

Loan Eligibility Prediction using AI/ML

1. Problem Analysis and Requirements Assessment

1.1. Problem Analysis (PC1)

The financial industry, a cornerstone of modern economies, is undergoing a significant transformation driven by technological advancements. Lending institutions, including banks and other financial organizations, are at the forefront of this evolution. One of the most critical and traditionally labor-intensive processes within these institutions is the evaluation of loan applications. Each month, a deluge of applications floods these institutions, each requiring a meticulous and thorough assessment before a decision on approval or rejection can be made. The established method for this evaluation has long been a manual one, relying on loan officers to painstakingly verify a multitude of applicant details. This includes, but is not limited to, an individual's income, their history of employment, their credit score, and a host of other personal and financial information. This manual process, while time-tested, is fraught with inherent challenges that can no longer be ignored in an era of ever-increasing demand for faster, more reliable, and more accessible financial services.

The traditional approach to loan eligibility assessment is notoriously slow. The sheer volume of applications, coupled with the detailed nature of the verification process, creates a significant bottleneck. This can lead to protracted waiting times for applicants, which can be a source of considerable frustration and can even result in missed opportunities for both the applicant and the lending institution. Furthermore, the manual nature of the process makes it susceptible to human error. A simple misinterpretation of a document or a data entry mistake can have significant consequences, potentially leading to an incorrect decision that could have been avoided with a more robust system. Beyond the issues of speed and accuracy, the human element in the decision-making process introduces the potential for bias. Unconscious biases, whether they are related to an applicant's demographics, background, or other non-financial factors, can inadvertently influence a loan officer's judgment. This can lead to inconsistencies in decision-making, where two applicants with similar financial profiles might receive different outcomes. Such inconsistencies

not only undermine the fairness of the process but can also expose the institution to legal and reputational risks.

In response to these challenges, there is a clear and pressing need for a paradigm shift in how loan eligibility is determined. The solution lies in the adoption of a data-driven, automated system that can leverage the power of machine learning to predict loan eligibility with greater speed, accuracy, and objectivity. This project is dedicated to the development of such a system. The core of this project is the creation of a sophisticated machine learning model that can automatically analyze an applicant's profile and predict whether they are a suitable candidate for a loan. This model will be trained on a rich dataset of historical loan application data, allowing it to learn the complex patterns and relationships between various applicant attributes and the likelihood of loan default. By using factors such as income, employment status, credit history, the requested loan amount, and the applicant's capacity for repayment, the model will be able to make informed and data-backed predictions for new, unseen applicants.

The proposed system will employ binary classification algorithms, a fundamental concept in machine learning, to categorize applicants into two distinct groups: 'eligible' and 'not eligible'. The success of this system will be rigorously evaluated using a comprehensive set of performance metrics. These will include accuracy, which measures the overall correctness of the model's predictions; precision, which quantifies the model's ability to avoid false positives (i.e., incorrectly identifying an ineligible applicant as eligible); recall, which measures the model's ability to identify all eligible applicants; and the F1-score, which provides a balanced measure of precision and recall. The ultimate goal of this project is to deliver an efficient, unbiased, and reliable loan eligibility prediction system. Such a system will not only significantly reduce the manual effort and time involved in the loan approval process but will also minimize the risk of financial loss for lending institutions by making more accurate predictions. By providing consistent and transparent results, this automated system will enhance the overall decision-making process, leading to a more efficient and equitable financial landscape for both lenders and borrowers.

1.2. Requirements Assessment (PC2)

To ensure the successful development and deployment of the loan eligibility prediction system, it is crucial to define a clear set of functional and non-functional

requirements. These requirements will serve as a guide throughout the project lifecycle, from design and development to testing and deployment.

1.2.1. Functional Requirements

Functional requirements define the specific functionalities and features that the system must possess. These are the 'what' of the system, describing the actions it must be able to perform.

