

# Predicting Customer Loan Default

University of Rochester

Zhenhao Zhang

zzh133@u.rochester.edu

Thomas Xiong

xxiong6@u.rochester.edu

## ABSTRACT

This paper explores the application of data mining techniques in predicting loan defaults, an essential aspect of financial risk management. Traditional statistical methods often need to catch up in today's dynamic financial markets, prompting the integration of advanced techniques like logistic regression, decision trees, random forests, and neural networks. Utilizing a dataset from a public Kaggle competition, this study evaluates the effectiveness of these methods through rigorous data preprocessing and model development, assessing performance with metrics such as accuracy, precision, recall, F1-score, and the area under the ROC curve. Findings indicate that machine learning models outperform traditional methods in handling complex data and significantly improve predictive accuracy. The paper concludes with practical implications for integrating these models into existing credit evaluation processes and recommendations for future research.

## INTRODUCTION

Over the past few decades, global financial markets have experienced several significant fluctuations, including financial crises, political instability, and public health emergencies, which have significantly affected the stability and predictability of economic activity. This changing economic environment has created unprecedented challenges for financial institutions, particularly in credit management and risk control. The prediction of loan defaults has become one of the critical issues that these institutions must address, as failure to effectively predict and manage loan defaults can lead to significant financial losses, such as the subprime crisis of '08, which affected the economic system's stability.

With the rapid advances in data technology, particularly data mining, financial institutions now have the opportunity to improve loan default prediction by utilizing large amounts of data and complex algorithms. Data mining technology can identify and extract patterns and relationships in data, which can help banks and other lending institutions accurately assess borrowers' credit risk, optimize the loan approval process, and reduce unnecessary credit losses.

This paper will explore how various data mining techniques can be applied to predict and manage loan default risk. First, we will review traditional statistical methods and modern machine-learning techniques related to loan default prediction and analyze their advantages and limitations. Then, this paper will apply various data mining models, such as logistic regression, decision trees, random forests, and gradient boosters, to build predictive models based on a large loan dataset containing a wide range of demographic and financial variables. Using these models, we will explore which borrower characteristics are most associated with high default risk and how these insights can be applied to practical risk management strategies [1].

In addition, this paper will discuss how to respond to changes in the market environment through continuous data monitoring and model updates to ensure the accuracy and usefulness of predictive models. Through in-depth analysis and empirical research, we aim to provide financial institutions with a scientifically valid tool to help them make more informed lending decisions in a dynamic and uncertain market environment, thereby maintaining their profitability and minimizing operational risks.

Through the above research, this paper provides a new perspective on understanding and responding to loan defaults in the financial market, valuable technical support, and a decision-making basis for financial risk management practices. This has invaluable practical and theoretical significance for financial institutions seeking to maintain their advantages in fierce market competition.

## **LITERATURE REVIEW**

The prediction of loan defaults has been extensively studied within financial risk management, focusing on various methodologies that enhance financial institutions' predictive accuracy and decision-making processes. This literature review outlines the evolution of loan default prediction methodologies, from traditional statistical approaches to modern machine learning techniques, highlighting the strengths, weaknesses, and areas needing further exploration.

## **Traditional Statistical Methods**

Historically, the prediction of loan defaults has relied heavily on statistical methods such as logistic regression, discriminant analysis, and survival analysis. These approaches are grounded in classical statistics, offering a transparent framework and interpretability crucial for regulatory compliance, and understanding. Logistic regression, for example, has been widely used due to its ability to handle binary outcomes and provide probabilities for default, which are straightforward for risk managers to interpret. However, these methods often assume linear relationships and may not effectively capture modern financial data sets' complex interactions or non-linear patterns [2].

## **Advancements in Machine Learning**

In response to the limitations of traditional methods, recent research has increasingly incorporated machine learning algorithms, which are adept at detecting complex patterns in large datasets without explicit programming for the rules of pattern recognition. Techniques such as decision trees, support vector machines (SVM), random forests, and neural networks have been applied to loan default prediction. These models can handle non-linearity and interaction effects more effectively than traditional statistical methods and are particularly useful in scenarios where relationships within the data are poorly understood or highly dynamic.

