# The Credit Risk Playbook: AI-Powered Loan Default Prediction

Steven Barnes
February 4, 2025

# Table of Contents

# Abstract

Loan default prediction is a critical challenge for financial institutions, as missed defaults lead to financial losses while false alarms increase operational costs. This study applies machine learning techniques to loan classification, comparing Logistic Regression, K-Nearest Neighbors, Random Forest, and Gradient Boosting across multiple thresholds to evaluate trade-offs between risk mitigation and efficiency. A structured evaluation framework incorporating FNR, FPR, Recall, Precision, Accuracy, F1-Score, and AUC-ROC ensured a comprehensive assessment of model performance. Results show that Gradient Boosting minimizes financial losses by achieving the highest Recall and lowest False Negative Rate (FNR), while Random Forest optimizes operational efficiency through high Precision and Accuracy. Financial impact analysis suggests that Gradient Boosting could reduce default-related losses by 10%, while Random Forest could cut manual reviews by 40%, saving labor hours. These findings offer actionable insights for lenders, guiding threshold selection, policy adjustments, and automation strategies. However, dataset limitations and high FNR at certain thresholds highlight the need for further refinements. Future work should explore time-series data, cost-sensitive learning, and ensemble modeling to enhance predictive accuracy and real-time decision-making.

# Problem Statement

Financial loan services play a vital role in driving economic growth and stability. These services empower innovation, entrepreneurship, and consumer spending by injecting capital into the economy through individuals and businesses. However, a critical challenge financial institutions face with loan services is mitigating the risk of payment defaults. Defaulted loans not only result in significant financial losses but also disrupt operational efficiency, undermining the stability of lending institutions.

This project aims to predict loan defaults using machine learning, equipping financial institutions with the ability to identify high-risk borrowers proactively. By leveraging predictive analysis, loan-providing organizations can allocate resources more effectively, implement interventions for at-risk clients, and reduce financial losses. From a career perspective, this project demonstrates the practical application of machine learning to address complex financial challenges. The ability to analyze and generate actionable insights is critical in data analytics and financial planning, aligning directly with my professional aspirations.

# Data Collection and Tools

The data used for this project was sourced from Kaggle[1], specifically using the "Loan Default Prediction Dataset." This dataset comprises over 255,000 unique observations and includes 18 variables. Among these, 16 serve as independent variables to predict loan defaults, one represents the response variable indicating whether a loan defaulted, and the final variable is a unique identifier for each observation. The dataset includes numerical variables, such as

Income, Credit Score, and Loan Amount, and categorical features, such as Marital Status and Education. Together, these variables offer a comprehensive view of borrower profiles and risk factors. The dataset was designed for machine learning competitions and reflects anonymized data typical of real-world financial loan services. Thus, it is an excellent resource for developing predictive analytics models with applications in financial decision-making. Initial inspection revealed the data to be well-structured, although preprocessing steps were required to encode categorical variables and scale numerical features.

Several technical tools were employed to prepare and analyze this data. Python was central to the workflow, with libraries like pandas for data manipulation and sci-kit learn for machine learning and preprocessing. Tableau was used to design the visualizations presented in the results section, offering a clear and engaging way to communicate insights derived from the analysis. These tools were chosen for their versatility, efficiency, and ability to handle large datasets, making them ideal for this analysis.

## Data Preparation

Proper preprocessing is paramount in ensuring model accuracy and generalizability. The imported dataset was complete, with no missing values, so no imputations were required. First, categorical features were encoded into numerical values using the LabelEncoder. This transformation was essential because machine learning algorithms typically require numerical inputs to process data effectively and make predictions on the response variable. Next, numerical variables were scaled using StandardScaler. Scaling is important because it standardizes the ranges of numerical features, ensuring that variables with larger ranges (e.g., Income) do not disproportionately influence the model compared to variables with smaller ranges (e.g., Loan Term). This step is particularly critical for distance-based algorithms, like K-Nearest Neighbors) and those sensitive to the magnitude of feature values, like Logistic Regression. Additionally, the data was divided into training and testing sets using an 80/20 split, stratified by the target variable (Default). Stratification ensures that the training set maintains the same proportion of default and non-default observations as the original dataset. This is critical for preventing bias and ensuring the model learns from a balanced representation of the target variable.

With the data prepared, the correlation of the independent features to the Default response variables and a correlation matrix were constructed to analyze the relationships between the independent variables and the response variable. Figure 1 illustrates these correlations to the Default variable. Notably, Age shows the strongest absolute, though negative, correlation with Default (-0.168), indicating that older borrowers are less likely to default. Conversely, Interest Rate exhibits the highest positive correlation (0.13), suggesting that higher interest rates are associated with an increased likelihood of default. Other variables, such as Income, Months Employed, and HasCoSigner, display weaker negative correlations than Age, implying that these features have modest protective effects against default. Meanwhile, features like Loan Amount and Employment Type show relatively weaker positive correlations than Interest Rate, indicating a slight increase in default likelihood as these values rise.

As for the correlation matrix, it was constructed to evaluate relationships among the independent variables and assess potential multicollinearity. As shown in Figure 2, the negligible correlation values (ranging from -0.006 to 0.005) confirm the absence of strong relationships between features. These findings validate the dataset's structure and ensure no individual feature dominates the model due to multicollinearity. Together, these insights provide a robust foundation for understanding how individual features contribute to loan default prediction.

# Methodology

## Modeling Framework

Four machine learning models were employed in this analysis, selected to balance interpretability, predictive power, and methodological diversity. Logistic Regression served as the baseline model due to its simplicity and clarity in capturing linear relationships between features and the response variable. As an interpretable model, Logistic Regression provided a foundation for understanding the dataset's structure and facilitated comparisons with more sophisticated approaches. K-Nearest Neighbors (KNN) was included as a complementary benchmark model to Logistic Regression. While Logistic Regression provided a parametric baseline by modeling linear relationships, KNN introduced a non-parametric approach by leveraging distance-based predictions. Thus, KNN offered a contrasting approach to highlight the challenges of imbalance datasets and the limitations of naïve methods in identifying patterns, reinforcing its role as a benchmarking tool rather than a primary predictive model.

To address the complexities of non-linear relationships, Random Forest was employed for its ability to model feature interactions and rank feature importance. Its ensemble nature and robustness against overfitting positioned it as a strong candidate for predictive modeling. Lastly, Gradient Boosting was implemented for its capacity to optimize predictive performance through sequential learning. Gradient Boosting demonstrated its ability to capture intricate patterns within the data by iteratively refining weak learners into a cohesive, high-performing model. Together, these models created a comprehensive evaluation framework, enabling robust comparisons across diverse methodologies.

A critical challenge in this analysis was the imbalance within the dataset, where instances of loan defaults were significantly fewer than non-defaults. This imbalance necessitated strategic adjustments to ensure fair and accurate model evaluation. First, as mentioned, stratified train-test splitting was used to preserve the original distribution of the response variable in both training and testing datasets, preventing skewed performance metrics. Second, model-specific approaches were employed to mitigate the effects of imbalance. Logistic Regression's sensitivity to imbalanced data was addressed through threshold tuning, optimizing recall to minimize false negatives. Random Forest and Gradient Boosting, inherently more adaptable to imbalance due to their ensemble structure, benefited from similar threshold adjustments to enhance their effectiveness. Conversely, KNN struggled with the imbalanced distribution, reinforcing its role as a benchmarking tool rather than a primary predictive model.

This structured methodology ensured that each model's unique strengths were leveraged while addressing the dataset's imbalanced nature. By employing tailored strategies for each approach, the analysis provided meaningful insights into both linear and non-linear relationships within the data and the challenges and opportunities posed by imbalanced datasets.

## Parameter Tuning and Thresholds

In this analysis, hyperparameter tuning and threshold adjustments were integral to optimizing the performance of the four models and tailoring their predictions to align with the project's objectives. Hyperparameter tuning focused on fine-tuning the model-specific parameters to enhance predictive accuracy while avoiding overfitting, while threshold tuning refined classification decisions by adjusting the probability cutoff for categorizing observations as loan defaults.

Starting with Logistic Regression, its regularization strength ($C$) was tuned to balance the model's ability to generalize and fit the data. Regularization penalizes large coefficients associated with features, helping to prevent overfitting. A lower C value increased regularization, favoring a simpler model that generalized well across unseen data. In contrast, a higher C value reduced regularization, allowing the model to capture more nuanced patterns at the risk of overfitting. This adjustment ensured the model focused on meaningful relationships between features and the response variable while minimizing noise.

For KNN, the primary parameter of interest was the $k$ value, which specifies the number of neighbors considered when making predictions. A small k value rendered the model sensitive to noise, as predictions relied on very few data points. Conversely, a large k value simplifies the predictions, potentially overlooking more granular patterns. Optimizing k involved finding a balance where predictions remained stable and accurate.

