

# A Threshold-Aware TabNet Framework for Imbalanced Bank Loan Default Prediction

Pallavi Ukey<sup>1</sup>, Shobha Rajak<sup>2</sup>, Prateek gupta<sup>3</sup>

<sup>1</sup>Research Scholar Department of CSE, Shriram Institute of Science & Technology, Jabalpur, M.P.

<sup>2</sup>Prof., Department of CSE, Shriram Institute of Science & Technology, Jabalpur, M.P.

<sup>3</sup>Prof., Department of CSE, Shriram Institute of Science & Technology, Jabalpur, M.P.

\*\*\*

**Abstract:** Accurate loan default prediction is a critical challenge for modern financial institutions, where class imbalance, non-linear feature interactions, and the high cost of misclassification demand sophisticated machine learning approaches. This paper proposes a Threshold-Aware TabNet framework for imbalanced binary classification of bank loan defaults, evaluated on a large-scale structured dataset of 252,000 applicant records containing 13 mixed-type features. The proposed framework integrates TabNet's sequential attention mechanism for automatic feature selection and interpretability with systematic decision threshold optimization to balance precision and recall according to real-world financial risk considerations. A comprehensive comparison is conducted against nine baseline models—including Logistic Regression, Random Forest, XGBoost, LightGBM, AdaBoost, Gradient Boosting, KNN, and ANN—under both raw and resampled training conditions. Performance is evaluated using a composite score weighted across ROC-AUC (30%), F1-score (30%), Recall (20%), and Precision (20%), intentionally excluding accuracy to avoid majority-class bias. Results demonstrate that the proposed TabNet model achieves the highest composite score (0.72), outperforming all baselines with a precision of 0.58, recall of 0.78, F1-score of 0.61, and ROC-AUC of 0.89 under threshold tuning. The study further demonstrates that conventional accuracy-based evaluation is insufficient for imbalanced credit risk datasets and that threshold optimization provides a practically deployable mechanism for risk-sensitive banking applications.

**Keywords:** TabNet, loan default prediction, class imbalance, threshold tuning, credit risk, deep learning for tabular data, XGBoost, feature attention, ROC-AUC, financial risk.

## I. INTRODUCTION

The rapid expansion of digital banking and the increasing volume of loan applications have created an urgent need for automated, data-driven credit evaluation systems. Non-performing loans (NPLs) represent one of the most significant sources of financial risk for banks: studies indicate that even a 5% increase in default misclassifications can lead to a 10% rise in NPLs, translating into substantial capital losses [1]. In the Indian banking context, the loan approval process involves multi-step assessment of credit scores, income stability, employment type, debt-to-income ratios, and collateral — a process increasingly suited to machine learning augmentation.

Despite significant research into machine learning for credit scoring, three persistent challenges remain inadequately addressed. First, real-world loan datasets are highly imbalanced: non-defaulters typically account for 85–90% of records, causing most classifiers to optimize for majority-class accuracy while failing to detect high-risk borrowers — the very cases most critical to financial institutions. Second, standard neural networks generally underperform on structured tabular data compared to tree-based ensemble methods, limiting the applicability of deep learning in this domain. Third, evaluation protocols

that rely on overall accuracy produce misleading performance claims, obscuring poor minority-class detection.

This paper addresses all three challenges through a unified Threshold-Aware TabNet framework. TabNet [4], a sequential attention-based architecture specifically designed for tabular data, is selected as the core model due to its ability to perform automatic feature selection, capture non-linear interactions without manual feature engineering, and provide interpretable decision masks. Combined with systematic decision threshold tuning, the proposed framework enables flexible control over the precision-recall trade-off in deployment, adapting to the asymmetric misclassification costs characteristic of credit risk scenarios.

The contributions of this paper are: (1) a comprehensive preprocessing and feature encoding pipeline for a large mixed-type loan dataset (252,000 records, 13 features); (2) a threshold-tuned TabNet classification framework achieving the best composite performance score among nine baselines; (3) a weighted composite evaluation metric (ROC-AUC, F1, Recall, Precision) that avoids accuracy bias; and (4) practical insights on class imbalance handling strategies across model families.

