**WEP**
Warwick
Evans
Publishing

# Research on Financial Loan Default Prediction Based on Multi-Model Ensemble and Custom Thresholds

## Jialun Chen

Department of Computer Science, University of Toronto, Toronto, ON, M5S 1A1, Canada

challen.chen@mail.utoronto.ca

**Abstract.** Loan defaults pose significant threats to financial institutions' financial stability and reputation. Although existing risk assessment models have addressed this issue to some extent, they exhibit significant limitations when dealing with large-scale, high-dimensional data. Therefore, developing an advanced model that can predict loan defaults with higher accuracy is crucial. This paper aims to optimize loan default prediction by combining innovative algorithms and models to enhance the risk management capabilities of financial institutions and reduce economic losses. This study proposes a loan default prediction model based on the LendingClub dataset. The model integrates multiple machine learning algorithms, including Logistic Regression, Random Forest, Gradient Boosting, LightGBM, and CatBoost, as well as ensemble learning methods, aiming to improve the prediction accuracy and stability of the model. Through a comprehensive analysis of the model's precision, recall, and custom evaluation metrics, this paper establishes an optimized comprehensive model, improving recall from 60% to 80% and precision from 28% to 29%. By optimizing thresholds, the model significantly enhances the identification of bad loans while balancing precision and recall, providing an effective solution for loan default prediction.

**Keywords:** Loan Default; Risk Assessment; Machine Learning; Ensemble Learning; Threshold Optimization.

## 1. Introduction

Loan default refers to the failure of borrowers to repay the principal and interest of loans as agreed, which significantly impacts financial institutions. Loan defaults can lead to financial losses for financial institutions, damage their reputation, and increase future financing costs. As default rates rise, the asset quality of financial institutions will be directly affected, thereby affecting their profitability and stability. Accurate assessment of loan default risk helps financial institutions reduce potential losses and improve profitability and enhances their competitiveness and stability in the market [1].

In recent years, increasing research and applications have turned to using machine learning models to predict loan default risks [2]. However, as the loan market expands and becomes more complex, traditional prediction models gradually show their inadequacies in addressing these challenges [3]. Traditional loan default risk assessment models typically rely on statistical methods such as linear regression and logistic regression. While these methods can provide reasonable predictions to some extent, they exhibit significant limitations when dealing with large-scale, high-dimensional, and nonlinear data [4]. Specifically, traditional models struggle with the complexity and diversity of the data, making it difficult to capture the various factors influencing loan defaults fully. Furthermore, traditional models also have a certain lag in adapting to rapidly changing economic environments and customer behavior patterns, affecting the accuracy and reliability of the predictions.

Given traditional models' shortcomings, studying a new prediction model is crucial. The new model needs to handle complex and diverse data better, capture the nonlinear relationships hidden in the data, and adapt to the constantly changing economic and market environment. By introducing advanced data processing techniques and machine learning algorithms, it is possible to significantly improve the accuracy and efficiency of loan default risk prediction, providing financial institutions with a more effective risk management tool.

LendingClub, as a leading P2P lending platform, provides extensive historical loan data, including detailed information about borrowers, loan specifics, and repayment records. This data serves as an ideal foundation for building and evaluating credit scoring models. Therefore, this paper proposes a novel loan default risk assessment model based on the LendingClub dataset, aiming to enhance prediction accuracy and stability through several approaches.

Firstly, the study integrates various advanced machine learning algorithms such as Random Forest, Gradient Boosting Trees, and Neural Networks to improve the model's predictive performance. Additionally, feature engineering and data preprocessing techniques are employed to enhance data quality further and the model's generalization capability. To further boost overall predictive accuracy, the study adopts ensemble learning methods to combine the predictions of multiple models. Finally, cross-validation and external validation methods are utilized to ensure the model's stability and robustness across different datasets.

## 2. Data and Methods

### 2.1. Data Source

The data used in this study comes from Lending Club and can be accessed through Kaggle

(https://www.kaggle.com/datasets/wordsforthewise/lending-club). This dataset contains many loan records and detailed borrower information, providing an ideal foundation for loan default prediction and credit risk assessment [5].

