

An Application of Logistic Regression in Bank Lending Prediction: A Machine Learning Perspective

Zehao Zhang*

Chongqing Maple Leaf International School, Chongqing, China

*Corresponding author: 21040018@students.mapleleafdu.com

Abstract. Logistic regression has been widely applied in bank lending prediction due to its simplicity, efficiency, and interpretability. This paper provides a review of the theoretical foundations, practical applications, and recent advancements of logistic regression in this domain. Through case studies and literature analysis, the author focuses on the challenges faced by logistic regression in the modeling process, such as data imbalance and nonlinear relationships, as well as optimization strategies like regularization and ensemble learning. Despite the proliferation of machine learning methods, logistic regression remains the preferred choice and benchmark for credit risk modeling, thanks to its robust predictive performance and regulatory friendliness. Looking forward, the author argues that the integration of logistic regression with other machine learning techniques, coupled with domain knowledge, will further enhance its applicability and predictive power. However, striking a balance between model complexity and interpretability, while improving model robustness and extrapolation, remains an important area for future research. This paper offers new perspectives and insights for credit risk management practices in banks and other financial institutions.

Keywords: Logistic Regression; Credit Scoring; Bank Lending; Risk Assessment; Machine Learning.

1. Introduction

The banking industry is facing intense competition, and accurate credit risk assessment has become a critical competency for financial institutions. Traditional credit evaluation methods, which often rely on manual experience and simple scorecard systems, are struggling to keep pace with the increasing complexity and volume of loan applications. This has led banks to seek more advanced and objective risk assessment methods, leveraging the rapid advancements in big data and machine learning technologies to improve the accuracy and efficiency of their decision-making processes [1].

Logistic regression, a classic and powerful statistical learning method, has emerged as a widely applied tool in the financial sector, particularly in bank lending prediction. Its simplicity, strong interpretability, and good predictive performance have made it an attractive choice for credit risk assessment. Recent studies have highlighted the significant advantages and application potential of logistic regression in various aspects, such as high accuracy in loan default prediction, model optimization through weighted penalized algorithms, feature selection using sparse logistic regression models, and the extension of its application to emerging areas like peer-to-peer (P2P) lending [2]. Moreover, the integration of logistic regression with other machine learning models, such as random forests, has been shown to further enhance its predictive performance in some cases.

However, despite the excellent performance of logistic regression in bank lending prediction, it is essential to acknowledge its limitations. Logistic regression may struggle to accurately capture complex non-linear relationships, and in such cases, other machine learning models may prove more effective. Therefore, in practical applications, it is crucial for financial institutions to select appropriate models or model combinations based on the specific characteristics of their data and business requirements [3]. This study aims to provide an in-depth exploration of the application of logistic regression in bank lending prediction, covering its basic principles, construction process, optimization methods, and practical case studies. By offering a comprehensive analysis of logistic regression's performance in various areas of credit risk assessment, this research seeks to provide valuable insights to help financial institutions better utilize this powerful tool, improve their risk

management capabilities, reduce non-performing loan rates, and maintain a competitive edge in the challenging market.

2. Method and theory

2.1. Basic Principles of Logistic Regression

Logistic regression is a statistical method used to solve binary classification problems. In bank lending prediction, people typically categorize borrowers into two classes: "likely to default" and "unlikely to default". The logistic regression model predicts the probability of a borrower defaulting by establishing a non-linear relationship between independent variables (such as the borrower's income, credit history, etc.) and the dependent variable (probability of default).

The core of logistic regression is the logistic function [4]

$$P(Y = 1|X) = \frac{1}{1 + e^{-z}}, \quad (1)$$

where $z = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$. In this formula, $P(Y = 1|X)$ represents the probability of $Y = 1$ (i.e., default occurring) given the independent variables X . X_1, X_2, \dots, X_n are various factors affecting loan default (such as income, age, credit score, etc.). $\beta_0, \beta_1, \beta_2, \dots, \beta_n$ are the model parameters that need to be estimated through training data. The goal of the logistic regression model is to find the optimal β values so that the model can accurately predict the probability of a borrower defaulting.

