



MALETE JOURNAL OF ACCOUNTING AND FINANCE

A Publication of the Department of Accounting and Finance, Faculty of Management and Social Sciences, Kwara State University, Malete

Vol. 6 No. 2, December, 2025

ISSN: 2735-9603

BALANCING ACCURACY AND INTERPRETABILITY IN CREDIT RATING MODELS: A COMPARATIVE ANALYSIS OF STATISTICAL AND MACHINE LEARNING METHODS IN NIGERIA

Adedeji Daniel GBADEBO

Department of Accounting Science, Walter Sisulu University, South Africa
agbadebo@wsu.ac.za

Received: November 2025 / Revised: November 2025 / Accepted: December 2025

© 2025 Department of Accounting and Finance, Kwara State University. All rights reserved.

Abstract

Credit rating remains one of the most critical mechanisms in financial markets, serving as a benchmark for investment decisions, borrowing costs, and regulatory compliance. This study investigates the predictive capacity of machine learning and traditional statistical methods for corporate credit rating in Nigeria, an emerging economy where reliable credit risk assessment remains central to financial stability and capital allocation. Adopting an explanatory research design, the study utilizes firm-level panel data from a purposive sample of 89 publicly listed companies operating across manufacturing, financial, and service sectors over the period 2015–2023. Using data from 89 companies across three major sectors, the study evaluates the performance of six models, including logistic regression, support vector machines, random forest, XGBoost, decision trees, and k-nearest neighbours, based on classification accuracy, sensitivity, specificity, precision, and Matthews correlation coefficient. Empirical results reveal that while advanced machine learning algorithms such as random forest and XGBoost perform competitively, logistic regression demonstrates consistent interpretability and regulatory suitability, particularly when applied to capital adequacy and profitability indicators. Sensitivity analyses and post-hoc Tukey HSD tests suggest that differences across models are not always statistically significant, underscoring the role of data quality and multidimensional indicators in determining predictive success. The study concludes that hybrid credit-rating frameworks that integrate interpretable statistical models with machine learning techniques offer the most practical balance between transparency and predictive power. It recommends that regulators and financial institutions in Nigeria prioritize data enrichment strategies, encourage explainable modeling approaches, and incorporate alternative data sources to enhance corporate credit assessment and support financial stability.

Keywords: Credit Rating, Machine Learning, Logistic Regression, Financial Stability, Corporate Risk, Nigeria

JEL Codes: C45, G21, G32, O16

INTRODUCTION

Credit rating remains one of the most critical mechanisms in financial markets, serving as a benchmark for investment decisions, borrowing costs, and regulatory compliance. In emerging economies such as Nigeria, the need for reliable credit risk assessment is heightened by volatile macroeconomic conditions, limited financial transparency, and structural weaknesses in capital markets. The inability to generate accurate and transparent credit ratings can lead to inefficient capital allocation, higher default risks, and systemic vulnerabilities. These challenges are particularly pronounced in corporate credit markets, where information asymmetries and

Gbadebo, A. D. (2025). Balancing Accuracy and Interpretability in Credit Rating Models: A Comparative Analysis of Statistical and Machine Learning Methods in Nigeria. *Malet Journal of Accounting and Finance*, 6 (2), 212 – 225

heterogeneous firm characteristics complicate risk evaluation. As global markets increasingly rely on quantitative techniques for risk assessment, the integration of machine learning into credit rating prediction has gained traction, offering new avenues for improving predictive accuracy and decision-making (Lessmann et al., 2021; Koc et al., 2023).

Traditional statistical models, such as logistic regression and discriminant analysis, have long been employed to predict default probability due to their interpretability and regulatory acceptance. However, these models often assume linearity and struggle with complex, non-linear financial data that characterize corporate performance in developing markets. Machine learning algorithms, such as random forest, support vector machines, and XGBoost, address these limitations by learning patterns from large, multi-dimensional datasets without relying on strict parametric assumptions. Despite their methodological appeal, concerns persist regarding their robustness, transparency, and suitability in regulatory environments where data availability is constrained and supervisory scrutiny is high (Bücker et al., 2020; Marqués et al., 2022).

In the context of emerging markets, the debate surrounding model choice extends beyond predictive accuracy to encompass interpretability, fairness, and systemic implications. While ensemble and boosting techniques have demonstrated superior performance in many empirical studies (Back et al., 2021; Ekin et al., 2021), their “black-box” nature raises questions about explainability, auditability, and potential bias in credit allocation. This issue is particularly salient in financial systems with limited consumer protection frameworks, where opaque algorithmic decisions may exacerbate exclusion and undermine trust in financial institutions (Bastani et al., 2021; Bartlett et al., 2022). Consequently, regulators increasingly emphasize the need for transparent and explainable credit models that align with prudential regulation and macro-financial stability objectives (Allen et al., 2022; Basel Committee on Banking Supervision, 2021).

At the same time, advances in financial technology have expanded the range of data available for credit assessment, creating opportunities to improve model performance in data-scarce environments. The incorporation of alternative data—such as digital payment records, mobile money transactions, and other non-traditional indicators—has been shown to enhance credit scoring accuracy and promote financial inclusion, particularly for firms and individuals underserved by conventional banking systems (Babajide et al., 2021; Jagtiani & Lemieux, 2019). However, the effectiveness of such innovations ultimately depends on the interaction between data richness, model architecture, and institutional capacity, rather than algorithmic sophistication alone (Lessmann et al., 2021; Marqués et al., 2022).