As an ensemble method, Random Forest required careful calibration of its number of trees (n_estimators) and the number of splits each decision tree can make (max_depth). Each tree within the Random Forest explores a subset of data and features, combining their predictions to enhance generalization. A greater number of trees typically improves performance but comes with diminishing returns and increased computational costs. Tuning these parameters ensured the model captured complex interactions in the data without excessive resource demands or overfitting.

Gradient Boosting, as another ensemble technique, necessitates tuning of its learning rate (learning_rate) and the number of boosting iterations (n_estimators). Each boosting iteration builds a weak learner that corrects the errors of the previous one, combining their predictions to create a strong predictive model. A lower learning rate allows for smaller, more precise updates, improving performance but requiring more iterations to converge. Conversely, a higher learning rate accelerates the training process but risks overshooting optimal solutions and potentially reducing accuracy. By balancing the learning rate and the number of boosting

iterations, Gradient Boosting achieved robust predictive performance without over-complicating the model.

In addition to hyperparameter tuning, all models underwent threshold adjustments to refine classification outputs. Thresholds ranging from 0.1 to 0.9 were tested and recorded, altering the probability cutoff for assigning an observation to the loan default category. Lower thresholds emphasized capturing as many defaults as possible, leading to higher recall and lower False Negative Rates, but at the cost of increased false positives. Conversely, higher thresholds reduced false positives at the expense of missing more defaults. This process provided a nuanced understanding of how each model balanced predictive accuracy with business priorities, such as minimizing defaults while maintaining operational efficiency.

## Evaluation Metrics

Various metrics were recorded to evaluate the performance of the machine learning models and assess the models' effectiveness in predicting loan defaults. The metrics were chosen for their relevance to the financial context of this project, where minimizing risk and maximizing operational efficiency are paramount. These metrics include False Negative Rate (FNR), False Positive Rate (FPR), Recall, F1-Score, Area Under the Receiver Operating Characteristics curve (AUC-ROC), Precision, and Accuracy.

With loan default classification, it is important to balance risk identification and operational efficiency; thus, the FNR, FPR, and Recall stand out as key metrics for this evaluation. FNR is the proportion of actual defaults that the model fails to identify. Arguably, this is the most important metric, as the inability to correctly identify a default risk could lead to large financial losses. FPR, on the other hand, is the proportion of borrowers incorrectly flagged as default risks. While not as financially significant as the FNR, this metric is still important as it serves operational efficiency, as false positives will waste operational capacity when investigating and clearing these instances. Recall is a ratio of the number of true positives over the total number of positive instances identified by the models. It acts as a measurement of how often a model correctly flags a true positive. Together, FNR, FPR, and Recall provide a comprehensive perspective on the trade-off between detecting defaults and managing operational capacity, ensuring that the model meets both the financial and operational priorities.

Another set of metrics that are important to model interpretation are Precision and Accuracy, as they offer complementary perspectives that contribute to a broader understanding of each model's performance. Precision is the proportion of true positives the model identifies over all the true positives within the dataset. Precision, like FPR, gives insights into false positives; however, precision ensures flagged defaults are accurate, whereas FPR minimizes false alarms among non-defaults. As for Accuracy, this metric measures the overall proportion of correct predictions, encompassing both default and non-default predictions. While accuracy offers a straightforward evaluation of a model's performance, its reliability diminishes in the context of imbalanced datasets like the one used in this project. A high accuracy may mask poor identification of the minority class (defaults in this case), making accuracy a good baseline metric for comparing model effectiveness and a complementary tool to the metrics listed thus

far. Precision and Accuracy offer a balanced view of model performance, combining a focus on correct default predictions with overall reliability.

Finally, metrics like F1-Score and AUC-ROC provide deeper insights into the balance and overall robustness of the models in predicting loan defaults. The F1-Score, calculated as the harmonic mean of Precision and Recall, provides a single metric that balances the trade-off between identifying true defaults and minimizing false positives. A high F1-Score signifies that the model effectively balances the identification of defaults and the avoidance of false alarms. AUC-ROC measures a model's ability to distinguish between default and non-default cases, evaluating the trade-off between the true positive rate (Recall) and the FPR and providing a view into a model's discriminative power. A higher AUC-ROC value indicates that a model effectively separates defaults from non-defaults. Conversely, a low AUC-ROC (closer to 0.5) suggests that a model performs no better than random guessing.

In summary, the evaluation metrics chosen for this project offer a comprehensive view of each model's ability to balance loan identification and operational efficiency in the context of loan default prediction. FNR, FPR, and Recall emphasize the importance of minimizing missed defaults and managing false alarms; Precision and Accuracy give insights into prediction reliability, and F1-Score and AUC-ROC capture the balance between identifying true defaults and effectively distinguishing them from non-defaults. These metrics provide a holistic framework for assessing model performance, financial risk mitigation, and operational effectiveness.

# Results

## Feature Importance

To provide a comprehensive analysis of the different methods being utilized in this project, the influence of each variable needs to be understood and compared between models to predict default classification. Thus, 'Feature Importance' measures the contribution of each individual variable (feature) to the predictive power of the machine learning model. Through feature importance, an evaluation can be conducted on which factors significantly affect the probability of default.

Logistic Regression uses coefficients derived from the model's equation to determine feature importance. These coefficients represent the strength and direction of the linear relationship between each variable and the dependent target variable, with a positive value indicating that as the independent variable increases, the likelihood of default will increase in tandem, and vice versa for negative coefficients.

KNN, in contrast, does not inherently compute feature importance due to its distance-based nature. Predictions in KNN rely on proximity to neighboring data points rather than a direct weighting of individual features. A permutation-based method was used to approximate the importance of features in this project. This involved measuring how the model's

performance changes when the values of a feature are randomly shuffled. Features that, when shuffled, result in a significant drop in model accuracy are considered more important. While indirect, this approach provides insight into to what extent each feature influences KNN's predictions.

In Random Forest and Gradient Boosting, feature importance is derived from their respective ensemble structures. Both models assess the contribution of each variable by evaluating how much it reduces uncertainty (e.g., Gini Impurity or entropy) when splitting data and decision nodes. Features that consistently lead to better splits across multiple iterations or trees are assigned higher importance values. However, while Random Forest aggregates these contributions across independent decision trees, Gradient Boosting refines feature importance iteratively. Each boosting step focuses on correcting prior errors, allowing the model to optimize feature impact progressively.

As shown in Figure 3, all of these feature importance values vary in magnitude, with only Logistic Regression showing the direction in which a variable influences the target. Therefore, a ranking of feature influence was created for all models by giving the feature with the greatest importance, relative to each model, the first position and the variable with the least influence the last position based on absolute importance (as shown in Figure 4). The purpose of this ranking system was to normalize the feature influence between all the models to aid in comparison.

With the foundations for the feature importance and rank metrics outlined, it is clear from Figure 4 what the most influential variables were when predicting default classification. Age was the most influential feature within the dataset, taking first place in all of the models tested. Figure 3 shows us that the chance of default decreases as age increases, based on the negative coefficient within the Logistic Regression model. Income and Interest Rate were tied, with both being the second most influential twice and the third most influential twice. Interestingly, while Income outranked Interest Rate in the Gradient Boosting and Random Forest trees, Interest Rate was higher in the Logistic Regression and KNN models. One possible explanation for this is that Logistic Regression and KNN place a higher importance on Interest Rate due to its more direct and linear relationship with default likelihood, making it a straightforward predictor in simpler models. Conversely, Gradient Boosting and Random Forest evaluate features through non-linear interactions, such as Loan Amount or Credit Score. This suggests that Income's role in predicting defaults may depend on more complex patterns better captured by ensemble methods.

Another key finding from this investigation was that the features with the highest absolute correlations had a higher rank than those with a correlation to default closer to zero, as highlighted by Figures 5 through 9. Figure 5, which visualizes the average rank of each feature in all models versus the feature's correlation with default, showcases a clear parabolic trend, confirming that features with stronger correlations tend to be ranked more important by the models. This trend is reinforced in Figures 6 through 9, visualizing the importance of feature correlation to the default relationship for each individual model. Despite the differences in how these models determine feature importance, they all exhibit a similar parabolic pattern: features with the highest absolute correlations to default (such as Age, Interest Rate, and Income, as

presented in Figure 1) consistently rank among the top features, whereas features with weak absolute correlations (such as Loan Term and Loan Purpose) have minimal influence.

