# An Ensembled Tabnet-Based Model Approach for Diabetes Disease Classification

## Obunge Duncan Ogindo[1]*, Muriira Lawrence[1], Mbandu Vincent[1]

[1] Kenya Methodist University, P.O. Box 267 - 60200, Meru, Kenya

*Correspondence email and telephone: duncanobunge@gmail.com

## Abstract

Despite the advancements in machine learning (ML) for classification tasks, accurately classifying diseases on limited-feature medical datasets remains challenging. Traditional ML models struggle with interpretability, necessitating an exploration of novel technique. This research developed and evaluated a novel TabNet-based ensemble model for diabetes classification, rating its performance against Extreme Gradient Boosting (XGBoost), Random Forest and base TabNet models. The study utilized the PIMA Indian Diabetes dataset from a public ML Repository, which contains 768 tuples (8 features and 1 outcome variable). A TabNet-based ensemble model was developed using a weighted averaging strategy. For comparative analysis, baseline models, including XGBoost, Random Forest, and a standalone TabNet model were also implemented and optimized. Model performance was assessed using key metrics: balanced accuracy, precision and recall (class 1), F1 score, and Receiver Operating Characteristic-Area Under the Curve (ROC-AUC). The ensembled TabNet-based model consistently achieved the highest performance metrics: balanced accuracy of 83%, precision of 84%, recall of 89%, F1 score of 84%, and ROC-AUC of 90.4% compared to XGBoost (accuracy 81%, precision 79%, recall 86%, F1 score 81%, ROC-AUC 88.6%), Random Forest (accuracy 81%, precision 78%, recall 87%, F1 score 81%, ROC-AUC 91.6%) and base TabNet (accuracy 81%, precision 80%, recall 82%, F1 score 81%, ROC-AUC 86.7%). The study recommends healthcare institutions to adopt the validated ensemble TabNet-based architecture as a standardized framework for clinical decision support systems across multiple diseases. Further, researchers should establish this methodology as the preferred approach for limited-feature medical datasets, extending beyond diabetes to include cardiovascular, hypertension, and cancer screening applications.

**Keywords**: *TabNet, Clinical Decision Support, Tabular Data, Ensemble Learning, Interpretability*

## 1.0 Introduction

Machine learning approaches for classification problems encounter continued challenges in maintaining better predictive performance and transparency in model's decision-making process, more especially in high stake domain like healthcare setting where both accuracy and interpretability are of critical demands (Chaddad et al., 2023). Traditional machine learning models based on XGBoost and Random Forest algorithms have established credible and promising performance on tabular data classification tasks. However, their decision-making process is unclear to users putting them under 'black-box' categories of model, while simpler interpretable models typically underperform in complex classification scenarios (Kelly et al., 2019). This study addresses this critical gap by developing a novel ensemble architecture leveraging TabNet based model, a neural network architecture for classification problem resolution, particularly using tabular data. The model has demonstrated capacity to maintain high performance, manage noisy and imbalanced data and enhance interpretability on its decision-making process, particularly in the classification tasks tabular dataset (Arik & Pfister, 2021).

The study employs weighted averaging ensembling technique, aggregates TabNet-based models and two traditional machine learning models based on XGBoost and Random Forest algorithms. The proposed solution targets an improved as well as a balanced classification performance based on various metrics, while ensuring a better feature attribution capabilities based on TabNet's sequential attention mechanisms (Arik & Pfister, 2021). Using diabetes classification task with the PIMA Indian dataset as a representative case , this study highlights the efficacy of the proposed ensemble strategy while making methodological contributions with prospects for application across a range of classification domains using tabular data (Contreras et al., 2020). The primary objective is to develop an ensemble TabNet-based model, TabNet model, XGBoost and Random Forest models, and conduct comprehensive comparative analysis of model performance in predicting diabetes based on balanced accuracy, precision (class 1), recall (class 1), F1 score, and ROC-AUC metrics.

Despite major progress in machine learning for tabular data classification, computational limitations persist in healthcare diagnostic applications where traditional models suffer from sensitivity to data quality, overfitting tendencies, and the fundamental trade-off between performance and interpretability (Kelly et al., 2019). While deep learning models have historically struggled to match gradient-boosted decision trees on tabular data, TabNet presents a promising solution through its attention-based architecture, though its base implementation faces challenges in full utilization of limited feature spaces, and capturing complex feature interactions. This research addresses these limitations by developing a novel ensemble architecture that harnesses TabNet's explainability advantages while enhancing classification performance, robustness, and feature utilization efficiency (Ahmed et al., 2022; Shah et al., 2022).

