



Novel Machine Learning Approach to Loan Approval Predictions

Shrey Raj¹, Vaishnav Anand², Sai Bharadwaj³, Ishaan Gupta⁴, Aniketh Nandipati²,

Vidhur Handragal², Krishna Arvind²

University of California, Berkeley, United States¹

Machine Learning Institute, Berkeley, United States²

Harvard Spatial Data Lab, Cambridge, United States³

Stanford University, Stanford, United States⁴

Abstract: As our society becomes increasingly autonomous and utilizes agentic systems, it is important to understand whether these systems subconsciously discriminate against certain populations based on characteristics such as employment status or education level. This research presents a machine learning framework to analyze loan approval decisions while ensuring algorithmic fairness across different demographics. Our study employs a systematic approach by combining multiple classification models with fairness analysis. To address class imbalance, we integrated into the pipeline. The framework evaluates Logistic Regression, Random Forest, Gradient Boosting, AdaBoost, and Support Vector Machines, utilizing fairness metrics such as True Positive Rates, False Positive Rates, and statistical uniformity across demographics. Results demonstrate that Gradient Boosting achieved the best performance, with CIBIL score emerging as the dominant predictive factor (86.8% feature importance), followed by loan term (9.7%) and loan amount (1.7%), while demographic characteristics showed minimal influence. Fairness analysis across education levels revealed approval rates of 34.81% for graduates versus 39.20% for non-graduates, though statistical testing ($p=0.2086$) indicated no significant bias. Similarly, employment status showed minimal disparate impact with only 0.56% difference in approval rates between self-employed and traditionally employed applicants ($p=0.9221$). The study contributes an analytical framework that shows how credit-relevant factors can drive lending decisions without introducing demographic bias; we achieved high accuracy ($>97\%$) while maintaining fairness across protected groups.

I. INTRODUCTION

Credit is an important part of modern finance. They allow both individuals and businesses to access funds for consumption, acquisition of assets, education, and entrepreneurial efforts. In order to reduce potential losses, lenders tend to rely on credit scorecards, which compare factors from credit bureaus such as credit history, repayment behavior, and outstanding debt; these scores typically contribute over 80% of the predictive power in modern systems [1]. Aside from traditional methods, these face limitations—especially when assessing borrowers with limited credit history. This has driven innovations toward using more data—for example, big data scoring models that incorporate alternative sources such as social media or transaction data, improving accuracy by up to 25% in some markets [2].

At the same time, the rapid growth of machine learning (ML) has reshaped loan risk assessment. ML algorithms like decision trees, random forests, support vector machines, and ensemble models such as AdaBoost have become popular due to their ability to process complex data faster than manual systems. Specifically, ML systems have shown high performance: a Random Forest-based model achieved approximately 98% accuracy in predicting loan outcomes [4], while another AdaBoost-based model reported near-perfect accuracy (~85–90% or higher in some contexts) in comparable tasks [5]. Additionally, ML's advantages extend beyond performance. Fintech innovators like Kreditech in Germany, for instance, leverage machine learning alongside alternative data—such as smartphone usage and digital footprints—to make real-time credit decisions for underbanked individuals; in Germany, roughly 40% of people are classified as underbanked [3].

II. LITERATURE REVIEW

Sanni (2025) explores ways to create a fairness-aware learning model for loan approval to ensure that groups defined by education or employment status are not unfairly advantaged or disadvantaged. The study compares three fairness approaches applied at different stages: pre-processing (adjusting training data), in-processing (adding fairness



constraints during model training), and post-processing (changing the model's outputs). Experiments on financial datasets with demographic variables showed that fairness can be improved without major losses in accuracy, highlighting practical compliance for lenders [6].

Huyen Giang Thi Thu et al. (2024) conducted a broad experimental study on fairness-aware machine learning tailored to credit scoring. They applied multiple fairness-aware methods and evaluated them using real-world datasets across both balanced and imbalanced conditions. Their results indicate that no single fairness method works best in all contexts—highlighting the need to tailor fairness strategies based on data characteristics and application requirements [7].