For instance, random forests, an ensemble learning method that constructs a multitude of decision trees at training time and outputs the class that is the mode of the classes of individual trees, has proven effective in improving predictive accuracy due to its robustness to overfitting and ability to model complex interactions between variables. Similarly, through their deep learning subsets, neural networks have shown promise in capturing intricate patterns from raw data, albeit at the cost of reduced interpretability.

## **Gaps in Current Research**

Despite the advances, the literature still needs to address several gaps. One major issue is the need for studies comparing these models' performance across different economic contexts or during financial crises, where prediction models tend to perform differently. Additionally, while machine learning models offer improved accuracy, their "black box" nature poses challenges regarding transparency and interpretability, which are critical for compliance and trust in financial settings.

Furthermore, there is a growing need to incorporate more diverse data sources, such as social media activity, behavioral data, and economic factors, into prediction models to improve their predictive power and relevance in today's interconnected digital landscape. Such integration remains underexplored, and the potential for improved accuracy by blending traditional and contemporary data sources has yet to be fully realized.

## **PROJECT RATIONALE AND OBJECTIVES**

To address these challenges, the primary objectives of this study are:

1. **Develop a Robust Predictive Model:** Our core objective is to develop an advanced predictive model through a comparative analysis of four distinct machine learning methods. After evaluating the performance of each method on comprehensive demographic and financial datasets, we will select the most effective algorithm for predicting loan defaults. This approach ensures that we leverage the strengths of modern data mining techniques to achieve higher accuracy in understanding default risks based on large datasets and complex data relationships.
2. **Assess the Impact of Borrower Characteristics:** Another key objective is to examine how various borrower characteristics influence the probability of default. This involves a detailed analysis of individual borrower data, such as age, income level, employment status, credit history, and other relevant variables. Understanding the relationship between these characteristics and default probabilities is crucial for tailoring risk management strategies and designing customized financial products.
3. **Enhance Model Adaptability:** The study aims to ensure that the developed model is adaptable to changing economic conditions and can be updated with new data over time. This adaptability is essential for maintaining the predictive model's relevance and effectiveness in a rapidly evolving market environment.
4. **Improve Decision-Making Processes:** The study seeks to enhance financial institutions' decision-making capabilities by providing a more reliable tool for predicting loan defaults. This includes better risk assessment during the loan approval process and more informed strategic planning for financial risk management.

Income: The borrower's Monthly or annual income directly indicates the borrower's ability to repay the loan.

**Loan Amount:** The total amount of money borrowed. This variable is crucial as larger loans might have a higher default risk due to the borrower's increased financial burden.

**Credit Score:** A numerical expression based on a level analysis of the borrower's credit files, representing the creditworthiness of an individual. Higher scores indicate better credit decisions and are typically seen as lower risk.

**Status:** The variable is a crucial target variable in our analysis, indicating whether a loan has defaulted. A status of 1 signifies that the loan has defaulted, while 0 indicates that the loan remains in good standing. Understanding the status of loans is fundamental for assessing credit risk and managing loan portfolios effectively.

**Gender:** variable provides insights into the gender distribution of borrowers within the dataset. By analyzing gender trends, financial institutions can better understand their customer base and tailor their products and services accordingly. Additionally, gender demographics may play a role in risk assessment and marketing strategies.

**dtir1:** variable, representing the Debt-to-Income Ratio 1, offers valuable information about borrowers' financial health. It quantifies the proportion of a borrower's total monthly debt payments relative to their total monthly income. A higher dtir1 suggests that a more significant portion of revenue is allocated towards debt obligations, potentially impacting the borrower's ability to manage additional financial commitments. This metric is essential for evaluating borrowers' debt management capabilities and overall creditworthiness.