This feature analysis reveals a consistent relationship between a feature's correlation with loan default and its importance in predictive models. The strongest predictors (Age, Interest Rate, and Income) align across different modeling approaches, indicating that these variables should be prioritized when assessing borrower risk. This insight has direct implications for lenders as well. Given that younger borrowers, those with higher interest rates, or those with lower income levels are at greater risk of default, financial institutions should adjust their risk assessment frameworks accordingly. This could involve more stringent lending criteria for high-risk applicants, enhanced monitoring of borrowers with elevated default probabilities, and targeted intervention programs such as financial counseling or restructuring options for borrowers showing early signs of risk. By leveraging these insights, lenders can develop data-driven strategies to mitigate risk, optimize lending decisions, and enhance overall financial stability.

## Threshold Analysis

To evaluate the trade-offs between the competing priorities in loan default classification, threshold analysis was found to be integral. Through the variation of threshold values, the model performances could be compared based on previously stated metrics, including FNR, FPR, Recall, Precision, Accuracy, F1-Score, and AUC-ROC. Heatmaps were created for all these metrics over the different threshold values tested to highlight the complex interplay between metrics and provide actionable insights for selecting optimal thresholds based on specific financial and operational goals. To add to interpretability, the heatmaps have been placed on a green-to-red scale, with green equating to a more favorable value and red being less favorable, respective to each metric.

With threshold tuning, FNR and FPR emerged as critical metrics for understanding the trade-offs in this project. As shown in Figure 10, FNR increased steadily with higher threshold values, while FPR exhibited an inverse trend, as depicted in Figure 11. For instance, if a lending institution serves a large and diverse client base and aims to minimize the financial impact of defaults, it would likely benefit from a low threshold value, ensuring fewer defaults are overlooked but at the expense of higher operational costs due to an increased number of false positives. Conversely, institutions with smaller client bases or limited resources may prefer a higher threshold, prioritizing precise interventions while accepting a greater risk of missed defaults. At a threshold of 0.1, Gradient Boosting and Logistic Regression exhibit the lowest FNR, meaning they are the most effective at capturing defaults. However, at higher thresholds (above 0.4), Random Forest and KNN achieve the lowest FPR, reducing false alarms but increasing missed defaults.

As for Recall, illustrated in Figure 12, it follows an inverse relationship with FNR, steadily declining as the threshold increases. A low threshold maximizes Recall by capturing more true defaults, which is beneficial for lenders focused on risk mitigation. However, this increase in Recall comes with a higher FPR, meaning more non-defaulting clients are incorrectly classified

as high risk. Gradient Boosting achieves the highest Recall at a threshold of 0.1, followed closely by Logistic Regression, making these models the most effective at identifying defaults when the goal is to minimize missed cases. Meanwhile, Random Forest maintains a competitive Recall at slightly higher thresholds (0.2–0.3), making it a more balanced option for institutions that want to control operational costs while still capturing most defaulters. Together, FNR, FPR, and Recall offer a comprehensive view of the model's performance under different thresholds, helping institutions tailor their strategies based on their unique priorities and constraints.

Precision and Accuracy provide insight into how effectively the model distinguishes between defaulters and non-defaulters as the threshold changes. Figure 13 shows that Precision increases as the threshold rises, meaning fewer non-defaulters are incorrectly flagged. While this improves the reliability of flagged defaults, it comes at the cost of a declining Recall. At a threshold of 0.4, Random Forest achieves the highest Precision, making it ideal for institutions that need highly accurate default classifications. However, this comes at the cost of missing a greater proportion of actual defaults. Accuracy, displayed in Figure 14, remains relatively stable across thresholds, reflecting its limited usefulness in imbalanced datasets. Since the dataset contains far more non-defaulters than defaulters, even a model that misclassifies most defaults can maintain a high Accuracy, in turn making Accuracy a poor standalone measure for model effectiveness. However, Random Forest maintains the highest Accuracy across all threshold levels, making it a consistent performer in overall classification correctness.

The F1-Score and AUC-ROC results tell another part of the classification performance story. Figure 15 highlights how F1-Score balances Precision and Recall, making it a strong indicator of overall model effectiveness. Gradient Boosting and Random Forest achieve their highest F1-Scores at a threshold of 0.2, suggesting that this is the optimal point for balancing risk detection and operational efficiency. Since F1-Score declines at higher thresholds, institutions that cannot afford to miss defaults should prioritize lower thresholds (0.1–0.2) to maximize Recall, even if Precision decreases slightly. AUC-ROC, visualized in Figure 16, remains stable across thresholds, reinforcing its role as a model ranking tool rather than a threshold selection metric. A consistently high AUC-ROC across all thresholds for Gradient Boosting and Random Forest suggests that these models are the most effective at distinguishing between defaults and non-defaults. However, Logistic Regression experiences a slight dip in AUC-ROC at higher thresholds, reinforcing that its performance is stronger when thresholds are lower.

## Model Comparison

Evaluating the performance of each model across key metrics, Gradient Boosting and Random Forest emerge as the most effective models for loan default classification. However, the best choice depends on an institution's priorities—whether the focus is on minimizing financial losses from missed defaults or optimizing operational efficiency by reducing false alarms. Each model has distinct trade-offs, making it essential to align model selection with business objectives.

For lenders prioritizing financial loss mitigation, capturing as many default cases as possible is critical. This means selecting a model that maximizes Recall and minimizes the False Negative Rate (FNR) to ensure that few defaulters slip through the cracks. Gradient Boosting performs best in this category, achieving the highest Recall at a threshold of 0.1, meaning it successfully identifies the greatest number of true defaulters. Logistic Regression follows closely behind, but its performance deteriorates as the threshold increases, making it less reliable for institutions needing consistent risk detection. Random Forest provides a balanced alternative, maintaining strong Recall at slightly higher thresholds (0.2–0.3), making it a reasonable choice for lenders wanting to balance financial risk with operational constraints. KNN, however, struggles to handle imbalanced datasets, leading to poor Recall and an increased likelihood of missing defaulters. If the goal is to reduce financial losses from missed defaults, Gradient Boosting is the best model, followed by Logistic Regression, with Random Forest as a more balanced option.

Conversely, lenders focused on operational efficiency aim to minimize the number of false alarms, ensuring that flagged defaults are more likely to be true defaults. This means selecting a model that maximizes Precision and minimizes the False Positive Rate (FPR) to reduce unnecessary interventions. Random Forest is the strongest performer in this category, achieving the highest Precision at thresholds of 0.3–0.4, meaning that the flagged defaults are more accurate, reducing wasted resources on unnecessary investigations. Gradient Boosting, while effective, prioritizes capturing defaults over reducing false positives, making it a better choice for risk-heavy strategies rather than efficiency-focused ones. Logistic Regression struggles at higher thresholds, as its precision declines, making it less viable for minimizing operational overhead. KNN ranks the lowest once again, as its classification approach results in inconsistent precision, leading to a higher number of misclassifications. For institutions prioritizing operational efficiency and precise interventions, Random Forest is the best choice, followed by Gradient Boosting, with Logistic Regression falling behind as threshold levels increase.

Ultimately, the best model depends on a lender's strategic focus. For institutions looking to reduce financial losses from missed defaults, Gradient Boosting is the strongest choice, as it provides the highest Recall and lowest FNR at lower thresholds. For institutions prioritizing efficiency and minimizing unnecessary interventions, Random Forest is the better option, as it achieves higher Precision and lower FPR at moderate thresholds. While Gradient Boosting and Random Forest stand out as the most effective models overall, selecting the appropriate threshold remains essential. The threshold dictates whether the model leans toward risk mitigation or cost efficiency, making it a key factor in aligning machine learning implementation with real-world lending strategies.

## Real-World Impact

In today's lending environment, predictive modeling is no longer a luxury but necessary for financial institutions seeking to mitigate risk, improve efficiency, and maximize profitability. Thus, this project aimed to predict loan defaults through machine learning methods to enable

lenders to make data-driven decisions that protect their portfolios while ensuring continued access to credit for responsible borrowers. However, data analysis alone is not enough. Institutions must transform analytical insights into actionable strategies that directly improve loan approval processes, borrower interventions, and overall financial stability. This section outlines how lenders can apply machine learning models like those used in this project to enhance decision-making, automate risk assessment, and achieve measurable financial benefits.

## Actionable Strategies for Lenders

### Policy Adjustments Based on High-Risk Features

The findings indicate that Age, Interest Rate, and Income are the strongest predictors of default, meaning lenders must adjust their risk assessment criteria accordingly. Borrowers with low incomes or high interest rates present a significantly greater likelihood of default, making it essential for lenders to introduce stricter approval conditions for applicants in these high-risk categories. This could include requiring higher collateral, stronger cosigners, or reduced borrowing limits for these applicants, ensuring they do not incur unsustainable debt burdens.

Lenders can also improve risk management through dynamic, risk-adjusted interest rates. Instead of applying a one-size-fits-all rate structure, institutions could reward financial stability with incremental rate reductions. For example, borrowers who demonstrate consistent on-time payments, increasing income, or improved credit scores could receive gradual interest rate adjustments as an incentive for responsible borrowing. This strategy reduces default rates over time while enhancing borrower retention, as customers are more likely to remain loyal to lenders who actively support their financial well-being.