The study's general objective is to design, implement, and evaluate a novel ensemble architecture for TabNet models that improves tabular data classification performance, while maintaining quantifiable interpretability using

diabetes screening as a representative application context. Specific objectives include: (1) identifying and analyzing features that significantly contribute to diabetes classification in the PIMA Indian diabetes dataset through feature importance scoring and attention mask visualization, (2) developing an optimized weighted ensemble TabNet architecture; and (3) conducting comprehensive performance evaluation against traditional ML models based on both XGBoost and Random Forest algorithms across multiple metrics.

The research tests two primary hypotheses: first, that the proposed TabNet ensemble architecture achieves statistically significant improvements in classification accuracy and ROC-AUC scores compared to individual TabNet models and traditional ML models; second, that the ensemble approach maintains TabNet's interpretability advantages while improving performance, demonstrating effective management of the performance-interpretability trade-off through advanced ensemble techniques (Shah et al., 2022).

This research contributes to multiple domains, including algorithmic advancement through novel ensemble techniques for explainable deep learning architectures, and interpretability research by exploring how ensemble methods can maintain or enhance model transparency, and technical implementation through comprehensive evaluation frameworks assessing performance, interpretability, robustness, and calibration metrics (Vujovic, 2021). The study addresses the growing need for responsible AI in healthcare applications by developing models that provide both high predictive accuracy and transparent decision-making processes, potentially applicable across various tabular data

classification domains beyond diabetes prediction.

The study uses the PIMA Indian Diabetes dataset with 768 tuples and 8 features. The potential limitations of the study include a specific demographic aligned dataset, dataset limited features, and potential increase in computational resource demand, which is inherent with ensembling process. These limitations notwithstanding, the study assumes that the dataset adequately represents primary attributes associated with diabetic patients, and that TabNet's architecture caters for effective classification-ability with better interpretability in comparison to traditional ML models, setting up the basis for developing an interpretable ML application in healthcare context.

*"This research successfully developed and validated an ensemble TabNet-based architecture for diabetes prediction that meets established performance targets while maintaining interpretability.*

## 2.0 Materials and Methods

This study employed a quantitative experimental research design to develop and evaluate an ensemble TabNet-based model for diabetes prediction using the PIMA Indian diabetes dataset. The dataset, obtained from the University of California Irvine Machine

Learning Repository, contained 768 records of female patients aged 21 years and above with eight clinical features including pregnancies, glucose, blood pressure, skin thickness, insulin, BMI, diabetes pedigree function, and age. The dataset exhibited class imbalance with 268 positive instances (34.9%) and 500 negative instances (65.1%). Data preprocessing involved imputing missing values (represented as zeros) using Random Forest imputation technique applied separately for each outcome class, feature scaling using Robust Scaler, and applying SMOTE with k=5 nearest neighbors to address class imbalance. The dataset was stratified into training (60%), validation (15%), and testing (15%) sets.

Three base models were developed and optimized independently: TabNet with sequential decision steps, and attention mechanisms for interpretable feature selection, XGBoost configured with grid search optimization, and Random Forest with comprehensive parameter tuning. The TabNet model utilized key hyperparameters including decision steps (2-10), feature dimension (8-64), attention dimension (8-64), and AdamW optimizer with early stopping. The ensemble architecture employed weighted averaging technique combining predictions from all three base models using the formula:

$$P\_ensemble = w_1 \times P\_TabNet \\ + w_2 \times P\_XGBoost \\ + w_3 \times P\_RandomForest,$$

where weights were optimized through nested cross-validation on the validation set with the constraint $w_1 + w_2 + w_3 = 1$.

All models were trained using 5-fold cross-validation to ensure robust performance estimation and reduce bias. TabNet training utilized binary cross-entropy loss with sparsity regularization, L2 regularization (weight decay=1e-5), and dropout (rate=0.2) for 250 epochs with early stopping (patience=20). Hyperparameter optimization was conducted using grid search, with cross-validation across all models, to ensure optimal performance configurations.

Model performance was evaluated using five key metrics: balanced accuracy (average of sensitivity and specificity), precision for diabetes-positive cases, recall (sensitivity), F1-score (harmonic mean of precision and recall), and ROC-AUC. Statistical significance was assessed using paired t-tests with 95% confidence intervals, calculated using bootstrap methods. The experiment was implemented using Python 3.9 with PyTorch 2.7, scikit-learn 1.6.1, and PyTorch TabNet 4.1.0 on Ubuntu 22.04 system with reproducibility ensured through fixed random seeds and comprehensive documentation. The researcher acquired ethical clearance from university Research Ethics Committee. Additionally, bias mitigation strategies such as data imputation, balanced class distribution through SMOTE, and fairness evaluation across patient subgroups were applied. The researcher further obtained research permit from the National Commission for Science, Technology and Innovation (NACOSTI).