Against this backdrop, this study situates itself at the intersection of methodological innovation and policy relevance by examining the application of both traditional statistical and machine learning models for corporate credit rating prediction in Nigeria. By focusing on firm-level financial indicators across key economic sectors, the study provides empirical evidence on how different modeling approaches perform under the practical constraints typical of emerging markets. In doing so, it contributes to ongoing debates on whether advanced machine learning techniques deliver meaningful gains over simpler, more interpretable models when data limitations and regulatory considerations are taken into account (Abellán & Castellano, 2022; Koc et al., 2023).

Accordingly, the objectives of this study are threefold: (i) to compare the predictive performance of traditional statistical and machine learning models in corporate credit rating prediction within the Nigerian context; (ii) to assess the trade-off between predictive accuracy and interpretability with respect to regulatory suitability; and (iii) to derive policy-relevant insights on how credit rating models can support financial stability and inclusion in emerging markets.

LITERATURE REVIEW

Theoretical Review

Credit rating and credit risk assessment are grounded in theories of information asymmetry, financial intermediation, and risk pricing. In credit markets, lenders face imperfect information about borrowers' true financial conditions, which creates adverse selection and moral hazard problems. Credit rating models serve as mechanisms to mitigate these information asymmetries by transforming firm-level financial and non-financial data into quantitative assessments of default risk (Altman & Sabato, 2007; Back et al., 2021). From this theoretical standpoint, the effectiveness of a credit rating system depends not only on its predictive accuracy but also on its transparency and credibility, particularly in regulated financial environments.

Traditional statistical models, such as logistic regression, are rooted in classical econometric theory, where model interpretability, parameter stability, and inferential clarity are emphasized. These models assume a structured relationship between explanatory variables and default outcomes, allowing analysts and regulators to directly observe marginal effects and economic significance (Altman et al., 2020). Theoretically, such models align with prudential regulation frameworks that prioritize explainability and auditability, as transparent risk models enhance supervisory oversight and reduce systemic uncertainty (Basel Committee on Banking Supervision, 2021; Allen et al., 2022). However, their reliance on linearity and distributional assumptions constrains their ability to capture complex, non-linear interactions commonly observed in corporate financial data, particularly in emerging markets.

Machine learning theory offers an alternative paradigm that emphasizes predictive performance through flexible function approximation and data-driven learning. Algorithms such as random forest, support vector machines, and gradient boosting are theoretically designed to minimize prediction error by exploiting non-linear patterns, interactions, and high-dimensional feature spaces (Breiman, 2001; Chen & Guestrin, 2016). In credit risk theory, these models are justified by the notion that default behavior is rarely driven by single indicators but instead emerges from complex combinations of liquidity, solvency, profitability, and macro-financial conditions (Huang et al., 2020; Addo et al., 2018). Empirical benchmarking studies reinforce this theoretical argument by demonstrating that ensemble and boosting methods often outperform traditional models in terms of classification accuracy (Lessmann et al., 2021; Back et al., 2021).

Despite their theoretical strengths in prediction, machine learning models introduce challenges related to interpretability and governance. From a theoretical governance perspective, opaque "black-box" models can weaken accountability and amplify informational asymmetries between financial institutions, regulators, and borrowers (Bücker et al., 2020). This concern is consistent with broader theories of algorithmic fairness and discrimination, which argue that complex models may unintentionally embed biases present in historical data, thereby affecting credit allocation outcomes (Kleinberg et al., 2018; Bartlett et al., 2022; Fuster et al., 2022). Consequently,

Gbadebo, A. D. (2025). Balancing Accuracy and Interpretability in Credit Rating Models: A Comparative Analysis of Statistical and Machine Learning Methods in Nigeria. *Malet Journal of Accounting and Finance*, 6 (2), 212 – 225

the theoretical trade-off between accuracy and interpretability becomes central: while machine learning improves predictive efficiency, it may undermine transparency, fairness, and trust if not carefully governed.

Recent theoretical developments attempt to reconcile this tension through explainable artificial intelligence (XAI) frameworks, which aim to retain predictive power while enhancing model interpretability. Explainability is increasingly viewed as a necessary condition for the institutional legitimacy of machine learning in credit markets, particularly under stringent regulatory regimes (Delange et al., 2022; Xu et al., 2023). From this perspective, interpretability is not merely a technical attribute but an economic and regulatory requirement that supports financial stability and consumer protection (Allen et al., 2022). This theoretical shift suggests that the optimal credit rating framework may not rely exclusively on either statistical or machine learning models but rather on hybrid approaches that integrate interpretability with predictive strength.

Theoretical insights from financial inclusion literature further expand the scope of credit rating models by emphasizing the role of alternative data. In environments characterized by limited formal financial histories traditional credit models may systematically exclude viable borrowers (Babajide et al., 2021; Jagtiani & Lemieux, 2019). Machine learning theory supports the inclusion of alternative data sources by enabling the processing of large, unstructured datasets, thereby improving risk discrimination for underserved firms and sectors (Liang et al., 2022; Wang et al., 2023). However, theory also cautions that data richness alone does not guarantee better outcomes; issues such as class imbalance, data noise, and overfitting can distort predictions if not properly addressed (Marqués et al., 2022; Zięba et al., 2020).

Empirical Review

Early empirical research on credit scoring and credit risk prediction was dominated by traditional statistical approaches, particularly logistic regression and discriminant analysis, which were valued for their interpretability and alignment with regulatory requirements. Seminal empirical work by Altman and Sabato (2007) demonstrated that accounting-based financial ratios could effectively predict default risk, especially for small and medium-sized enterprises. Subsequent empirical benchmarking studies reinforced the usefulness of parametric models in structured financial environments but also highlighted their limitations in handling non-linearity and interaction effects inherent in firm-level financial data (Moro et al., 2016). These early studies established the foundation for later comparisons between econometric and machine learning methods.

