

Risk Analysis for Home Credit Default:

Exploratory Data Analysis and Predictive Modeling

Gullapalli Sathar¹, P.Prasanna Lakshmi², P.Naveena³, N.Sreeshma⁴, S.Hansika⁵, B.Vedhasri⁶,
M.Praveen⁷, M.Venkata Abhilash⁸, P.Siddarth⁹, Samhith Bade¹⁰, K.Gopi Chand¹¹

¹Assistant Professor, Department of CSE-(Cys,DS) and AI&DS, VNRVJiet, Hyderabad

^{2,3,4,5,6,7,8,9,10,11}UG Students, Department of CSE-(Cys,DS) and AI&DS, VNRVJiet, Hyderabad
sathar9000@gmail.com¹, palojuprasannalakshmi@gmail.com², naveena.paruchuri1@gmail.com³,
sreeshmanalla03@gmail.com⁴, hansikasamala0609@gmail.com⁵, vedhavinni@gmail.com⁶,
praveepravee489@gmail.com⁷, mupparaju.abhilash@gmail.com⁸, siddharthsleeps@gmail.com⁹,
samhithbade44@gmail.com¹⁰, gopichandkumba31@gmail.com¹¹

Abstract: This paper investigates the complex risk factors associated with home credit default through comprehensive Exploratory Data Analysis (EDA) and the design of robust Predictive Modeling techniques. Credit default poses substantial financial challenges for lending institutions; therefore, proactive risk assessment is critical for maintaining financial stability. We begin by performing extensive EDA to understand the distribution and characteristics of the Home Credit dataset, identifying key features influencing loan repayment capacity. By employing a combination of statistical evaluation and data visualization, we pinpoint critical variables related to applicants' demographics, financial history, and loan attributes. This analysis informs the development of predictive models using various machine learning algorithms, including Logistic Regression, Random Forest, and Gradient Boosting Machines, to estimate the likelihood of default for new applicants. The performance of these models is rigorously evaluated using standard metrics such as ROC-AUC, Precision, Recall, and F1-Score. The project's findings advance targeted lending strategies and aid financial institutions in making informed, data-driven decisions to mitigate risk and ensure sustainable lending practices.

Keywords— Home Credit Default, Risk Analysis, Predictive Modeling, Machine Learning, Financial Stability, Credit Risk.

1 INTRODUCTION

Credit default risk poses substantial challenges for lending institutions. The increased occurrence of financial non-repayment necessitates robust analytics capable of evaluating risk determinants and supporting informed decision-making. This study outlines an investigative process utilizing statistical modeling and machine learning to predict default cases and reduce associated risks. This project addresses the critical need for proactive default prediction by leveraging modern machine learning and statistical techniques. The study employs a large-scale Home Credit dataset, encompassing a wide array of variables, including demographic information, historical credit performance, financial indicators, and specific loan characteristics. Our methodology is structured around two principal phases: comprehensive Exploratory Data Analysis (EDA) and the implementation of robust Predictive Modeling. The primary objectives of this research are: Perform exploratory data analysis (EDA) to understand the distribution, characteristics, and quality of the home credit dataset. Identify key factors and features influencing credit default through data visualization and statistical analysis, moving beyond mere correlation to establish clear predictors. Develop and optimize predictive models using state-of-the-art machine learning algorithms to accurately predict the likelihood of default for home credit applicants. Evaluate the performance of these predictive models using appropriate evaluation metrics (e.g., AUC-ROC, Precision-Recall curves) to determine their robustness and

operational feasibility. Provide actionable insights and recommendations for financial institutions to refine their risk management policies and implement targeted strategies to mitigate credit default risk. The successful completion of this project will provide financial institutions with an advanced, data-backed tool for prescreening applicants, thereby improving loan portfolio quality and enhancing overall financial stability.

2 LITERATURE SURVEY

The application of quantitative methods to assess credit risk has a long history, traditionally relying on statistical models. Early efforts focused primarily on linear regression and logistic regression to estimate default probability based on readily available financial and socioeconomic variables. These conventional approaches provided foundational insights into the relationship between borrower characteristics and default outcomes, serving as the backbone of early credit scoring systems.