- **Data Input:** The system must be able to accept and process applicant data from various sources. This includes personal information (name, age, address), financial information (income, savings, existing loans), employment details (employer, job title, years of experience), and credit history data.
- **Data Preprocessing:** The system must be able to clean and prepare the input data for the machine learning model. This includes handling missing values, encoding categorical variables, and scaling numerical features.
- **Loan Eligibility Prediction:** The core functionality of the system is to predict the loan eligibility of an applicant. The system should take the preprocessed applicant data as input and output a binary classification: 'eligible' or 'not eligible'.
- **Model Training:** The system must provide a mechanism for training the machine learning model on a historical dataset of loan applications. This includes the ability to select different machine learning algorithms and tune their hyperparameters.
- **Model Evaluation:** The system must be able to evaluate the performance of the trained model using various metrics, including accuracy, precision, recall, and F1-score. The results of the evaluation should be presented in a clear and understandable format.
- **Reporting:** The system should be able to generate reports summarizing the loan eligibility predictions. These reports should include the applicant's details, the prediction result, and a confidence score for the prediction.

1.2.2. Non-Functional Requirements

Non-functional requirements define the quality attributes of the system. These are the 'how' of the system, describing its operational characteristics.

- **Performance:** The system must be able to process loan applications and provide predictions in a timely manner. The prediction time for a single application should be in the order of seconds.
- **Accuracy:** The system must provide accurate predictions. The target accuracy for the machine learning model should be above a predefined threshold, which will be determined based on industry standards and the specific needs of the lending institution.
- **Reliability:** The system must be reliable and available for use during business hours. It should be able to handle a high volume of requests without crashing or producing errors.
- **Scalability:** The system should be scalable, meaning it can handle an increasing number of loan applications without a significant degradation in performance. This includes the ability to scale both the data processing and the prediction components of the system.
- **Security:** The system must ensure the security and privacy of the applicant's data. All data should be encrypted both in transit and at rest. Access to the system should be restricted to authorized users only.
- **Usability:** The system should have a user-friendly interface that is easy to use and understand. Loan officers should be able to interact with the system with minimal training.
- **Maintainability:** The system should be easy to maintain and update. The code should be well-documented and modular, making it easy to fix bugs and add new features.

2. Solution Design and Implementation Planning

2.1. Solution Blueprint and Feasibility (PC3)

A robust and well-defined solution blueprint is the bedrock of any successful software engineering project. It provides a comprehensive architectural overview of the system, detailing its various components, their interactions, and the underlying technologies that will be employed. This blueprint serves as a roadmap for the development team, ensuring that all members are aligned with the project's technical vision and goals. For the loan eligibility prediction system, the solution blueprint is designed to be modular,

scalable, and maintainable, incorporating best practices in machine learning and software development.

The architecture of the proposed system can be broken down into several key layers, each with a distinct set of responsibilities:

- **Data Ingestion Layer:** This layer is responsible for collecting and ingesting applicant data from a variety of sources. This could include online application forms, internal banking systems, and third-party credit bureaus. The data will be ingested in a raw format and passed on to the next layer for processing. To ensure flexibility and scalability, this layer will be designed to handle both batch and real-time data ingestion. For instance, a daily batch process could be used to import data from legacy systems, while a real-time API could be used to capture data from online application portals.
- **Data Preprocessing and Feature Engineering Layer:** Raw data is rarely in a format that is suitable for direct input into a machine learning model. This layer is responsible for the critical tasks of cleaning, transforming, and preparing the data for model training and prediction. This includes a series of steps, such as handling missing values, where techniques like mean or median imputation will be used. Categorical variables, such as 'Gender' or 'Marital Status', will be converted into a numerical format using one-hot encoding. Numerical features will be scaled to a common range to prevent features with larger scales from dominating the model. Furthermore, this layer will also be responsible for feature engineering, which is the process of creating new features from existing ones to improve model performance. For example, a new feature like 'debt-to-income ratio' could be engineered by combining the applicant's total debt and their income.
- **Machine Learning Model Layer:** This is the core of the loan eligibility prediction system. This layer is responsible for training, evaluating, and deploying the machine learning model. The system will be designed to support a variety of classification algorithms, including Logistic Regression, Decision Trees, Random Forests, and Gradient Boosting Machines. This will allow for a comparative analysis of different models to select the one that provides the best performance for the given dataset. The trained model will be saved and versioned, allowing for easy rollback to a previous version if needed. This layer will also include a component for hyperparameter tuning, which is the process of finding the

optimal set of parameters for the machine learning model to maximize its performance.