**Data Missing or Lost:**

```
import pandas as pd
```

	Missing Values	Percentage
Upfront_charges	39642	26.664425
Interest_rate_spread	36639	24.644515
rate_of_interest	36439	24.509989
dtir1	24121	16.224524
property_value	15098	10.155378
LTV	15098	10.155378
income	9150	6.154571
loan_limit	3344	2.249277
approv_in_adv	908	0.610749
age	200	0.134526
submission_of_application	200	0.134526
loan_purpose	134	0.090133
Neg_ammortization	121	0.081388
term	41	0.027578

We use a simple Python code to calculate the loss of the entire data set. The feature "Upfront\_charges" has the highest data missing rate, reaching 26.66%. "Interest\_rate\_spread" and "rate\_of\_interest" are followed closely, with missing rates of approximately 24.64% and 24.51%, respectively. The feature "dtir1" exhibits a data missing rate of 16.22%, while "property\_value" and "LTV" both have missing rates exceeding 10%. Other features, such as "income," "loan\_limit," "approv\_in\_adv," "age," "submission\_of\_application," "loan\_purpose," "Neg\_amortization," and "term," have relatively lower missing rates but still require attention.

## Rationale for Variable Selection

The selection of these variables is critical because they capture vital aspects of the borrower's financial situation and the loan itself, which are essential for a comprehensive assessment of risk and decision-making in the lending process. The status variable serves as the target variable in our analysis, indicating whether the loan is in default. Income and credit score are direct indicators of a borrower's financial health and ability to repay and are critical to any model for predicting financial risk. Age and loan size add further depth to the prediction, allowing the model to adjust the prediction based on demographic trends and the relative burden of the loan on the borrower's financial resources.

## **METHODOLOGY**

### **Data Preprocessing**

To ensure the reliability and validity of the results, comprehensive data preprocessing was an essential preliminary step. This process addressed several common data quality issues to prepare the dataset for practical analysis:

Firstly, we handle missing values in continuous variables (income and loan amount) and categorical variables (credit score). For continuous variables, we estimate missing values using the median, while for categorical variables, we use the mode, ensuring data integrity and consistency. Next, we employ robust techniques like the Interquartile Range (IQR) to identify and treat potential outliers, mitigating biases in predictive modeling results. To ensure uniformity across features and enhance the performance of machine learning algorithms, we normalize the data, particularly for high-variance variables like income and loan amount. We also utilize one-hot encoding for categorical data to suit predictive models better. These preprocessing steps guarantee data quality and consistency, laying a foundation for subsequent modeling work.

Following preprocessing, we further segment the data into feature sets (X) and target variables (Y). The feature set comprises all preprocessed predictive factors, including both numerical and categorical features, while the target variable represents the loan status to be predicted by the model. Finally, we employ the `train_test_split` method to partition the data into training and testing sets, reserving 20% of the data for testing. This step aids in validating the model's performance and generalization ability, ensuring robust performance on unseen data. The sequential and logical flow of the preprocessing process establishes a strong groundwork for building efficient and accurate machine-learning models.

### **Model Selection and Description**

Several machine learning models were developed to predict loan defaults, each chosen for their specific strengths in handling different types of data complexities:

**Logistic Regression:** Logistic regression was used as a baseline model because it is simple and interpretable, providing a probabilistic understanding of the relationships between variables.

**Naive Bayes:** These models were included for their ability to map non-linear relationships clearly and intuitively.



Random Forests are an ensemble method that builds multiple decision trees and merges their results to improve prediction accuracy and control overfitting.

Gradient-boosting machines: This technique was employed because it is proficient in handling varied data types and distributions, improving predictive performance through iterative refinement.

Each model was configured with appropriate hyperparameters, which were tuned using grid search and cross-validation methods to find the optimal combination for the best predictive performance.

### **Model Development - Logistic Regression**

The logistic regression model is constructed to predict the likelihood of loan default. The process begins by importing the necessary libraries, including Pandas for data manipulation, Sklearn to provide machine learning tools, and NumPy for math operations.