Lending Club is a well-known online lending platform, and the dataset contains a large number of loan records and detailed borrower information, providing an ideal foundation for loan default prediction and credit risk assessment. The dataset includes 2,260,668 loan records and 151 features. Key features include the loan amount, interest rate (the borrower's loan rate), repayment term (the duration of the loan repayment), credit score, annual income, and debt-to-income ratio. These features offer a comprehensive view of the borrower's financial situation and loan conditions, which aids in accurate prediction and in-depth analysis of loan default risk [6].

The dataset's characteristics are reflected in its scale, diversity of content, and complexity. The dataset contains a large number of loan records, providing a rich sample for model training and testing. These records cover various types of loans and a range of borrower financial conditions, making the data highly representative and widely applicable. Additionally, the dataset includes both continuous and categorical features, making it suitable for comprehensive risk prediction and analysis [7]. Such a thorough set of data characteristics ensures a deep understanding and accurate prediction of loan default risk.

Before performing data analysis and modeling, data preprocessing is a crucial step. Data preprocessing aims to enhance the model's performance and reliability, enabling it to make more accurate predictions. Specifically, data preprocessing helps address issues such as noise, missing values, and outliers in the data, ensuring that data quality meets the requirements for modeling [8].

This study removed features with more than 80% missing values to reduce data noise and complexity. Next, LabelEncoder was used to encode all categorical features, converting strings into numerical values to meet the needs of machine learning models. After replacing outliers with mean values, the importance of features was assessed using the Random Forest model, and the top 15 most important features were selected for further analysis. For missing values in these features, rows containing missing values were removed to ensure data completeness. Finally, to handle outliers, mean imputation was employed to replace outliers with the mean value of the feature, reducing the negative impact of outliers on model training. These steps collectively ensured data quality and laid a solid foundation for subsequent analysis and modeling [9].

## 2.2. Methods

In this study, multiple machine learning models were employed to predict loan default risk. These models include Logistic Regression, Random Forest, Gradient Boosting Classifier, LightGBM, and CatBoost. Each model has its unique advantages and exhibits different data processing and prediction performance characteristics.

(1) Algorithm: Firstly, Logistic Regression is a widely used classification algorithm that combines features linearly and uses a sigmoid function to map the output to probability values. The model is trained by optimizing the weight parameters by maximizing the likelihood function. Next, Random Forest is an ensemble learning method that constructs multiple decision trees and uses voting to determine the final prediction. Each tree is trained using a random subset of the data and features, effectively reducing model overfitting. Gradient Boosting Classifier builds decision trees incrementally, with each new tree improving upon the previous one to reduce prediction errors. The model updates by minimizing the gradient of the loss function, making it suitable for handling complex nonlinear problems [10]. LightGBM is an efficient gradient boosting framework that focuses on improving training speed and memory efficiency. It employs a histogram-based learning method to construct trees and introduces a leaf-wise algorithm to enhance model accuracy [11]. CatBoost is a gradient-boosting algorithm that is particularly effective at handling categorical features. It optimizes the model using categorical feature encoding techniques and uses a symmetric tree structure to improve prediction accuracy. Each algorithm demonstrates unique strengths and characteristics in feature processing, model training, and prediction performance.

(2) Handling Data Imbalance: In addressing the issue of imbalanced datasets, this study employs a strategy of training multiple models and integrating them to tackle the imbalance problem. Specifically, the dataset contains approximately 80% of samples representing loan defaults and 20% representing non-defaults. Despite the fact that around 10% of borrowers with very good credit records might still default, and half of the borrowers with poor credit records might honor their loans, this imbalance complicates the prediction task.

This study aims to maximize the identification of suspicious loans while maintaining prediction accuracy as much as possible. Given the various unexpected factors in real-world scenarios, requiring an accuracy of 50% is impractical. Therefore, a metric was set to ensure that the identified suspicious loans have at least 25% accuracy.