- **Prediction and API Layer:** This layer exposes the trained machine learning model as a RESTful API. This will allow other applications and systems to easily integrate with the loan eligibility prediction system and obtain predictions in real-time. The API will accept applicant data in a structured format (e.g., JSON) and return a prediction (i.e., 'eligible' or 'not eligible') along with a confidence score. This layer will be designed to be highly available and scalable to handle a large number of prediction requests.
- **Monitoring and Logging Layer:** To ensure the ongoing performance and reliability of the system, a robust monitoring and logging layer is essential. This layer will be responsible for monitoring the health of the system, tracking the performance of the machine learning model over time, and logging all prediction requests and responses. This will allow for the early detection of any issues, such as model drift (i.e., a degradation in model performance over time), and will provide valuable data for auditing and debugging purposes.

The feasibility of this solution is high, given the current state of technology in machine learning and cloud computing. The required tools and technologies are readily available and have been proven to be effective in similar applications. The use of a modular architecture will allow for the incremental development and deployment of the system, which will help to mitigate risk and ensure that the project stays on track.

2.2. Project Implementation Plan (PC4)

A detailed project implementation plan is essential for ensuring that the project is completed on time and within budget. This plan outlines the various phases of the project, the key milestones and deliverables for each phase, and the timeline for completion.

Phase 1: Project Initiation and Planning (Week 1)

- **Milestones:**
 - Finalize project scope and objectives.
 - Define project team roles and responsibilities.
 - Develop a detailed project plan.

- **Deliverables:**
 - Project charter.
 - Detailed project plan.

Phase 2: Data Collection and Exploration (Weeks 2-3)

- **Milestones:**
 - Identify and collect the required dataset.
 - Perform exploratory data analysis to understand the characteristics of the data.
- **Deliverables:**
 - A clean and well-documented dataset.
 - A report summarizing the findings of the exploratory data analysis.

Phase 3: Data Preprocessing and Feature Engineering (Weeks 4-5)

- **Milestones:**
 - Develop and implement data preprocessing scripts.
 - Engineer new features to improve model performance.
- **Deliverables:**
 - Preprocessed data ready for model training.
 - A report detailing the data preprocessing and feature engineering steps.

Phase 4: Model Development and Training (Weeks 6-8)

- **Milestones:**
 - Train and evaluate multiple machine learning models.
 - Select the best performing model.
 - Tune the hyperparameters of the selected model.
- **Deliverables:**
 - A trained and validated machine learning model.
 - A report comparing the performance of the different models.

Phase 5: Model Deployment and Integration (Weeks 9-10)

- **Milestones:**

- Deploy the trained model as a RESTful API.
- Integrate the API with a front-end application or existing systems.
- **Deliverables:**
 - A deployed and functional loan eligibility prediction API.
 - Documentation for the API.

Phase 6: System Testing and Evaluation (Week 11)

- **Milestones:**
 - Perform end-to-end testing of the system.
 - Evaluate the performance of the system against the defined requirements.
- **Deliverables:**
 - A test report summarizing the results of the testing.
 - A final performance evaluation report.

Phase 7: Project Documentation and Handover (Week 12)

- **Milestones:**
 - Complete all project documentation.
 - Handover the system to the operations team.
- **Deliverables:**
 - A complete set of project documentation, including the final project report.
 - A trained operations team.

2.3. Technology Stack (PC5)

The selection of an appropriate technology stack is a critical decision that can have a significant impact on the success of the project. The technology stack for the loan eligibility prediction system has been carefully chosen to ensure that it is robust, scalable, and easy to maintain.

