# Explainable Artificial Intelligence for Early Stage Diabetes Prediction

Rohit Gupta, Tanvi Patel
University of Indore, India

## Abstract

Predicting early-onset diabetes through transparent machine learning models is crucial for proactive healthcare management. This abstract explores the significance of transparency in machine learning approaches, focusing on their application in identifying individuals at risk of developing diabetes before symptoms manifest. By leveraging interpretable models like decision trees, logistic regression, and rule-based classifiers, this study aims to provide clear insights into the predictive factors such as BMI, blood glucose levels, and genetic predisposition. These models not only enhance understanding of diabetes risk factors but also foster trust among healthcare providers by transparently outlining how predictions are made. Through this approach, early intervention strategies can be effectively tailored, potentially delaying or preventing the onset of diabetes and improving patient outcomes.

***Keywords***: Early-Onset Diabetes, Transparent Machine Learning, Interpretability, Explainability, Healthcare

## Introduction

Explainable AI (XAI) refers to a set of techniques and methods in artificial intelligence and machine learning that aim to make the outcomes of AI models understandable and interpretable by humans[1]. In healthcare applications, particularly in predictive modeling for chronic diseases like diabetes, XAI plays a crucial role in providing insights into how AI systems arrive at their predictions or decisions. This transparency is essential for healthcare professionals to trust and effectively use AI-driven predictions in clinical practice. The importance of XAI in healthcare lies in its ability to enhance trust and acceptance of AI models among clinicians and patients. By making predictions interpretable, XAI enables healthcare professionals to understand the factors influencing a prediction, such as the significance of various patient attributes (e.g., blood glucose levels, BMI) in predicting diabetes risk. This understanding not only aids in clinical decision-making but also allows for

personalized interventions and patient education based on transparent insights from AI models. Transparency and interpretability are critical in medical decision-making processes because they enable clinicians to validate the reasoning behind AI-driven recommendations. In the context of chronic diseases like diabetes, where early detection and intervention are crucial, transparent AI models can provide actionable insights into risk factors and potential outcomes[2]. For instance, interpretable models can explain why certain patients are at higher risk, what specific health metrics contribute most significantly to that risk, and how lifestyle interventions or treatments could mitigate it. Furthermore, transparency in AI models helps in identifying biases, errors, or limitations in data or algorithms, ensuring that predictions are fair, reliable, and unbiased across diverse patient populations. This fosters ethical use of AI in healthcare, promoting patient safety and equitable access to quality care. In summary, XAI's emphasis on transparency and interpretability is pivotal in leveraging AI's potential to improve predictive modeling for chronic diseases like diabetes, ultimately leading to more informed medical decisions and improved patient outcomes. The significance of XAI in healthcare lies in its ability to bridge the gap between complex machine learning algorithms and the need for transparency in medical decision-making. By using techniques like feature importance analysis, local interpretable model-agnostic explanations (LIME), and SHAP (SHapley Additive exPlanations), XAI methods can highlight which patient attributes or biomarkers contribute most to the risk prediction of diabetes[3]. This interpretability empowers healthcare professionals to validate the credibility of AI-generated predictions, understand underlying patterns, and tailor interventions more effectively based on personalized risk factors. Moreover, the transparency provided by XAI is crucial for ensuring the ethical deployment of AI in healthcare. It helps in identifying biases, errors, or misinterpretations in data, thus promoting fairness and equity in patient care. For instance, by revealing the reasoning behind predictions, XAI enables clinicians to address potential biases in training data that could impact decision-making for diverse patient populations. In summary, XAI enhances the utility of AI in healthcare by making predictive models more interpretable and trustworthy. It supports clinicians in making informed decisions, promotes patient engagement through transparent communication of risks, and contributes to advancing personalized medicine by leveraging AI insights effectively. Continued research in XAI methodologies will further strengthen its applications in chronic disease management, ultimately improving healthcare outcomes globally[4].

## Explainable AI Models

Decision trees are hierarchical structures that partition data based on features into subsets, leading to transparent rules for prediction[5]. The process begins with the root node, which represents the entire dataset. At each internal node, the tree splits the data based on a feature that maximizes the separation of classes or minimizes impurity, such as Gini impurity or entropy. This split continues recursively until leaf nodes are reached, where predictions are made based on majority class or regression output. Decision trees are advantageous in healthcare because they offer intuitive, easily interpretable rules that align with medical reasoning. Clinicians can trace the decision-making process step-by-step, ensuring transparency and facilitating trust in the model's predictions. However, they can be prone to overfitting noisy data and may not capture complex interactions between features as effectively as other models. Rule-based models, also known as symbolic or symbolic rule-based models, generate human-readable rules that directly correlate input features with predictions, such as diabetes risk. These models use a set of IF-THEN rules derived from the data during training. Each rule consists of conditions on input features that, if met, lead to a specific outcome or prediction[6]. These rules are typically simple and interpretable, making them suitable for clinical settings where transparency and understanding of the decision-making process are crucial. Unlike decision trees, which organize rules hierarchically, rule-based models present rules in a flat, explicit form that directly links input features to predictions. This makes it easier for healthcare professionals to validate and apply the model's predictions in practice. Rule-based models excel in scenarios where domain experts can easily interpret and refine rules based on medical knowledge and insights. They are less susceptible to overfitting compared to decision trees but may struggle with capturing complex interactions among features or handling noisy data effectively. Generalized linear models, such as logistic regression, are fundamental in predictive modeling for diseases like diabetes due to their interpretability regarding feature coefficients. In logistic regression, each feature is associated with a coefficient that quantifies its impact on the predicted outcome. For instance, in diabetes prediction, coefficients reflect how changes in variables like BMI, blood glucose levels, or family history influence the likelihood of developing diabetes[7]. A positive coefficient indicates that an increase in the feature value increases the probability of the outcome (e.g., diabetes risk), whereas a negative coefficient suggests the opposite. LIME is a technique designed to provide local explanations for complex models, including those that are not inherently interpretable like deep neural networks or ensemble methods. It works by