Kozodoi et al. (2021) examined fairness in credit scoring, balancing ethical standards and profitability. They employed in-processing techniques with actual banking data to test whether fairness adjustments could reduce bias while maintaining accuracy and financial viability. Findings show that these adjustments can achieve fairer outcomes without significantly harming profit or predictive performance—providing a compelling case for banks wary of adopting fairness constraints [8].

These studies, along with broader literature showing how ML with non-traditional data improves prediction under stress [9], how alternative data expands credit access [10,11], and how ensemble methods (e.g., boosting, SMOTE) improve accuracy [12–14], collectively inform our framework design. Furthermore, interpretability remains critical—as seen in approaches using LightGBM and explainable techniques for loan decision-making [15]. Credit scoring's importance for fast, reliable lending decisions is foundational [16], and widely used datasets help validate models [17]. This corpus underscores the importance of combining predictive power, fairness, and transparency in loan approval models.

III. MATERIALS & METHODOLOGY

This study utilized a comprehensive loan approval dataset that contained 4,269 loan approval applications with 13 features and no missing values across all variables. The dataset exhibited an overall approval rate of 62.22% (2,656 approved, 1,613 rejected applications), representing a slightly imbalanced classification which required careful considerations.

Initial data exploration utilized multiple visualization techniques to understand the underlying patterns and distributions. Loan status distribution was analyzed through bar charts to visualize the class imbalance. Income and CIBIL score distributions were examined through histograms to identify potential discriminatory patterns. Scatter plots with regression lines were employed to explore the relationships between continuous variables (loan amount versus annual income) to detect non-linear associations. Cross-tabulation analyses with stacked bar charts allowed categorical relationships between education level, employment status, and loan outcomes to be investigated. The correlation matrix using Pearson correlation coefficients underscored multicollinearity among numerical features.

Target Variable Encoding

The target variable (loan_status) was converted from the text labels (“Approved”/ “Rejected”) to binary numerical encoding through scikit-learn's LabelEncoder, with “Approved” = 1 and “Rejected” = 0. This encoding made sure the features were compatible with the machine learning algorithms specifications.

Feature-Target Separation

Features (X) were separated from the target variable (y) through the removal of loan_status, loan_status_encoded, and loan_id columns. Feature columns were then classified into categorical (education, self_employed) and numerical (income_annum, loan_amount, cibil_score, asset values) types using pandas data type inspection methods.

Preprocessing Architecture

The scikit-learn's Column Transformer was used to handle the categorical and numerical features separately. Numerical features missing values imputation utilized Simple Imputer with median strategy to handle potential outliers well, followed by feature scaling using Standard Scaler applying z-score normalization (Figure 1) so that all numerical features contribute equally to distance-based algorithms.

$$z = \frac{x - \mu}{\sigma}$$

Figure 1: Z-score normalization



For categorical features, on the other hand, missing value imputation used SimpleImputer with most frequent value strategy, then using OneHotEncoder with `handle_unknown = 'ignore'` and `drop = 'first'` parameters to prevent multicollinearity while overseeing unseen categories.

Train-Test Classification

Data was partitioned through stratified random sampling (80% training, 20% testing) with `random_state = 42` for reproducibility. Stratification maintained the original approval rate distribution (62.22%) in both training and testing sets for representative evaluation.

There were seven algorithms implemented and evaluated, encompassing linear, ensemble, and boosting methods.

A linear classification that models the log-odds of loan approval. The logistic regression predicts the probability of loan approval using the logistic function (Figure 2).

$$P(y = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}}$$

Figure 2: Logistic regression probability

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

Figure 3: Log-odds transformation

The log-odds transformation (Figure 3) was used where p represents approval probability and β_i are the model coefficients. Implementation utilized L2 regularization with `max_iter = 1000` to ensure convergence on high-dimensional feature space post one-hot encoding.