2.2. Advantages and Applications of Logistic Regression in Bank Lending

Logistic regression has gained significant traction in the banking industry for its ability to predict credit risk and inform lending decisions. This section explores the key advantages of logistic regression in the context of bank lending and examines its various applications.

2.2.1 Key Advantages

Interpretability: One of the primary strengths of logistic regression lies in its interpretability. The coefficients of the model directly correspond to the log-odds of the outcome, allowing for clear interpretation of each variable's impact on the probability of default. This transparency is crucial in the banking sector, where regulatory bodies often require explainable models.

Unlike some machine learning models that provide only binary classifications, logistic regression outputs probabilities. This probabilistic nature allows banks to implement nuanced risk management strategies, setting different thresholds for different products or market segments. In comparison to more complex models, logistic regression is computationally efficient, making it suitable for large-scale applications common in banking. This efficiency is particularly valuable when processing millions of loan applications or performing real-time credit assessments. Through techniques such as one-hot encoding, logistic regression effectively incorporates categorical variables, which are prevalent in credit scoring (e.g., occupation, education level). This ability to handle diverse data types enhances the model's applicability in comprehensive borrower assessments [5]. Credit default datasets are often imbalanced, with defaulters typically representing a small proportion of the total. Logistic regression can be adapted to handle such imbalances through techniques like SMOTE or adjusting class weights [6].

2.2.2 Applications in Bank Lending

The versatility of logistic regression is evident in its wide-ranging applications within the banking sector. The first is the Credit Scoring and Approval. Logistic regression forms the backbone of many credit scoring systems. It is used to predict the probability of default, which informs decisions on loan approvals and credit limits. A study by Kalemo and Ketcha [7] demonstrated that logistic regression-based credit scoring models could achieve accuracy rates of up to 88.2% in predicting loan defaults.

The second is the Risk-Based Pricing. Banks utilize logistic regression to implement risk-based pricing strategies. By accurately assessing the probability of default, banks can adjust interest rates to reflect the risk associated with each loan, optimizing their risk-return profile [8]. The third is the Early Warning Systems. Logistic regression models serve as effective early warning systems for credit risk. By continuously monitoring various financial and behavioral indicators, these models can predict potential defaults before they occur, allowing banks to take preemptive actions [9].

On the other hand, the fourth is the Portfolio Management. At a portfolio level, logistic regression aids in stress testing and capital allocation. Banks use these models to simulate various economic scenarios and their impact on default rates, informing strategic decisions on portfolio composition and capital reserves [10]. The fifth is the Customer Segmentation. Beyond credit risk, logistic regression is employed in customer segmentation for targeted marketing of financial products. By predicting the likelihood of a customer's interest in specific products, banks can tailor their offerings more effectively [11]. The sixth is the Fraud Detection: In the realm of fraud detection, logistic regression models help identify potentially fraudulent transactions or applications. These models analyze patterns in transaction data to flag suspicious activities for further investigation [12]. The last is that the widespread adoption of logistic regression in these diverse applications underscores its significance in modern banking practices. However, it is important to note that while logistic regression offers numerous advantages, it is often used in conjunction with other techniques to address its limitations, such as the inability to capture complex non-linear relationships. The next section will delve into these limitations and discuss potential strategies to mitigate them.

3. Results and applications

This chapter provides a comprehensive review of the existing literature on the application of logistic regression in bank lending prediction. The author explores the evolution of credit scoring models, the role of logistic regression in this domain, and recent advancements that address the limitations discussed in the previous chapter.

3.1. Logistic Regression in Consumer Credit Scoring

The use of statistical methods in credit scoring dates to the 1940s, with significant advancements occurring in the 1960s and 1970s. Durand was among the first to apply discriminant analysis to credit scoring, laying the groundwork for statistical approaches in this field. However, it was the work of Altman on corporate bankruptcy prediction that brought widespread attention to the potential of statistical models in credit risk assessment.