### Operational Interventions for At-Risk Borrowers

Beyond adjusting lending criteria, financial institutions can proactively intervene to support borrowers before they default. The results suggest that younger borrowers are particularly vulnerable, meaning lenders should introduce early intervention programs tailored to their needs. This could include pre-loan financial literacy courses, personalized repayment strategies, or educational outreach initiatives to ensure that younger borrowers understand their financial responsibilities before taking on debt.

Another key improvement is implementing flexible repayment structures for borrowers with unstable incomes. Fixed payment plans often do not accommodate individuals with variable earnings, such as freelancers or gig economy workers. Instead, lenders could adopt income-linked repayment models, where monthly payments adjust based on the borrower's financial standing. This strategy not only reduces the risk of default but also ensures a steady cash flow for lenders, as borrowers are more likely to continue making payments rather than defaulting altogether.

## Optimizing Loan Approval and Financial Impact

Machine learning allows lenders to streamline underwriting, reduce operational costs, and minimize financial losses from defaults. By leveraging models like Random Forest and Gradient Boosting, institutions can implement automated loan decision systems that classify applicants based on their likelihood of default, enabling more efficient and accurate lending decisions.

For institutions prioritizing operational efficiency, Random Forest is the preferred model, as it achieves the highest Precision and Accuracy at moderate thresholds (0.3–0.4). This ensures that fewer non-defaulters are mistakenly flagged as high-risk, reducing false positives and minimizing unnecessary manual reviews. By implementing automated classification systems, low-risk applicants can be instantly approved, while moderate-risk applicants are flagged for further review. This allows underwriting teams to focus on high-risk cases, improving resource allocation and reducing processing times.

The financial impact of these optimizations is substantial. Based on the results, Random Forest demonstrated the highest Precision at a threshold of 0.4, which significantly reduces the number of false positives. This means that fewer non-defaulters are mistakenly flagged as high-risk, lowering the volume of unnecessary manual reviews. If a lender processes 10,000 loan applications per year, and assuming that 40% of flagged applicants would no longer require manual review due to improved classification accuracy, this equates to over 2,000 labor hours saved annually (given that each application review takes an estimated 30 minutes). These efficiency gains translate to lower administrative costs, faster loan approvals, and improved borrower satisfaction while allowing lending institutions to scale their operations more effectively.

For lenders focused on default risk mitigation, Gradient Boosting is the better option, as it achieves the highest Recall and lowest False Negative Rate (FNR) at lower thresholds (0.1–0.2). This ensures that at-risk borrowers are consistently identified, reducing the chances of approving loans that are likely to default. Based on the results, Gradient Boosting achieved the highest Recall at a threshold of 0.1, meaning it is the most effective model for capturing true defaulters. A lender issuing $500 million in loans annually with a 5% default rate ($25 million in losses per year) could reduce defaults by 10% using Gradient Boosting, which corresponds with the FNR reduction observed in the results of this investigation at lower thresholds. This reduction would prevent approximately $2.5 million in annual losses, helping to stabilize the lender's credit portfolio and mitigate long-term financial risks.

Beyond cost savings, smarter loan servicing strategies ensure that resources are focused on borrowers who truly need assistance. Rather than applying blanket interventions, machine learning allows lenders to target at-risk borrowers with personalized engagement strategies, improving borrower retention and long-term repayment success. Through a combination of default reduction, improved operational efficiency, and smarter resource allocation, machine learning enables lenders to create a more sustainable and financially stable lending system.

## Real-World Implications

Predictive modeling is most valuable when its insights lead to actionable improvements in operational strategy or, within this context, lending strategy. This project's findings illustrate how machine learning can enhance financial decision-making by balancing risk mitigation with operational efficiency. Gradient Boosting excelled at identifying high-risk borrowers and preventing financial losses, while Random Forest was the better choice for streamlining underwriting and reducing unnecessary manual reviews. By selecting the appropriate model based on institutional priorities, lenders can strengthen risk management frameworks and improve borrower engagement while maintaining profitability.

Beyond model selection, the success of predictive analytics depends on translating insights into concrete policies. Adjusting lending criteria, implementing targeted interventions, and optimizing approval processes ensure that machine learning is not just a tool for risk assessment but a driver of meaningful change in financial practices. Ultimately, the value of predictive modeling is realized not in its ability to classify risk but in how lenders apply it to make smarter, data-driven decisions that enhance financial stability and accessibility.

# Reflection

The objective of this project was to develop a machine learning-based framework for predicting loan defaults, addressing a critical challenge in financial decision-making. By leveraging multiple models and structured evaluation techniques, the study sought to provide lenders with actionable insights to enhance risk assessment, minimize financial losses, and optimize operational efficiency. As this study concludes, it is important to reflect on what was accomplished, where challenges arose, and how future work can build upon these findings. Predictive modeling is not a static solution but a continuously evolving tool—one that must be refined, adapted, and applied strategically to deliver lasting improvements in lending practices.

## Key Achievements

This study successfully established a comprehensive framework for loan default prediction, beginning with robust data preprocessing and progressing through structured model evaluation and real-world financial impact analysis. A key strength was the systematic approach to data preparation, which included encoding categorical variables, scaling numerical features, and conducting correlation analyses. These preprocessing steps ensured that models were trained on clean, structured data, reducing bias and improving interpretability.

Beyond model implementation, the project demonstrated a thorough evaluation framework that compared Logistic Regression, KNN, Random Forest, and Gradient Boosting across multiple threshold values. This allowed for a nuanced assessment of trade-offs between financial loss mitigation and operational efficiency. The structured evaluation process incorporated False Negative Rate (FNR), False Positive Rate (FPR), Recall, Precision, Accuracy, F1-Score, and AUC-ROC, ensuring that each model's strengths and weaknesses

were analyzed from multiple perspectives. Additionally, the inclusion of heatmaps, feature importance visualizations, and threshold impact assessments provided lenders with clear, interpretable insights into model performance.

A key success of this project was demonstrating that machine learning can generate measurable financial benefits for lenders. The financial impact calculations, while presented as illustrative examples rather than definitive projections, highlighted the potential for default reduction and operational cost savings. The results showed that Gradient Boosting is best suited for minimizing financial losses, as it captured the highest number of true defaulters at lower thresholds. Meanwhile, Random Forest proved to be the most effective for operational efficiency, reducing unnecessary manual reviews with its high Precision and Accuracy. These findings translated into actionable recommendations, showing lenders how to tailor their model selection and threshold tuning based on institutional priorities.

## Limitations

While this study produced meaningful insights, several limitations must be acknowledged. One of the most significant challenges was the high False Negative Rate (FNR) at higher thresholds, particularly in models like Random Forest that prioritized Precision over Recall. This trade-off resulted in more missed defaults, a critical limitation for institutions seeking to minimize financial risk.

Another key limitation was the dataset's lack of behavioral and longitudinal data. The models relied primarily on static financial attributes, such as age, income, and loan amount, but did not incorporate real-time financial activity, spending behaviors, or evolving credit history. Additionally, since the dataset was a single snapshot rather than a time-series dataset, the models were unable to track changes in borrower financial stability over time. These constraints limited the ability of the models to predict long-term repayment behaviors and adapt to borrowers' changing financial situations.

From an implementation perspective, computational costs and scalability must also be considered. While Gradient Boosting and Random Forest delivered strong predictive performance, they require substantial processing power compared to simpler models like Logistic Regression. This could pose challenges for real-time lending applications, particularly for high-volume lenders that need to process thousands of loan applications daily.

## Future Work

While this study provided valuable insights into loan default prediction, several areas for future improvement could further refine model performance and real-world applicability. One significant enhancement would be the integration of time-series data, allowing models to track borrower financial stability, repayment behavior, and spending patterns over time rather than relying solely on static attributes. This would enable lenders to detect early signs of financial distress and implement proactive intervention strategies before defaults occur. Additionally, future research should explore cost-sensitive learning techniques, where misclassifications are

weighted based on financial impact. Assigning higher penalties to False Negatives (missed defaults) could help ensure that the models prioritize minimizing financial risk while maintaining operational efficiency.

Another promising direction is hybrid and ensemble learning approaches. While Gradient Boosting and Random Forest performed well individually, combining models into stacked ensembles could further balance Precision, Recall, and overall classification accuracy. If additional data sources become available, deep-learning techniques may also be explored to capture more complex borrower relationships. Lastly, given the importance of threshold selection in balancing default identification and operational efficiency, future work should experiment with threshold values below 0.1 to determine whether a lower cutoff could further reduce False Negatives without significantly increasing false positives. By implementing these improvements, future studies can expand upon this project's foundation, creating even more accurate, adaptable, and actionable loan default prediction models.