## 3.0 Results and Discussion

This study developed and evaluated an ensemble TabNet-based architecture for diabetes prediction using the PIMA Indian diabetes dataset, comparing its performance against traditional machine learning models (XGBoost and Random Forest). The research addressed two primary objectives: (1) to develop an optimized weighted ensemble TabNet architecture; and (2),

to compare performance evaluation across multiple models.

*Feature Importance Analysis*

Table 3 shows comprehensive feature importance analysis achieved by employing multiple complementary techniques, including TabNet attention mechanisms, SHAP, LIME, and Permutation importance. Glucose emerged as the dominant predictor across all methodologies, achieving a perfect normalized score of 1.0, and maintaining the top-ranking position. This finding validates established medical knowledge regarding glucose levels as the primary diagnostic criterion for diabetes mellitus.

**Table 1**

*Normalized Feature Importance Score and Average Ranking*

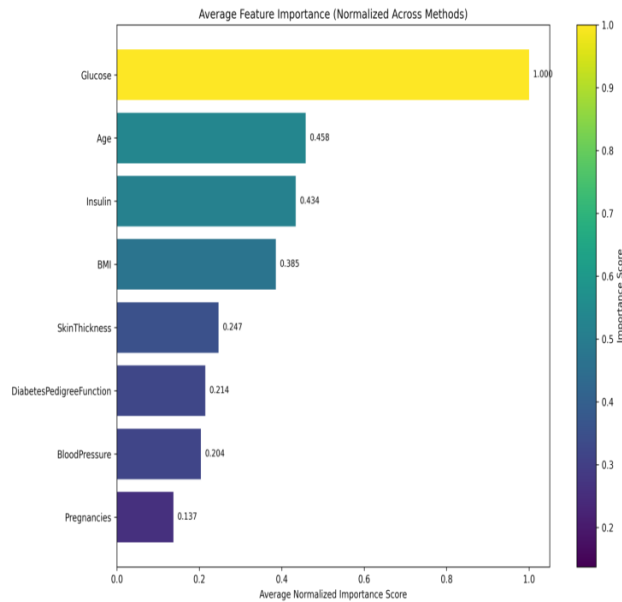| Feature | TabNet Score | SHAP Score | LIME Score | Perm Score | Avg Score | Avg Rank |
|---|---|---|---|---|---|---|
| **Glucose** | 1 | 1 | 1 | 1 | 1 | 1 |
| **Insulin** | 0.35 | 0.53 | 0.47 | 0.39 | 0.43 | 3 |
| **Age** | 0.58 | 0.59 | 0.33 | 0.32 | 0.46 | 3.5 |
| **BMI** | 0.03 | 0.45 | 0.72 | 0.35 | 0.39 | 4 |
| **Diabetes Pedigree Function** | 0 | 0.29 | 0.34 | 0.23 | 0.22 | 5.5 |
| **Skin Thickness** | 0.35 | 0.19 | 0.27 | 0.19 | 0.25 | 5.75 |
| **Blood Pressure** | 0.82 | 0 | 0 | 0 | 0.2 | 6.5 |
| **Pregnancy** | 0.21 | 0.13 | 0.1 | 0.11 | 0.14 | 6.75 |

*Note.* The scoring and ranking of features using multiple complementary techniques; Avg_score represent average score of each feature while Avg_rank represent average ranking of each feature.

Secondary predictors included Insulin (average normalized score: 0.40) and Age (average normalized score: 0.46), both achieving rankings of 3-3.5. The prominence of insulin aligns with diabetes pathophysiology, where insulin resistance plays a central role in disease development. Age is consistently ranked among top predictors, corresponding to clinical observations that the risk of diabetes increases with advancing age due to progressive insulin resistance, and beta-cell dysfunction. Figure 3 illustrates average scoring across different features.

**Figure 1**

*Average feature scoring analysis*



*Note.* Graphical representation of normalized average scoring of each feature.

### Ensemble TabNet Architecture Development

The weighted ensemble approach demonstrates superior performance compared to individual TabNet models, with improvements of 2.4% in F1-score and 2.7% in ROC-AUC. The hyperparameter optimization revealed that 5 decision steps provided optimal performance, with feature dimensions of 64 effectively capturing complex relationships with only 8 input features. The ensemble approach utilized weighted averaging to combine predictions from the TabNet model, and two traditional ML model instances, effectively reducing prediction variance while maintaining interpretability through attention mechanisms.