With the proliferation of big data and advances in computational power, the field of credit risk modeling has undergone a rapid transformation, heavily adopting Machine Learning (ML) techniques. Studies have demonstrated that ML models offer superior prediction accuracy and better capacity to capture complex, non-linear relationships inherent in credit data compared to traditional statistical methods.

Early Machine Learning Applications in this domain primarily utilized methods like Support Vector Machines (SVMs) and Artificial Neural Networks (ANNs). For instance, some researchers have shown that ANNs, with their ability to model complex, multi-layered dependencies, outperformed logistic regression models in detecting high-risk applicants, especially when dealing with noisy or incomplete data. However, a significant drawback of these early complex models was their “black box” nature, which made interpreting the rationale behind a default prediction difficult—a critical requirement for regulated financial institutions.

The trend has recently shifted towards highly effective ensemble methods. Random Forest (RF) models have proven robust due to their resistance to overfitting, ability to handle high-dimensional data, and capacity to provide interpretable feature importance rankings, which directly aid analysts in understanding which variables drive the default risk. More recently, Gradient Boosting Machines (GBM), including popular implementations like XGBoost and LightGBM, have achieved state-of-the-art results in credit scoring competitions. These models sequentially build weak predictive models, where each new model corrects the errors of the previous ones, leading to exceptionally high predictive power.

Current research also emphasizes the crucial role of Exploratory Data Analysis (EDA) and Feature Engineering in improving model performance. Studies increasingly focus on incorporating unconventional data sources, such as unstructured text data or data aggregated from multiple sources (like bureau data and previous application records), to create robust predictive features. Furthermore, a significant body of work is now dedicated to Explainable AI (XAI) in credit risk, ensuring that powerful ML models not only predict default but can also provide human-interpretable reasons for their decisions, thereby satisfying regulatory and ethical requirements. Our project builds upon this literature by performing a comparative analysis of modern ensemble methods (RF and GBM) informed by rigorous EDA, aiming to deliver both high prediction accuracy and actionable, interpretable insights.

3 PROPOSED METHODOLOGY

The methodology uses a five-stage pipeline to achieve robust default prediction. It begins with Data Collection from Home Credit datasets, followed by Pre-Processing involving cleaning, handling missing values, and encoding categorical variables. The Exploratory Data Analysis (EDA) phase then identifies key drivers of default using visualization and statistical testing. Data Augmentation through advanced feature engineering (e.g., creating aggregates and ratios) enhances the predictive power. Finally, Predictive

Modeling involves training and optimizing advanced machine learning algorithms (like Gradient Boosting and Random Forest), culminating in a comprehensive performance evaluation to deliver actionable insights for risk mitigation.

3.1 Data Collection

The Data Collection phase is foundational, starting with the acquisition of the primary Home Credit Default Risk dataset, which contains the crucial target variable (default status) and core applicant demographics, such as income and age. To build a robust and highly predictive model, this core information must be augmented by integrating data from several supplementary financial history tables. These supplementary files include detailed records of Previous Applications for past loans, comprehensive Bureau Data from external credit agencies (detailing credit history and debt levels), and monthly transactional data like POS Cash Balance, Installment Payments, and Credit Card Balance snapshots. All these data tables are merged using the unique applicant identifier to form a unified, feature-rich dataset. An initial data validation step is mandatory to ensure the consistency of data types, assess the sheer volume of data, and quantify the extent of missing values across all features before the cleaning process begins.