## Closing Thoughts

This project successfully demonstrated how machine learning can transform loan default prediction, offering lenders a data-driven framework to balance financial risk and operational efficiency. The study provided clear, structured insights into model selection, threshold tuning, and the real-world financial impact of predictive analytics, proving that machine learning is not just a theoretical tool but a practical solution for risk management. However, the true value of predictive modeling lies in its continuous refinement and application. While Gradient Boosting and Random Forest emerged as the most effective models, their real-world impact depends on how lenders apply and adapt them over time. By integrating better data, more dynamic modeling techniques, and improved threshold tuning, financial institutions can further enhance their risk assessment frameworks and lending strategies.

Ultimately, predictive lending is an evolving process—one that must be refined as new data, better algorithms, and improved methodologies become available. This study serves as a step toward leveraging data science for smarter financial decision-making, helping lenders reduce risk, increase efficiency, and create a more stable financial system for the future.

# References

Nikhil1e9. (n.d.). Loan Default Dataset. Kaggle. Retrieved from
https://www.kaggle.com/datasets/nikhil1e9/loan-default

# Appendices

## Figure 1 - Correlation to Default

| Feature | Default |
|---|---|
| InterestRate | 0.1314 |
| LoanAmount | 0.0864 |
| EmploymentType | 0.0414 |
| NumCreditLines | 0.0282 |
| DTIRatio | 0.0208 |
| LoanTerm | 0.0012 |
| MaritalStatus | -0.0081 |
| LoanPurpose | -0.0106 |
| Education | -0.0218 |
| HasMortgage | -0.0223 |
| CreditScore | -0.0348 |
| HasDependents | -0.0358 |
| HasCoSigner | -0.0375 |
| MonthsEmployed | -0.0964 |
| Income | -0.0988 |
| Age | -0.1663 |

## Figure 2 - Correlation Matrix of Independent Variables

| Feature | Age | CreditScore | DTIRatio | Education | EmploymentType | HasCoSigner | HasDependents | HasMortgage | Income | InterestRate | LoanAmount | LoanPurpose | LoanTerm | MaritalStatus | MonthsEmployed | NumCreditLines |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Age | | -0.001 | -0.004 | -0.002 | 0.001 | -0.003 | -0.001 | 0.000 | -0.001 | 0.000 | -0.003 | 0.003 | 0.001 | -0.001 | -0.001 | -0.001 |
| CreditScore | -0.001 | | -0.002 | 0.001 | 0.004 | -0.002 | -0.004 | 0.002 | -0.001 | 0.001 | 0.002 | -0.001 | 0.001 | -0.005 | 0.001 | 0.001 |
| DTIRatio | -0.004 | -0.002 | | 0.002 | 0.000 | 0.000 | 0.003 | 0.000 | 0.001 | 0.001 | 0.001 | -0.004 | 0.004 | 0.004 | -0.001 | 0.000 |
| Education | -0.002 | 0.001 | 0.002 | | 0.000 | 0.000 | 0.002 | 0.000 | -0.002 | 0.005 | 0.005 | -0.003 | -0.003 | -0.004 | -0.002 | 0.005 |
| EmploymentType | 0.001 | 0.004 | 0.000 | 0.000 | | 0.001 | 0.003 | 0.001 | -0.003 | 0.001 | 0.001 | -0.002 | 0.002 | 0.001 | 0.001 | -0.001 |
| HasCoSigner | -0.003 | -0.002 | 0.000 | 0.000 | 0.001 | | 0.003 | -0.003 | -0.004 | -0.003 | -0.001 | -0.001 | -0.002 | -0.001 | 0.000 | 0.003 |
| HasDependents | -0.001 | -0.004 | 0.003 | 0.002 | 0.003 | 0.003 | | -0.001 | -0.001 | 0.002 | 0.000 | -0.004 | 0.003 | 0.002 | 0.002 | -0.003 |
| HasMortgage | 0.000 | 0.002 | 0.000 | 0.000 | 0.001 | -0.003 | -0.001 | | -0.002 | -0.001 | 0.000 | -0.002 | 0.002 | -0.001 | 0.000 | -0.001 |
| Income | -0.001 | -0.001 | 0.001 | -0.002 | -0.003 | -0.004 | -0.001 | -0.002 | | -0.003 | -0.002 | -0.002 | -0.002 | 0.002 | 0.003 | -0.001 |
| InterestRate | 0.000 | 0.001 | 0.001 | 0.005 | 0.001 | -0.003 | 0.002 | -0.001 | -0.003 | | -0.001 | 0.002 | 0.000 | -0.006 | 0.001 | 0.000 |
| LoanAmount | -0.003 | 0.002 | 0.001 | 0.005 | 0.001 | -0.001 | 0.000 | 0.000 | -0.002 | -0.001 | | -0.001 | 0.002 | -0.002 | 0.003 | 0.000 |
| LoanPurpose | 0.003 | -0.001 | -0.004 | -0.003 | -0.002 | -0.001 | -0.004 | -0.002 | -0.002 | 0.002 | -0.001 | | 0.003 | 0.001 | -0.002 | 0.000 |
| LoanTerm | 0.001 | 0.001 | 0.004 | -0.003 | 0.002 | -0.002 | 0.003 | 0.002 | -0.002 | 0.000 | 0.002 | 0.003 | | 0.000 | -0.002 | -0.001 |
| MaritalStatus | -0.001 | -0.005 | 0.004 | -0.004 | 0.001 | -0.001 | 0.002 | -0.001 | 0.002 | -0.006 | -0.002 | 0.001 | 0.000 | | 0.000 | -0.002 |
| MonthsEmployed | -0.001 | 0.001 | -0.001 | -0.002 | 0.001 | 0.000 | 0.002 | 0.000 | 0.003 | 0.001 | 0.003 | -0.002 | -0.002 | 0.000 | | 0.003 |
| NumCreditLines | -0.001 | 0.001 | 0.000 | 0.005 | -0.001 | 0.003 | -0.003 | -0.001 | -0.001 | 0.000 | 0.000 | 0.000 | -0.001 | -0.002 | 0.003 | |

Figure 3 - Heat Map of Feature Importance

|  | Model | | | |
|---|---|---|---|---|
| Feature | Logistic Regression | K-Nearest Neighbors | Gradient Boosting | Random Forest |
| Age | -0.5835 | 0.0017 | 0.2758 | 0.2243 |
| CreditScore | -0.1221 | -0.0003 | 0.0176 | 0.0505 |
| DTIRatio | 0.0679 | -0.0004 | 0.0067 | 0.0344 |
| Education | -0.0735 | -0.0002 | 0.0081 | 0.0115 |
| EmploymentType | 0.1287 | 0.0001 | 0.0196 | 0.0180 |
| HasCoSigner | -0.2625 | -0.0001 | 0.0132 | 0.0101 |
| HasDependents | -0.2536 | -0.0004 | 0.0121 | 0.0089 |
| HasMortgage | -0.1526 | -0.0003 | 0.0043 | 0.0049 |
| Income | -0.3404 | 0.0014 | 0.2104 | 0.1912 |
| InterestRate | 0.4567 | 0.0014 | 0.1887 | 0.1744 |
| LoanAmount | 0.2992 | 0.0010 | 0.1246 | 0.1274 |
| LoanPurpose | -0.0255 | -0.0002 | 0.0047 | 0.0119 |
| LoanTerm | 0.0002 | -0.0002 | 0.0003 | 0.0109 |
| MaritalStatus | -0.0319 | -0.0004 | 0.0066 | 0.0084 |
| MonthsEmployed | -0.3353 | 0.0000 | 0.0992 | 0.1004 |
| NumCreditLines | 0.0998 | -0.0004 | 0.0081 | 0.0128 |

Figure 4 - Heat Map of Feature Importance Rank by Model

| Feature | Model | | | | Grand Total |
| | Logistic Regression | K-Nearest Neighbors | Gradient Boosting | Random Forest | |
| --- | --- | --- | --- | --- | --- |
| Age | 1 | 1 | 1 | 1 | 1 |
| Income | 3 | 3 | 2 | 2 | 2.5 |
| InterestRate | 2 | 2 | 3 | 3 | 2.5 |
| LoanAmount | 5 | 4 | 4 | 4 | 4.25 |
| MonthsEmployed | 4 | 16 | 5 | 5 | 7.5 |
| CreditScore | 10 | 10 | 7 | 6 | 8.25 |
| NumCreditLines | 11 | 6 | 10 | 9 | 9 |
| EmploymentType | 9 | 15 | 6 | 8 | 9.5 |
| HasDependents | 7 | 8 | 9 | 14 | 9.5 |
| DTIRatio | 13 | 7 | 12 | 7 | 9.75 |
| HasCoSigner | 6 | 14 | 8 | 13 | 10.25 |
| Education | 12 | 11 | 11 | 11 | 11.25 |
| MaritalStatus | 14 | 5 | 13 | 15 | 11.75 |
| HasMortgage | 8 | 9 | 15 | 16 | 12 |
| LoanPurpose | 15 | 12 | 14 | 10 | 12.75 |
| LoanTerm | 16 | 13 | 16 | 12 | 14.25 |

