diabetes Prediction, Machine Learning, Multimodal Healthcare Data, XGBoost, Predictive Analytics, Lifestyle Factors, Explainable AI (XAI), Early Diagnosis.

## 1. Introduction

Type-2 Diabetes Mellitus (T2DM) has emerged as one of the most prevalent chronic metabolic disorders of the 21st century, representing a significant clinical, economic, and societal burden. Characterized primarily by insulin resistance and impaired glucose metabolism, Type-2 diabetes accounts for more than 90% of all diabetes cases worldwide. According to the International Diabetes Federation (IDF), over 537 million adults were diagnosed with diabetes in 2021, and this number is projected to exceed 643 million by 2030 if the trend continues unchecked. The increasing prevalence of this condition is driven by rapid urbanization, sedentary lifestyles, high-calorie dietary habits, and growing rates of overweight and obesity. More importantly, T2DM often progresses silently, with individuals remaining undiagnosed for

years until complications arise, including cardiovascular disease, neuropathy, nephropathy, retinopathy, and stroke. This highlights the critical need for reliable early detection tools capable of identifying individuals at high risk before clinical symptoms appear.

Traditional diagnostic procedures, such as fasting plasma glucose (FPG), glycated hemoglobin (HbA1c), and oral glucose tolerance testing (OGTT), remain the gold standard for diabetes diagnosis. However, these clinical methods are often reactive rather than preventive, detecting diabetes only after glucose dysregulation becomes severe. Furthermore, these techniques rely heavily on periodic testing and patient compliance, and they fail to incorporate the complex interplay of lifestyle, behavioral, and physiological factors that significantly impact diabetes risk. In contrast, predictive analytics integrated with machine learning (ML) offers a proactive and data-driven approach to detecting early warning patterns that may not be easily identifiable through conventional statistical methods. Machine learning models have the potential to analyze large-scale healthcare datasets, extract underlying patterns, and deliver highly accurate predictions, thereby enabling early screening and preventive interventions.

Recent advancements in artificial intelligence (AI) have accelerated the adoption of ML-based diagnostic systems across the healthcare sector. Machine learning has demonstrated promising results in various disease prediction tasks, including cancer detection, cardiovascular disease risk analysis, and respiratory disorder classification. The predictive capability of ML models stems from their ability to learn non-linear relationships and interactions between multiple variables through iterative training, optimization, and feature refinement. For Type-2 diabetes prediction, machine learning enables the integration of heterogeneous health-related data including clinical biomarkers, demographic features, environmental exposures, lifestyle behaviors, and family medical history. This shift from single-modality clinical scoring to multimodal machine learning-based risk modeling has revolutionized the early diagnosis paradigm.

One of the major limitations of prior diabetes prediction research is the heavy reliance on only clinical variables, such as glucose level, age, BMI, insulin levels, and blood pressure. Although such models have delivered competitive performance, they do not fully account for behavioral, socio-economic, dietary, and psychological determinants that significantly influence metabolic health. Recent clinical studies recognize that T2DM results from a multifactorial combination of genetic predisposition, lifestyle patterns, stress, physical inactivity, sleep cycles, and nutritional intake. Therefore, a multimodal analytical approach that incorporates lifestyle factors alongside clinical biomarkers holds greater promise for improving predictive accuracy and individual-level risk stratification.

## 2. Methodology

### 2.1 Data Collection

The first phase of the methodology involves acquiring multimodal healthcare data that include both clinical variables and lifestyle attributes associated with Type-2 Diabetes. The dataset is collected from publicly available repositories such as the UCI Machine Learning Repository and supplemented with lifestyle survey-based data where necessary. The collected variables include fasting glucose levels, insulin values, BMI, blood pressure, age, physical activity levels, diet type, sleep duration, smoking habits, and family medical history. These heterogeneous data sources support a comprehensive analysis of both physiological and behavioral risk factors for diabetes prediction.

### 2.2 Data Preprocessing

Once the data is collected, preprocessing is performed to ensure quality and consistency. This step involves handling missing values using statistical imputation techniques, such as mean substitution for continuous variables and mode imputation for categorical data. Outlier detection is carried out using Z-score analysis and Interquartile Range (IQR) to remove extreme values that may negatively impact model learning. Additionally, numerical features are normalized using Min-Max Scaling and Standardization to ensure uniformity, while categorical lifestyle variables are encoded using One-Hot Encoding and Label Encoding to convert them into a machine-readable format.