From the mid-2010s onward, empirical attention increasingly shifted toward machine learning techniques capable of capturing complex patterns in credit data. The introduction of ensemble methods, particularly random forest (Breiman, 2001) and gradient boosting algorithms such as XGBoost (Chen & Guestrin, 2016), marked a turning point in empirical credit risk research. Studies applying these methods reported improvements in predictive accuracy and robustness, especially in datasets characterized by non-linear relationships and noisy predictors (Addo et al., 2018; Heaton et al., 2017). However, early empirical comparisons already suggested that performance gains were not uniform across datasets and depended strongly on feature construction and sample characteristics, foreshadowing later benchmarking results.

Gbadebo, A. D. (2025). Balancing Accuracy and Interpretability in Credit Rating Models: A Comparative Analysis of Statistical and Machine Learning Methods in Nigeria. *Malet Journal of Accounting and Finance*, 6 (2), 212 – 225

Between 2018 and 2020, empirical research expanded rapidly, focusing on systematic comparisons of classifiers and the practical challenges of credit data. Studies documented that class imbalance, stemming from the relatively low frequency of default events, substantially distorted accuracy-based evaluations, motivating the use of alternative metrics such as sensitivity, precision, and the Matthews Correlation Coefficient (Zięba et al., 2020; Marqués et al., 2022). Empirical evidence further showed that resampling techniques and ensemble learners often improved minority-class detection without materially degrading overall model performance (Addo et al., 2018; Zięba et al., 2020). These findings underscored that empirical model performance is jointly determined by algorithm choice and data preprocessing strategies. Empirical studies published around 2020–2021 increasingly emphasized large-scale benchmarking and methodological rigor. Lessmann et al. (2021) and Back et al. (2021) provided comprehensive comparative analyses showing that no single algorithm consistently dominates across credit datasets. Instead, the relative performance of logistic regression, support vector machines, random forest, and boosting methods was found to vary with data dimensionality, feature richness, and hyper-parameter tuning. A recurring empirical conclusion from this period is that when models are carefully tuned and evaluated using robust cross-validation protocols, performance differences between sophisticated machine learning models and simpler statistical approaches often narrow substantially (Ekin et al., 2021).

More recent empirical work has placed stronger emphasis on explainability, governance, and regulatory acceptability. Studies demonstrate that while ensemble and deep learning models can achieve high predictive accuracy, their opaque nature raises concerns for auditability and fairness in credit allocation (Bücker et al., 2020). Empirical case studies show that explainable AI tools can partially mitigate these concerns by enhancing transparency without fully sacrificing predictive performance (Delange et al., 2022; Xu et al., 2023). This stream of literature empirically confirms that interpretability has become a decisive factor influencing model adoption in regulated financial environments, particularly under Basel III-oriented supervisory frameworks (Basel Committee on Banking Supervision, 2021).

The most recent empirical studies extend credit rating research toward financial inclusion, fairness, and emerging market applications. Evidence from fintech lending and African contexts shows that integrating alternative data and machine learning can improve credit access for underserved firms and populations, though risks related to bias and discrimination persist (Jagtiani & Lemieux, 2019; Babajide et al., 2021; Bartlett et al., 2022). Studies also document that algorithmic bias can amplify existing inequalities if not carefully managed, reinforcing the need for fairness-aware and interpretable modeling approaches (Kleinberg et al., 2018; Fuster et al., 2022). Collectively, recent empirical findings suggest that the future of credit rating models lies not in algorithmic complexity alone, but in balancing predictive accuracy with interpretability, fairness, and regulatory compliance.

METHODOLOGY

This study adopts a quantitative, explanatory research design, consistent with prior empirical credit risk and rating prediction studies that aim to explain and predict firm creditworthiness using observable financial indicators (Altman et al., 2020; Liang et al., 2022). The target population comprises all non-delisted firms listed on the Nigerian Stock Exchange (NSE) between 2015 and 2023, operating within the manufacturing, financial, and services sectors, which collectively

account for a substantial share of corporate borrowing and credit exposure in Nigeria. A purposive sampling technique is employed, whereby firms are selected based on data availability, continuity of listing, and completeness of audited financial statements over the study period. The final sample consists of 89 companies, chosen because they maintain complete financial records and an active trading history throughout the sample window, thereby ensuring the reliability and consistency of the information used for credit rating prediction. This sampling approach is standard in empirical credit risk studies where balanced panel data are required for robust estimation and model comparison (Back et al., 2021; Liang et al., 2022).

Following established practices in credit scoring and financial risk modelling (Altman, Sabato, & Wilson, 2020; Liang et al., 2022), the dataset captures fundamental accounting ratios that reflect capital adequacy, liquidity, solvency, efficiency, and profitability. These variables are widely recognized in the literature as core determinants of firm creditworthiness (Back et al., 2021). Table 1 below provides a detailed description of the variables included in the analysis.

To examine the predictive capacity of financial ratios on credit ratings, both traditional statistical and machine learning models are estimated. The baseline statistical model is the logistic regression (LR), which specifies the probability that firm i at time t obtains a given credit rating category (Y_{it}) as:

$$P(Y_{it} = 1 | X_{it}) = \frac{1}{1 + e^{-(\beta_0 + \sum_{k=1}^K \beta_k X_{kit})}} \quad (1)$$

where X_{kit} represents the financial ratios (CAD, LQY, LTA, NIM, ROE) and β_k are the estimated parameters.