## II. RELATED WORK

### A. Machine Learning for Loan Default Prediction

Classical machine learning approaches including Logistic Regression, Decision Trees, SVM, and Random Forest have been widely applied to credit risk assessment. Bhargav and Sashirekha demonstrated that Random Forest (79.44% precision) outperforms Decision Tree (67.28%) for loan approval prediction [43]. Rath et al. confirmed Random Forest superiority (80% accuracy) over LR (73%), DT (79%), and SVM (75%) [53]. Ensemble methods consistently outperform single classifiers, with XGBoost, LightGBM, and Gradient Boosting achieving strong performance on imbalanced financial datasets [5, 6].

### B. Deep Learning Approaches

Deep learning models—CNN, LSTM, and MLP—have been applied to loan risk assessment but typically require extensive tuning and may underperform tree-based methods on tabular data. Wang et al. proposed a stacking-based deep learning model achieving a 6% improvement in joint loan approval on imbalanced data [44]. Khashman (2009) developed a backpropagation neural network for credit risk evaluation, achieving competitive results but limited interpretability [49]. ANN-based approaches often suffer from majority-class dominance without explicit imbalance handling, resulting in near-zero minority-class recall [61].

### C. Research Gap

Despite extensive research, key limitations persist across existing work: (i) most studies focus on binary accuracy without macro-level precision-recall analysis; (ii) SMOTE is frequently applied before train-test splitting, inflating results through data leakage; (iii) deep learning models with attention mechanisms for tabular data (e.g., TabNet) remain underexplored in the bank loan context; and (iv) systematic threshold optimization as a deployment strategy is largely absent. This work directly addresses these gaps.

## III. DATASET DESCRIPTION

The dataset used in this study is a publicly available structured loan application dataset containing 252,000 records and 13 features per applicant, sourced from Kaggle. The target variable Risk Flag is binary: 0 denotes non-defaulter (low risk) and 1 denotes defaulter (high risk). The dataset exhibits severe class imbalance — approximately 176,803 non-defaulters (~88%) versus 24,797 defaulters (~12%) — closely reflecting real-world credit portfolios.

Features include: Income (annual), Age, Experience (years), Married/Single (marital status), House Ownership (rented/owned), Car Ownership (yes/no), Profession (20+ categories), CITY (50+ categories), STATE (20+ categories), CURRENT\_JOB\_YRS, and CURRENT\_HOUSE\_YRS. The mixture of low-cardinality categorical features (marital status, car ownership), high-

cardinality categorical features (profession, city, state), and continuous numerical features makes this dataset an ideal tested for tabular deep learning evaluation.

Key dataset characteristics: mixed feature types (numerical and categorical), presence of missing values, high-cardinality categorical columns requiring label encoding, strong class imbalance, and right-skewed distributions in financial attributes (income, loan-related features). These characteristics directly motivated the preprocessing pipeline and model selection described in Section IV.

## IV. PROPOSED METHODOLOGY

### A. Preprocessing Pipeline

A systematic preprocessing pipeline was applied to transform raw records into model-ready feature vectors. Missing numerical values were imputed using column-wise medians, selected for outlier robustness. Missing categorical values were replaced with an explicit 'Unknown' category. Low-cardinality categorical features (Married/Single, House\_Ownership, Car\_Ownership) were encoded via One-Hot Encoding to preserve category independence. High-cardinality features (Profession, CITY, and STATE) were encoded using Label Encoding to avoid dimensionality explosion while retaining categorical information.

Numerical features (Income, Age, Experience, CURRENT\_JOB\_YRS, CURRENT\_HOUSE\_YRS) were standardized using Z-score normalization (zero mean, unit variance), essential for gradient-based optimization in TabNet and neural models. One-Hot encoded boolean columns were converted from True/False to 0/1. The preprocessed dataset was stratified-split into 80% training and 20% testing subsets, ensuring equal class representation in both partitions. SMOTE was applied exclusively within the training set for models lacking native imbalance handling (KNN, basic ANN), preventing synthetic sample leakage into evaluation data.

### B. Imbalance Handling Strategy

Imbalance was handled through model-appropriate mechanisms: (i) `class_weight='balanced'` for Logistic Regression, SVM, Random Forest, and LightGBM; (ii) `scale_pos_weight=ratio` for XGBoost; (iii) manual oversampling via SMOTE for KNN and ANN. TabNet was evaluated under threshold tuning without resampling, relying on its attention mechanism and the optimization of the decision boundary to handle the imbalanced distribution. All imbalance-handling techniques were applied strictly within the training partition.