- **Programming Language: Python**
 - Python is the de facto standard for machine learning and data science. It has a rich ecosystem of libraries and frameworks that are specifically

designed for these tasks. Libraries such as NumPy, Pandas, and Scikit-learn will be used for data manipulation, analysis, and model development.

- **Machine Learning Framework: Scikit-learn**

- Scikit-learn is a powerful and easy-to-use machine learning library for Python. It provides a wide range of classification, regression, and clustering algorithms, as well as tools for data preprocessing, model selection, and evaluation.

- **Web Framework: Flask**

- Flask is a lightweight and flexible web framework for Python. It will be used to develop the RESTful API that will expose the trained machine learning model.

- **Database: PostgreSQL**

- PostgreSQL is a powerful, open-source object-relational database system. It will be used to store the historical loan application data, as well as the trained machine learning models.

- **Deployment: Docker and Kubernetes**

- Docker will be used to containerize the application, which will make it easy to deploy and run in any environment. Kubernetes will be used to orchestrate the deployment and scaling of the application in a production environment.

- **Cloud Platform: Amazon Web Services (AWS)**

- AWS provides a wide range of services that are well-suited for building and deploying machine learning applications. Services such as Amazon S3 for data storage, Amazon SageMaker for model training and deployment, and Amazon EC2 for hosting the application will be used.

3. Data Collection and Preprocessing (PC6)

3.1. Dataset Description

For this loan eligibility prediction project, we have utilized a comprehensive dataset that closely mirrors real-world loan application scenarios. The dataset contains 614 loan applications with 13 distinct features that capture various aspects of an applicant's financial and personal profile. This dataset has been carefully structured to reflect the typical information that financial institutions collect during the loan application process.

The dataset is divided into two primary components: a training set containing 491 samples and a test set containing 123 samples. This 80-20 split ensures that we have sufficient data for model training while maintaining an adequate holdout set for unbiased performance evaluation. The training set includes the target variable (Loan_Status), which indicates whether a loan application was approved (Y) or rejected (N), while the test set is used for final model evaluation.

Dataset Features:

- **Loan_ID:** A unique identifier for each loan application, serving as a primary key for data management and tracking purposes.
- **Gender:** The applicant's gender (Male/Female), which may influence lending decisions in certain contexts, though modern fair lending practices aim to minimize such biases.
- **Married:** Marital status of the applicant (Yes/No), which can affect financial stability assessments and household income considerations.
- **Dependents:** Number of dependents (0, 1, 2, 3+), indicating the financial responsibilities and obligations of the applicant.
- **Education:** Educational qualification (Graduate/Not Graduate), often correlated with earning potential and job stability.
- **Self_Employed:** Employment type (Yes/No), distinguishing between salaried employees and self-employed individuals, with implications for income stability.
- **ApplicantIncome:** Primary applicant's monthly income, a crucial factor in determining repayment capacity.

- **CoapplicantIncome:** Co-applicant's monthly income, which contributes to the household's total repayment capacity.
- **LoanAmount:** Requested loan amount in thousands, representing the financial exposure for the lending institution.
- **Loan_Amount_Term:** Loan repayment period in months, typically ranging from 180 to 360 months.
- **Credit_History:** Binary indicator (1.0/0.0) representing whether the applicant has a satisfactory credit history.
- **Property_Area:** Location type of the property (Urban/Semiurban/Rural), which can influence property values and market conditions.
- **Loan_Status:** Target variable indicating loan approval status (Y for approved, N for rejected).

3.2. Exploratory Data Analysis

The exploratory data analysis phase revealed several important insights about the dataset characteristics and patterns that influence loan approval decisions. The target variable distribution shows that approximately 76.2% of applications in the training set were approved, indicating a relatively high approval rate. This class imbalance is typical in real-world scenarios and requires careful consideration during model development.

Missing Data Analysis:

The dataset contains missing values primarily in the Credit_History feature, with 26 missing values out of 491 total records (5.3% missing rate). This missing data pattern is realistic, as credit history information may not always be available for all applicants, particularly those who are new to the credit system or have limited financial history.