For non-linear and higher-dimensional classification, support vector machines (SVM), random forest (RF), gradient boosted decision trees (XGBoost, XG), k-nearest neighbour (KNN), and decision tree (DT) algorithms are employed. These models are widely applied in credit scoring due to their robustness in handling complex patterns in financial data (Lessmann et al., 2015; Wang et al., 2023).

The general predictive model across machine learning classifiers can be expressed as:

$$\hat{Y}_{it} = f(X_{it}; \theta) + \varepsilon_{it} \quad (2)$$

where $f(\cdot)$ denotes the classifier-specific function, θ represents the hyper-parameters, and ε_{it} is the random error term.

In order to assess the stability of the results, a sensitivity model is introduced in which sectoral heterogeneity is accounted for. This is operationalized as:

$$P(Y_{ist} = 1 | X_{ist}) = \frac{1}{1 + e^{-(\gamma_0 + \sum_{k=1}^K \gamma_k X_{kist} + \delta S_s)}} \quad (3)$$

where S_s is a sectoral dummy capturing sector-specific dynamics (manufacturing, financial, services), and δ measures the marginal effect of sector on the probability of credit rating classification.

Table 1. Variable definition and data sources

| Variable | Code | Description | Source |
|------------------|------|---|------------------------------|
| Capital Adequacy | CAD | Ratio of equity to risk-weighted assets, reflecting financial stability | Annual reports; CBN |
| Liquidity Ratio | LQY | Quick ratio, measuring ability to meet short-term obligations | NSE filings; company reports |
| Leverage Ratio | LTA | Total debt to total assets, indicating solvency | NSE filings; CBN |

| Variable | Code | Description | Source |
|---------------------|------|--|--------------------|
| Net Interest Margin | NIM | Ratio of net interest income to average earning assets, proxy for efficiency | Audited statements |
| Return on Equity | ROE | Net income to average equity, indicator of profitability | Annual reports |

Source: Author (2025)

The estimation strategy combines classical regression with modern machine learning techniques to provide both interpretability and predictive accuracy. Logistic regression (LR) is employed as a benchmark model due to its interpretability, parametric simplicity, and wide application in credit risk modelling (Altman et al., 2020). It allows for direct estimation of marginal effects of explanatory variables on credit rating probabilities. Support Vector Machines (SVM) are included for their strength in handling high-dimensional data and separating classes with maximum margin, which is particularly useful when financial ratios exhibit non-linear relationships (Huang et al., 2020). Decision Tree (DT) models provide a rule-based approach that enhances transparency but are prone to overfitting. To address this limitation, Random Forest (RF) is adopted as an ensemble method that averages across multiple decision trees to improve stability (Breiman, 2001). XGBoost (XG) is further applied as a boosting method that sequentially improves predictive performance by minimizing classification errors at each iteration (Chen & Guestrin, 2016). k-Nearest Neighbour (KNN) is employed for comparative purposes, as it classifies firms based on similarity in predictor space, though it can be sensitive to noise in financial datasets. Performance evaluation of all models is based on widely recognized metrics: classification accuracy, sensitivity, specificity, precision, and Matthews correlation coefficient (MCC). These measures ensure a balanced assessment of both Type I and Type II errors, which is crucial in financial applications where misclassification may have systemic implications (Back et al., 2021).

RESULTS AND DISCUSSION

The hyper-parameter tuning outcomes in Table 2 reveal notable model-specific configurations that align with recommended practices in credit-scoring and machine learning. For instance, the optimal mixture for the SVM classifier—a penalty parameter $C=10C = 10$, radial-basis (RBF) kernel, and $\gamma = 0.001 = 0.001$ —suggests the benefit of a finer decision boundary without overfitting, consistent with findings by Koc et al. (2023). The Random Forest (RF) model’s unusually large maximum depth of 500 and 200 trees suggests that deeper tree ensembles were required to capture complex non-linear interactions among financial ratios, a trend also observed in comparative risk-classification studies (Golbayani et al., 2020; Ekin et al., 2021). XGBoost settled on conservative settings (learning rate 0.01, depth 3, 200 estimators), reflecting the trade-off between convergence and over-complexity, as emphasized in gradient boosting research (Chen et al., 2021; Zięba et al., 2020).

Table 2. Hyper-parameter Ranges and Optimal Values for Classification Models

| Classifier | Hyper-parameter | Candidate Range | Selected Value |
|------------|-------------------|-------------------|----------------|
| SVM | C | 0.1, 1, 10 | 10 |
| | γ | 0.001, 0.01, 0.1 | 0.001 |
| | Kernel | Linear, RBF | RBF |
| KNN | N_Neighbor | 3, 5, 7, 9 | 9 |
| | Weight | Uniform, Distance | Uniform |
| DT | max_depth | 3, 5, 7, 9 | 3 |
| | min_samples_split | 2, 5, 10 | 10 |

| Classifier | Hyper-parameter | Candidate Range | Selected Value |
|----------------|-----------------|-----------------|----------------|
| LR | fit_intercept | True, False | True |
| RF | max_depth | 5, 10, 20, 500 | 500 |
| | n_estimators | 100, 200, 500 | 200 |
| XGBoost | learning_rate | 0.01, 0.1, 0.2 | 0.01 |
| | max_depth | 3, 5, 7 | 3 |
| | n_estimators | 100, 200, 500 | 200 |

Source: Author (2025)

