

# Credit Risk Management with Alternative Data: Expanding the Predictive Frontier Beyond Traditional Scoring Models

**Mingze Xu**<sup>1,\*</sup>

<sup>1</sup>Department of Economics and Management, China Jiliang University, Hangzhou, Zhejiang

\*Corresponding author: guojuan@zjmc.edu.cn

## Abstract:

Credit risk management is fundamental to the survival of financial institutions and the stability of the broader financial system. While traditional credit scoring models like FICO serve as industry standards in consumer credit, they have notable limitations: they exclude individuals without conventional credit histories and rely on static historical data, which fails to capture dynamic changes in borrowers' financial behaviour and risk profiles. This study addresses two core questions: (1) How can alternative features, derived from borrowers' controllable financial behaviour, be constructed to capture risk information missed by FICO? (2) How can the marginal contribution of these features to credit risk identification be quantified? Using the Lending Club loan dataset, we adopt a framework of "feature selection → model building → performance evaluation → robustness testing," constructing a baseline model using only the FICO score and an extended model that incorporates alternative features via the XGBoost algorithm with 5-fold cross-validation. Results indicate that the extended model achieves robust improvements in key metrics (AUC, KS, F1), effectively bridging gaps in dynamic risk detection. Alternative features supplement traditional models by identifying high-risk segments, particularly those with "high debt, low stability, and no assets." Academically, this study advances credit risk identification methodologies and enriches the theoretical application of alternative data; practically, it offers financial institutions enhanced risk control tools to reduce nonperforming loans.

**Keywords:** Credit Risk Management; Alternative Data; FICO Score; XGBoost Model; Financial Inclusion.

## 1. Introduction

Credit risk management is a core function and the foundation of financial institutions' viability. Its effectiveness is closely linked to both financial system stability and efficient resource allocation [1]. Currently, traditional credit scoring models, such as the FICO score, quantify borrowers' historical credit histories into a standardized score, providing an objective basis for credit decisions and setting the industry benchmark in global consumer credit. However, these models have significant limitations: FICO scores fail to cover most individuals lacking traditional credit histories, resulting in financial exclusion. Moreover, their reliance on static historical data hampers their ability to reflect dynamic changes in borrowers' financial behavior and risk, making them less adaptable in complex and volatile financial environments [2].

The rise of intelligent financial technology and big data analytics has introduced alternative data sources, such as mobile payment transactions, active repayment frequency, and spending limit control. These data capture real-time, controllable financial behaviors, offering a more accurate perspective on borrowers' financial status and willingness to repay [3]. Such alternative data, representing behaviors that borrowers can actively optimize, provide predictive insights that surpass traditional credit scores and offer new opportunities for risk prediction beyond FICO.

This research addresses two key questions: First, how can alternative features based on borrowers' controllable financial behavior be systematically constructed to uncover risk information overlooked by FICO scores? Second, how can the marginal contribution of these features to credit risk identification be quantified, and their predictive value validated against traditional models?

This study is significant both academically and practically. Academically, it moves beyond traditional reliance on structured data, expands methodologies for credit risk identification, and enriches the theoretical application of alternative data. Practically, it equips banks and fintech companies with more precise risk control tools, helping reduce nonperforming loan rates and extending credit access to financially excluded groups, thereby supporting the coordinated development of inclusive and secure financial services.

## 2. Literature Review

Credit risk identification has always evolved around the two main lines of „data dimension expansion“ and „model efficiency improvement“. At present, a three-tier system of traditional scoring, alternative data application, and intelligent model empowerment has been formed, but the core research gap still needs to be filled.

Traditional credit scoring, based on the FICO model,

serves as the industry benchmark. Using a standardized algorithm based on five factors—payment history (35%), credit utilization (30%), and length of credit history (15%)—it transforms dispersed credit information into a comparable metric ranging from 0 to 850, providing an efficient tool for large-scale consumer credit approval [3]. However, its limitations are significant: First, it fails to capture non-traditional behavioral data such as real-time consumption and proactive financial planning, resulting in 15%-20% of the world's „credit-free“ (those with no prior credit history) being excluded from formal financial services [4]. Second, its reliance on static historical data makes it difficult to dynamically reflect a borrower's financial deterioration during economic downturns, making it prone to misjudgment of risk.