### C. TabNet Architecture

TabNet implements a sequential attention mechanism that processes features through  $N$  decision steps, each step selecting a sparse subset of features via an Attentive Transformer using Sparsemax activation. The feature transformer at each step applies shared and step-specific

fully connected layers with Batch Normalization and Gated Linear Units (GLU), enabling efficient gradient flow. The final output aggregates representations across all decision steps:

$$y = \sum_t f(M_t \odot X),$$

Where  $M_t$  denotes the sparse attention mask at step  $t$  and  $f(\cdot)$  is the feature transformer. A sigmoid output layer produces class probabilities  $P(y=1|X)$  for binary classification. The model was trained using the Adam optimizer with binary cross-entropy loss, early stopping on validation ROC-AUC, and a maximum of 100 epochs. This architecture eliminates manual feature engineering while providing inherent interpretability through the learned attention masks.

#### D. Threshold Optimization

Due to the asymmetric misclassification costs in credit risk (false negatives — missed defaulters — carry higher financial penalties than false positives), the default classification threshold of 0.5 is suboptimal. The proposed framework evaluates probability thresholds ranging from 0.10 to 0.55, computing Accuracy, Precision, Recall, F1-score, and ROC-AUC at each threshold. The optimal threshold is selected based on application-specific requirements: high-recall configuration for risk-sensitive deployment, or balanced F1 for general-purpose use.

#### E. Evaluation Framework

Models are ranked using a weighted composite score that excludes accuracy to avoid majority-class bias:

$$\text{Score} = 0.30 \times \text{ROC-AUC} + 0.30 \times \text{F1} + 0.20 \times \text{Recall} + 0.20 \times \text{Precision}$$

ROC-AUC (30%) captures overall class separation independent of threshold. F1-score (30%) balances precision and recall. Recall (20%) prioritizes detection of actual defaulters. Precision (20%) controls false alarm rates. All models are evaluated on identical stratified test sets using probability-based predictions, enabling threshold-consistent comparison.

### V. EXPERIMENTAL RESULTS

#### A. Accuracy Comparison

Table I presents best-accuracy results across all evaluated models. Most models converge near 88% accuracy due to majority-class dominance — a model predicting all non-defaults would achieve 88% without learning any discriminative patterns. The proposed TabNet model achieves a marginally superior accuracy of 89.07%, confirming that it maintains competitive overall correctness while simultaneously improving minority-class detection.

TABLE I: ACCURACY COMPARISON ACROSS MODELS

Model	Strategy	Best Acc.
Logistic Regression	Baseline (no balancing)	~0.88
Gradient Boosting	No oversampling	0.88
ANN	No balancing	0.88
XGBoost	Cost-sensitive + tuning	0.88
TabNet (Proposed)	Threshold tuning (0.55)	0.89 ★

#### B. Precision and ROC-AUC Analysis

Beyond accuracy, the critical distinction between models emerges in precision and ROC-AUC. Gradient Boosting achieves the highest individual precision (0.72) but at the cost of near-zero recall (0.17), making it operationally ineffective for defaulter detection despite high precision. Logistic Regression achieves perfect recall (1.00) but precision collapses to 0.14, indicating excessive false positives. The proposed TabNet achieves a balanced precision of 0.58 and recall of 0.78, indicating that it correctly identifies 78% of all actual defaulters while maintaining acceptable false-positive rates.

ROC-AUC analysis reveals that Random Forest achieves the highest single-metric discrimination (0.94), followed by XGBoost (0.90) and TabNet (0.89). ANN without balancing collapses to ~0.50 ROC-AUC — equivalent to random guessing — confirming that unbalanced neural network training fails entirely under this class distribution.

#### C. Composite Ranking

Table II presents the composite ranking across all models using the weighted scoring formula. TabNet ranks first (score: 0.72) due to its balanced combination of high recall, competitive precision, and strong ROC-AUC without requiring synthetic data generation. Random Forest ranks second (0.71) with the highest ROC-AUC but slightly lower recall and precision than TabNet. XGBoost ranks third (0.70), achieving the highest recall but slightly lower precision balance.