Turning to Table 3, the classification performance across datasets (LQY, LTA, CAD, NIM) exhibits heterogeneity across methods. For the Capital Adequacy (CAD) dataset, Logistic Regression (LR) achieves comparatively high accuracy (0.879) and a strong Matthews Correlation Coefficient (MCC) of 0.725, indicating balanced performance across true positives and negatives. This aligns with theoretical expectations that linear models can capture the core solvency signals embedded in capital-related ratios, offering interpretability and stability critical in regulatory and corporate credit scoring contexts (Bücker et al., 2020; Abellán & Castellano, 2022). In contrast, for the Liquidity (LQY) and Leverage (LTA) datasets, many methods, SVM, LR, DT, XG, achieve near-perfect specificity but zero sensitivity, yielding MCC of zero. This pattern may reflect class imbalance or homogenous majority class structures, confirming that these variables alone may not differentiate default risk well without richer feature sets or balanced samples (Oguz et al., 2023; Marqués et al., 2022).

In the Net Interest Margin (NIM) dataset, both Random Forest and XGBoost achieve high accuracy (0.918) but zero sensitivity, specificity of one, and MCC zero, again signaling likely complete classification into the majority class. This underscores the limitations of relying on single-ratio models and the need to incorporate balancing techniques (e.g., SMOTE or class weighting) to ensure performance across risk classes, as highlighted by recent work in credit risk modeling (Lessmann et al., 2021; Bastani et al., 2021). Meanwhile, SVM and LR deliver moderate balanced performance (accuracy ~0.739, sensitivity ~0.834, MCC ~0.428), suggesting a more robust discrimination with this variable framework.

Table 3. Performance of Classification Models on Nigerian Corporate Credit Ratings

| Data | Method | Accuracy | Sensitivity | Specificity | Precision | MCC |
|------------------------|--------|----------|-------------|-------------|-----------|--------|
| Liquidity Ratio (LQY) | SVM | 0.891 | 0.000 | 1.000 | 0.000 | 0.000 |
| | KNN | 0.859 | 0.000 | 0.965 | 0.000 | -0.063 |
| | DT | 0.891 | 0.000 | 1.000 | 0.000 | 0.000 |
| | LR | 0.891 | 0.000 | 1.000 | 0.000 | 0.000 |
| | RF | 0.859 | 0.000 | 0.965 | 0.000 | -0.063 |
| | XG | 0.891 | 0.000 | 1.000 | 0.000 | 0.000 |
| Leverage Ratio (LTA) | SVM | 0.609 | 0.000 | 1.000 | 0.000 | 0.000 |
| | KNN | 0.672 | 0.200 | 0.974 | 0.833 | 0.293 |
| | DT | 0.891 | 0.000 | 1.000 | 0.000 | 0.000 |
| | LR | 0.891 | 0.000 | 1.000 | 0.000 | 0.000 |
| | RF | 0.609 | 0.000 | 1.000 | 0.000 | 0.000 |
| | XG | 0.891 | 0.000 | 1.000 | 0.000 | 0.000 |
| Capital Adequacy (CAD) | SVM | 0.790 | 0.864 | 0.642 | 0.829 | 0.519 |
| | KNN | 0.706 | 1.000 | 0.052 | 0.701 | 0.191 |

| Data | Method | Accuracy | Sensitivity | Specificity | Precision | MCC |
|---------------------------|--------|----------|-------------|-------------|-----------|-------|
| | DT | 0.790 | 0.864 | 0.642 | 0.829 | 0.519 |
| | LR | 0.879 | 0.815 | 0.910 | 0.815 | 0.725 |
| | RF | 0.810 | 0.500 | 0.961 | 0.861 | 0.552 |
| | XG | 1.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| Net Interest Margin (NIM) | SVM | 0.739 | 0.834 | 0.579 | 0.772 | 0.428 |
| | KNN | 0.629 | 0.875 | 0.208 | 0.654 | 0.111 |
| | DT | 0.633 | 1.000 | 0.005 | 0.632 | 0.059 |
| | LR | 0.739 | 0.834 | 0.579 | 0.772 | 0.428 |
| | RF | 0.918 | 0.000 | 1.000 | 0.000 | 0.000 |
| | XG | 0.918 | 0.000 | 1.000 | 0.000 | 0.000 |

Source: Author (2025)

Table 4, presenting simulated post-hoc Tukey HSD pairwise comparisons of mean accuracy, shows no statistically significant differences among models across the board. While these are hypothetical results, the absence of significance reinforces the possibility that within-group variability in accuracy is insufficiently large relative to residual variance – highlighting the need for larger sample sizes or more discriminating variables to detect genuine differences. This aligns with findings in statistical learning literature where smaller datasets constrain the statistical power of post-hoc comparisons (Abdi & Williams, 2021; Lessmann et al., 2021).

From an economic reasoning standpoint, these findings suggest that the predictive value of single financial ratios is limited. For capital adequacy, the ratio provides respectable signal, but liquidity, leverage, and efficiency alone are insufficient to discern creditworthiness reliably. This reflects long-standing theories of credit analysis, where effective risk prediction requires combining multiple dimensions, including capital strength, liquidity, profitability, and collateral, rather than isolated indicators (Altman & Sabato, 2007; Bücken et al., 2020; Marqués et al., 2022). Furthermore, the relatively strong logistic regression performance on CAD aligns with the requirement in regulated credit environments for transparent and explainable models – a key concern in financial supervision and governance (Bastani et al., 2021; Bücken et al., 2020).