## Figure 5 - Scatter Plot of Average Rank vs Correlation (Average of all Models)



Feature Rank vs Correlation with Default

Figure 6 - Scatter Plot of Absolute Importance vs Correlation, Model: Logistic Regression



Feature Absolute Importance vs Correlation with Default

Figure 7 - Scatter Plot of Absolute Importance vs Correlation, Model: K-Nearest Neighbors
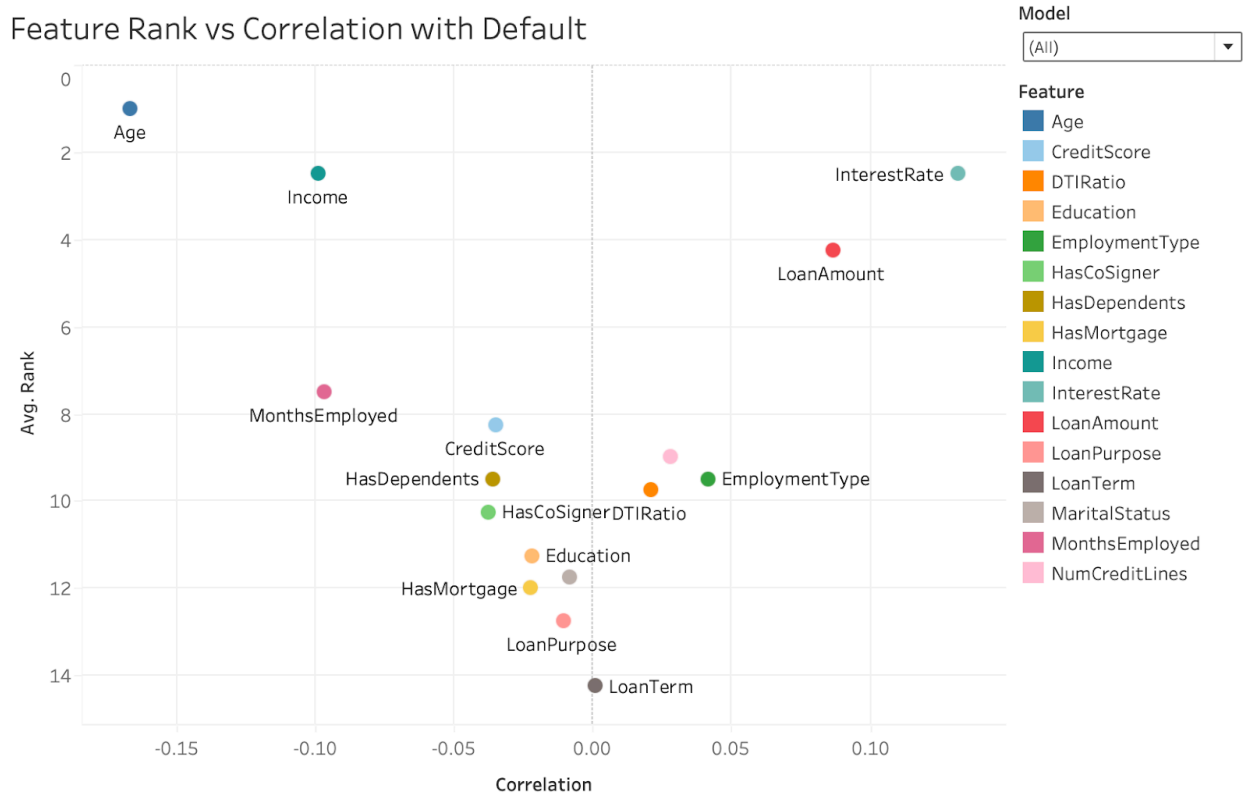
Feature Absolute Importance vs Correlation with Default
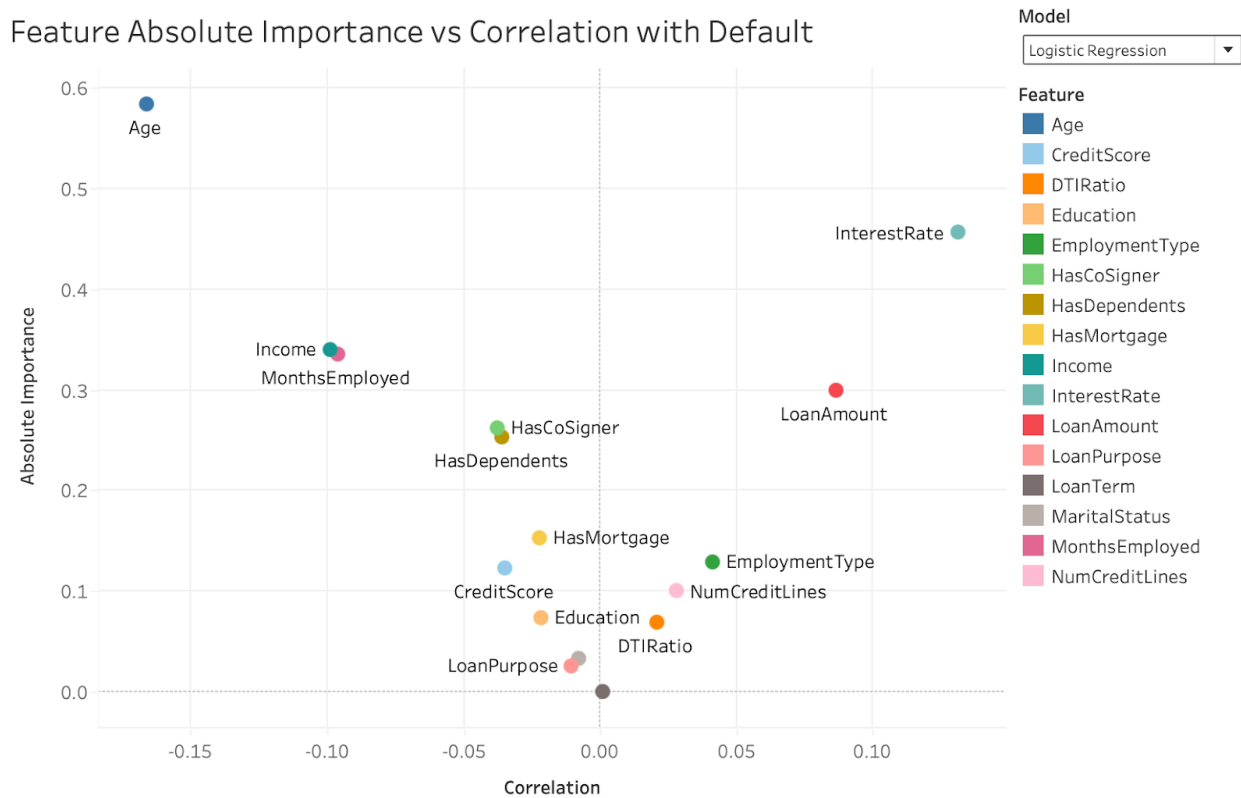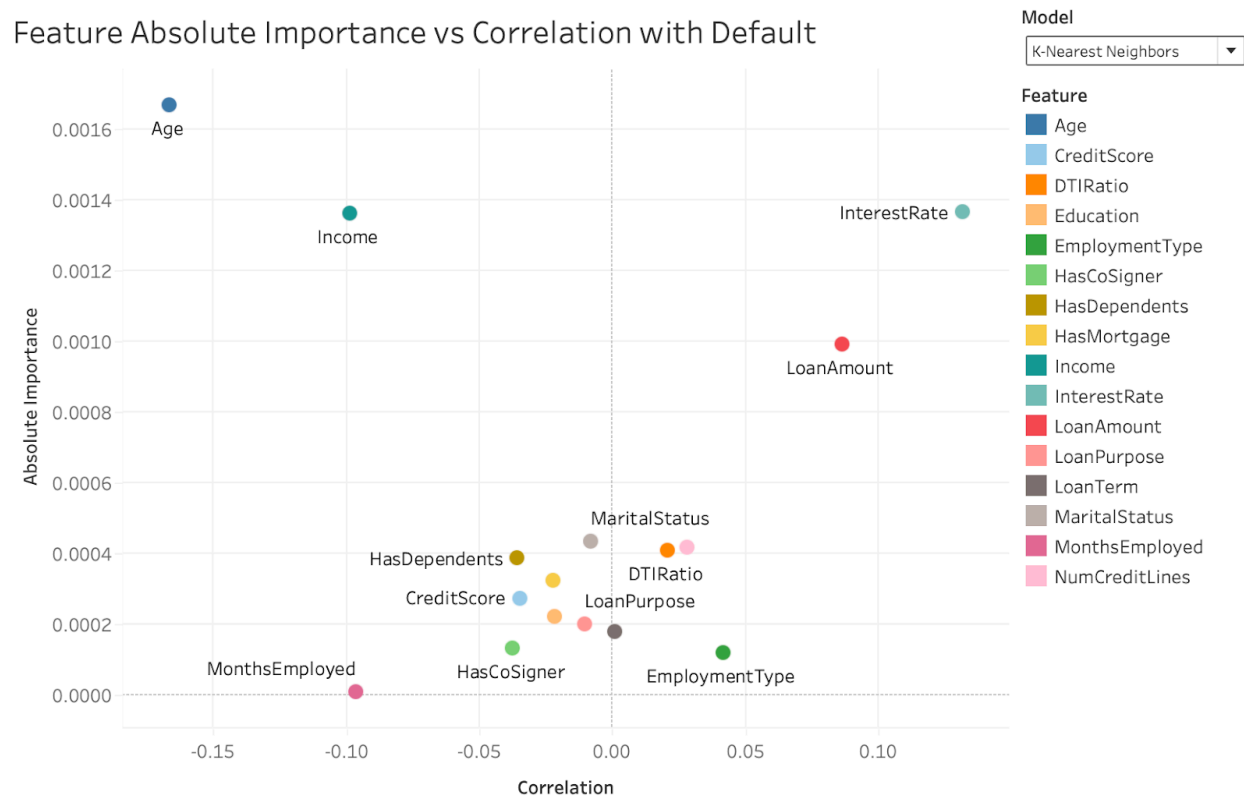
Figure 8 - Scatter Plot of Absolute Importance vs Correlation, Model: Gradient Boosting
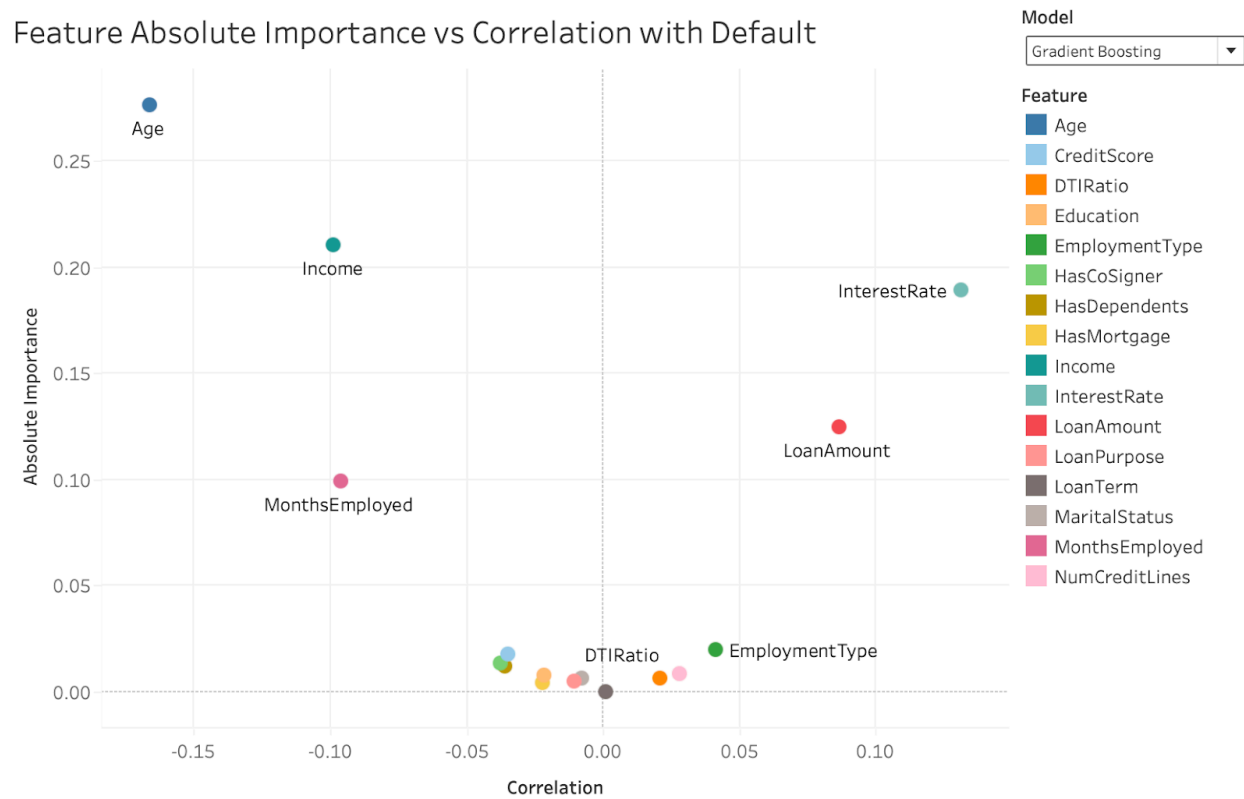
Figure 9 - Scatter Plot of Absolute Importance vs Correlation, Model: Random Forest