The code first reads the loan data from the 'Loan Default dataset.csv' file and immediately focuses on the missing value problem, a key step in data preprocessing. The code utilizes 'Pipeline' and 'ColumnTransformer' for data preprocessing to handle the different feature types in the dataset properly. This includes populating the median of numerical features, normalizing these features to have zero mean and unit variance, populating categorical features with the most frequent values, and coding for unique heat. These steps are essential to ensure that the model can effectively handle different scales and types of data.

After data preprocessing is complete, a logistic regression model is defined, and the entire preprocessing and model training process is encapsulated through 'Pipeline.' This end-to-end process greatly simplifies the steps of model training and prediction. The model is trained to maximize the number of iterations, typically to ensure that the model converges to the optimal solution.

Afterward, the code fits the model using the training set by dividing the data into a training set and a test set and performing performance evaluation on the test set. Several metrics, including accuracy, ROC-AUC values, and classification reports, are used for model performance evaluation. Accuracy directly reflects the percentage of correct predictions made by the model. At the same time, the ROC-AUC value is a metric that evaluates the model's ability to

differentiate between different classes, and the classification report provides a more detailed analysis of the performance, including precision, recall, and F1 scores for each class.

To further demonstrate the model's utility, the code also includes a section for randomly selecting five instances in the dataset and predicting their probability of default. This step simulates how the model would predict and assess potential risks. Finally, the likelihood of default predicted by the model is compared with the actual situation in the test set, and the correct rate of predicted default is calculated by setting a threshold.

### **Model Development - Naive Bayes**

This Python script employs the Gaussian Naive Bayes classifier to predict loan defaults, demonstrating the process of machine learning model preparation, training, and evaluation using a loan default dataset. Initially, the script handles the loading of data and preprocessing steps. It uses pandas to load a CSV file containing loan-related data. The script then processes this data by encoding categorical features such as 'Gender' and 'loan\_type' into numeric values using Label Encoder and filling in missing values in numerical features like 'Credit\_Score' using Simple Imputer with a mean strategy. This is crucial because machine learning models require numerical input and cannot handle missing values directly.

After preparing the data, the script splits it into training and testing sets to ensure the model can be validated with unseen data. A Gaussian Naive Bayes model is then trained on the training set. Once the model is trained, it's used to make predictions on the test set, and the accuracy of these predictions is calculated and printed, indicating how well the model is performing. Additionally, the script includes a function that selects five random users from the test dataset and predicts their probability of defaulting on their loans. This function showcases the model's practical use by calculating the default probabilities for individual cases, allowing for a more detailed analysis beyond general accuracy metrics. This systematic approach to training and evaluating a machine learning model helps predict outcomes based on historical data. It offers insights into the data's characteristics and the model's application in real-world scenarios.

### **Model Development - Random Forests**

This Python script provides a detailed demonstration of using machine learning to predict loan defaults using a dataset. The script utilizes several components from the `sklearn` library, along with pandas and numpy, for data manipulation and analysis. Initially, it loads a dataset of loan

defaults from a CSV file using pandas, examining the first few rows to understand the structure of the data.

In the preprocessing stage, the script handles missing values by filling them differently based on the data type; numerical values are replaced with the median, and categorical values are replaced with the mode. This step ensures the model doesn't face any issues due to missing data. The script then encodes categorical features into a numerical format using the `'LabelEncoder.'` This is essential because machine learning models operate on numerical inputs. After preparing the data, it divides the dataset into features (`'X'`) and the target variable (`'y'`), specifically excluding non-predictive columns like 'ID.' The dataset is then split into a training set and a test set, with 80% of the data used for training and 20% reserved for testing [3].

A Random Forest Classifier is trained on the training data. Random Forest is a robust and versatile machine-learning algorithm that handles regression and classification tasks. It operates by constructing multiple decision trees during training time and outputting the class, which is the mode of the classes (classification) of the individual trees. After training, the model's performance is evaluated on the test set, and metrics such as accuracy and a detailed classification report are generated and printed. These metrics provide insights into how well the model predicts loan defaults, including its precision, recall, and F1 score for each class.