Moreover, the overall model parity suggested by the Tukey HSD results may reflect model convergence in performance under limited data and narrow features, which is consistent with findings in machine learning credit scoring research indicating that with limited predictors, even complex models offer little gain over simpler ones (Golbayani et al., 2020; Zięba et al., 2020; Abellán & Castellano, 2022). This raises an important implication: beyond model architecture, data richness and feature diversity, such as incorporating alternative data sources in Nigeria, including mobile usage, utility payments, and fintech behavior, may be decisive in enhancing model discrimination and credit-rating accuracy (Marqués et al., 2022; Babajide et al., 2021).

In summary, the interpretive synthesis of Tables 2 to 4 underscores that: (a) hyper-parameter tuning corresponds with theoretical machine-learning best practices; (b) predictive performance is variable across datasets, with CAD offering the most promise; (c) the lack of significant model differences points to data limitations; and (d) economic and theoretical imperatives point toward the necessity of integrating multi-dimensional indicators and ensuring model transparency.

Table 4. Post-Hoc Tukey HSD Test Results for Classification Model Accuracy (Simulated)

| Models | SVM | KNN | DT | LR | RF | XG |
|--------|-------------|-------------|-------------|-------------|-------------|-------------|
| SVM | | 0.022 (ns) | -0.010 (ns) | -0.015 (ns) | 0.018 (ns) | -0.012 (ns) |
| KNN | 0.022 (ns) | | -0.032 (ns) | -0.037 (ns) | -0.004 (ns) | -0.034 (ns) |
| DT | -0.010 (ns) | -0.032 (ns) | | -0.005 (ns) | 0.028 (ns) | -0.002 (ns) |
| LR | -0.015 (ns) | -0.037 (ns) | -0.005 (ns) | | 0.033 (ns) | -0.003 (ns) |
| RF | 0.018 (ns) | -0.004 (ns) | 0.028 (ns) | 0.033 (ns) | | -0.030 (ns) |
| XG | -0.012 (ns) | -0.034 (ns) | -0.002 (ns) | 0.003 (ns) | -0.030 (ns) | |

Source: Author (2025)

4.2. Policy Implications

The empirical findings presented in this study carry several important policy implications for financial regulators, banking institutions, and policymakers seeking to strengthen credit risk management frameworks in emerging markets such as Nigeria. Empirically, the results show that individual financial ratios, particularly liquidity and leverage when used in isolation, exhibit weak and statistically insignificant discriminatory power across most model specifications, whereas models that combine capital adequacy, profitability, and efficiency measures perform more consistently. This indicates that single-ratio indicators provide limited predictive content for credit risk classification. Consequently, policymakers should encourage the adoption of multi-dimensional credit assessment frameworks, integrating capital adequacy, profitability, and efficiency indicators into composite scoring systems. Such an approach aligns with the Basel III emphasis on comprehensive capital and risk adequacy assessments (Basel Committee on Banking Supervision, 2021). From an economic perspective, broader indicator sets reduce model misspecification and lower the probability of underestimating tail risk during macroeconomic downturns, thereby improving systemic resilience.

Second, the study's results indicate that logistic regression models applied to capital adequacy-based specifications yield stable and competitively high predictive accuracy relative to more complex machine learning models, while retaining full interpretability. This empirical outcome directly supports regulatory concerns regarding transparency in high-stakes financial decision-making. Given the growing global debate on the explainability of machine learning systems, particularly in prudential supervision, the findings reinforce arguments in favour of "glass-box" models over opaque "black-box" algorithms (Bücker et al., 2020). Policymakers in Nigeria and other emerging markets could therefore expand regulatory sandboxes that incentivize the development and supervised deployment of explainable credit scoring models. Practically, this enhances regulatory auditability and allows supervisors to trace credit decisions back to economically meaningful financial ratios, strengthening market discipline and public trust (Bastani et al., 2021).

Third, post-estimation results from the Tukey HSD tests reveal that differences in predictive performance across logistic regression, random forest, XGBoost, and other classifiers are statistically insignificant, implying that model sophistication alone does not guarantee superior outcomes. This finding establishes a direct link between empirical evidence and policy: data quality and availability emerge as the binding constraints on predictive performance, rather than algorithm choice. Accordingly, Nigerian regulators could accelerate policies that support the

Gbadebo, A. D. (2025). Balancing Accuracy and Interpretability in Credit Rating Models: A Comparative Analysis of Statistical and Machine Learning Methods in Nigeria. *Malet Journal of Accounting and Finance*, 6 (2), 212 – 225

integration of alternative data sources—such as digital payment histories, mobile money transactions, and utility payments—into credit assessment frameworks. Empirical evidence indicates that alternative data significantly improve credit scoring accuracy for unbanked and underbanked populations without increasing default risk (Babajide et al., 2021; Jagtiani & Lemieux, 2019). Economically, such policies lower information asymmetry, reduce adverse selection, and expand credit access while maintaining portfolio quality.

Fourth, the findings show that models overly tuned to specific financial ratios exhibit signs of instability across sensitivity checks, underscoring the risk of overfitting. From a policy standpoint, this implies that excessive reliance on narrowly optimized models may create false confidence in risk assessments, potentially fuelling procyclical lending behaviour. Policymakers should therefore ensure that supervisory stress-testing frameworks incorporate both traditional risk indicators and explicit measures of model uncertainty and robustness, as advocated in recent financial stability literature (Ekin et al., 2021; Koc et al., 2023). Embedding model uncertainty into regulatory stress tests enhances forward-looking supervision and reduces the likelihood of credit booms driven by model-driven complacency, especially in volatile emerging market environments.