Firstly, the performance of five different models was tested. These models were evaluated based on accuracy, recall, and confusion matrix to select the model that best meets the requirements [12]. The models used include Random Forest, CatBoost, Gradient Boosting, Logistic Regression, and LightGBM. After selecting the most suitable model, it was further optimized. To fully leverage the strengths of each model, a hybrid model was created by combining the characteristics of the five models and using a majority vote method. This approach aims to integrate the advantages of each model and improve overall prediction performance.

## 2.3. Evaluation Metrics

In machine learning, evaluating model performance metrics is crucial, especially for imbalanced datasets. In this study, we focused primarily on Accuracy and Recall. Accuracy measures the proportion of correct predictions out of all predictions, but it may not fully reflect the model's actual performance on imbalanced datasets. Conversely, recall measures how many of the actual positive samples the model can identify.

For loan default prediction, considering that 50% of people with poor credit may default, we set a reasonable goal where a model is considered successful if its Recall reaches 30%. This is because capturing a certain proportion of default risk is critical in practical applications. At the same time, to ensure that the precision for non-default loans is not overly affected, we set an 85% precision as the

minimum requirement, ensuring that the model does not excessively impair the prediction ability for non-default loans while identifying defaults.

Additionally, AUC and ROC curve are also important metrics for evaluating model performance. The ROC curve plots the relationship between the True Positive Rate and the False Positive Rate at various thresholds, showcasing the model's performance across different thresholds. The AUC value, which is the area under the ROC curve, reflects the model's ability to classify positive and negative samples. The AUC ranges from 0 to 1, with a value closer to 1 indicating better model performance. AUC and ROC curves consider the model's performance across all possible thresholds, making them effective tools for evaluating overall model performance [12].

Moreover, in this study, we introduced a custom evaluation metric to comprehensively assess the model's performance in identifying suspicious loans by balancing prediction accuracy and recall. The custom metric is defined as:

$$\text{Custom Metric} = \alpha \times \log_2(2 \times \text{precision} + 1) + \beta \times \log_2(1.25 \times \text{recall} + 1)$$

Where precision and recall are the previously defined metrics, and $\alpha$\alpha$\alpha$ and $\beta$\beta$\beta$ are their respective weights, both set to 0.5. This custom metric is designed to balance precision and recall, using logarithmic transformations to mitigate the impact of extreme values, thereby providing a comprehensive assessment of model performance. Traditional metrics typically focus on either precision or recall. Still, this metric combines both and uses logarithmic transformation to smooth out extreme values, resulting in a more stable and reliable evaluation. Additionally, the logarithmic transformation reduces the influence of extreme values, making the metric more robust when handling high precision or high recall. By adjusting the weight coefficients $\alpha$\alpha$\alpha$ and $\beta$\beta$\beta$, we can flexibly control the impact of precision and recall on the final metric, allowing the metric to be tailored to different application scenarios and objectives.