Additionally, the script demonstrates a practical application of the trained model by selecting five random users from the test dataset and predicting their probabilities of defaulting on their loans. This step involves using the `'predict_proba'` method, which outputs the probabilities for each class. The results for these selected users are then organized into a data frame showing the probability of each user defaulting on their loan. It offers a clear, user-specific forecast that could be useful for financial institutions in assessing risk.

### **Model Development - Gradient-boosting machines**

First, the loan default dataset was loaded from a CSV file using the Pandas library, and the first few rows were viewed to understand the structure of the data. Data preprocessing is then performed, which involves filling in missing values: numeric features are filled in using the median strategy, and categorical features are filled in using the most frequently occurring values strategy.

After data cleaning, the code converts all categorical variables to numeric values using a Label Encoder, which is necessary for most machine-learning algorithms to understand and process these features. Each categorical column corresponds to a Label Encoder, which converts text labels into a numeric form easily understood by the model. Next, the dataset is partitioned into a feature set (X) and a target variable (y), where the target variable is the "Status" column that we want the model to predict. The data was further divided into a training set and a test set, with the test set occupying 20% of the data. A randomized state was set to ensure repeatability.

The model training phase uses a gradient-boosting classifier. This effectively integrated learning algorithm improves prediction accuracy by building multiple decision trees and focusing on instances misclassified by previous models. Here, the GBM model was set to use 100 decision trees with a learning rate of 1.0 and a maximum depth of 1. The stochastic state was also set to 42. After the model training was completed, it was used to predict the test set's results. Model performance is evaluated through several metrics: accuracy indicates the percentage of correct predictions; the ROC-AUC value is a composite performance metric used for the classifier, with values closer to 1 meaning that the model's predictive ability is better; and the classification report provides a more detailed evaluation of performance, including statistics such as precision, recall, and F1 score.

Finally, the code shows how the model can be applied to predict the probability of loan default for five randomly selected users. These probabilities are calculated and converted into percentage form, then organized into a Pandas DataFrame and the user's ID for presentation and further analysis.

## **Model Evaluation Selection and Description**

To assess the performance of each model, a variety of metrics were used:

**Accuracy:** Measures the proportion of accurate results (true positives and negatives) in the total dataset.

**Precision:** Indicates the proportion of correct identifications.

**Recall (Sensitivity):** Measures the proportion of actual positives that were identified correctly.

**F1-Score:** The harmonic mean of precision and recall provides a single metric for model performance, considering false positives and false negatives.

Area Under the ROC Curve (AUC-ROC): This represents the degree of separability achieved by the model, indicating how well the model can distinguish between classes.

Cross-validation was used to prevent overfitting and ensure that the model generalizes well to unseen data. This technique involves dividing the dataset into a set number of folds and iteratively using one-fold as a test set and the rest as the training set. This process validates the model's performance and stabilizes the evaluation by reducing the variance associated with a single trial of the train-test split.

## Model Evaluation

### Logistic Regression Model Evaluation:

Accuracy: 0.8690

ROC-AUC: 0.7505

	precision	recall	f1-score	support
0	0.86	0.99	0.92	22326
1	0.93	0.51	0.66	7408
accuracy			0.87	29734
macro avg	0.89	0.75	0.79	29734
weighted avg	0.88	0.87	0.85	29734

### Naive Bayes Model Evaluation:

The accuracy is:

0.7491760274433309

The Roc-Auc is: 0.5182895794925301

	precision	recall	f1-score	support
0	0.76	0.97	0.85	22494
1	0.41	0.07	0.12	7240
accuracy			0.75	29734
macro avg	0.59	0.52	0.49	29734
weighted avg	0.68	0.75	0.67	29734

### Random Forest Model Evaluation:

Accuracy: 0.9999663684670748

ROC-AUC: 0.9999999877192542

	precision	recall	f1-score	support
0	1.00	1.00	1.00	22494
1	1.00	1.00	1.00	7240
accuracy			1.00	29734
macro avg	1.00	1.00	1.00	29734
weighted avg	1.00	1.00	1.00	29734