TABLE II: FINAL COMPOSITE RANKING (ALL MODELS)

#	Model	Prec.	Recall	F1	AUC	Score
1	TabNet (Proposed)	0.58	0.78	0.61	0.89	0.72
2	Random Forest	0.55	0.76	0.64	0.94	0.71
3	XGBoost	0.51	0.81	0.62	0.90	0.70
4	LightGBM	0.51	0.79	0.59	0.87	0.67
5	KNN	0.46	0.87	0.60	0.87	0.66
6	Gradient Boosting	0.72	0.17	0.25	0.70	0.47
7	ANN (Oversampled)	0.16	0.78	0.27	0.61	0.43

#	Model	Prec.	Recall	F1	AUC	Score
8	ANN (Baseline)	0.64	0.00	0.00	0.50	0.32
9	Logistic Reg.	0.14	1.00	0.22	0.55	0.39

$$Score = 0.30 \times AUC + 0.30 \times F1 + 0.20 \times Recall + 0.20 \times Precision$$

#### D. Threshold Sensitivity Analysis

TabNet probability outputs were evaluated across thresholds from 0.10 to 0.55. At threshold 0.55, the model achieves peak accuracy (89.07%) while maintaining the best F1-score (0.61) and recall (0.78) combination. Lowering the threshold below 0.30 increases recall toward 0.90+ but reduces precision to below 0.40, significantly increasing false alarms. For high-risk banking environments where missed defaults carry greater penalties, a threshold of 0.35–0.40 offers superior recall (>0.85) with acceptable precision (~0.48). This threshold flexibility is a practical advantage of the proposed framework over models relying on fixed 0.5 boundaries.

#### E. Why Accuracy is Misleading

The fundamental limitation of accuracy-based evaluation on imbalanced datasets is illustrated by comparing ANN Baseline and TabNet: ANN Baseline achieves 88% accuracy and 0.64 precision but delivers 0.00 recall and 0.00 F1 — it never predicts a single defaulter. TabNet also achieves ~88–89% accuracy but with 0.78 recall and 0.61 F1, representing a categorically different and practically useful model. This analysis confirms that accuracy is an unreliable metric for imbalanced credit risk classification and that composite metrics weighted toward ROC-AUC and F1 are essential for honest evaluation.

### VI. DISCUSSION

The experimental results yield several practically significant insights for credit risk modeling. First, tree-based ensemble methods (Random Forest, XGBoost, and LightGBM) provide strong baseline performance under class weighting, benefiting from their inherent resistance to class imbalance through split-criterion adjustments. However, they exhibit a fixed precision-recall trade-off that cannot be dynamically adjusted post-training without threshold manipulation, which may not always be straightforward.

Second, TabNet's sequential attention mechanism provides a distinct advantage in feature utilization for high-cardinality mixed-type tabular datasets. The model dynamically prioritizes the most informative features at each decision step — such as Income, CURRENT\_JOB\_YRS, and Experience — while attenuating noise from high-cardinality categorical features like CITY and STATE. This selective attention capability reduces overfitting risk and improves generalization without requiring explicit feature selection preprocessing.

Third, threshold tuning represents a practical deployment mechanism that transforms a probability estimator into a risk-calibrated decision system. The proposed framework

explicitly separates model training from deployment configuration, allowing banking professionals to select threshold levels that align with institutional risk tolerance — a significant operational advantage over fixed-threshold classifiers.

Fourth, the study reveals that SMOTE, while beneficial for models without native imbalance handling, can be counterproductive for tree-based models that already implement class weighting. Gradient Boosting and ANN with SMOTE showed lower performance than their class-weighted counterparts, confirming that oversampling introduces synthetic noise that degrades decision boundaries when class weighting already provides adequate minority-class focus.

### VII. CONCLUSION

This paper presented a Threshold-Aware TabNet framework for imbalanced bank loan default prediction, evaluated on a 252,000-record mixed-type dataset with severe class imbalance (~88%/12%). The proposed framework integrates TabNet's sequential attention architecture with systematic threshold optimization, achieving the highest composite performance score (0.72) among nine evaluated baselines — including Random Forest, XGBoost, LightGBM, Gradient Boosting, AdaBoost, KNN, and ANN variants.

Key findings are: (1) TabNet achieves the best balance of Precision (0.58), Recall (0.78), F1 (0.61), and ROC-AUC (0.89) among all models; (2) accuracy is a misleading metric for imbalanced credit datasets — ANN Baseline achieves 88% accuracy with 0.00 recall; (3) threshold tuning provides a practically deployable mechanism for risk-sensitive classification; (4) class weighting is generally preferable to SMOTE for tree-based models; and (5) a composite evaluation metric excluding accuracy is essential for meaningful model comparison in credit risk research.