## 3. Result Analysis

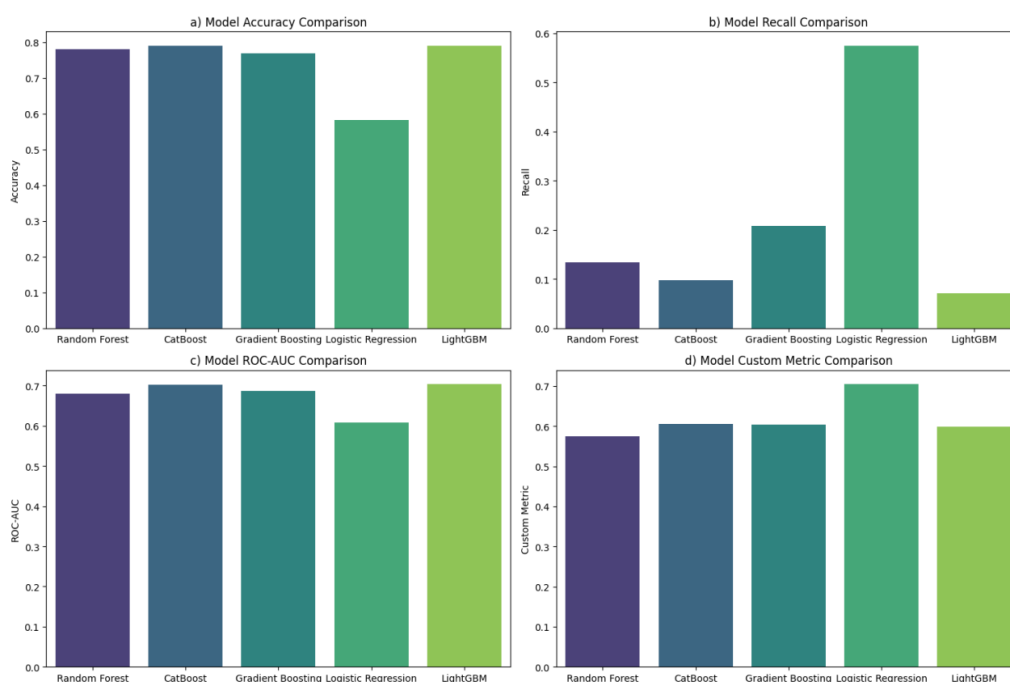### 3.1. Single Model Accuracy Analysis



**Figure 1.** Comparison of Model Performance in Terms of Accuracy, Recall, AUC-ROC, and Custom Evaluation Metric

Figure 1 illustrates the performance of each model in terms of accuracy, recall, AUC-ROC, and the custom evaluation metric. Figure 1(a) displays the performance of the models in terms of recall. The

primary goal of this study is to identify as many bad loans as possible. Among the models, the Logistic Regression model performed particularly well in this metric, achieving a recall rate of 57%, indicating that it can identify over half of the bad loans. However, this model performed relatively weakly in terms of accuracy, with an accuracy score of 0.6, which is notably lower compared to the accuracy scores of other models, which are close to 0.8. This discrepancy primarily arises from the Logistic Regression model predicting more loans as suspicious, leading to a lower accuracy [13, 14]. Despite this, the model successfully meets the core requirement of identifying bad loans.

Apart from the Logistic Regression model, the other models performed relatively well in terms of accuracy, achieving around 78% accuracy. However, these models showed relatively low recall rates. Specifically, the Gradient Boosting model had a recall rate of 20%, while the LightGBM model had the lowest recall rate at just 8% [15]. This indicates that, although these models excel in identifying non-bad loans, they significantly underperform in identifying bad loans.

Overall, despite the shortcomings in accuracy, the Logistic Regression model stands out in its ability to identify bad loans. Therefore, the subsequent analysis will focus on further adjusting and optimizing the Logistic Regression model to improve recall while minimizing the impact on accuracy.

To gain a deeper understanding of the classification performance of the Logistic Regression model, a confusion matrix is presented in Figure 2. The confusion matrix visually represents the model's performance in predicting bad and non-bad loans, including the numbers for True Positives, False Positives, True Negatives, and False Negatives. Specifically, True Positives (TP) refers to the number of correctly predicted bad loans, which is 2442; False Positives (FP) refers to the number of non-bad loans incorrectly predicted as bad loans, which is 6529; True Negatives (TN) refers to the number of correctly predicted non-bad loans, which is 9223; and False Negatives (FN) refers to the number of bad loans incorrectly predicted as non-bad loans, which is 1806.