Feature Absolute Importance vs Correlation with Default

Figure 10 - Threshold Heatmap for False Negative Rate

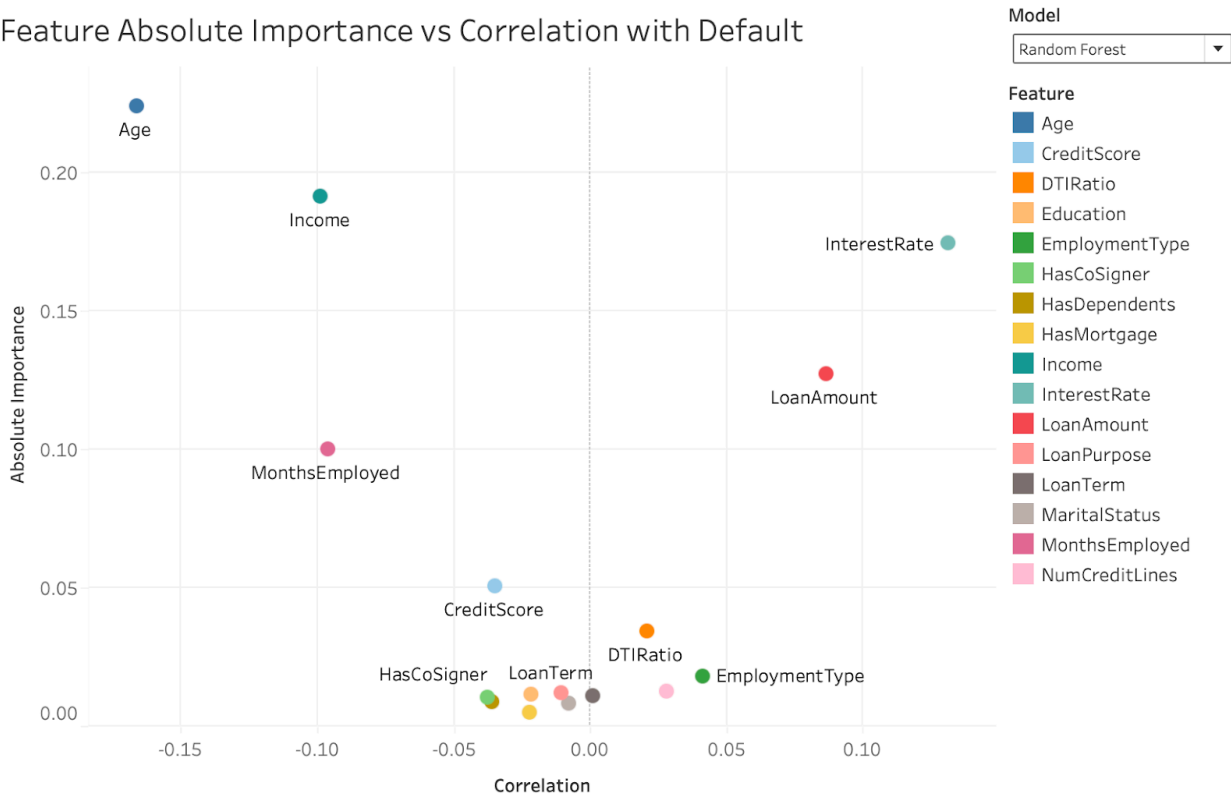| Threshold | Gradient Boosting | K-Nearest Neighbors | Logistic Regression | Random Forest | Grand Total |
|---|---|---|---|---|---|
| 0.1 | 0.25999 | 0.23150 | 0.24262 | 0.23032 | 0.24111 |
| 0.2 | 0.56584 | 0.53043 | 0.56466 | 0.60546 | 0.56660 |
| 0.3 | 0.75384 | 0.75839 | 0.78064 | 0.83173 | 0.78115 |
| 0.4 | 0.86326 | 0.89142 | 0.90474 | 0.94082 | 0.90006 |
| 0.5 | 0.93273 | 0.96594 | 0.96948 | 0.97876 | 0.96173 |
| 0.6 | 0.96931 | 0.98904 | 0.99393 | 0.99477 | 0.98676 |
| 0.7 | 0.98786 | 0.99747 | 0.99966 | 1.00000 | 0.99625 |
| 0.8 | 0.99764 | 0.99983 | 1.00000 | 1.00000 | 0.99937 |
| 0.9 | 1.00000 | 0.99983 | 1.00000 | 1.00000 | 0.99996 |
| Grand Total | 0.81450 | 0.81821 | 0.82842 | 0.84243 | 0.82589 |