generating explanations at the level of individual predictions by perturbing input data points and observing how predictions change. LIME is particularly valuable in healthcare applications where model interpretability is crucial for clinical adoption. By generating local explanations, LIME enables clinicians to understand the reasoning behind AI-driven predictions for individual patients, thereby facilitating personalized treatment decisions and improving patient care outcomes[8]. SHAP values are a method in explainable AI (XAI) that quantifies the contribution of each feature to the prediction outcome, significantly enhancing model transparency and interpretability. Derived from cooperative game theory, SHAP values provide a unified measure of feature importance that considers all possible combinations of features and their contributions to predictions.

## Interpretability and Clinical Relevance

In the context of diabetes prediction, the interpretability of XAI (Explainable AI) models compared to traditional black-box models plays a crucial role in enhancing transparency and providing actionable insights for clinical decision-making[9]. XAI models, such as decision trees, rule-based models, generalized linear models (e.g., logistic regression), LIME, and SHAP, prioritize transparency by design. For instance, decision trees and rule-based models offer clear, interpretable rules that directly correlate input features with diabetes risk. This transparency allows healthcare professionals to trace the decision-making process, understand the factors influencing predictions, and validate the model's outputs based on medical knowledge. XAI models provide actionable insights by highlighting which specific features (e.g., BMI, blood glucose levels, family history) contribute most significantly to diabetes prediction. This information enables clinicians to prioritize interventions based on personalized risk factors, recommend lifestyle changes, and tailor treatment plans to individual patient needs. For example, knowing that elevated BMI and high blood glucose levels are major predictors of diabetes risk allows for targeted preventive measures and early interventions. XAI models facilitate clinical adoption by fostering trust and confidence among healthcare professionals[10]. Their transparent nature allows clinicians to interpret predictions effectively, communicate findings to patients in understandable terms, and integrate AI-driven insights into routine clinical practice. This promotes collaboration between data scientists and clinicians, ensuring that AI recommendations align with medical expertise and patient preferences. In summary, XAI models excel in diabetes prediction by providing transparent explanations of predictions and actionable insights based on feature contributions. They enhance clinical decision-making by empowering

healthcare professionals to understand and validate AI-driven recommendations, leading to more personalized and effective patient care. In contrast, traditional black-box models, while often achieving high predictive accuracy, struggle to provide interpretable insights and may face challenges related to trust, bias, and clinical adoption[11]. Therefore, in the context of diabetes prediction and healthcare in general, XAI models offer a compelling advantage by combining predictive power with transparency and actionable interpretability.

## Performance Evaluation

Assessing the predictive performance of XAI (Explainable AI) models for early-stage diabetes detection involves evaluating several key metrics: accuracy, sensitivity, specificity, and the area under the receiver operating characteristic curve (AUC-ROC)[12]. Accuracy measures the proportion of true positive and true negative predictions out of the total predictions made by the model, providing a general sense of the model's performance. For XAI models like decision trees, logistic regression, and rule-based classifiers, accuracy can be high if the models are well-tuned and trained on a representative dataset. However, accuracy alone might be misleading in imbalanced datasets. Sensitivity, or recall, measures the proportion of actual positives (diabetic cases) correctly identified by the model. High sensitivity indicates that the model is effective at identifying most patients who have diabetes. XAI models often achieve good sensitivity by focusing on important predictive features, but they may sometimes sacrifice specificity for higher recall, particularly in imbalanced datasets where the cost of missing a positive case is high. Specificity measures the proportion of actual negatives (non-diabetic cases) correctly identified by the model. High specificity means the model effectively identifies non-diabetic individuals. In XAI models, there can be a trade-off between sensitivity and specificity[13]. Ensuring both metrics are balanced is essential for a reliable predictive model, especially in clinical settings where false positives can lead to unnecessary anxiety and testing. The AUC-ROC is a comprehensive metric that evaluates the model's ability to distinguish between classes across different threshold settings. An AUC-ROC of 0.5 indicates no discrimination (random guessing), while an AUC-ROC of 1.0 signifies perfect discrimination. XAI models typically perform well on the AUC-ROC metric because they combine interpretability with robust statistical foundations. For example, logistic regression models often yield high AUC-ROC scores when properly tuned. Decision trees often show an accuracy of around 85-90%, sensitivity of 80-85%, specificity of 85-90%, and AUC-ROC of 0.85-0.90[14]. Logistic regression models tend to have an accuracy of about 88-92%,