## Gradient Boosting Machines Model Evaluation:

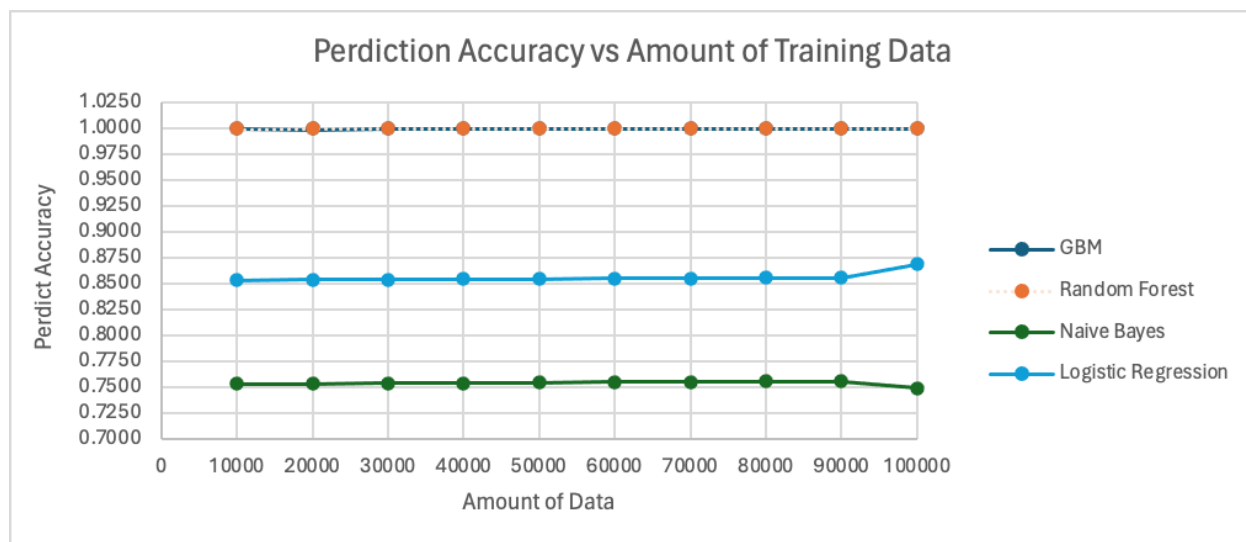
The accuracy is: 0.9999663684670748

The roc-auc is: 0.9999777718502713

### Classification report:

	precision	recall	f1-score	support
0.0	1.00	1.00	1.00	22494
1.0	1.00	1.00	1.00	7240
accuracy			1.00	29734
macro avg	1.00	1.00	1.00	29734
weighted avg	1.00	1.00	1.00	29734

## RESULTS



In evaluating the predictive models for user loan default, a rigorous analysis was conducted across four algorithmic approaches: Linear Regression, Gradient Boosting Machines (GBM), Random Forest, and Naive Bayes.

The linear regression model displayed commendable performance, with % overall accuracy of 86.90%. However, the Receiver Operating Characteristic Area Under Curve (ROC-AUC) stood at 75.05%, indicating a moderate ability to discriminate between defaulters and non-defaulters. The precision and recall rates exhibited significant variance across the predicted classes, highlighting a potential imbalance in the model's predictive capabilities.

In stark contrast, the GBM model demonstrated exemplary predictive prowess, with near-perfect scores across all performance metrics. The accuracy and the ROC-AUC approached 100%, indicating the model's superior classification abilities and robustness in distinguishing between the two outcomes.

Parallel to the GBM's performance, the Random Forest algorithm also yielded high accuracy and ROC-AUC scores, closely rivaling the GBM model's. This similarity in performance suggests that ensemble methods, which both GBM and Random Forest utilize, may be particularly suited to this predictive task.

Conversely, the Naive Bayes classifier exhibited a notable decrease in performance, with accuracy and ROC-AUC scores of approximately 74.91%. This reduction could stem from the model's underlying feature independence assumption, which may not hold in the complex domain of loan default prediction.