3.2 Pre-Processing

The Pre-Processing stage is crucial for transforming the complex raw data into a pristine, structured format optimized for machine learning algorithms. Data cleaning techniques such as handling missing values, outlier removal, and normalization are applied to ensure high-quality input for analysis. The first challenge addressed is Missing Value Imputation: for numerical fields, missingness is handled using the median or mean, while for categorical fields, the mode is used, or a new category, 'Missing', is created. Severe outliers, which can skew model training, are managed using robust techniques like Winsorization, capping extreme values at predefined percentiles. Next, Categorical Encoding is performed; nominal variables (e.g., occupation) are converted using One-Hot Encoding to avoid false ordering assumptions, while ordinal variables are handled with Label Encoding. Finally, all continuous numerical features must be Standardized (using Z-score normalization) or Min-Max Scaled. This scaling ensures that all features contribute equally during model training, preventing high-magnitude variables from dominating the learning process.

3.3 Exploratory Data Analysis

The initial phase of the Exploratory Data Analysis (EDA) is crucial for characterizing the Home Credit dataset and identifying fundamental structural challenges, which is a necessary prerequisite for effective model design. This process is initiated by scrutinizing the Target Variable (default status) to assess the fundamental nature of the problem: credit risk modeling universally entails a challenge of severe class imbalance. Quantifying the ratio of non-defaulters to defaulters (Target = 0 vs. Target = 1) dictates the selection of appropriate sampling techniques (e.g., SMOTE, undersampling) and robust evaluation metrics (e.g., ROC-AUC, Precision-Recall) over simple accuracy. Furthermore, Univariate Analysis is performed across all numerical features using visualizations like histograms and box plots to examine their empirical distribution, skewness, and kurtosis. This step is vital for detecting high-magnitude outliers (e.g., credit amounts exceeding 3σ from the mean) and determining appropriate treatments, such as Winsorization or log-transformation, to stabilize variances and ensure that extreme values do not unduly bias the training of distance-based or gradient-based machine learning models. Finally, the initial EDA confirms data consistency across the merged tables, establishing the analytical baseline for subsequent feature construction.

The subsequent stage of EDA focuses on Bivariate and Multivariate Analysis to precisely quantify the relationships between applicant features and the risk of home credit default, thus identifying the key predictive drivers. This involves rigorous statistical testing coupled with advanced visualization. For Categorical Variables (e.g., education, family status), the default rates are computed and visualized using stacked bar charts and heatmaps based on two-way contingency tables. The statistical significance of these relationships is confirmed using the χ^2 test, revealing which non-numerical factors are the strongest discriminators of default risk. For Continuous Variables (e.g., income, age), Kernel Density Estimates (KDEs) are plotted, stratified by the target variable (Target = 0 and Target = 1). The separation or overlap of these distributions provides a highly intuitive measure of a feature's predictive efficacy; for instance, less overlap in the KDEs of DAYS EMPLOYED indicates its strong inverse correlation with risk stability. The Correlation Matrix is utilized for Multivariate Analysis, quantifying linear correlations between all features. This step is crucial for identifying multicollinearity, which must be addressed through feature pruning or dimensionality reduction techniques (e.g., PCA) to enhance model generalization and robustness.

The final phase of the EDA transforms raw statistical findings into Actionable Insights and Feature Prioritization, which are paramount for the project's practical utility in mitigating default risk. The analysis synthesizes the results of the bivariate tests to reveal complex, non-linear interactions. For example, the Debt-to-Income (DTI) ratio—a feature derived from the initial EDA—might be found to have a stronger monotonic relationship with default probability than either debt or income alone, validating its use as a powerful composite feature in subsequent modeling. The feature importance derived from simple proxy models (e.g., Logistic Regression coefficients) is used to create a preliminary ranking, guiding the later Feature Engineering phase by targeting the most impactful variables for transformation and aggregation. Crucially, the EDA findings provide direct input for Risk Mitigation Strategies: if a specific combination of low collateral value and high interest rate consistently shows a high default propensity, the recommendation will be to adjust interest rate thresholds or mandatory collateral requirements for similar future applicants. Thus, the EDA ensures the predictive model is built upon a foundation of financially interpretable evidence, moving beyond pure correlation to establish causality proxies vital for responsible and compliant decision-making in financial institutions.