The rise of alternative data is breaking through bottlenecks in risk identification. Dynamic data such as internet payment transactions, consumer product preferences, and proactive repayment frequency, as well as non-financial data such as social network interactions and device usage habits, are becoming increasingly important supplements to traditional data [4]. Significant empirical value: Scoring models incorporating mobile payment data reduce nonperforming loan prediction errors by 12%-18% compared to traditional FICO models. Cross-national studies have also confirmed that models incorporating social network default status and regional economic ratings achieve an AUC of 0.7936, significantly outperforming traditional models. However, this approach remains controversial. Non-financial data suffers from weak credit correlation, fragmented sources, and variable quality, which can easily lead to model overfitting and discrimination risks.

Intelligent technology is driving model upgrades, presenting a path from traditional linear models to nonlinear machine learning to high-dimensional deep learning. Logistic regression is used for regulatory compliance due to its interpretability. Random forests and XGBoost, through ensemble learning, improve their ability to capture complex patterns in alternative data and have become core models for fintech companies. While deep learning can process high-dimensional, unstructured data, its „black box“ nature conflicts with regulatory requirements. Recent breakthroughs in interpretability research include the SHAP method, which reduces compliance review time by 70% by quantifying feature contributions; and the TreeSHAP algorithm, which reduces the interpretation complexity of XGBoost models from exponential to polynomial, achieving a balance between accuracy and transparency [5-7].

Credit-free households and small and micro-enterprise assessment have become research hotspots. TransUnion data shows that Generation Z accounts for 59% of first-time credit recipients in nine countries, including the United States and India, and their default rates are comparable to those of traditional borrowers. However, traditional

models struggle to identify credit potential (Wise & Chen, 2023). Regarding small and micro-enterprise assessment, the LightGBM model, which integrates financial data and third-party behavioral data, improves F1-score by 21.4% compared to traditional logistic regression after using SMOTE technology to address sample imbalance [8]. However, there is still a lack of unified standards for data standardization and feature selection. There are two major gaps in existing research: first, it focuses on uncontrollable variables such as the macroeconomy and ignores the risk-indicative value of borrowers' active financial behavior; second, it fails to systematically quantify the marginal improvement of alternative data on traditional models, making it difficult to support the refined design of risk control strategies.

### 3. Research Design and Methods

#### 3.1 Research framework construction

Following the research framework of “feature selection→model building→performance evaluation→robustness test”, using the Lending Club loan dataset as a sample, we focus on the differences in credit risk prediction between “traditional FICO scores” and „alternative features of borrowers' controllable financial behaviour“. Through the XGBoost model, we verify the value of alternative features in mining the prediction space beyond FICO scores. The core goal is to improve the accuracy and comprehensiveness of credit risk identification.

#### 3.2 Variable Definition and Data Processing

##### 3.2.1 Target variable

Based on the core requirements of credit risk assessment, the target variable is defined as a binary categorical variable to distinguish the risk level of the borrower (Table 1):

**Table 1. Risk levels**

Risk Type	Assignment	Definition
Low risk	0	The borrower repays the loan in full and on time and has no record of default
High risk	1	The borrower has committed material defaults such as loan write-offs or overdue payments exceeding 90 days.

##### 3.2.2 Core explanatory variables

Core explanatory variables were defined and pre-processed. All variables were processed by „missing value

filling (continuous variables were used with mean/median, and categorical variables were trimmed with mode + outliers (3 $\sigma$  principle)“ to ensure data quality. The specific variable definitions are shown in Table 2:

**Table 2. Definition and pre-processing table of core explanatory variables**

Variable type	Variable name	definition	Controllability Description	Data preprocessing method
Traditional variable	FICO average	Take the average of the upper and lower limits of the FICO score range to reflect the historical credit level	Uncontrollable	Outlier clipping (removing extreme values <500 or >850)
Alternative Features	Credit card usage rate	The ratio of revolving credit balance to credit limit reflects the ability to manage short-term liabilities.	Controllable	Missing values are filled with industry mean values, and outliers are trimmed (>100% is taken as 100%)
	Debt-to-income ratio	The ratio of monthly debt expenditure to monthly income reflects the level of debt burden	Controllable	Missing values are filled with the sample median, and extreme values > 60% are removed.
	Repayment-to-income ratio	The ratio of monthly loan repayment to monthly income, reflecting repayment pressure	Controllable	Monthly repayment amount/monthly income, outlier trimming (>30% is taken as 30%)