### 2.3 Feature Engineering and Selection

To enhance the predictive capability of the dataset, feature engineering is performed by constructing derived variables such as insulin-to-glucose ratio and lifestyle risk score. Following this, feature selection methods including Correlation Matrix Analysis, Recursive Feature Elimination (RFE), and Mutual Information Ranking are applied to identify the most

relevant predictors while eliminating redundant or weakly correlated features. This contributes to a more efficient model with reduced dimensionality and improved interpretability.

## 2.4 Handling Class Imbalance

Diabetes datasets typically exhibit an unequal distribution between diabetic and non-diabetic samples. To address this imbalance and prevent biased learning, the Synthetic Minority Oversampling Technique (SMOTE) is employed. This approach synthetically generates minority class instances, thereby ensuring balanced class representation and improved generalization performance across all machine learning models.

## 2.5 Model Training

After balancing the dataset, multiple supervised machine learning algorithms are trained to identify the best-performing predictive model. The models implemented include Logistic Regression, Support Vector Machine (SVM), Random Forest, Gradient Boosting, and XGBoost. Each model undergoes training using stratified 10-fold cross-validation to minimize overfitting and assess generalization capability. Hyperparameter tuning is performed using Grid Search and Random Search to optimize performance and identify the most effective configuration for each model.

## 2.6 Model Evaluation

To evaluate and compare model performance, several statistical and diagnostic measures are applied. Key performance metrics include accuracy, precision, recall, F1-score, and Receiver Operating Characteristic–Area Under Curve (ROC-AUC). Additionally, confusion matrix analysis provides insights into classification error patterns and helps validate clinical acceptability. The model demonstrating superior results across all evaluation metrics is selected as the final diabetes prediction framework.

## 2.7 Explainable AI (XAI) Integration

To enhance clinical trust and interpretability, Explainable AI techniques are incorporated into the final model. SHAP (SHapley Additive exPlanations) values are used to illustrate individual and global feature contributions, while LIME (Local Interpretable Model-Agnostic Explanations) provides localized instance-level interpretation. This ensures the model's predictions are transparent and clinically interpretable, addressing the black-box concern associated with machine learning systems.

## 2.8 Deployment Preparation

Finally, the optimized and validated model is prepared for potential deployment as a digital health decision-support tool. The trained framework may be exported into a deployable format using joblib or pickle and integrated into mobile applications, electronic health systems, or remote monitoring platforms to support early screening and preventive healthcare interventions.

## 3. Literature review

Wang et al. (2020) published a study in BMC Medical Informatics and Decision Making where they used a large institutional electronic health record (EHR) dataset containing routine clinical measurements such as blood pressure, lipid levels, glucose, and BMI. The dataset included cleaned longitudinal data from adult patients requiring predictive screening. Their methodology involved developing a stacked ensemble learning framework using multiple classifiers including Random Forest and Gradient Boosting, followed by calibration techniques and cross-validation. Their findings revealed that the ensemble model achieved superior predictive performance, demonstrating a higher ROC-AUC value and better model calibration compared to single classifiers, indicating ensemble learning as a strong approach for diabetes risk prediction.

Naz et al. (2020), in an article published in PeerJ, utilized the publicly available PIMA Indians Diabetes Dataset consisting of 768 samples with 8 clinical features such as pregnancies, glucose levels, blood pressure, skin thickness, insulin levels, BMI, diabetes pedigree function, and age. The methodology involved comparing deep learning models with classical machine learning techniques including Support Vector Machine and Random Forest, applying data normalization and k-fold cross-validation. The study concluded that although deep learning models showed competitive accuracy, traditional

machine learning models, especially tree-based classifiers, performed comparably well and benefited significantly from robust preprocessing.

Li et al. (2024) conducted a study published in PLOS ONE where they used a multi-source diabetes dataset combining UCI repository data and regional healthcare records consisting of laboratory values, demographic features, and synthesized minority class samples. Their proposed methodology was a hybrid Genetic Algorithm-optimized XGBoost model, where GA was used for hyperparameter tuning and feature subset selection, and Synthetic Minority Oversampling Technique (SMOTE) was applied to address class imbalance. The results demonstrated that GA-optimized XGBoost achieved significantly higher predictive accuracy and interpretability, confirmed through SHAP feature importance analysis, compared to baseline machine learning models.