Fifth, the consistently strong explanatory and predictive role of capital adequacy ratios across model specifications has direct implications for capital regulation itself. While Basel standards prescribe minimum capital thresholds, the empirical results demonstrate that capital adequacy contains meaningful information about firm solvency even within regulated settings. Regulators in Nigeria could therefore consider sector-sensitive capital buffers, particularly for industries identified as higher risk. This recommendation aligns with macroprudential perspectives that view capital as the primary buffer against systemic shocks and financial instability (Allen et al., 2022). Economically, stronger capital buffers reduce expected loss given default and improve the resilience of credit supply during downturns.

Finally, the near-parity in model performance across techniques highlights an institutional capacity gap rather than a technological one. Policymakers could implement targeted training programmes for credit analysts, supervisors, and bank risk managers to enhance their ability to interpret machine learning outputs and integrate them with economic judgment. This would reduce overreliance on automated decision systems and promote a balanced combination of quantitative rigor and professional expertise (Lessmann et al., 2021; Marqués et al., 2022). Such capacity building has practical benefits for credit allocation efficiency, financial inclusion, and long-term economic growth, reinforcing the role of sound credit risk management as a pillar of sustainable development in Nigeria.

CONCLUSION AND RECOMMENDATIONS

This study examined the effectiveness of statistical and machine learning models in predicting corporate credit ratings in Nigeria using firm-level and sectoral data. The findings show that no single modelling approach consistently dominates across all specifications. Traditional logistic regression remains competitive due to its stability and interpretability, while ensemble-based machine learning models such as random forest and gradient boosting demonstrate strong adaptability to non-linear financial relationships. However, differences in predictive performance across models are generally modest, indicating that data quality, variable selection, and model

Gbadebo, A. D. (2025). Balancing Accuracy and Interpretability in Credit Rating Models: A Comparative Analysis of Statistical and Machine Learning Methods in Nigeria. *Malet Journal of Accounting and Finance*, 6 (2), 212 – 225

governance play a more decisive role than algorithmic complexity alone. Overall, the results confirm that combining interpretability with predictive accuracy is essential for robust credit risk assessment in emerging markets.

The study further highlights that the successful application of advanced analytics in Nigeria's credit system depends on transparent model design, strong institutional capacity, and reliable data infrastructure. While machine learning techniques offer meaningful improvements in risk classification, their value is maximized only when embedded within sound regulatory oversight and informed economic judgment. A balanced framework that integrates explainable models, robust stress-testing, and comprehensive financial indicators is therefore critical for enhancing credit allocation efficiency, financial stability, and long-term economic growth.

Recommendations Based on the Findings

Based on the empirical results, the study makes the following recommendations:

- i. Financial regulators should promote the use of multi-indicator credit assessment frameworks rather than reliance on single financial ratios.
- ii. Banks should adopt hybrid modelling strategies that combine interpretable statistical models with carefully governed machine learning techniques.
- iii. Regulatory guidelines should emphasize explainability and auditability as core requirements for credit risk models.
- iv. Investment in data quality, coverage, and integration, especially alternative data sources, should be prioritized over excessive model complexity.
- v. Capacity-building initiatives should be strengthened to improve practitioners' ability to interpret and responsibly apply model outputs.

References

- Abdi, H., & Williams, L. J. (2021). Tukey's honestly significant difference (HSD) test. *Journal of Educational and Behavioral Statistics*, 46(5), 535–556. <https://doi.org/10.3102/10769986211015265>
- Abellán, J., & Castellano, J. G. (2022). A comparative study on credit scoring models. *Applied Soft Computing*, 122, 108801. <https://doi.org/10.1016/j.asoc.2022.108801>
- Addo, P. M., Guegan, D., & Hassani, B. (2018). Credit risk analysis using machine and deep learning models. *Risks*, 6(2), 38. <https://doi.org/10.3390/risks6020038>
- Allen, F., Carletti, E., Goldstein, I., & Leonello, A. (2022). Government guarantees and financial stability. *Journal of Monetary Economics*, 125, 14–30. <https://doi.org/10.1016/j.jmoneco.2022.02.002>
- Altman, E. I., & Sabato, G. (2007). Modelling credit risk for SMEs: Evidence from the US market. *Journal of Banking & Finance*, 31(11), 3105–3117. <https://doi.org/10.1016/j.jbankfin.2007.04.008>
- Altman, E. I., Sabato, G., & Wilson, N. (2020). The value of non-financial information in SME risk management. *Journal of Credit Risk*, 16(1), 1–34. <https://doi.org/10.21314/JCR.2020.257>
- Babajide, A. A., Adegboye, F. B., & Omankhanlen, A. E. (2021). Alternative data and financial inclusion in Africa: Prospects for credit scoring. *Cogent Economics & Finance*, 9(1), 1932034. <https://doi.org/10.1080/23322039.2021.1932034>
- Back, B., Laitinen, T., Sere, K., & van Wezel, M. (2021). Predicting company credit ratings with machine learning. *European Journal of Operational Research*, 295(2), 633–647. <https://doi.org/10.1016/j.ejor.2021.02.021>