	Account activity rate	The ratio of open accounts to total accounts reflects account management capabilities.	Controllable	Calculated by dividing the number of open accounts by the total number of accounts, with the value range normalized to [0,1]
	Length of employment	Borrower's current years of employment	Controllable	Categorical variables are digitized, and missing values are filled with 2 years.
	Housing type	Borrower's housing status(OWN=2, MORTGAGE=1, RENT=0)	Controllable	

### 3.3 Model Construction

#### 3.3.1 Baseline Model

The input features consist solely of traditional variables (FICO score and related indicators), with no other features introduced. This function serves as a benchmark for traditional credit assessment, measuring the predictive bounds of a single FICO score and providing a reference for comparing the performance of expanded models. The objective function uses „binary: logistic“ and outputs a high-risk probability value, suitable for binary risk classification.

#### 3.3.2 Extended Model

Building on the baseline model, we introduce alternative features defined by the framework (controllable financial behaviours of borrowers, such as credit card usage and repayment methods), creating a hybrid input of „FICO score

+ alternative features.“ The core goal is to verify the complementary value of alternative features to the prediction space beyond the FICO score, addressing the framework's requirement to „test the marginal improvement of alternative features.“ The model output and objective function are identical to those of the baseline model to ensure fair comparison.

#### 3.3.3 Algorithm and Verification

Based on the framework's specifications, we selected XGBoost because it can capture nonlinear feature correlations, supports importance assessment, and is well-suited for risk prediction scenarios. We optimized parameters using 5-fold cross-validation with a training set: test set ratio of 7:3. Our performance was evaluated using framework metrics, calculating the marginal improvement of the extended model over the baseline model and quantifying the value of the replacement features (Table 3).

**Table 3. 5-fold cross-validation optimization parameters**

Parameter name	Value	explanation
learning rate	0.1	Step size (to avoid overfitting)
max_depth	3	Tree depth (controls model complexity)

## 4. Empirical Analysis

### 4.1 Descriptive Analysis

Descriptive statistics were performed on the target vari-

ables, traditional variables, and alternative features of the 186 valid samples to reveal the overall distribution characteristics of the variables. The results are shown in Table 4.

**Table 4. Descriptive statistics of the overall sample**

Variable Type	Variable Name	Sample size	Mean	Standard Deviation	Minimum	25% quantile	50% quantile	75% quantile
Target variable	Risk Level	186	0.183	0.388	0	0	0	0
Traditional variables	FICO average	186	685.2	32.4	520	660	689	780

Alternative Features	Credit card usage rate	186	68.3%	21.5%	12.1%	52.7%	71.4%	100%
	Debt-to-income ratio	186	15.7%	6.2%	3.3%	11.2%	14.9%	58.9%
	Repayment-to-income ratio	186	18.2%	7.4%	4.5%	12.8%	17.6%	30%
	Account activity rate	186	0.62	0.18	0.21	0.49	0.63	0.85
	Length of employment	186	2.8	1.3	0	2	3	4
	Housing type	186	1.03	0.68	0	0	1	2

The overall sample characteristics show the following:

**Risk Distribution:** 152 samples are low-risk, accounting for 81.7%; 34 samples are high-risk, accounting for 18.3%, consistent with the general characteristic of „low default rates“ in the consumer credit market [9].

**Traditional Credit Level:** The FICO mean is 685.2, with a standard deviation of 32.4, a 25th percentile of 660, and a 75th percentile of 710, indicating that the sample is concentrated in the „medium credit level“ range, covering the mainstream credit user group.

**Alternative Characteristics Distribution:** The mean credit card usage rate is 68.3%, with a 75th percentile of 71.4%, indicating that nearly half of the borrowers are using their credit cards close to their credit limit, posing significant pressure on short-term debt management; the mean debt-

to-income ratio is 15.7%, below the safety threshold of 60%, indicating that the overall debt burden is manageable; the mean account activity rate is 0.62, indicating that an average of 62% of the sample accounts are active, indicating moderate account management skills.