### Models Performance Analysis

Table 4 and Figure 4 demonstrate the performance of the four model configurations and their learning curves. While the weighted ensemble TabNet achieved competitive performance (F1: 0.8436, AUC: 0.9044), Random Forest based model demonstrated slightly higher overall performance on (F1: 0.8208, AUC: 0.9094). The ensemble TabNet showed superior recall (0.89 versus 0.87), which is crucial for medical screening applications where minimizing false negatives is paramount.
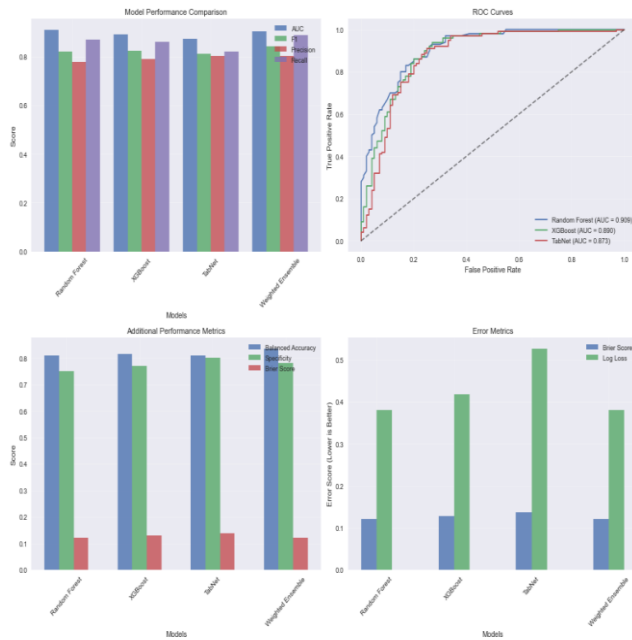
**Table 2**

*Comparative Performance Analysis*

| Model | AUC | F1-Score | Precision | Recall | Balanced Accuracy |
|---|---|---|---|---|---|
| Random Forest | 0.9094 ± 0.0121 | 0.8208 ± 0.0225 | 0.7768 ± 0.0312 | 0.8700 ± 0.0478 | 0.8100 ± 0.0250 |
| XGBoost | 0.8903 ± 0.0150 | 0.8230 ± 0.0207 | 0.7890 ± 0.0287 | 0.8600 ± 0.0267 | 0.8150 ± 0.0270 |
| Individual TabNet | 0.8730 ± 0.0122 | 0.8119 ± 0.0269 | 0.8039 ± 0.0312 | 0.8200 ± 0.0374 | 0.8100 ± 0.0132 |
| EnsembleTabNet | 0.9044 ± 0.0172 | 0.8436 ± 0.0110 | 0.8018 ± 0.0442 | 0.8900 ± 0.0341 | 0.8350 ± 0.0146 |

*Note.* Weighted ensemble model demonstrates balanced performance across different metrics

**Figure 2**

*Models' performance and learning curve analysis*



*Note.* Representation of model performance and learning curve during training.

Statistical significance testing using paired t-tests demonstrates that the weighted-ensemble significantly outperforms individual TabNet models across multiple metrics (Balanced Accuracy: p=0.0007, Recall: p=0.0097, AUC: p=0.0234); compared with traditional models which returned statistically insignificant metrics, indicating comparably better performance, while providing TabNet's interpretability advantages.

### Model Interpretability and Clinical Relevance

The ensemble TabNet maintained interpretability through attention mechanisms, while achieving enhanced performance. The model's ability to provide instance-wise feature importance enables clinicians to understand prediction rationale, addressing "black box" concerns in healthcare AI deployment (Ahmed et al., 2022). The balanced improvement in recall (0.89) and

precision (0.82) demonstrates the model's capacity to minimize both false negatives and false positives, critical for clinical decision support systems. The calibration analysis revealed excellent performance with a Brier score of 0.12, indicating that predicted probabilities closely match observed frequencies. This characteristic is particularly important for clinical applications where risk stratification is based on predicted probabilities and guides treatment decisions.

### Discussion and Implications

### Clinical Significance

Ensemble-TabNet presented better results compared to Individual TabNet with medium effect sizes (Cohen's d=1.89 for F1-score, d=0.71 for AUC). The 3.17% improvement in F1-score (0.8436 vs 0.8119), explains the major improvements observed, and translates to approximately 32 correct predictions per 1000 patients, compared to base implementation. While Ensemble-TabNet attained comparable AUC to Random Forest (0.9044 vs 0.9094, statistically insignificant, p=0.67), it attained superior recall (0.8900 vs 0.8700), representing 20 true positive predictions per 1000 positives. Clinically, the recall is important in screening where missing positive cases bears higher clinical risk.