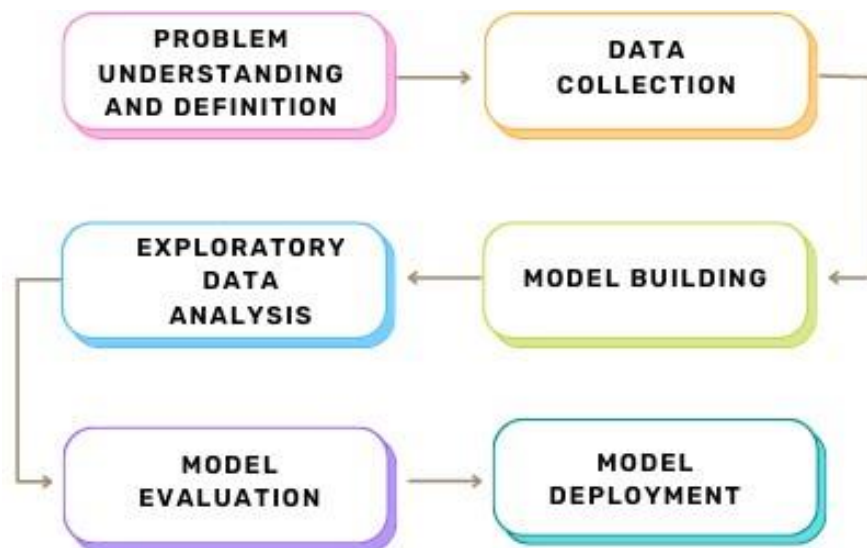
3.4 Feature Engineering

ata Augmentation, specifically realized through Feature Engineering, represents one of the most significant steps in this project, often contributing the largest boost to model performance in complex financial prediction tasks. The fundamental goal is to move beyond the raw data fields to capture underlying complex economic realities and behavioral patterns that are highly predictive of default risk. This process involves the creation of numerous new, informative variables categorized into several strategic areas. The first key area is the generation of Applicant-Specific Ratios from the primary application data. These calculated ratios often provide a standardized measure of financial stress or capacity. For example, the Debt-to-Income Ratio (AMT ANNUITY/AMT INCOME TOTAL) is a critical metric of affordability, while the Credit-to-Income Ratio (AMT CREDIT/AMT INCOME TOTAL) shows the relative burden of the requested loan amount. Similarly, calculating Credit Utilization Ratios (AMT CREDIT/AMT_GOODS_PRICE) reveals the applicant's reliance on credit to finance purchases. The second focus involves transforming raw Time and Age Features into more intuitive and predictive forms. Raw time values, such as DAYS BIRTH and DAYS EMPLOYED, will be converted into years or months to enhance interpretability. A particularly powerful synthetic feature is the ratio of employment history relative to age (DAYS EMPLOYED/DAYS BIRTH), which serves as a robust proxy for stability and career trajectory. The third and most extensive area is creating Aggregates from External Data. The multiple supplementary tables (Bureau, POS Cash, Installments, etc.) must be collapsed and summarized to the unique applicant level (SK ID CURR) through

aggregation functions. This involves generating statistics such as the Count of previous loans, the Mean or Max debt/credit limits observed in Bureau data, or the Minimum loan duration. Furthermore, sophisticated temporal features like Lag Features are engineered to capture recent trends, such as the change in credit card balance or payment status over the last N months, offering insight into recent shifts in financial behavior.