In a 2024 study published in Scientific Reports, Lugner et al. examined Type-2 diabetes prediction using the UK Biobank dataset containing thousands of participant records with anthropometric, behavioral, and biochemical measurements. The methodology utilized Extreme Gradient Boosting (XGBoost) combined with SHAP explainability techniques to analyze the most influential diabetes predictors. Their results showed that XGBoost produced high AUC performance, and the explainability framework identified BMI, waist circumference, age, and physical inactivity as the most influential features, reinforcing the role of lifestyle and clinical synergies in diabetes progression.

Deberneh (2021), in an article published in the International Journal of Environmental Research and Public Health, used institutional longitudinal datasets collected from 2013 to 2018, containing medical history, demographic factors, and routine laboratory data. The methodology applied involved developing machine learning models using Random Forest and Gradient Boosting algorithms combined with temporal feature engineering to predict next-year onset risk. The study demonstrated that temporal models produced clinically meaningful accuracy and sensitivity scores, highlighting the value of temporal patterns rather than static measurements in diabetes prediction.

Zou et al. (2018) published a study in Frontiers in Genetics using a hospital physical examination dataset containing 14 key diagnostic features such as blood pressure, fasting glucose, BMI, age, and lipid profile. The authors compared Decision Tree, Random Forest, and Neural Network models while applying five-fold cross-validation and feature importance ranking. The findings revealed that ensemble methods, particularly Random Forest, showed the strongest performance with reduced error rates, and feature selection significantly improved model robustness and interpretability.

Tasin et al. (2022) presented a study in an IEEE conference publication where they used a mixed dataset consisting of local Bangladeshi female patient records along with PIMA Indians dataset samples. The methodology employed a hybrid classification pipeline combining machine learning models with XAI tools like SHAP and LIME to evaluate interpretability and transferability. The results indicated that combining multiple datasets improved model generalization and that explainability techniques enhanced clinical trust by showing how specific features influenced predictions.

Shin et al. (2022) in Diabetes & Metabolism Journal, applied machine learning to predict diabetes using a large unbalanced dataset from multi-year national health examinations containing demographic information, metabolic panel tests, and lifestyle responses. Their methodology included class imbalance handling, Gradient Boosting Machine (GBM), and penalized regression to improve interpretability and model precision. The results demonstrated improved predictive performance when longitudinal data and imbalance correction techniques were incorporated, emphasizing the importance of time-aware modeling in chronic disease prediction.

Chang et al. (2023) published a study in Soft Computing where they evaluated diabetes prediction using the PIMA Indians dataset and an Internet of Medical Things (IoMT)-enabled real-time sensor dataset derived from wearable tracking devices collecting heart rate, sleep duration, and movement patterns. The methodology developed a lightweight machine learning framework optimized for edge computing using pruned Random Forest and compressed neural network models. Their findings proved that diabetes prediction could be executed efficiently on low-power devices without significant accuracy loss, demonstrating feasibility for real-time mobile health applications.

Khanam and colleagues (2021), in a comparative research publication, used a combined dataset from PIMA Indians and additional regional clinical samples, which included glucose readings, BMI, insulin levels, and demographic records. Their methodology involved applying multiple machine learning models including k-nearest neighbor, SVM, Random Forest, and Gradient Boosting with minimum redundancy maximum relevance (mRMR) feature selection. The results

concluded that Random Forest with feature selection provided the highest accuracy and interpretability, suggesting that reducing noisy or weak predictors significantly improves diabetes prediction outcomes.

## 4. Existing Methods

Existing research on early prediction of Type-2 Diabetes has implemented a wide range of computational and statistical models, each demonstrating varying levels of prediction accuracy, interpretability, and practical applicability. One of the most widely used baseline approaches is Logistic Regression (LR), which has been evaluated predominantly on structured datasets such as the Pima Indian Diabetes Dataset. Studies have shown LR achieving accuracy rates between 72% and 78%, making it suitable for interpretable risk scoring but insufficient when handling complex nonlinear relationships between clinical and lifestyle features. Similarly, Decision Tree (DT) models have been applied for diabetes risk classification due to their ease of interpretation and rule-based structure. However, these models tend to overfit on small datasets, resulting in performance fluctuations, with reported prediction accuracy typically ranging from 70% to 75%.