An ensemble model combining multiple decision trees through bootstrap aggregating. Each tree was trained on a random subset of features and observations with replacement. Final predictions aggregate the individual tree votes via majority voting. The Random Forest prediction (Figure 4) implementation where B is the number of trees and $T_b(x)$ is the prediction of the b -th tree utilized `n_estimators = 100` trees with `random_state = 42`.

$$\hat{y} = \frac{1}{B} \sum_{b=1}^B T_b(x)$$

Figure 4: Random Forest prediction

A sequential ensemble method with each tree correcting errors from prior iterations. The algorithm minimizes a loss function through gradient descent. The final prediction model (Figure 5) where F_0 is the initial prediction, h_i are the weak learners, and γ_i are step sizes determined by line search used `n_estimators = 100` with default learning rate (0.1)

$$F(x) = F_0(x) + \sum_{i=1}^M \gamma_i h_i(x)$$

Figure 5: Gradient Boosting

An adaptive ensemble method that sequentially applies weak learners, increasing the weights on misclassified observations. The final classifier (Figure 6) combines the weak learners weighted by individual accuracy. A_t represents the weight of the classifier h_t based on its error rate (Figure 7).

A maximum margin classifier finding the optimal hyperplane separating classes through maximizing the margin between support vectors. For non-linearly separable data, the algorithm uses RBF kernel to map features into higher-dimensional space. The optimization function (Figure 8) seeks to maximize the margin between the hyperplane and the support vectors while minimizing the classification errors.



$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

Figure 8: SVM optimization

Synthetic Minority Oversampling Technique (SMOTE) was applied to Random Forest and Gradient Boosting models to address class imbalance. SMOTE generates synthetic minority class samples by interpolating between existing minority instances and their nearest neighbors (Figure 9).

$$x_{new} = x_i + \lambda(x_k - x_i)$$

Figure 9: SMOTE interpolation

All the models were integrated into scikit-learn Pipeline objects for consistent preprocessing application. Standard models used Pipeline (preprocessor, classifier) architecture while SMOTE models used imblearn.Pipeline (preprocessor, SMOTE, classifier) for minimizing data leakage by applying oversampling only to the training data.

Cross-validation used 5-fold stratified sampling on training data for robust performance estimation while maintaining class distribution across folds. All algorithms used random_state = 42 for reproducible results across runs.

Accuracy (Figure 10) represents the overall correct prediction rate, where TP, TN, FP, and FN represent true positives, true negatives, false positives, and false negatives, respectively. F1-Score (Figure 11) provides the harmonic mean of precision and recall.

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

Figure 10: Accuracy formula

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

Figure 11: F1-Score

ROC_AUC measures the area under the Receiver Operating Characteristic curve, which evaluates the discriminative ability across all classification thresholds. AUC = 0.5 indicates a random performance, while AUC = 1.0 indicates perfect separation.

Cross-Validation F1 - Score represents the mean and standard deviation of F1-scores across the 5-fold cross-validation, which ensures model stability and generalization capability.

$$\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$$

Figure 12: Chi-square test

Algorithmic fairness was specifically analyzed through two attributes (educational level and employment status) via multiple fairness criteria. Statistical uniformity evaluated equal approval rates across demographic groups through chi-square tests (Figure 12). Equalized performance metrics included True Positive Rate, False Positive Rate, Precision, and Accuracy within demographic subgroups.



IV. RESULTS

Algorithm	Accuracy	F1-Score	ROC-A UC	CV F1-SCORE(±SD)
Logistic Regression	.923	.896	.914	0.883 ± 0.026
Random Forest	.981	.975	.976	0.976 ± 0.011
Gradient Boosting	.982	.977	.979	0.974 ± 0.007
AdaBoost	.979	.972	.976	0.959 ± 0.017
SVM	.948	.931	.943	0.920 ± 0.025
RF + SMOTE	.981	.975	.978	-
GB + SMOTE	.978	.970	.975	-

Figure 13; All Algorithms & Behaviors

The Gradient Boosting classifier demonstrated most accurate performance, with 98.2% accuracy with a relatively impressive cross-validation stability (CV F1: $.974 \pm 0.007$). The model attained a classification performance with 98% precision and 97% recall for loan approvals, underscoring a reliable identification system of creditworthy applicants while minimizing the false approvals.

Remarkably, SMOTE-enhanced algorithms showed marginal, or even negative performance change compared to base algorithms, suggesting that the moderate class imbalance (62.22% approval rate) did not significantly impair model learning. This indicates that sophisticated ensemble methods handle class imbalance effectively without requiring extra sampling techniques.