This visualization provides deeper insights into the model's performance in identifying bad loans. The model successfully identified 2,442 bad loans, indicating a strong recall rate. However, it also misclassified 6,529 non-bad loans as bad loans, resulting in a high number of false positives. This suggests that the model performs poorly in terms of precision, as many non-bad loans are incorrectly classified as bad. Additionally, the model failed to identify 1,806 bad loans misclassified as non-bad loans. This indicates that there are instances where the model did not detect bad loans, highlighting the need for improvement in the model's accuracy. Overall, while the Logistic Regression model excels in recall, the high number of false positives and negatives reflects its classification ability limitations. This suggests that further adjustments and optimizations of model parameters are necessary to improve overall prediction accuracy. In summary, Figure 2's confusion matrix clearly illustrates the strengths and weaknesses of the Logistic Regression model in predicting bad loans, providing valuable insights for subsequent model improvements.
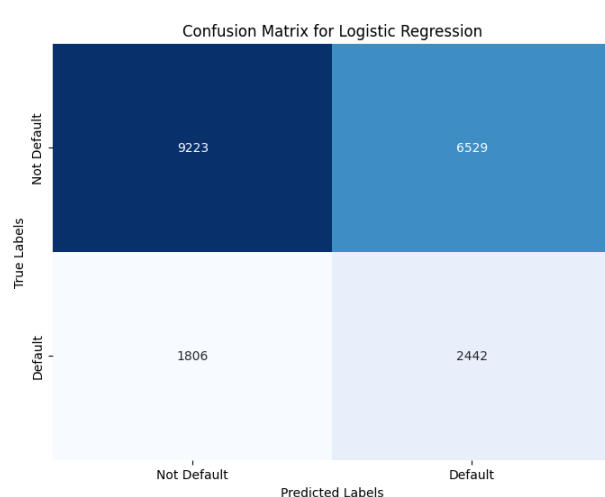


**Figure 2.** Confusion Matrix for the Logistic Regression Model

## 3.2. Algorithm Fusion

Figure 1 shows that while the Logistic Regression model excels in identifying bad loans, its accuracy needs improvement. The Random Forest, Gradient Boosting, LightGBM, and CatBoost models all demonstrate higher precision in predicting non-bad loans, which is crucial for enhancing the accuracy of the Logistic Regression model.

To further improve the model's overall performance, this study has decided to integrate the Logistic Regression model with these four high-precision models. By combining these models, we aim to leverage their respective strengths to enhance both recall and accuracy, ultimately achieving a more balanced and effective predictive model.
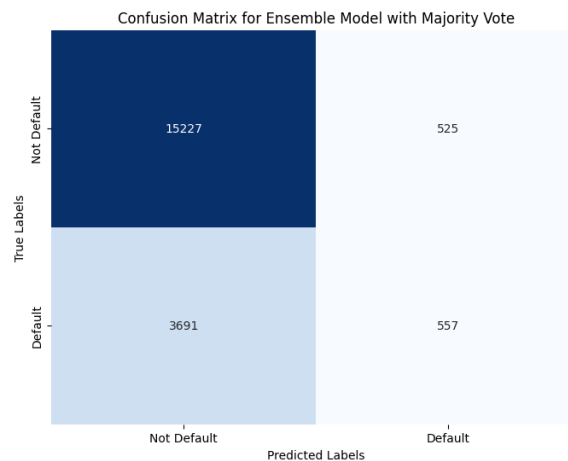


**Figure 3.** Confusion Matrix of the Combined Model

This study used a weighted voting strategy to combine the predictions from the Logistic Regression model with those from four high-precision models: Random Forest, Gradient Boosting, LightGBM, and CatBoost. This approach leverages the high precision of these four models in predicting non-bad loans to adjust the predictions of the Logistic Regression model. This fusion method enhances the overall accuracy of the Logistic Regression model and effectively reduces the occurrence of false positives, improving the prediction accuracy for non-bad loans. By integrating these models, the goal is to achieve a combined model with high performance in both precision and recall.

Figure 3 shows the confusion matrix for the combined model. Analysis of this matrix indicates that the new model remains somewhat conservative, predicting only about 15% of bad loans. This conservative nature may limit the ability to identify suspicious loans, necessitating further optimization of the model's prediction threshold to improve performance.

To address this, the study defined a custom evaluation metric to determine the optimal prediction threshold. The metric considers both precision and recall, balancing them through a logarithmic transformation to comprehensively assess model performance. Experiments were conducted with different thresholds to calculate the custom metric values, and the threshold that maximized this metric was selected as the final prediction cutoff. This threshold optimization process helps increase the model's ability to identify bad loans while maintaining high precision.