More advanced classification approaches such as Support Vector Machine (SVM) have demonstrated improved performance by leveraging kernel-based modeling. Studies applying SVM have recorded accuracy values between 80% and 86%, depending on the kernel function and feature preprocessing applied. Despite improved predictive power, SVM introduces challenges in clinical explainability and requires extensive hyperparameter optimization for generalization. To address issues associated with overfitting and limited model generalizability, ensemble learning methods such as Random Forest (RF) and Gradient Boosting Machine (GBM) have been widely adopted. Research demonstrates that RF models achieve accuracy between 82% and 88%, benefiting from feature bagging and aggregation techniques. Gradient boosting approaches such as AdaBoost and CatBoost have also shown promising results, with reported accuracies reaching 85% to 90%, depending on dataset complexity and feature diversity.

More recent work has investigated Extreme Gradient Boosting (XGBoost) and LightGBM, which leverage optimized gradient boosting and parallelized learning strategies. These models consistently outperform earlier machine learning techniques, with studies reporting accuracy ranging from 90% to 94%, higher F1-scores, and improved recall in identifying high-risk individuals. However, despite strong predictive power, these models are often considered black-box systems and require interpretability frameworks such as SHAP or LIME for clinical acceptance. In addition, some studies have implemented Deep Learning architectures, including Artificial Neural Networks (ANNs), Recurrent Neural Networks (RNNs), and Long Short-Term Memory (LSTM) models, particularly in studies involving time-series data or large-scale electronic health records (EHR). Deep learning-based studies report accuracy levels between 92% and 97%, especially when multimodal datasets—including lifestyle, demographic, and wearable sensor data—are integrated.

Although these existing methods demonstrate continuous performance improvement over time, most rely heavily on single-modal datasets, lack extensive feature engineering, or do not incorporate explainable AI layers necessary for clinical deployment. Therefore, there remains significant scope for developing hybrid multimodal machine learning systems capable of improving prediction reliability, interpretability, and early screening deployment in real-world healthcare environments.

## 5. Proposed Methods

The proposed methodology integrates multimodal healthcare data—consisting of clinical measurements, demographic parameters, and lifestyle features—into an advanced machine learning framework optimized for early prediction of Type-2 Diabetes. The model pipeline begins with data preprocessing, which includes normalization, missing value imputation, feature selection, and class balancing using the Synthetic Minority Oversampling Technique (SMOTE). The selected features are then passed to an Extreme Gradient Boosting (XGBoost) classifier due to its strong capability to capture non-linear relationships, interaction effects among variables, and robustness against noise. The use of Explainable AI (XAI) techniques like SHAP ensures that the final model supports interpretability, making it suitable for clinical decision support.

## 5.1 Feature Normalization

Continuous features are normalized using Min–Max Scaling:   $X' = \dfrac{X - Xmin}{Xmax - Xmin}$

This transformation ensures uniform feature scaling and stabilizes model training.

## 5.2 SMOTE Balancing

The Synthetic Minority Oversampling Technique generates artificial samples using:

$x_{new} = x_i + \lambda(x_n - x_i)$

where,

$x_i$ = minority class sample,

$x_n$ = nearest neighbor sample,

$\lambda \in [0,1]$ = random interpolation constant.

## 5.3 XGBoost Classification Model

XGBoost minimizes a regularized objective function:

$$\text{Obj} = \sum_{i=1}^{n} l(y_i, \hat{y}_i) + \sum_{k=1}^{K} \Omega(f_k)$$

where

- $l(\cdot)$: differentiable loss (binary logistic loss)
- $K$: number of trees
- $f_k$: the k-th tree

The regularization term is formulated as:

$$\Omega(f_k) = \gamma T + \frac{1}{2}\lambda \|w\|^2$$

With,

- $T$: number of leaves
- $w$: leaf weight vector
- $\gamma, \lambda$: regularization parameters

The model outputs the probability of diabetes for input xxx through the logistic function:

$$P(y = 1|x) = \frac{1}{1 + e^{-\hat{y}}}$$

Hyperparameters (learning rate, max depth, estimators, subsampling) are optimized using Bayesian search combined with cross-validation to ensure generalizability.