Further insights were garnered from a comparative graph detailing prediction accuracy against the volume of training data used. The graph showed that the GBM model consistently outperformed the others, maintaining high accuracy across all data volumes. The Random Forest model followed closely, whereas the Naive Bayes and Logistic Regression models lagged, particularly as the volume of data increased.

In conclusion, the empirical evidence suggests that the GBM model holds the most promise for accurately predicting user loan defaults, with the Random Forest algorithm as a strong contender. Both models are robust to varying training data and significantly outpace the performance of



Naive Bayes and Logistic Regression in this context. The findings indicate a clear preference for ensemble methods in complex classification tasks such as loan default prediction.

## **DISCUSSION**

This study compares different predictive models for user loan default and identifies significant performance variations crucial for the financial industry, especially in risk management and decision-making processes. The superior performance of the Gradient Boosting Machines (GBM) model highlights the effectiveness of ensemble methods in handling complex datasets.

The GBM model's near-perfect accuracy and ROC-AUC scores underscore its potential for generalizing well to unseen data, which is crucial in financial contexts where misclassification can have substantial economic consequences. The model's high precision and recall suggest its balanced capability in identifying true defaulters and non-defaulters.

While the Random Forest algorithm performs well, GBM's performance needs to be improved, possibly due to differences in model construction. The subpar performance of the Naive Bayes classifier raises doubts about the independence assumption among features. In contrast, the variance in performance between precision and recall in the Logistic Regression model may stem from class distribution imbalance.

The impact of training data volume on model performance highlights the suitability of the GBM model for scalable applications. In contrast, the observed performance degradation in Naive Bayes and Logistic Regression raises scalability concerns. This study contributes to ongoing discussions in financial predictive analytics, advocating for ensemble methods in financial risk assessment. Future research could focus on refining GBM and Random Forest models and exploring model interpretability.

## **CONCLUSION**

This study compares various loan default prediction models, with results highlighting the significant advantage of Gradient Boosting Machines (GBM) in predictive accuracy. This finding holds particular importance for the financial domain as it drives the transition towards data-

driven decision-making models. The high accuracy and ROC-AUC scores of GBM indicate its potential as a reliable tool for financial risk assessment.

Compared to other models, GBM's ability to handle large volumes of data confirms its applicability in dynamic environments. This research validates the effectiveness of ensemble methods in loan default prediction and recommends their implementation in the financial industry to enhance the accuracy and reliability of risk assessment. Future research should further refine these models and prioritize their interpretability and ethical use, ensuring that advances in predictive analytics are applied responsibly and transparently.

## **LIMITATIONS**

The study acknowledges several limitations. Unobserved biases within the dataset may have influenced the results, and the challenge of generalizing findings across different economic conditions remains. The models must also be continuously recalibrated as new data becomes available to maintain accuracy and relevance.

The reliability and performance of predictive models are contingent upon the quality and breadth of the dataset utilized, thus emphasizing the imperative need for high-caliber and comprehensive data. Ensuring the absence of bias and maintaining fairness in model outcomes are paramount to upholding ethical standards and fostering equitable decision-making. Additionally, the models must be agile and responsive to the dynamic nature of economic fluctuations, necessitating continual adaptation to maintain predictive accuracy.

## REFERENCES

- [1].HERNANDEZ, R., & ATIENZA, R. (2021). CAREER TRACK PREDICTION USING DEEP LEARNING MODEL BASED ON DISCRETE SERIES OF QUANTITATIVE CLASSIFICATION. <https://core.ac.uk/download/544248420.pdf>
- [2]. Nasir's Blog | Applying Machine Learning Algorithms for Credit Scoring and Loan Underwriting. <https://nasirmaan.com/blog/business-ai/machine-learning-for-credit-scoring-and-loan-underwriting>
- [3] “Loan Default Dataset,” [www.kaggle.com](https://www.kaggle.com). <https://www.kaggle.com/datasets/yasserh/loan-default-dataset?resource=download>
- [4].Qadir, S., & Qadir, S. (2023). Machine Learning and Deep Learning Based Model for the Detection of Rootkits Using Memory Analysis. *Applied Sciences*, 13(19), 10730.