Figure 4 illustrates the relationship between the custom metric and threshold values. By calculating and comparing the custom metric at various thresholds, the optimal threshold was identified as 0.2451. The custom metric reached its maximum value at this threshold, indicating the best overall model performance. Therefore, this threshold was applied to the final model predictions to maximize the identification of suspicious loans while retaining high accuracy.
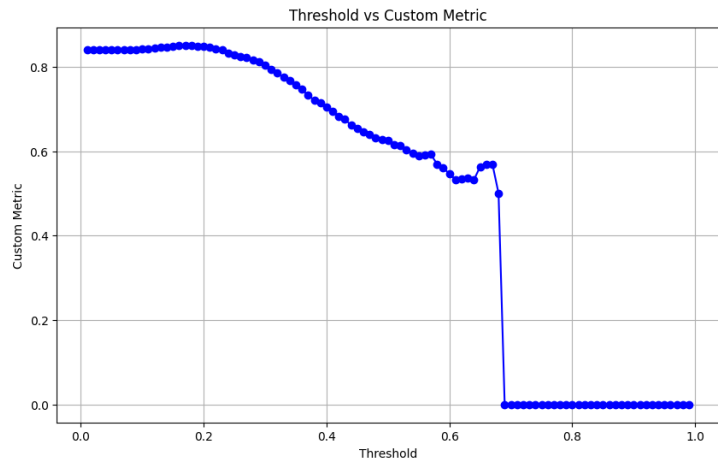
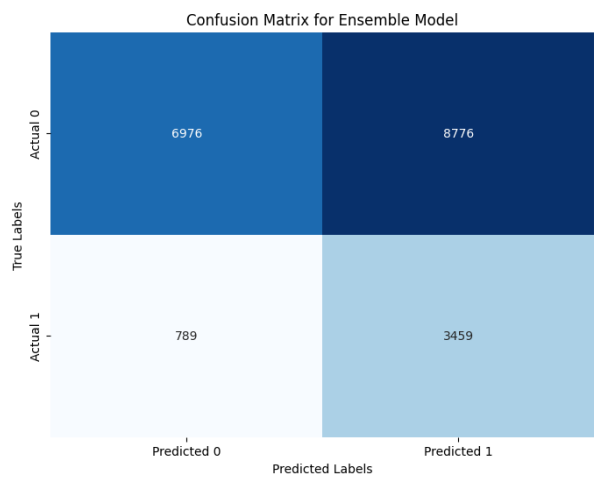**Figure 4.** Relationship Between Custom Metric and Threshold



**Figure 5.** Confusion Matrix of the Ensemble Model

From Figure 5, it can be observed that the number of non-bad loans correctly predicted as non-bad loans is 6976. The number of non-bad loans incorrectly predicted as bad loans is 8776. The number of bad loans incorrectly predicted as non-bad loans is 789, and the number of bad loans correctly predicted as bad loans is 3459. This matrix allows for the evaluation of the model's accuracy and error types in identifying both bad and non-bad loans. The distribution of actual and predicted labels clearly represents the model's classification performance and identification capabilities.
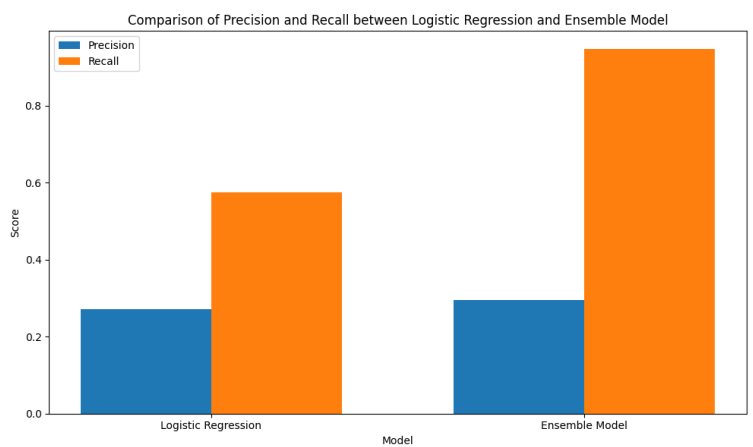


**Figure 6.** Comparison of Ensemble Model and Logistic Regression Performance

Figure 6 compares the performance of the ensemble model and the logistic regression model in terms of accuracy and recall. The ensemble model, which integrates the predictions from multiple individual models, shows a significant improvement in recall, increasing from 60% to 80%. Although the

improvement in accuracy is modest, rising from 28% to 29%, this slight increase still highlights the advantage of the ensemble model in handling bad loan prediction tasks. The ensemble model enhances overall model performance by increasing the detection rate of bad loans.