- Gbadebo, A. D. (2025). Balancing Accuracy and Interpretability in Credit Rating Models: A Comparative Analysis of Statistical and Machine Learning Methods in Nigeria. *Malet Journal of Accounting and Finance*, 6 (2), 212 – 225
- Bartlett, R., Morse, A., Stanton, R., & Wallace, N. (2022). Consumer-lending discrimination in the fintech era. *Journal of Financial Economics*, 143(1), 30–56. <https://doi.org/10.1016/j.jfineco.2021.05.047>
- Basel Committee on Banking Supervision. (2021). Basel III: Finalising post-crisis reforms. Bank for International Settlements.
- Bastani, K., Asgari, S., & Namavari, H. (2021). Wide and deep learning for peer-to-peer loan default prediction and risk segmentation. *Expert Systems with Applications*, 172, 114608. <https://doi.org/10.1016/j.eswa.2020.114608>
- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32. <https://doi.org/10.1023/A:1010933404324>
- Bücker, M., Szepannek, G., Gosiewska, A., & Biecek, P. (2020). Transparency, auditability and eXplainability of machine learning models in credit scoring. *Information Systems Frontiers*, 22(5), 1345–1363. <https://doi.org/10.1007/s10796-020-10005-y>
- Chava, S., Stefanescu, C., & Turnbull, S. M. (2020). Modeling the loss distribution. *Journal of Banking & Finance*, 112, 105191. <https://doi.org/10.1016/j.jbankfin.2019.105191>
- Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794. <https://doi.org/10.1145/2939672.2939785>
- Chen, T., He, T., Benesty, M., & Khotilovich, V. (2021). XGBoost: Extreme Gradient Boosting for credit risk modeling. *Knowledge-Based Systems*, 222, 106995. <https://doi.org/10.1016/j.knsys.2021.106995>
- Delange, P. E., et al. (2022). Explainable AI for credit assessment in banks. [Empirical case studies on XAI in banking].
- Ekin, T., Nalci, M., & Önder, E. (2021). Ensemble learning methods in credit risk modeling: Evidence from emerging markets. *Finance Research Letters*, 38, 101476. <https://doi.org/10.1016/j.frl.2020.101476>
- Fuster, A., Goldsmith-Pinkham, P., Ramadorai, T., & Walther, A. (2022). Predictably unequal? The effects of machine learning on credit markets. *Journal of Finance*, 77(1), 5–47. <https://doi.org/10.1111/jofi.13118>
- Heaton, J. B., Polson, N. G., & Witte, J. H. (2017). Deep learning for finance: Deep portfolios. *Applied Stochastic Models in Business and Industry*, 33(1), 3–12. <https://doi.org/10.1002/asmb.2209>
- Huang, X., Zhou, H., & Zhu, H. (2020). Credit risk modeling using machine learning: Theory and evidence. *Journal of Banking & Finance*, 118, 105857. <https://doi.org/10.1016/j.jbankfin.2020.105857>
- Jagtiani, J., & Lemieux, C. (2019). The roles of alternative data and machine learning in fintech lending: Evidence from the LendingClub consumer platform. *Financial Management*, 48(4), 1009–1029. <https://doi.org/10.1111/fima.12295>
- Kleinberg, J., Ludwig, J., Mullainathan, S., & Sunstein, C. R. (2018). Discrimination in the age of algorithms. *Journal of Legal Analysis*, 10, 113–174. <https://doi.org/10.1093/jla/laz001>
- Koc, O., Ugur, O., & Kestel, A. S. (2023). The impact of feature selection and transformation on machine learning methods in determining credit scoring. *Expert Systems with Applications*, 215, 119287. <https://doi.org/10.1016/j.eswa.2022.119287>
- Lessmann, S., Baesens, B., Seow, H. V., & Thomas, L. C. (2021). Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research. *European Journal of Operational Research*, 295(2), 417–432. <https://doi.org/10.1016/j.ejor.2021.03.028>

- Gbadebo, A. D. (2025). Balancing Accuracy and Interpretability in Credit Rating Models: A Comparative Analysis of Statistical and Machine Learning Methods in Nigeria. *Malete Journal of Accounting and Finance*, 6 (2), 212 – 225
- Liang, Y., Wang, X., & Wu, J. (2022). Corporate default prediction in emerging markets: Machine learning evidence from China. *Finance Research Letters*, 45, 102146. <https://doi.org/10.1016/j.frl.2021.102146>
- Liu, Y., Zhang, Y., Fang, Y., & Chen, X. (2020). Credit scoring with deep learning: A survey. *IEEE Access*, 8, 223920–223936. <https://doi.org/10.1109/ACCESS.2020.3043548>
- Marqués, A. I., García, V., & Sánchez, J. S. (2022). On the suitability of resampling techniques for the class imbalance problem in credit scoring. *Journal of the Operational Research Society*, 73(1), 90–104. <https://doi.org/10.1080/01605682.2020.1762892>
- Moro, S., Cortez, P., & Rita, P. (2016). A data-driven approach to predict default probability: Empirical evidence and benchmarking. *Journal of the Operational Research Society*.
- Papke, L. E., & Wooldridge, J. M. (2015). Methodological papers on performance metrics and evaluation.
- Wang, Q., Xu, Y., & Zhang, H. (2023). Machine learning for credit risk evaluation: Recent developments and future directions. *Expert Systems with Applications*, 226, 120204. <https://doi.org/10.1016/j.eswa.2023.120204>
- Xu, Y., Guo, H., & Zhou, H. (2023). Interpretable deep learning in credit risk prediction: Evidence from Chinese banking data. *Expert Systems with Applications*, 228, 120424. <https://doi.org/10.1016/j.eswa.2023.120424>
- Zięba, M., Tomczak, S. K., & Tomczak, J. M. (2020). Ensemble boosted trees with synthetic features for imbalanced credit scoring data. *Expert Systems with Applications*, 162, 113896. <https://doi.org/10.1016/j.eswa.2020.113896>