Future work will investigate cost-sensitive TabNet training with domain-specific misclassification penalties, integration of SHAP-based counterfactual explanations for regulatory compliance, temporal modeling of borrower behavior using Temporal TabNet or LSTM-TabNet hybrids, and cross-institutional dataset validation to confirm generalizability across diverse lending environments.

### REFERENCES

- [1] M. Vasheghani, E. N. Farokhi, B. Dolatshah, "Forecasting loan, deferred rate and customer segmentation in banking industry: a computational intelligence approach," *Array*, vol. 27, 2025.
- [2] N. Uddin et al., "An ensemble machine learning based bank loan approval predictions system with a smart application," *Int. J. Cogn. Comput. Eng.*, vol. 4, pp. 327–339, 2023.
- [3] L. Hota, P. K. Jain, A. Kumar, "A comparative performance assessment for prediction of loan

approval in financial sector," *Procedia Comput. Sci.*, vol. 258, pp. 298–307, 2025.

- [4] S. O. Arik and T. Pfister, "TabNet: Attentive interpretable tabular learning," *Proc. AAAI Conf. Artif. Intell.*, vol. 35, pp. 6679–6687, 2021.
- [5] H. A. Alamri and V. Thayananthan, "Bandwidth control and extreme gradient boosting for DDoS in SDN," *IEEE Access*, vol. 8, 2020.
- [6] C. N. Nwafor, O. Z. Nwafor, "Determinants of non-performing loans: an explainable ensemble and deep neural network approach," *Finance Res. Lett.*, vol. 56, 2023.
- [7] D. P. Singh et al., "Predictive modeling for bank loan approval: from data to decisions," *Procedia Comput. Sci.*, vol. 259, pp. 1426–1431, 2025.
- [8] S. Zandi et al., "Attention-based dynamic multilayer graph neural networks for loan default prediction," *Eur. J. Oper. Res.*, vol. 321, no. 2, pp. 586–599, 2025.
- [9] V. L. Hess and B. Damasio, "Machine learning in banking risk management: mapping a decade of evolution," *Int. J. Inf. Manag. Data Insights*, vol. 5, no. 1, 2025.
- [10] M. Tavakoli et al., "Multi-modal deep learning for credit rating prediction using text and numerical data streams," *Appl. Soft Comput.*, vol. 171, 2025.
- [11] M. Madaan et al., "Loan default prediction using decision trees and random forest: a comparative study," *IOP Conf. Ser. Mater. Sci. Eng.*, vol. 1022, 2021.
- [12] G. B. Rath, D. Das, and B. Acharya, "Modern approach for loan sanctioning in banks using machine learning," *Adv. Mach. Learn. Comput. Intell.*, pp. 179–188, Springer, 2021.
- [13] P. Bhargav and K. Sashirekha, "A machine learning method for predicting loan approval by comparing the random forest and decision tree algorithms," *J. Surv. Fish. Sci.*, vol. 10, no. 1S, pp. 1803–1813, 2023.
- [14] A. Khashman, "A neural network model for credit risk evaluation," *Int. J. Neural Syst.*, vol. 19, no. 4, pp. 285–294, 2009.
- [15] M. Abdullah et al., "Forecasting nonperforming loans using machine learning," *J. Forecast.*, 2023.

right now she is looking to upgrade herself in a field of teaching.

Author2- Shobha Rajak

Description- Prof. Shobha Rajak is an Assistant Professor in Computer science and engineering department at the Shri Ram Institute of Science and Technology, with 15 years of teaching experience in higher education. She specializes in Machine Learning and Deep Learning, with a strong interest in research, academic mentoring, and emerging technologies in Artificial Intelligence. She is actively involved in teaching, project guidance, and research activities in the field of computer science and intelligent systems.



## BIOGRAPHIES

Author1- Pallavi ukey

Description - Pallavi Ukey is a Mtech final year student of Computer science and engineering department of the Shri Ram Institute of Science and Technology Jabalpur. She is actively involved in teaching as a lecturer of diploma's students from last 6 years in the field of computer science dependent. And