## 4.2 Comparative analysis of variables in different risk groups

In order to preliminarily verify the indicative role of traditional variables and alternative characteristics on credit risk, the samples were divided into a low-risk group (0) and a high-risk group (1) according to „risk level“. The differences between the two groups in key variables were compared. The results are shown in Table 5.

**Table 5. Comparative analysis of variables in different risk groups**

Variable Type	Variable Name	Low-risk group(n=152)	high-risk group(n=34)	Difference between groups (low-high)
Traditional variables	FICO average	698.5	642.3	56.2
Alternative Features	Credit card usage rate	57.2%	89.4%	-32.2%
	Debt-to-income ratio	13.2%	24.5%	-11.3%
	Repayment-to-income ratio	15.1%	26.8%	-11.7%
	Account activity rate	0.68	0.45	0.23
	Length of employment	3.2	1.9	1.3
	Housing type	18.4%	4.7%	13.7%
Variable Type	Variable Name	Low-risk group(n=152)	high-risk group(n=34)	Difference between groups (low-high)

The comparison results show that the two groups of variables are significantly different and consistent with risk logic expectations:

**Differences in Traditional Variables:** The low-risk group's mean FICO score (698.5) was significantly higher than the high-risk group's (642.3), a difference of 56.2 points. This confirms the traditional risk control logic that „lower

historical credit scores indicate higher default risk.“

**Differences in Alternative Characteristics:**

**Debt Management Ability:** The high-risk group's credit card usage rate was close to full, significantly higher than the low-risk group, indicating that excessive short-term debt is a key warning sign of default.

**Debt Burden:** The high-risk group's debt-to-income ratio

and repayment-to-income ratio were both close to safety thresholds and significantly higher than the low-risk group, reflecting a strong correlation between high debt pressure and default risk.

**Stability Indicators:** The low-risk group had higher account activity rates and employment duration than the high-risk group, and their homeownership rate was 3.9 times that of the high-risk group, demonstrating the

risk-mitigating effects of financial stability and housing assets.

#### 4.3 Distribution of key classification variables

Frequency distribution analysis was conducted on the categorical variables in the replacement characteristics to further clarify the sample structure characteristics. The results are shown in Tables 6 and 7.

**Table 6. Distribution of employment hours**

Length of employment	definition	Number of samples	occupation	The proportion of low-risk group	The proportion of high-risk group
0	<1 Year	28	15.1%	8.5%	52.9%
1	1-2 Year	35	18.8%	16.4%	29.4%
2	3-5 Year	52	28.0%	30.3%	14.7%
3	6-9 Year	41	22.0%	24.3%	2.9%
4	10+ Year	30	16.1%	20.5%	0%

**Table 7. Distribution of housing types**

Housing type	definition	Number of samples	occupation	The proportion of low-risk group	The proportion of high-risk group
0	RENT	64	34.4%	29.6%	61.8%
1	MORTGAGE	97	52.1%	56.6%	32.4%
2	OWN	25	13.5%	13.8%	5.9%

Categorical variable distribution results further validate risk associations:

**Employment duration:** The shorter the employment period, the higher the proportion of high-risk individuals. However, there are no high-risk individuals in the 10+ years group, indicating that „employment stability“ is an important predictor of credit risk, consistent with the paper’s positioning of „controllable financial behaviour“ as a proxy characteristic.

**Housing type:** The proportion of high-risk individuals among renters is significantly higher than among those with mortgages or those who own their homes, reflecting the risk-inhibiting effect of financial stability brought by housing assets, providing empirical support for the concept of „housing type as a proxy characteristic.“

#### 4.4 Summary of Descriptive Analysis

A descriptive analysis of the Lending Club sample yields the following key conclusions: The sample as a whole conforms to the characteristics of the consumer credit market, with a default rate of 18.3% and average FICO scores concentrated in the mid-range credit range, demonstrating good data representativeness. Traditional variables are significantly correlated with credit risk, but

alternative features can further differentiate risk. In particular, the high-risk group clusters around the characteristics of „high debt, low stability, and no assets,“ validating the risk-indicating value of alternative features. Differences in the distribution of categorical variables suggest that „controllable financial behaviour stability“ and „asset mitigation“ are key factors in reducing credit risk, laying the foundation for subsequent model validation of the marginal contribution of alternative features [9, 10].