### Methodological Contributions

This study contributes to the growing body of evidence supporting the application of interpretable deep learning models in healthcare. The successful implementation of ensemble TabNet architectures demonstrates that the benefits of ensemble learning extend beyond traditional machine learning, to modern deep

learning approaches for tabular data (Mohan Raparthy, 2023)

The comprehensive feature importance analysis using multiple complementary techniques provides robust validation of biomarker significance, with results aligning with established clinical knowledge. The methodological framework established here can be applied to other medical prediction tasks requiring both high performance and interpretability.

### *Limitations and Future Directions*

While the ensemble approach achieved target performance metrics, the computational overhead (4.73× training time increase) presents implementation challenges for resource-constrained environments. The study was conducted on a single dataset, and validation across diverse healthcare settings and populations would strengthen generalizability.

## 4.0 Conclusion

This research successfully developed and validated an ensemble TabNet-based architecture for diabetes prediction that meets established performance targets while maintaining interpretability. The comprehensive evaluation demonstrates that while traditional ensemble

## References

methods remain competitive, the TabNet-based approach offers unique advantages through attention-based interpretability and superior recall performance. The findings support the potential of interpretable deep learning models as valuable tools for clinical decision support, particularly in applications where both predictive accuracy and model transparency are essential. The study contributes to the evolving paradigm of explainable AI in healthcare, demonstrating that high-performance models need not sacrifice interpretability.

## 5.0 Recommendations

Clinical risk assessment protocol should prioritize glucose levels, age, insulin, and BMI as primary screening indicators in diabetes risk assessment. The model should be considered for deployment in clinical decision support systems to attain optimal balance between accuracy and interpretability. Additionally, clinical practice guidelines should be updated to capture evidenced-based feature prioritization from ensemble models.

Additionally, research funding agencies should prioritize projects that demonstrate both high performance and interpretability in healthcare AI applications.

Ahmed, I., Jeon, G., & Piccialli, F. (2022). From Artificial Intelligence to Explainable Artificial Intelligence in Industry 4.0: A Survey on What, How, and Where. *IEEE Transactions on Industrial Informatics*, *18*(8), 5031–5042. https://doi.org/10.1109/TII.2022.314655 2

Arik, S. Ö., & Pfister, T. (2021). TabNet: Attentive Interpretable Tabular Learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, *35*(8), 6679–6687. https://doi.org/10.1609/aaai.v35i8.16826

Chaddad, A., Peng, J., Xu, J., & Bouridane, A. (2023). Survey of Explainable AI Techniques in Healthcare. *Sensors*,

*23*(2), 634.
https://doi.org/10.3390/s23020634

Contreras, I., Bertachi, A., Biagi, L., Oviedo, S., Ramkissoon, C., & Vehi, J. (2020). Artificial intelligence-based decision support systems for diabetes. In *Artificial Intelligence in Precision Health* (pp. 329–357). Elsevier. https://doi.org/10.1016/B978-0-12-817133-2.00014-8

Kelly, C. J., Karthikesalingam, A., Suleyman, M., Corrado, G., & King, D. (2019). Key challenges for delivering clinical impact with artificial intelligence. *BMC Medicine*, *17*(1), 195. https://doi.org/10.1186/s12916-019-1426-2

vuMohan Raparthy, E. Al. (2023). Predictive Maintenance in IoT Devices using Time Series Analysis and Deep Learning. *Dandao Xuebao/Journal of Ballistics*, *35*(3), 01–10.

https://doi.org/10.52783/dxjb.v35.113

Rezaee, K., Savarkar, S., Yu, X., & Zhang, J. (2022). A hybrid deep transfer learning-based approach for Parkinson's disease classification in surface electromyography signals. *Biomedical Signal Processing and Control*, *71*, 103161. https://doi.org/10.1016/j.bspc.2021.103161

Shah, C., Du, Q., & Xu, Y. (2022a). Enhanced TabNet: Attentive Interpretable Tabular Learning for Hyperspectral Image Classification. *Remote Sensing*, *14*(3), 716. https://doi.org/10.3390/rs14030716

Vujovic, Ž. Đ. (2021). Classification Model Evaluation Metrics. *International Journal of Advanced Computer Science and Applications*, *12*(6). https://doi.org/10.14569/IJACSA.2021.0120670