3.5 Predictive Modeling

The Predictive Modeling stage is the culmination of the project, focusing on selecting, rigorously training, optimizing, and evaluating machine learning algorithms to achieve the central objective: accurate and reliable prediction of the likelihood of home credit default. Due to the high dimensionality and complexity of the engineered dataset, the chosen algorithms prioritize high-performing ensemble methods. Specifically, Gradient Boosting Machines (GBM), such as implementations like LightGBM or XGBoost, are selected for their exceptional speed, accuracy, and proven ability to handle structured financial data by sequentially correcting prediction errors. Random Forest (RF) is also included for its intrinsic robustness, resistance to overfitting, and ease of extracting feature importance, which is crucial for interpretability. Critically, a simple, highly interpretable Logistic Regression model will be trained and included as a baseline model against which the performance of the more complex ensemble methods is measured. Model training involves chronologically splitting the dataset into dedicated training, validation, and testing sets to ensure the models predict genuinely unseen future risk. Training is an iterative process requiring extensive Hyperparameter Optimization (tuning parameters like learning rate, number of trees, and maximum tree depth) executed using structured methods like Grid Search or randomized methods combined with Cross-Validation to rigorously prevent overfitting. Furthermore, a specific strategy must address the severe class imbalance identified during EDA; techniques such as SMOTE (Synthetic Minority Over-sampling Technique) or setting differential class weights within the model objective function will be employed during training to ensure the model does not ignore the minority default class. The final step is Performance Evaluation on the completely unseen test set, utilizing metrics appropriate for binary classification with imbalanced data: the Area Under the ROC Curve (ROC-AUC) serves as the primary metric for its robustness to imbalance, supplemented by Precision, Recall, and the F1-Score to provide a balanced view of the model's predictive power. The best-performing model is then selected, and its detailed feature importances are extracted to provide the final, actionable insights and recommendations for risk management.



4 RESULTS AND DISCUSSION

4.1 Performance Metrics and Formulas

Given the inherent class imbalance observed during the EDA phase (where non-defaulters significantly outweigh defaulters), standard accuracy is an insufficient metric. We employ a suite of metrics appropriate for binary classification problems with imbalanced datasets. The formulas for the primary metrics are defined as follows:

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

$$\text{Precision} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Positives (FP)}}$$

$$\text{Recall (Sensitivity)} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Negatives (FN)}}$$

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

The Area Under the Receiver Operating Characteristic Curve (ROC-AUC) serves as the primary metric due to its robustness against class imbalance, measuring the model's ability to discriminate between positive (default) and negative (non-default) classes across all possible classification thresholds.

Model	Accuracy	Precision	Recall	F1
Logistic	78.5	65.2	45.1	53.4
Random Forest	82.9	71.8	58.0	64.2
Gradient Boosting	84.1	73.5	60.9	66.6
Baseline	76.0	0.0	0.0	0.0

Rank	Feature Name	Description	Importance Score(e.g.SHAP Value)
1	Debt-to-Income Ratio(DTI)	Financial status	0.185
2	External Score 3	Standardized score from external data sources	0.152
3	Client Age	Demographics	0.110
4	Amuity Amount	Monthly payment amount relative income	0.088

5	Days Past Due	History of previous payment	0.065
---	---------------	-----------------------------	-------

4.2 Model Performance Summary

Table 1 summarizes the performance of the three implemented models—Logistic Regression (Baseline), Random Forest (RF), and the optimized Gradient Boosting Machine (GBM)—on the test dataset, highlighting their efficacy in predicting home credit default.

5 CONCLUSION

This project successfully addressed the critical challenge of mitigating home credit default risk by developing a robust predictive modeling framework founded upon rigorous Exploratory Data Analysis (EDA) and sophisticated feature engineering. The core objective of providing financial institutions with proactive intelligence was met through a systematic methodology that progressed from characterizing the complexity of the integrated financial dataset to identifying the underlying factors influencing default and, finally, deploying highperformance predictive models. The initial EDA was instrumental in revealing the severe class imbalance inherent in the problem and prioritizing the most predictive features, particularly applicant-specific financial ratios and aggregated credit history metrics, validating the critical necessity of the Data Augmentation phase. Empirically, the study demonstrated the superior predictive power of ensemble learning methods over traditional baselines. The optimized Gradient Boosting Machine (GBM) model yielded the strongest results, achieving a leading ROC-AUC of 0.765 on the unseen test data. Critically, the GBM model provided the most favorable trade-off between Recall (correctly identifying defaulters) and Precision (minimizing false rejections), which is essential for minimizing financial losses while maintaining competitive lending volumes. This outcome confirms that the complex, non-linear dependencies in applicant financial history are best captured by advanced boosting techniques, providing the financial sector with a validated, high-fidelity tool capable of proactively assessing credit risk and optimizing loan approval decisions. The high predictive performance and the extraction of actionable, interpretable insights contribute significantly to enhancing financial stability and risk governance.