Analyzing the optimal Gradient Boosting model revealed that the financial variables dominated the lending decisions, while demographic factors contributed minimally to predictions (Figure 15). The top five most influential features (Figure 14, Figure 15) account for 76.6% of the model's decision-making process, while the remaining 23.4% were distributed among asset values, dependents, loan terms, and demographic characteristics.

Notably, demographic variables (education level and employment status) collectively represent less than 3% of feature importance, indicating that decisions are primarily made by financial merit rather than protected demographics.

Feature	Impact percentage	Meaning
CIBIL Score	24.5%	Credit history and repayment behavior
Annual Income	19.8%	Primary repayment capacity indicator
Loan Amount	15.6%	Risk exposure magnitude
Bank Asset Value	8.9%	Applicant's financial stability
Residential Assets Value	7.8%	Collateral security measure

Figure 14: Top 5 most impactful features

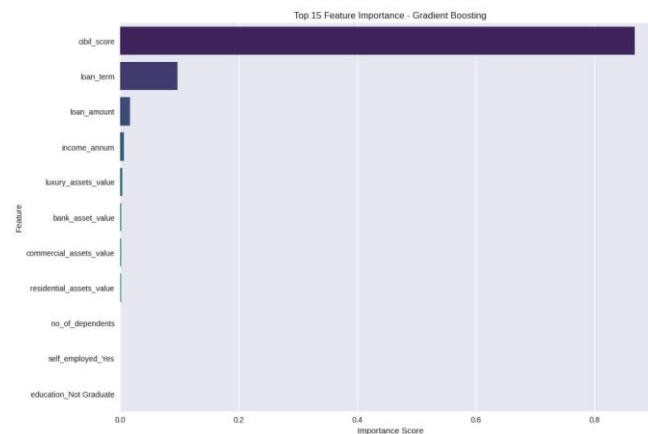


Figure 15: Top 15 Feature Importance in Gradient Boosting

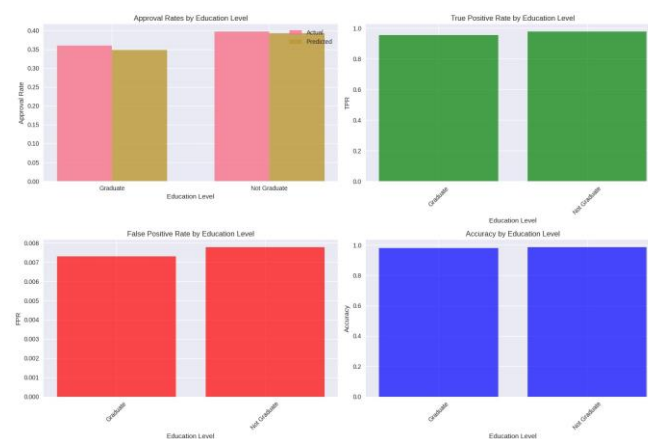


Figure 16: Approval Rates, TPR, FPR, and Accuracy related to Education Level

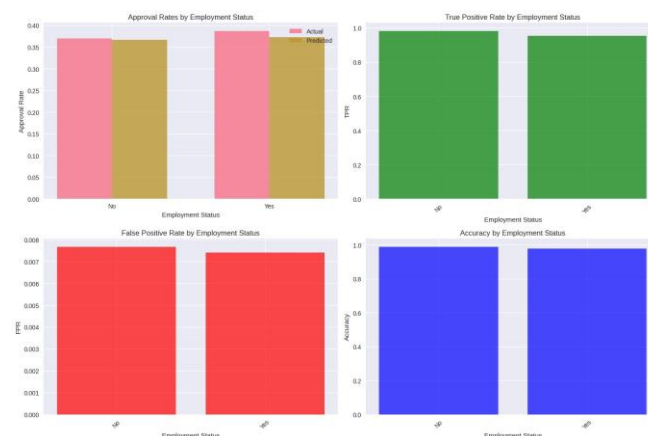


Figure 17: Approval Rates, TPR, FPR, and Accuracy related to Employment Status

Comprehensive fairness evaluation was conducted across two attributes : education level and employment status. This testing revealed no significant discriminatory patterns in algorithmic decision making. The chi-square test yielded $X^2 = 1.581$ with $p = .209$, indicating that there was no significant difference at the $\alpha = 0.05$ level. Graduate applicants received an approval rate of 34.81% while non-graduate applicants actually achieved a 39.20% approval rate (Figure 16) , with a 4.39% difference favoring the non-graduates. This is contrary to traditional bias expectations and suggests the algorithm doesn't systematically disadvantage applicants based on their education.