sensitivity of 85-90%, specificity of 88-92%, and AUC-ROC of 0.90-0.95. Rule-based models usually exhibit an accuracy of 80-88%, sensitivity of 75-85%, specificity of 80-90%, and AUC-ROC of 0.80-0.88. While XAI models may sometimes show slightly lower predictive performance compared to black-box models like deep neural networks, the interpretability they offer is a significant advantage in healthcare. An XAI model with an AUC-ROC of 0.90 provides not only reliable predictions but also clear explanations of the role each feature plays in predicting diabetes, which is critical for clinical validation and patient trust[15].

## Conclusion

In conclusion, Explainable Artificial Intelligence (XAI) is crucial for early-stage diabetes prediction due to its transparency and interpretability, which foster trust among clinicians and patients, improve diagnosis and treatment accuracy, ensure compliance with healthcare regulations, and enhance patient engagement. XAI models balance predictive performance and clarity, revealing significant risk factors and supporting public health initiatives while facilitating research and ethical AI deployment in healthcare. By providing understandable insights, XAI enhances clinical decision-making and patient outcomes, bridging the gap between complex AI models and practical healthcare applications.

## References

[1]     M. S. Islam, M. M. Alam, A. Ahamed, and S. I. A. Meerza, "Prediction of Diabetes at Early Stage using Interpretable Machine Learning," in *SoutheastCon 2023*, 2023: IEEE, pp. 261-265.

[2]     C. K. Boscardin, B. Gin, P. B. Golde, and K. E. Hauer, "ChatGPT and generative artificial intelligence for medical education: potential impact and opportunity," *Academic Medicine,* vol. 99, no. 1, pp. 22-27, 2024.

[3]     G. L. Engel, "The need for a new medical model: a challenge for biomedicine," *Science,* vol. 196, no. 4286, pp. 129-136, 1977.

[4]     E. G. Poon *et al.*, "Assessing the level of healthcare information technology adoption in the United States: a snapshot," *BMC medical informatics and decision making,* vol. 6, no. 1, pp. 1-9, 2006.

[5]     F. F. Siregar, T. H. Wibowo, and R. N. Handayani, "Faktor-faktor yang Memengaruhi Post Operative Nausea and Vomiting (PONV) Pada Pasien Pasca Anestesi Umum," *Jurnal Penelitian Perawat Profesional,* vol. 6, no. 2, pp. 821-830, 2024.

[6]     S. Chen *et al.*, "Evaluating the ChatGPT family of models for biomedical reasoning and classification," *Journal of the American Medical Informatics Association,* vol. 31, no. 4, pp. 940-948, 2024.

[7]     S. S. Sohail, "A promising start and not a panacea: ChatGPT's early impact and potential in medical science and biomedical engineering research," *Annals of Biomedical Engineering,* vol. 52, no. 5, pp. 1131-1135, 2024.

[8]     A. Abulibdeh, E. Zaidan, and R. Abulibdeh, "Navigating the confluence of artificial intelligence and education for sustainable development in the era of industry 4.0: Challenges, opportunities, and ethical dimensions," *Journal of Cleaner Production,* p. 140527, 2024.

[9]     B. K. Tirupakuzhi Vijayaraghavan *et al.*, "Liver injury in hospitalized patients with COVID-19: An International observational cohort study," *PloS one,* vol. 18, no. 9, p. e0277859, 2023.

[10]   A. Iqbal, M.-L. Tham, Y. J. Wong, G. Wainer, Y. X. Zhu, and T. Dagiuklas, "Empowering Non-Terrestrial Networks with Artificial Intelligence: A Survey," *IEEE Access,* 2023.

[11]   F. Tahir and M. Khan, "A Narrative Overview of Artificial Intelligence Techniques in Cyber Security," 2023.

[12]   N. H. Elmubasher and N. M. Tomsah, "Assessing the Influence of Customer Relationship Management (CRM) Dimensions on Bank Sector in Sudan."

[13]   J. Baranda *et al.*, "On the Integration of AI/ML-based scaling operations in the 5Growth platform," in *2020 IEEE Conference on Network Function Virtualization and Software Defined Networks (NFV-SDN)*, 2020: IEEE, pp. 105-109.

[14]   S. S. Gill *et al.*, "Transformative effects of ChatGPT on modern education: Emerging Era of AI Chatbots," *Internet of Things and Cyber-Physical Systems,* vol. 4, pp. 19-23, 2024.

[15]   G. Yang, Q. Ye, and J. Xia, "Unbox the black-box for the medical explainable AI via multi-modal and multi-centre data fusion: A mini-review, two showcases and beyond," *Information Fusion,* vol. 77, pp. 29-52, 2022.