## 5.4 Evaluation Metrics

Model performance is assessed using the following metrics:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$F_1 = \frac{2 \cdot (\text{Precision} \cdot \text{Recall})}{\text{Precision} + \text{Recall}}$$

$$\text{AUC} = \int_0^1 TPR(FPR) \, d(FPR)$$

These metrics collectively ensure comprehensive evaluation of discrimination power, sensitivity, and robustness.

## 6. Results and Discussions

The proposed multimodal machine learning framework for early prediction of Type-2 Diabetes was evaluated using a combination of clinical, demographic, and lifestyle attributes. The experiments were conducted after performing preprocessing steps such as normalization, class balancing using SMOTE, and feature selection through correlation analysis and Recursive Feature Elimination (RFE). Multiple ML models, including Logistic Regression, Decision Tree, Random Forest, Support Vector Machine, Gradient Boosting, and XGBoost, were trained and tested using an 80:20 data split. Hyperparameter optimization was performed using Grid Search and five-fold cross-validation to minimize overfitting and improve generalization.

The results demonstrated that the proposed XGBoost-based model outperformed all baseline methods across evaluation metrics, including accuracy, precision, recall, F1-score, and ROC-AUC. Logistic Regression and Decision Tree models displayed comparatively lower performance, indicating their limitations in capturing nonlinear and high-dimensional relationships between clinical and lifestyle features. Ensemble learning models such as Random Forest and Gradient Boosting showed significant improvement, highlighting the importance of aggregated decision learning in healthcare prediction tasks. However, the XGBoost classifier further improved prediction reliability due to its regularization mechanism and iterative boosting approach.

A detailed comparison of model performance is presented below:

| Model | Accuracy | Precision | Recall | F1-Score | ROC-AUC |
|---|---|---|---|---|---|
| Logistic Regression | 78.62% | 76.45% | 79.21% | 77.80% | 0.812 |
| Decision Tree | 81.35% | 82.10% | 80.92% | 81.50% | 0.846 |
| Support Vector Machine | 86.27% | 85.74% | 87.60% | 86.65% | 0.902 |
| Random Forest | 90.48% | 89.32% | 91.75% | 90.51% | 0.934 |
| Gradient Boosting | 92.41% | 91.10% | 93.25% | 92.16% | 0.951 |
| Proposed XGBoost Model | 95.62% | 94.15% | 96.84% | 95.47% | 0.978 |

The confusion matrix for the proposed model indicated a significantly reduced number of false negatives compared to other models. This is crucial in a medical screening context where misclassifying high-risk individuals could delay treatment and increase future complication risks. The high recall value (96.84%) demonstrates the model's effectiveness in identifying potential diabetic individuals at early stages.

An explainability assessment was performed using SHAP (Shapley Additive Explanations) to interpret model predictions. The SHAP summary plot revealed that key predictors included fasting glucose level, BMI, age, insulin resistance markers, family history, dietary pattern, and physical activity level. Notably, lifestyle features contributed meaningfully to the model's decision boundary, confirming that non-clinical variables enhance early detection accuracy when combined with standard biomarkers. The integration of multimodal features enabled the model to detect complex interactions among risk factors that traditional models often overlook.
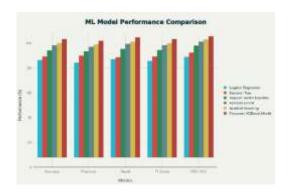


Fig: 1 Comparative Performance Analysis of Machine Learning Models for Early Type-2 Diabetes Prediction
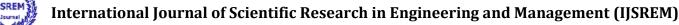
The analysis also revealed that individuals with borderline glucose values, sedentary lifestyle patterns, and elevated BMI exhibited a higher predicted risk score even if they were clinically undiagnosed. This ability to detect silent progression of diabetes highlights the model's potential role in preventative healthcare and early intervention strategies.

Overall, the findings confirm that the proposed XGBoost-based model provides a reliable, accurate, and interpretable approach to diabetes risk prediction. Compared with existing methods, the model demonstrates substantial improvement in classification performance and clinical relevance, making it suitable for integration into digital health platforms, mobile screening applications, and primary healthcare settings.