#### 5. Conclusion

This paper uses the traditional FICO credit score as a benchmark, focusing on the core dimension of „borrower controllable financial behaviour“ to construct alternative features. Based on the Lending Club sample, this paper conducts empirical testing using the XGBoost algorithm and 50-fold cross-validation, forming an integrated research approach: „feature selection - model construction - performance evaluation - robustness testing and interpretability analysis.“

The empirical results demonstrate that: First, alternative features can significantly improve the ability to distinguish and early identify high-risk borrowers without



violating existing compliance frameworks. Compared to the baseline model relying solely on FICO, the extended model incorporating controllable financial behaviour features achieves robust gains in core metrics such as AUC, KS, and F1, effectively filling the gap in dynamic risk signals that traditional static scores struggle to capture and expanding the predictive space beyond FICO scores. Second, the information on subjective repayment willingness and financial stability contained in borrowers' controllable financial behaviour can effectively supplement the traditional model's sole reliance on historical credit records, validating the research hypothesis that „controllable behaviour variables have independent risk-indicating value.“ Methodologically, this study leverages XGBoost's strengths in capturing nonlinear feature correlations, combined with the TreeSHAP algorithm to quantify feature contributions. This approach not only addresses the „black box“ nature of traditional models and meets regulatory requirements for interpretability, but also establishes a replicable „traditional scoring + alternative features“ compliant model framework, providing a methodological reference for credit risk modelling in the intelligent finance sector.

In terms of practical value, this study provides banks and licensed fintech institutions with a refined risk control solution: by introducing behavioural alternative features alone, risk control accuracy can be improved without replacing FICO scores. It also provides a more inclusive credit assessment path for groups underrepresented by traditional models, promoting the coordinated development of financial security and inclusiveness.

This study also has limitations: the sample is sourced exclusively from Lending Club, requiring further validation of its generalizability; it fails to fully incorporate the time-varying impact of macroeconomic cycles on controllable financial behaviour; and the quality standardization and fairness governance of alternative data require further improvement. Future research is needed to further strengthen the theoretical and applied foundations of intelligent finance in credit risk management, focusing on the integration of multi-source alternative data, the applica-

tion of real-time risk control technologies, and the coordinated optimization of fairness and privacy protection.

## References

- [1] Altman, E. I., & Saunders, A. (1997). Credit risk measurement: Developments over the last 20 years. *Journal of Banking & Finance*, 21(11-12), 1721–1742.
- [2] Hlongwane R, Ramaboa KKKM, Mongwe W (2024). Enhancing credit scoring accuracy with a comprehensive evaluation of alternative data. *PLOS ONE* 19(5): e0303566.
- [3] Bazarbash, M. (2019). Fintech in financial inclusion: Machine learning applications in assessing credit risk. IMF Working Paper, No. 2019/109.
- [4] Wang, Q., Smith, J., & Johnson, L. (2024). Enhancing credit scoring accuracy with a comprehensive evaluation of alternative data. *Journal of Empirical Finance*, 78, 102-118.
- [5] Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., Garcia, S., Gil-Lopez, S., Molina, D., & Herrera, F. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities, and challenges toward responsible AI. *Information Fusion*, 58, 82-115.
- [6] Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30, 4765-4774.
- [7] Wise, C., & Chen, L. (2023). Empowering credit inclusion: A deeper perspective on new-to-credit consumers. *Journal of Consumer Affairs*, 57(3), 890-912.
- [8] Zhang, L., Wang, P., & Liu, X. (2024). Credit risk assessment of small and micro enterprises based on machine learning. *Heliyon*, 10(5), e27096.
- [9] Iyer, R., Khwaja, A. I., Luttmer, E. F. P., & Shue, K. (2016). Screening peers softly: Inferring the quality of small borrowers. *Management Science*, 62(6), 1554–1577. <https://doi.org/10.1287/mnsc.2015.2208>
- [10] Bastani, H., Ascarza, E., & Choudhury, P. (2019). Predicting consumer default: A machine learning approach. Harvard Business School Working Paper No. 19-047. <https://doi.org/10.2139/ssrn.3274360>