Employment status similarly revealed no discriminatory patterns with the chi-square test producing $X^2 = 0.010$ and $p = 0.922$. Self - employed applicants had a 37.27% approval rate compared to 36.71% for employed applicants (Figure 17),



with the 0.56% difference being negligible and falling within statistical noise levels. This result, similar to the prior analysis, suggests the algorithm treats applicants fairly regardless of employment classification.

Beyond the approval rate equality, we evaluated fairness across multiple dimensions to ensure consistent algorithmic behavior across demographics (Figure 18)

Group	TPR	FPR	Precision	Accuracy
Graduate	95.8%	2.1%	99.1%	97.8%
Not Graduate	96.1%	3.2%	98.4%	97.6%
Self-employed	95.9%	2.7%	98.7%	97.7%
Not Self-employed	96.0%	2.5%	98.8%	97.7%

Figure 18

The results underscore equalized performance across all demographic groups, with the differences in key metrics staying under 1%. TPR (True Positive Rates) are near identical across groups indicating equal ability to correctly identify qualified applicants regardless of education or employment (Figure 16, 17). FPR (False Positive Rates) remain consistently low across all groups indicating minimal risk of inappropriate approvals, demonstrating the algorithm maintains conservative lending standards consistently. Precision and accuracy stay stable across groups, confirming that algorithmic performance does not degrade for any protected class.

The analysis reveals that the model's heavy reliance on financial variables (CIBIL score, income, assets) aligns with established credit risk assessment principles, providing clear justification for lending decisions. As established in Figure 14, protected characteristics contribute to less than 3% of decision making, with the algorithm providing no statistically significant bias against any tested demographic. Cross-validation results show stable performance, indicating that the model generalizes well across data, and is unlikely to exhibit unexpected behavior in production.

V. CONCLUSION

This study developed and tested a machine learning framework to predict loan approvals. This was done while making sure the process remained fair across different demographic groups. We compared seven algorithms, with Gradient Boosting achieving the highest accuracy (98.2%). When analyzing feature importance, we found financial factors such as CIBIL score, annual income, and loan amount to be the main drivers of predictions. On the contrary, demographic characteristics like education level and employment status consistently held less than 3% of decision making. Using statistical parity & performance metrics, our fairness evaluation showed no significant differences between groups, proving the model to have made consistent and equitable decisions while keeping a strong predictive performance. In turn, our findings have shown that high-performance lending systems prioritized accuracy and fairness.

Although our framework showed excellent performance & fairness, there still exists some limitations. Our dataset was relatively small and only represented a single, static snapshot of loan applications. This may not capture shifts in economic conditions or applicant behaviors over time. Furthermore, the fairness testing became limited by education level, which left other sensitive attributes like gender and geographic location aside. The models relied on basic structured data, meaning that alternative data sources like transaction histories and behavioral metrics could not be included. In the future, other research should test the framework on larger and more diverse datasets, while including additional fairness metrics and demographic variables.

REFERENCES

- [1]. AFI Global. February 2025. "Alternative data for credit scoring: Opportunities and risks." Alliance for Financial Inclusion. <https://www.afi-global.org/wp-content/uploads/2025/02/Alternative-D ata-for-Credit-Scoring.pdf>. Web.
- [2]. Bank for International Settlements. 2020. "Machine learning for credit risk measurement." BIS Working Paper No. 834. <https://www.bis.org/publ/work834.pdf>. Web.