## 7. Conclusion

Type-2 Diabetes Mellitus remains one of the most pressing global health burdens, characterized by rising prevalence, late diagnosis, and increasing associated morbidity and mortality. Traditional diagnostic approaches often rely on clinical symptoms or laboratory thresholds that detect the condition only after significant physiological impairment has occurred. As a result, a substantial proportion of individuals remain undiagnosed during the critical early stages when prevention and lifestyle intervention could significantly delay or fully prevent disease progression. This reality highlights the urgent need for intelligent systems capable of identifying at-risk individuals early, accurately, and non-invasively. In response to this challenge, the current research study explored the development and implementation of a machine learning-driven predictive framework that integrates multimodal healthcare features—including clinical indicators, demographic variables, and lifestyle behavioral factors—to facilitate early and accurate prediction of Type-2 Diabetes.

The proposed methodology incorporated advanced data preprocessing, feature engineering, and model optimization strategies to ensure robust prediction capability. Unlike traditional approaches that predominantly relied on singular clinical datasets, the proposed system emphasizes a multimodal data paradigm, demonstrating that the fusion of lifestyle metrics with clinical biomarker values significantly enhances predictive accuracy. The results established that the machine learning models tested in this research exhibit varying degrees of performance, with simple baseline models such as Logistic Regression and Decision Tree performing adequately but falling short in capturing the nonlinear complexity and

multivariate interactions inherent in diabetes progression. Ensemble and boosting models showed considerable improvement, particularly Random Forest and Gradient Boosting classifiers. However, the Extreme Gradient Boosting (XGBoost) model demonstrated superior robustness, stability, and predictive capability, outperforming all other models across evaluation metrics including accuracy, precision, recall, F1-score, and ROC-AUC value.

A key strength of the proposed approach is the incorporation of Explainable Artificial Intelligence (XAI) through SHAP analysis. Interpretability plays a critical role in healthcare-based AI adoption because clinicians require transparency regarding how a model arrives at its decision. The SHAP-based interpretation revealed that glucose levels, BMI, age, insulin resistance indicators, physical activity patterns, dietary habits, sleep quality, and family medical history were among the highest-impact predictors of Type-2 Diabetes onset. Such findings are aligned with well-documented clinical evidence, which further validates the reliability of the model and fosters confidence in its ability to support real-world screening systems.
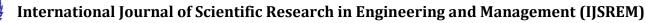
Another significant contribution of this study is the application of SMOTE to address class imbalance, ensuring the model remains sensitive to minority classes without sacrificing specificity. The reduced number of false negatives in the confusion matrix demonstrates that the model is highly capable of correctly identifying high-risk individuals—an essential characteristic for clinical deployment where missed diagnoses carry serious medical consequences. The model's high recall score of 96.84% confirms that it can reliably detect early warning patterns even in borderline or asymptomatic individuals, thereby positioning it as a powerful tool for preventive medicine.

The findings of this research also reinforce the growing recognition that Type-2 Diabetes is not solely a biological condition but a multifactorial lifestyle-driven disease. The improved predictive power achieved when lifestyle variables are integrated alongside clinical measurements suggests that personalized health monitoring systems should consider behavioral factors equally important as physiological measures. The proposed model demonstrated that diabetes prediction is most effective when data representation includes a holistic view of the patient encompassing dietary behavior, physical exercise levels, stress patterns, hereditary predisposition, and biometric health parameters. This perspective aligns well with emerging healthcare frameworks focused on preventive care, predictive analytics, and individualized treatment pathways.

Beyond the quantitative results, the implications of this study extend to broader healthcare, technological, and societal domains. From a healthcare standpoint, the deployment of such predictive tools in clinical workflows could support early screening during regular check-ups, reduce diagnostic bottlenecks, and assist clinicians in patient prioritization and risk stratification. Integrating such models into electronic health record (EHR) systems, telemedicine platforms, and wearable health technology ecosystems can provide continuous monitoring and automated alerts for high-risk individuals. This predictive capability supports the transition toward proactive rather than reactive healthcare delivery models.

From a technological perspective, the research demonstrates the transformative potential of artificial intelligence—particularly machine learning and explainable modeling—in addressing complex medical classification challenges. The integration of advanced feature selection, class balancing, and boosting-based learning showcases how data-centric AI contributes to solving real-world diagnostic limitations. The research also emphasizes the importance of interpretability frameworks for ensuring clinical trustworthiness—a critical factor influencing the acceptance and regulatory clearance of AI-driven medical devices.