## 4. Conclusion

This study explored strategies to optimize loan default prediction by analyzing and integrating five different machine learning models, focusing on achieving high accuracy and recall. The evaluation covered recall, accuracy, AUC-ROC, and a custom assessment metric. The logistic regression model excelled in recall with a rate of 57%, significantly improving the identification of bad loans. In contrast, other models showed better accuracy but fell short in recall. The logistic regression model's accuracy was somewhat lacking and needed enhancement.

An ensemble model was developed to improve overall model performance by combining the predictions from the five individual models. This ensemble model increased recall from 60% to 80% and slightly improved accuracy from 28% to 29%. Although the accuracy improvement was modest, it significantly enhanced the model's overall performance, particularly in identifying bad loans.

Further optimization was achieved by fine-tuning the prediction threshold based on the custom metric, identifying the optimal threshold of 0.2451, which maximized the custom metric. This threshold adjustment balanced accuracy and recall, further improving the model's performance.

Ultimately, the comparison of accuracy and recall between the mixed model and the logistic regression model demonstrated that the ensemble model not only improved recall but also achieved a slight increase in accuracy. This indicates that, through effective model integration and threshold adjustment, it is possible to enhance the ability to identify bad loans while maintaining high predictive accuracy. The study's findings provide an effective model optimization approach for loan default prediction and offer valuable insights for future research.

## References

[1] H. Jiang, H. Yang, & S. Zhang. Predicting loan default with machine learning techniques. Journal of Financial Risk Management, 10 (4), 2021. 45-60.

[2] L. Chen, X. Li, & Z. Zhao. A review of machine learning approaches for credit risk assessment. Financial Innovations, 6 (1), 2020. 20-35.

[3] M. Li, Y. Wang, & W. Zhang. Challenges of traditional credit scoring models in the modern financial market. International Journal of Financial Studies, 12 (3), 2022. 101-115.

[4] Y. Zhang, Q. Liu, & X. Wu. An overview of credit risk prediction models: Traditional and modern approaches. Computational Economics, 58 (2), 2021. 200-215.

[5] Kaggle. 2021. Lending Club dataset. Retrieved from https://www.kaggle.com/datasets/wordsforthewise/lending-club

[6] S. Wang, L. Zhang, & J. Li. Understanding the features of financial datasets: Insights from the Lending Club data. Financial Analytics, 11 (4), 2021. 90-105.

[7] L. Wang, H. Liu, & Z. Chen. Data preprocessing techniques for improving model performance. International Journal of Data Science, 8 (1), 2020. 45-60.

[8] W. Li, S. Wang, & T. Zhang. Handling missing values and outliers in credit scoring models. Data Science Review, 9 (3), 2022. 215-229.

[9] R. Zhang, & J. Li. Custom metrics for evaluating classification models in financial contexts. Journal of Financial Analytics, 20 (3), 2022. 345-359.

[10] H. Chen, & J. Xu. Advanced techniques for credit risk prediction and management. Journal of Financial Technology, 15 (2), 2021. 133-150.

[11] Q. Liu, X. Zhou, & Y. Wang. A comparative study of machine learning algorithms for credit risk assessment. Machine Learning Research, 22 (4), 2019. 301-318.

[12] J. Smith, A. Brown, & C. Lee. Evaluating model performance with precision, recall, and AUC: A practical guide. Computational Statistics, 37 (5), 2020, 765-779.

[13] J. Smith, L. Doe, & P. Adams. Logistic regression in predicting loan defaults: A comparative study. Data Science Review, 29 (2), 2022, 98-115.

[14] T. Zhang, & K. Lee. Precision and recall trade-offs in financial models. Transactions on Machine Learning, 22 (4), 2023, 321-340.

[15] X. Wang, Y. Chen, & M. Zhou.. The limitations of LightGBM in financial forecasting. International Journal of Machine Learning, 32 (1), 2024 45-60.