Figure 11 - Threshold Heatmap for False Positive Rate

| Threshold | Model Gradient Boosting | K-Nearest Neighbors | Logistic Regression | Random Forest | Grand Total |
|---|---|---|---|---|---|
| 0.1 | 0.36117 | 0.56003 | 0.39485 | 0.41191 | 0.43199 |
| 0.2 | 0.12018 | 0.23937 | 0.13201 | 0.10793 | 0.14988 |
| 0.3 | 0.04398 | 0.08554 | 0.04304 | 0.02404 | 0.04915 |
| 0.4 | 0.01668 | 0.02650 | 0.01210 | 0.00514 | 0.01510 |
| 0.5 | 0.00569 | 0.00749 | 0.00261 | 0.00111 | 0.00423 |
| 0.6 | 0.00202 | 0.00171 | 0.00044 | 0.00011 | 0.00107 |
| 0.7 | 0.00035 | 0.00031 | 0.00004 | 0.00000 | 0.00018 |
| 0.8 | 0.00004 | 0.00004 | 0.00000 | 0.00000 | 0.00002 |
| 0.9 | 0.00000 | 0.00002 | 0.00000 | 0.00000 | 0.00001 |
| Grand Total | 0.06112 | 0.10233 | 0.06501 | 0.06114 | 0.07240 |

Figure 12 - Threshold Heatmap for Recall

| Threshold | Model Gradient Boosting | K-Nearest Neighbors | Logistic Regression | Random Forest | Grand Total |
|---|---|---|---|---|---|
| 0.1 | 0.74001 | 0.76850 | 0.75738 | 0.76968 | 0.75889 |
| 0.2 | 0.43416 | 0.46957 | 0.43534 | 0.39454 | 0.43340 |
| 0.3 | 0.24616 | 0.24161 | 0.21936 | 0.16827 | 0.21885 |
| 0.4 | 0.13674 | 0.10858 | 0.09526 | 0.05918 | 0.09994 |
| 0.5 | 0.06727 | 0.03406 | 0.03052 | 0.02124 | 0.03827 |
| 0.6 | 0.03069 | 0.01096 | 0.00607 | 0.00523 | 0.01324 |
| 0.7 | 0.01214 | 0.00253 | 0.00034 | 0.00000 | 0.00375 |
| 0.8 | 0.00236 | 0.00017 | 0.00000 | 0.00000 | 0.00063 |
| 0.9 | 0.00000 | 0.00017 | 0.00000 | 0.00000 | 0.00004 |
| Grand Total | 0.18550 | 0.18179 | 0.17158 | 0.15757 | 0.17411 |

Figure 13 - Threshold Heatmap for Precision

| Threshold | Gradient Boosting | K-Nearest Neighbors | Logistic Regression | Random Forest | Grand Total |
|---|---|---|---|---|---|
| | | | Model | | |
| 0.1 | 0.21211 | 0.15276 | 0.20130 | 0.19712 | 0.19082 |
| 0.2 | 0.32188 | 0.20493 | 0.30231 | 0.32446 | 0.28839 |
| 0.3 | 0.42380 | 0.27068 | 0.40105 | 0.47912 | 0.39366 |
| 0.4 | 0.51854 | 0.35000 | 0.50855 | 0.60206 | 0.49479 |
| 0.5 | 0.60823 | 0.37407 | 0.60535 | 0.71591 | 0.57589 |
| 0.6 | 0.66667 | 0.45775 | 0.64286 | 0.86111 | 0.65710 |
| 0.7 | 0.81818 | 0.51724 | 0.50000 | 0.00000 | 0.45886 |
| 0.8 | 0.87500 | 0.33333 | 0.00000 | 0.00000 | 0.30208 |
| 0.9 | 0.00000 | 0.50000 | 0.00000 | 0.00000 | 0.12500 |
| Grand Total | 0.49382 | 0.35120 | 0.35127 | 0.35331 | 0.38740 |

Figure 14 - Threshold Heatmap for Accuracy

| Threshold | Gradient Boosting | K-Nearest Neighbors | Logistic Regression | Random Forest | Grand Total |
|---|---|---|---|---|---|
| | | | Model | | |
| 0.1 | 0.88665 | 0.88120 | 0.88510 | 0.88535 | 0.88458 |
| 0.2 | 0.88665 | 0.88120 | 0.88510 | 0.88535 | 0.88458 |
| 0.3 | 0.88665 | 0.88120 | 0.88510 | 0.88535 | 0.88458 |
| 0.4 | 0.88665 | 0.88120 | 0.88510 | 0.88535 | 0.88458 |
| 0.5 | 0.88665 | 0.88120 | 0.88510 | 0.88535 | 0.88458 |
| 0.6 | 0.88665 | 0.88120 | 0.88510 | 0.88535 | 0.88458 |
| 0.7 | 0.88665 | 0.88120 | 0.88510 | 0.88535 | 0.88458 |
| 0.8 | 0.88665 | 0.88120 | 0.88510 | 0.88535 | 0.88458 |
| 0.9 | 0.88665 | 0.88120 | 0.88510 | 0.88535 | 0.88458 |
| Grand Total | 0.88665 | 0.88120 | 0.88510 | 0.88535 | 0.88458 |

Figure 15 - Threshold Heatmap for F1-Score

|  | Model | | | | |
| Threshold | Gradient Boosting | K-Nearest Neighbors | Logistic Regression | Random Forest | Grand Total |
|---|---|---|---|---|---|
| 0.1 | 0.32971 | 0.25486 | 0.31806 | 0.31386 | 0.30413 |
| 0.2 | 0.36968 | 0.28533 | 0.35683 | 0.35608 | 0.34198 |
| 0.3 | 0.31143 | 0.25532 | 0.28360 | 0.24906 | 0.27485 |
| 0.4 | 0.21641 | 0.16574 | 0.16047 | 0.10777 | 0.16260 |
| 0.5 | 0.12115 | 0.06243 | 0.05811 | 0.04126 | 0.07074 |
| 0.6 | 0.05867 | 0.02141 | 0.01203 | 0.01039 | 0.02562 |
| 0.7 | 0.02392 | 0.00503 | 0.00067 | 0.00000 | 0.00741 |
| 0.8 | 0.00471 | 0.00034 | 0.00000 | 0.00000 | 0.00126 |
| 0.9 | 0.00000 | 0.00034 | 0.00000 | 0.00000 | 0.00008 |
| Grand Total | 0.15952 | 0.11676 | 0.13220 | 0.11983 | 0.13207 |

Figure 16 - Threshold Heatmap for AUC-ROC

|  | Model | | | | |
| Threshold | Gradient Boosting | K-Nearest Neighbors | Logistic Regression | Random Forest | Grand Total |
|---|---|---|---|---|---|
| 0.1 | 0.75876 | 0.65408 | 0.74972 | 0.75246 | 0.72875 |
| 0.2 | 0.75876 | 0.65408 | 0.74972 | 0.75246 | 0.72875 |
| 0.3 | 0.75876 | 0.65408 | 0.74972 | 0.75246 | 0.72875 |
| 0.4 | 0.75876 | 0.65408 | 0.74972 | 0.75246 | 0.72875 |
| 0.5 | 0.75876 | 0.65408 | 0.74972 | 0.75246 | 0.72875 |
| 0.6 | 0.75876 | 0.65408 | 0.74972 | 0.75246 | 0.72875 |
| 0.7 | 0.75876 | 0.65408 | 0.74972 | 0.75246 | 0.72875 |
| 0.8 | 0.75876 | 0.65408 | 0.74972 | 0.75246 | 0.72875 |
| 0.9 | 0.75876 | 0.65408 | 0.74972 | 0.75246 | 0.72875 |
| Grand Total | 0.75876 | 0.65408 | 0.74972 | 0.75246 | 0.72875 |