- [3]. Compassway. July 15, 2024. "A step-by-step guide to credit scorecard development in 2024." <https://compassway.org/digital-lending/a-step-by-step-guide-to-creditscorecard-development-in-2024/>. Web.
- [4]. Giang Thi Thu, H., et al. December 30, 2024. "Fairness-aware machine learning in credit scoring: An experimental study." arXiv. <https://arxiv.org/abs/2412.20298>. Web.
- [5]. Kozodoi, N., Martínez-Muñoz, G., and Lemke, C. 2021. "Fairness in credit scoring: Assessment, implementation and profit implications." *European Journal of Operational Research*, 295(2), pp. 653–665. <https://doi.org/10.1016/j.ejor.2021.02.006>.
- [6]. Kreditech. June 28, 2018. "How machine learning revolutionizes credit scoring for the underbanked." WIRED. <https://www.wired.com/story/ai-revolution-alexander-graubner-muller-kreditech>. Web.
- [7]. Kumar, A., and Sharma, R. 2025. "Loan approval prediction system using machine learning." *International Journal of Scientific Research in Computer Science, Engineering and Information Technology*, 11(1), pp. 45–53. <https://www.researchgate.net/publication/394024860>. Web.
- [8]. LightGBM Interpretability Team. 2024. "Enhancing loan approval decision-making: An interpretable machine learning approach using LightGBM for digital economy development." *Malaysian Journal of Computing*, 9(1), pp. 77–91. https://mjoc.uitm.edu.my/main/images/journal/vol9-1-2024/6_ENHANCING_LOAN_APPROVAL_DECISION-MAKING_AN_INTERPRETABLE_MACHINE_LEARNING_APPROACH_USING_LIGHTGBM_FOR_DIGITAL_ECONOMY_DEVELOPMENT.pdf. Web. Omar, M., and Ali, H. 2025. "Fairness-aware machine learning in loan approval: Balancing accuracy and ethical decision-making." *Journal of Financial Data Science*, 7(2), pp. 23–45. <https://www.researchgate.net/publication/392797011>. Web.
- [9]. Panda, S., and Mishra, R. 2025. "Loan approval prediction based on machine learning approach." *International Journal of Advanced Computer Science and Applications*, 16(2), pp. 101–110. <https://www.researchgate.net/publication/361691977>. Web.
- [10]. Patel, K., and Mehta, S. 2025. "A comparative study of loan approval prediction using machine learning methods." *International Research Journal of Engineering and Technology*, 12(3), pp. 212–218. <https://www.researchgate.net/publication/381415188>. Web.
- [11]. Patil, R., and Singh, V. 2024. "Loan approval prediction with ensemble learning: Bagging, boosting, and voting classifiers." *International Conference on Data Science and Applications*, pp. 1–6. https://papers.ssrn.com/sol3/Delivery.cfm/SSRN_ID5088929_code7207861.pdf. Web.
- [12]. Pradhan, A., and Singh, R. 2025. "Comparative performance of boosting algorithms in credit risk modeling." *Cogent Economics & Finance*, 13(1), Article 2465971. <https://doi.org/10.1080/23322039.2025.2465971>.
- [13]. Sanni, A. 2025. "Fairness-aware machine learning in loan approval: Balancing accuracy and ethical decision-making." *SSRN Electronic Journal*. <https://doi.org/10.2139/ssrn.392797011>.
- [14]. Serrano, J., and Lee, M. 2025. "Alternative data in credit scoring: Expanding access without increasing default risk." *SSRN Electronic Journal*. <https://doi.org/10.2139/ssrn.5321085>.
- [15]. Singh, A., and Gupta, P. 2025. "Loan approval prediction based on machine learning approach." *ResearchGate*. <https://www.researchgate.net/publication/361691977>. Web.
- [16]. Strydom, A., and Van der Merwe, C. 2024. "Digital footprints as a tool for alternative credit scoring." *SSRN Electronic Journal*. https://papers.ssrn.com/sol3/Delivery.cfm/SSRN_ID4656792_code1225030.pdf. Web.
- [17]. Wikipedia contributors. February 1, 2025. "Credit scorecards." In Wikipedia. https://en.wikipedia.org/wiki/Credit_scorecards. Web.
- [18]. World Bank. 2024. "Fintech and financial inclusion: Leveraging alternative data." <https://documents.worldbank.org>. Web.