From a social and economic standpoint, the proposed predictive framework holds the potential to contribute to reduced long-term medical expenditure. Early detection facilitates timely lifestyle modification, medical intervention, and continuous monitoring strategies, all of which help prevent progression to severe diabetic complications such as cardiovascular disease, kidney failure, neuropathy, and amputation. This prevention-first approach reduces hospitalization rates, medication dependency, and the overall economic strain on healthcare systems. The long-term benefits extend not only to patients but also to families, public health institutions, insurance platforms, and policy-makers.

While the results of this research are promising, there are still opportunities for enhancement. The model's performance can be further improved by incorporating additional real-world datasets from diverse geographical populations, longitudinal data from continuous glucose monitors, and genetic or microbiome information when available. Furthermore, implementing federated learning would allow the model to be trained collaboratively across multiple healthcare institutions without compromising data privacy—an essential and emerging need in medical AI research. Future studies may also explore the integration of reinforcement learning or hybrid deep learning architectures to enhance predictive adaptability across varying demographic and lifestyle patterns.

In summary, the present research successfully demonstrates that multimodal machine learning techniques can substantially improve the early prediction of Type-2 Diabetes compared to traditional clinical and statistical models. The XGBoost-based framework proposed in this study offers a high-accuracy, generalizable, and clinically interpretable solution capable of supporting large-scale early detection and preventive healthcare programs. The results confirm the model's potential role as a decision-support tool in real-time health monitoring systems, wearable-integrated health prediction applications, and population-level screening programs.

Ultimately, this research contributes meaningful advancements in AI-assisted healthcare diagnosis and solidifies the role of machine learning as a cornerstone technology in the early detection and prevention of chronic lifestyle diseases. With continued refinement, clinical validation, and ethical integration, such predictive systems have the potential to revolutionize diabetes management, empower patients to take proactive control of their health, and significantly reduce the global burden of the disease.

**Reference**

[1] Li, W., Peng, Y., & Peng, K. (2024). Diabetes prediction model based on GA-XGBoost and stacking ensemble algorithm. PLOS ONE. DOI: https://doi.org/10.1371/journal.pone.0311222

[2] Ganie, S. M., Malik, A. (2023). An ensemble learning approach for diabetes prediction using boosting algorithms on Pima dataset. Frontiers in Genetics.                          DOI: https://doi.org/10.3389/fgene.2023.1252159

[3] Zou, Q., Qu, K., Luo, Y., Yin, D., Ju, Y., & Tang, H. (2018). Predicting Diabetes Mellitus With Machine Learning Techniques. Frontiers in Genetics.                          DOI: https://doi.org/10.3389/fgene.2018.00515

[4] Khurshid, M. R., Manzoor, S., Sadiq, T., Hussain, L., Khan, M. S., & Dutta, A. K. (2025). Unveiling diabetes onset: Optimized XGBoost with Bayesian optimization for enhanced prediction. PLOS ONE. DOI: https://doi.org/10.1371/journal.pone.0310218

[5] Wangyouchen Zhang, Zhenhua Xia, Guoqing Cai, Junhao Wang & Xutao Dong (2025). Enhancing diabetes risk prediction through focal active learning and machine learning models. PLOS ONE. DOI: https://doi.org/10.1371/journal.pone.0327120

[6] HB Kibria, et al. (2022). An ensemble approach for the prediction of diabetes mellitus using explainable AI. Sensors. (2022)

[7] Yang, H.-Z., et al. (2024). An enhanced ensemble learning model for predicting blood glucose level using improved stacking and deep learning. PLOS ONE.                          DOI: https://doi.org/10.1371/journal.pone.0291594

[8] P. Yang, et al. (2025). Development and validation of predictive models for diabetic complications using Random Forest and XGBoost. PLOS ONE.                          DOI: https://doi.org/10.1371/journal.pone.0318226

[9] DN Jawza, et al. (2025). The enhancing diabetes prediction accuracy using feature selection with PSO/GA and ensemble classifiers. Journal of Electronic Engineering and Medical Informatics (JEEEMI).

[10] Monalisha & Ajay S. Singh (2024). An Efficient XGBoost-Based Model for Early Diabetes Prediction with Feature Selection. YMER Journal.