**Literature Review and Research Motivation**

Diabetes is a chronic metabolic condition defined by elevated blood glucose levels and is a major public health concern due to its link to serious complications such as cardiovascular disease, neuropathy, and kidney failure [1]. Early diagnosis plays a critical role in improving patient outcomes, as timely treatment can significantly reduce the risk of long-term damage [1]. In this context, The application of machine learning is particularly helpful. By analyzing patient data, Machine learning models can identify individuals at risk of developing diabetes more accurately than many traditional methods [2]. However, developing reliable predictive models is not without its challenges. The complexity of diabetes, the imbalanced nature of many medical datasets, and the multifaceted relationships between risk factors like glucose levels, BMI, and insulin resistance all make model development more difficult. Despite these challenges, I believe that with thoughtful feature engineering and careful model evaluation, machine learning holds strong potential for advancing early detection and personalized care in diabetes management.

**Prior Work on Diabetes Prediction Using ML**

Several studies have explored the use of machine learning for diabetes prediction, and I noticed that many of them rely on the Pima Indians Diabetes Dataset. It's a well-known benchmark dataset that includes data from 768 patients and 9 clinical features such as glucose tolerance, blood pressure, and BMI [3]. In [4], the authors used several machine learning algorithms for diabetes prediction, including K-Nearest Neighbors (K-NN), Decision Trees (DT), Logistic Regression (LR), Bernoulli Naive Bayes (BNB), and Support Vector Machines (SVM). To address class imbalance, the study applied the Synthetic Minority Over-sampling Technique (SMOTE), and used Principal Component Analysis (PCA) for dimensionality reduction. Interestingly, K-NN achieved the highest accuracy (79.6%) when SMOTE was applied, followed closely by BNB at 77.2%. However, performance dropped significantly when PCA was introduced—for instance, K-NN's accuracy fell to 44.4%—suggesting that dimensionality reduction may have removed features important for diagnosis. The study used standard evaluation metrics such as accuracy, precision, recall, and F1-score, but did not explore hyperparameter tuning or advanced ensemble methods, indicating potential areas for future improvement. Prior work focused on classical ML models (K-NN, BNB) but did not evaluate state-of-the-art ensemble methods (e.g., XGBoost, Random Forest), which often outperform traditional algorithms in medical diagnostics. While classification accuracy was assessed, no unsupervised learning techniques (e.g., clustering) were introduced to explore patient subgroups or risk levels.

**Bridging the gap**

In this study, we extend prior research by incorporating advanced machine learning models and systematic optimization techniques. Specifically, we used XGBoost, Random Forest, Support Vector Machines (SVM), and Logistic Regression, applying Optuna for hyperparameter tuning to maximize predictive performance. Unlike earlier studies that relied on default settings, our approach focuses on carefully optimizing parameters such as learning rate, tree depth, and regularization strength. This allows us to improve model generalizability and potentially discover stronger predictive patterns within the dataset.

Our optimized XGBoost model achieved an accuracy of 76.6%, which is comparable to the best-performing model from prior studies (K-NN at 79.6%). However, our approach offers several distinct advantages. First, XGBoost provides superior interpretability through built-in feature importance analysis, that allows for a deeper understanding of the factors driving predictions. Second, we include clustering insights into the modeling process, allowing more targeted interventions and enhanced patient profiling.

**Research questions**

This study is guided by the following research questions:

1. Can hyperparameter-optimized ensemble models (such as XGBoost and Random Forest) outperform traditional classifiers like K-NN and Bernoulli Naive Bayes (BNB) in predicting diabetes risk?

2. Does incorporating clustering-based risk stratification yield actionable insights that go beyond standard binary classification?

3. Which clinical features are most influential in distinguishing high-risk subgroups?

By combining supervised learning for prediction with unsupervised learning for patient segmentation, we aim to build a framework that goes beyond just improving accuracy. The goal is to support more personalized, data-driven healthcare by not only identifying who is at risk but also understanding how different groups of patients may require different types of intervention.

References

1. Improving care and promoting health in populations: *standards of medical care in diabetes—2022*. (2021). *Diabetes Care*, *45*(Supplement_1). https://doi.org/10.2337/dc22-s001

2. Shojaee-Mend, H., Velayati, F., Tayefi, B., & Babaee, E. (2024). Prediction of diabetes using data mining and Machine Learning Algorithms: A cross-sectional study. *Healthcare Informatics Research*, *30*(1), 73–82. https://doi.org/10.4258/hir.2024.30.1.73

3. Learning, U. M. (2016, October 6). *Pima Indians Diabetes Database*. Kaggle. https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database

4. Iparraguirre-Villanueva, O., Espinola-Linares, K., Flores Castañeda, R. O., & Cabanillas-Carbonell, M. (2023). Application of machine learning models for early detection and accurate classification of type 2 diabetes. *Diagnostics*, *13*(14), 2383. https://doi.org/10.3390/diagnostics13142383