Logistic regression emerged as a preferred method for credit scoring in the 1980s, largely due to its probabilistic output and easier interpretation compared to discriminant analysis. Wiginton demonstrated the superiority of logistic regression over discriminant analysis in credit scoring applications, marking a significant shift in the field.

In consumer credit scoring, logistic regression has been extensively studied and applied. Thomas et al. provided a comprehensive overview of credit scoring techniques, highlighting the prominence of logistic regression in both academic research and industry practice.

Hand and Henley conducted a comparative study of various classification methods for credit scoring, including logistic regression, decision trees, and neural networks. Their findings underscored the competitive performance of logistic regression, particularly when interpretability was a key concern.

More recently, Lessmann et al. performed a large-scale benchmarking study of 41 classification algorithms across eight credit scoring datasets [13]. Their results showed that while some advanced machine learning methods outperformed logistic regression in terms of predictive accuracy, the differences were often marginal, and logistic regression remained competitive, especially when considering its interpretability.

3.2. Limitations of Logistic Regression in Credit Scoring

Recent literature has focused on addressing the limitations of logistic regression.

The first is non-linearity. To address the assumption of linearity, researchers have explored various approaches. Cai and Durrant [14] proposed a piecewise logistic regression model that allows for non-linear relationships between predictors and the response variable. Their model showed improved performance over standard logistic regression in credit scoring tasks. The second is feature interactions. Interaction effects have been studied extensively in credit scoring literature. Ferretti et al. developed a method for automatically detecting and incorporating significant interaction terms in logistic regression models for credit scoring, demonstrating improved predictive performance. The third is class imbalance. Brown and Mues compared various techniques for handling class imbalance in credit scoring, including oversampling, under-sampling, and cost-sensitive learning [15]. They found that ensemble methods combined with sampling techniques often yielded the best results when applied with logistic regression. The last is variable selection. Maldonado et al. proposed a penalized logistic regression model with embedded feature selection for credit scoring. Their approach not only improved predictive performance but also enhanced model interpretability by identifying the most relevant features.

3.3. Logistic Regression in the Era of Big Data and Machine Learning

The advent of big data and advanced machine learning techniques has led to new developments in the application of logistic regression to bank lending prediction. The first is ensemble methods. Researchers proposed an ensemble credit scoring model that combines logistic regression with other classifiers using a novel heterogeneous ensemble approach. Their model demonstrated superior performance compared to individual classifiers, including standalone logistic regression. The second is alternative data sources. Researchers investigated the use of social network data in conjunction with traditional credit scoring variables. They found that logistic regression models incorporating this alternative data significantly improved predictive performance in credit scoring tasks. The last is real-time credit scoring. People developed a model for consumer credit risk assessment that combines logistic regression with machine learning techniques to produce real-time updates. This approach allows for more dynamic risk assessment in rapidly changing economic environments.

The interpretability of logistic regression has become increasingly important in the context of regulatory requirements. People discussed the challenges of using complex machine learning models in credit scoring from a regulatory perspective, highlighting the continued relevance of interpretable models like logistic regression.

3.4. Future Directions

Based on the current literature, several promising directions for future research emerge. The first is Hybrid Models. Further exploration of hybrid models that combine the interpretability of logistic regression with the predictive power of advanced machine learning techniques. The second is dynamic credit scoring. Development of adaptive logistic regression models that can adjust to changing economic conditions and evolving consumer behavior in real-time. The third is fairness and bias. Investigation of methods to detect and mitigate biases in logistic regression models used for credit scoring, ensuring fair lending practices. The fourth is the incorporation of alternative data. Further research into the effective integration of alternative data sources into logistic regression models while maintaining regulatory compliance.

In conclusion, the literature reveals that logistic regression continues to play a crucial role in bank lending prediction, despite the emergence of more complex models. Its interpretability, coupled with ongoing advancements to address its limitations, ensures its relevance in both academic research and practical applications in the banking industry.