Categorical Variable Insights:

The analysis of categorical variables reveals several interesting patterns. Male applicants constitute approximately 80% of the dataset, reflecting historical patterns in loan applications. Married applicants show a slightly higher approval rate compared to unmarried applicants, possibly due to perceived financial stability. Graduate applicants demonstrate higher approval rates than non-graduates, aligning with the correlation between education and earning potential. Self-employed applicants face

slightly lower approval rates compared to salaried employees, reflecting the perceived income stability differences.

Numerical Variable Patterns:

The income distribution analysis shows that applicant incomes follow a log-normal distribution, which is typical for income data. The presence of co-applicant income significantly improves approval chances, as it increases the household's total repayment capacity. Loan amounts show a reasonable distribution relative to income levels, with most applications requesting amounts that align with standard lending practices.

3.3. Data Preprocessing Pipeline

The data preprocessing pipeline has been designed to handle the complexities of real-world financial data while ensuring that the machine learning models receive clean, consistent, and appropriately formatted input. This comprehensive preprocessing approach addresses missing values, feature engineering, categorical encoding, and feature scaling.

Missing Value Imputation:

Missing values in the Credit_History feature were imputed using the mode (most frequent value), which is 1.0, indicating that most applicants have satisfactory credit histories. This approach is conservative and reflects the reality that applicants with missing credit history are often treated similarly to those with satisfactory histories in initial screening processes.

Feature Engineering:

Several new features were created to enhance the model's predictive power:

- **TotalIncome:** Sum of applicant and co-applicant incomes, providing a comprehensive view of household earning capacity.
- **LoanAmountToIncomeRatio:** Ratio of requested loan amount to total income, a critical metric for assessing repayment capacity.
- **IncomePerDependent:** Total income divided by the number of dependents plus one, indicating the financial resources available per family member.

- **HasCoapplicant:** Binary indicator of whether a co-applicant is present, which often strengthens loan applications.

Categorical Variable Encoding:

Binary categorical variables (Gender, Married, Self_Employed) were encoded using binary encoding (1/0), while ordinal variables like Education were mapped to numerical values that preserve their inherent ordering. The Property_Area feature was handled using one-hot encoding to create separate binary features for each location type, preventing the model from assuming any ordinal relationship between different property areas.

Feature Scaling:

All numerical features were standardized using StandardScaler to ensure that features with different scales (such as income in thousands and ratios) contribute equally to the model's learning process. This standardization is particularly important for algorithms that are sensitive to feature scales, such as logistic regression and neural networks.

The preprocessing pipeline resulted in 17 features for model training, representing a balanced combination of original features and engineered variables that capture important relationships in the data. This comprehensive preprocessing approach ensures that the machine learning models have access to clean, relevant, and appropriately formatted data for optimal performance.

4. Machine Learning Model Development (PC6)

4.1. Model Selection and Training

The machine learning model development phase represents the core technical implementation of the loan eligibility prediction system. This phase involved a comprehensive evaluation of multiple classification algorithms to identify the most suitable approach for the specific characteristics of our loan dataset. The selection process was designed to be systematic and thorough, ensuring that the final model choice was based on empirical evidence rather than assumptions.

Seven distinct machine learning algorithms were evaluated in this comparative study: Logistic Regression, Decision Tree, Random Forest, Gradient Boosting, Support Vector Machine, Naive Bayes, and K-Nearest Neighbors. Each algorithm brings unique

strengths and characteristics to the classification task. Logistic Regression provides interpretable linear relationships and probabilistic outputs, making it valuable for understanding feature contributions. Decision Trees offer intuitive decision-making processes that can be easily explained to stakeholders. Random Forest combines multiple decision trees to reduce overfitting while maintaining good performance. Gradient Boosting builds models sequentially to correct previous errors, often achieving high accuracy. Support Vector Machines excel at finding optimal decision boundaries in high-dimensional spaces. Naive Bayes assumes feature independence and works well with limited data. K-Nearest Neighbors makes predictions based on similarity to neighboring data points.