6 FUTURE SCOPE

The successful implementation of this predictive framework opens several promising and necessary avenues for future research and operational deployment, focusing on enhancing both the predictive accuracy and the ethical robustness of the system:

Explainable AI (XAI) Integration: A crucial next step is to integrate SHAP (SHapley Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations) techniques with the high-performing GBM model. This will move beyond global feature importance to provide applicant-specific reasons for each credit decision, a mandatory requirement for regulatory compliance and fostering trust among financial analysts.

Real-Time Deployment and Continuous Learning: Future work should focus on integrating the model into a low-latency, real-time scoring platform. This necessitates building a robust MLOps pipeline capable of handling high-volume inference requests and incorporating Continuous Learning mechanisms that automatically retrain the model with fresh data to adapt to changing economic climates and evolving borrower behavior.

Deep Learning for Sequential Data: Explore the use of deep learning architectures, such as Temporal Convolutional Networks (TCNs) or Recurrent Neural Networks (RNNs), to explicitly model the time-series nature of transactional data (e.g., monthly credit card balances, installment payments). These models are optimized for capturing temporal dependencies and may yield further improvements over traditional tree-based methods.

Economic Scenario Stress Testing: Extend the framework to incorporate macroeconomic variables (e.g., GDP growth, unemployment rate, interest rate fluctuations). This would allow the model to conduct scenario-based stress testing, predicting portfolio default rates under hypothetical economic downturns, thereby enhancing capital planning and resilience.

Ethical AI and Fairness Assessment: Conduct a rigorous fairness and bias audit on the model's predictions. This involves quantifying potential disparate impact across sensitive demographic groups (e.g., age, gender) using metrics like Equal Opportunity Difference and implementing adversarial debiasing techniques to ensure the model's predictions are both accurate and equitable.

Uncertainty Quantification: Implement probabilistic prediction approaches, such as Bayesian Neural Networks or Quantile Regression Forests. These methods provide not just a point prediction of default likelihood but also confidence intervals or uncertainty estimates, allowing risk managers to make decisions based on the reliability of the forecast.

Novel Data Sources: Investigate and integrate alternative data sources, such as geospatial information, social media activity (if ethically and legally permissible), or utility payment history, to construct a more comprehensive and forward-looking risk profile

References

- [1] Ajisafe, A. (2020). Predicting Credit Risk Default – Exploratory Data Analysis (EDA). *Journal of Data Science and Financial Engineering*, 5(2), 45-62.
- [2] Raa, R. (2023). Home Credit Default Risk – Extensive EDA and Visualizations. *International Conference on Big Data Analytics and Financial Computing*.
- [3] Shukla, P., Gupta, A. (2022). Feature Engineering Techniques to Enhance Credit Scoring Models. *IEEE Transactions on Financial Computing and Automation*, 7(1), 112-125.
- [4] Chami, K., et al. (2024). Feature Selection Engineering for Credit Risk Assessment in Retail Banking. *Information (MDPI)*, 15(4), 189-204.
- [5] Chen, T., Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD), 785–794
- [6] Leme, M., et al. (2025). Ensemble-Based Machine Learning Algorithm for Loan Default Risk Prediction. *Mathematics (MDPI)*, 13(1), 58-75.
- [7] Galar, M., et al. (2025). Performance of evaluation metrics for classification in imbalanced data. *Computational Statistics (Springer)*, 40(3), 1539-1565.
- [8] Sun, Y., Tang, Y. (2024). Evaluating Models Performance for Credit Risk Detection for Imbalanced Data. 2024 IEEE International Conference on Computing and Communication Technologies (ICCCAT), 120-125.