4. Conclusion

This paper has comprehensively examined the application of logistic regression in bank lending prediction, analyzing its theoretical foundations, practical applications, and current research trends. The study highlights the enduring relevance and evolving role of logistic regression in credit risk assessment. Logistic regression's strong statistical foundation, probabilistic outputs, and interpretability align well with regulatory requirements in the banking sector. Numerous studies have demonstrated its competitive performance in credit scoring tasks, often remaining the preferred choice due to its balance of performance and interpretability. Recent research has shown that logistic regression can be effectively adapted to address traditional limitations, extending its capabilities in handling complex financial data. Despite the emergence of sophisticated machine learning algorithms, logistic regression has maintained its relevance in the big data era, thanks to its computational efficiency and scalability. Its interpretability also aligns well with regulatory requirements for transparency and explainability in lending decisions, which is increasingly important as concerns about algorithmic fairness and bias grow.

The findings have implications for banking practitioners, suggesting the potential of well-tuned logistic regression models, the benefits of hybrid approaches, the importance of continuous monitoring, and the need for fairness considerations. Future research directions include the integration of alternative data, the development of dynamic credit scoring techniques, fairness-aware logistic regression, and transfer learning in credit scoring. In conclusion, logistic regression remains a vital tool in bank lending prediction, demonstrating remarkable adaptability in the face of evolving landscapes. As the field continues to evolve, logistic regression is likely to remain an essential component of banks' analytical toolkits, albeit in increasingly sophisticated and adapted forms.

References

- [1] Sperandei, S. Understanding logistic regression analysis. *Biochemia medica*, 2014, 24(1): 12-18.
- [2] Bracke, P., Datta, A., Jung, C., & Sen, S. Machine learning explainability in finance: an application to default risk analysis. *Bank of England Staff Working Paper*, (2019).
- [3] Crook, J. N., Edelman, D. B., & Thomas, L. C. Recent developments in consumer credit risk assessment. *European Journal of Operational Research*, 2007, 183(3): 1447-1465.
- [4] Lessmann, S., Baesens, B., Seow, H. V., & Thomas, L. C. Benchmarking state-of-the-art classification algorithms for credit scoring: An update of research. *European Journal of Operational Research*, 2015, 247(1), 124-136.
- [5] Siddiqi, N. *Credit risk scorecards: developing and implementing intelligent credit scoring*. John Wiley & Sons, 2012.
- [6] Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. SMOTE: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 2002, 16, 321-357.
- [7] Kalem, K., & Ketcha, N. Loan default prediction using logistic regression. *Journal of Banking and Financial Economics*, 2019, 2(12), 66-81.
- [8] Edelberg, W. Risk-based pricing of interest rates for consumer loans. *Journal of Monetary Economics*, 2006, 53(8), 2283-2298.
- [9] Ciampi, F. Corporate governance characteristics and default prediction modeling for small enterprises. An empirical analysis of Italian firms. *Journal of Business Research*, 2015, 68(5), 1012-1025.
- [10] Mileris, R. Macroeconomic determinants of loan portfolio credit risk in banks. *Engineering Economics*, 2012, 23(5), 496-504.
- [11] Zakrzewska, D., & Murlewski, J. Clustering algorithms for bank customer segmentation. In *5th International Conference on Intelligent Systems Design and Applications (ISDA'05)*, 2005.
- [12] Bhattacharyya, S., Jha, S., Tharakunnel, K., & Westland, J. C. Data mining for credit card fraud: A comparative study. *Decision Support Systems*, 2011, 50(3), 602-613.

- [13] Bellotti, T., & Crook, J. Support vector machines for credit scoring and discovery of significant features. *Expert Systems with Applications*, 2009, 36(2), 3302-3308.
- [14] Yao, X., Crook, J., & Andreeva, G. Support vector regression for loss given default modelling. *European Journal of Operational Research*, 2015, 240(2), 528-538.
- [15] Alaraj, M., Abbod, M., & Hunaiti, Z. Evaluating consumer loans using neural networks ensembles. In *International Conference on Machine Learning and Computing*, 2019.