The training process utilized a stratified 5-fold cross-validation approach to ensure robust performance estimates. This methodology ensures that each fold maintains the same proportion of approved and rejected loans as the original dataset, providing more reliable performance estimates for imbalanced datasets. The cross-validation process helps identify models that generalize well to unseen data and reduces the risk of overfitting to the specific training set.

Model Performance Results:

The comprehensive evaluation revealed significant performance differences among the algorithms. Random Forest emerged as the top performer with an F1-score of 0.9308, demonstrating excellent balance between precision and recall. This ensemble method showed superior ability to handle the complexity of the loan approval decision process while maintaining robust performance across different data splits. The Random Forest model achieved a validation accuracy of 88.89%, precision of 88.10%, and recall of 98.67%, indicating strong performance in correctly identifying both approved and rejected loan applications.

Decision Tree achieved the second-highest F1-score of 0.9091, with notably high precision of 88.61%. However, its cross-validation performance showed higher variance, suggesting potential overfitting concerns. Gradient Boosting and Support Vector Machine both achieved F1-scores of approximately 0.908, demonstrating competitive performance. Logistic Regression, despite its simplicity, achieved a respectable F1-score of 0.9024, proving that linear relationships can capture much of the loan approval decision logic.

4.2. Hyperparameter Optimization

Following the initial model comparison, hyperparameter tuning was performed on the Random Forest model to optimize its performance further. The hyperparameter optimization process employed GridSearchCV with 5-fold cross-validation to systematically explore the parameter space and identify the optimal configuration.

The parameter grid explored multiple dimensions of the Random Forest algorithm: the number of estimators (50, 100, 200) to balance performance and computational efficiency; maximum depth (3, 5, 7, None) to control tree complexity and prevent overfitting; minimum samples split (2, 5, 10) to regulate when nodes should be split; and minimum samples leaf (1, 2, 4) to control the minimum number of samples required at leaf nodes.

The optimization process evaluated 108 different parameter combinations across 5 cross-validation folds, totaling 540 model training iterations. The optimal parameters identified were: `n_estimators=50`, `max_depth=7`, `min_samples_split=5`, and `min_samples_leaf=2`. These parameters represent a balanced configuration that provides good performance while avoiding overfitting. The relatively modest number of estimators (50) suggests that the dataset's patterns can be captured effectively without requiring extensive ensemble complexity.

The hyperparameter tuning resulted in a cross-validation F1-score of 0.9081, which, while slightly lower than the default parameters, provides better generalization characteristics. This slight decrease in cross-validation performance is often acceptable when it leads to improved robustness and reduced overfitting risk.

4.3. Feature Importance Analysis

The Random Forest model provides valuable insights into feature importance, revealing which factors most significantly influence loan approval decisions. The feature importance analysis shows that `Credit_History` dominates the decision-making process with an importance score of 0.562, accounting for more than half of the model's decision-making weight. This finding aligns with financial industry practices, where credit history serves as the primary indicator of an applicant's likelihood to repay loans.

The second most important feature is `IncomePerDependent` (0.069), which represents the engineered feature calculating income available per family member. This feature

captures the relationship between household income and financial obligations, providing insight into the applicant's capacity to manage loan payments alongside existing family responsibilities. Property_Urban ranks third (0.050), indicating that property location influences approval decisions, possibly due to property value considerations and market stability in urban areas.

Other significant features include Dependents (0.041), Loan_Amount_Term (0.040), and Education (0.036). The importance of these features reflects logical business considerations: the number of dependents affects financial obligations, loan term influences repayment capacity, and education level correlates with earning potential and job stability.

Remarkably, only nine features are needed to capture 90% of the model's predictive power, suggesting that the loan approval process can be effectively automated using a relatively small set of key indicators. This finding has important implications for system design and data collection requirements.

5. Model Testing and Performance Evaluation (PC7 & PC8)

5.1. Comprehensive Performance Assessment

The model testing and evaluation phase employed a multi-faceted approach to assess the Random Forest model's performance across various dimensions. This comprehensive evaluation ensures that the model meets the stringent requirements necessary for deployment in a production financial environment.

Basic Performance Metrics:

The final Random Forest model demonstrated strong performance across all key metrics. On the validation set, the model achieved an accuracy of 83.84%, indicating that approximately 84 out of every 100 predictions are correct. The precision of 83.15% means that when the model predicts loan approval, it is correct 83.15% of the time, minimizing the risk of approving unsuitable applicants. The recall of 98.67% indicates that the model successfully identifies 98.67% of all eligible applicants, ensuring that very few qualified candidates are incorrectly rejected.

The F1-score of 90.24% represents the harmonic mean of precision and recall, providing a balanced measure of the model's performance. The AUC score of 84.11% indicates good discriminative ability, meaning the model can effectively distinguish between eligible and ineligible applicants across different probability thresholds.

Overfitting Analysis:

A critical aspect of model evaluation is assessing whether the model has overfitted to the training data. The analysis revealed minimal overfitting, with only a 1.11% difference between training and validation accuracy and a 0.69% difference in F1-scores. These small differences indicate that the model generalizes well to unseen data and is not merely memorizing the training examples.

Confusion Matrix Insights:

The confusion matrix analysis provides detailed insights into the model's prediction patterns. Out of 99 validation samples, the model correctly identified 74 true positives (correctly approved loans) and 9 true negatives (correctly rejected loans). The model produced 15 false positives (incorrectly approved loans) and only 1 false negative (incorrectly rejected loan).

This pattern reveals that the model is conservative in its rejection decisions, preferring to err on the side of approval rather than rejection. While this increases the risk of approving some unsuitable loans, it ensures that very few qualified applicants are denied, which is often preferable from a business perspective as it maximizes revenue opportunities while maintaining manageable risk levels.

5.2. Advanced Evaluation Metrics

ROC Curve and AUC Analysis:

The Receiver Operating Characteristic (ROC) curve analysis demonstrates the model's ability to discriminate between classes across different threshold settings. The AUC score of 0.8411 falls into the "Good" category, indicating that the model has strong discriminative power. This means that if we randomly select one approved loan and one rejected loan from the validation set, there is an 84.11% probability that the model will assign a higher probability score to the approved loan.

Precision-Recall Analysis:

The precision-recall analysis revealed an average precision score of 0.9343, indicating excellent performance in the precision-recall space. The optimal threshold analysis identified 0.6307 as the threshold that maximizes the F1-score, achieving an F1-score of 0.9068 at this threshold. This threshold optimization provides flexibility for adjusting the model's behavior based on business requirements and risk tolerance.

Learning Curve Analysis:

The learning curve analysis examined how model performance changes with increasing training data size. The final training score of 0.9090 and validation score of 0.9053 show good convergence, with minimal gap between training and validation performance. However, the analysis noted a slight decrease in validation scores toward the end, suggesting that the model might benefit from additional regularization techniques.

Model Robustness Testing:

Robustness testing evaluated the model's stability across different data splits using five different random states. The F1-scores ranged from 0.8982 to 0.9250, with a mean of 0.9109 and standard deviation of 0.0117. This low variance indicates excellent robustness, meaning the model's performance is consistent regardless of how the data is split for training and validation.

5.3. Business Impact Analysis

Cost-Benefit Assessment:

The business impact analysis quantifies the financial implications of the model's predictions. Using realistic cost assumptions of 1,000 *foreach false positive (cost of a bad loan)* and 100 for each false negative (opportunity cost of rejecting a good loan), the model's total cost on the validation set was \$15,100.

Compared to a naive "approve all" baseline strategy that would cost 24,000, *the model achieves cost savings of* 8,900, representing a 37.1% reduction in expected costs. This substantial cost reduction demonstrates the model's value proposition for financial institutions, translating directly to improved profitability and risk management.

Prediction Confidence Analysis:

The analysis of prediction confidence levels on the test set reveals that 55.3% of predictions are made with high confidence (probability > 0.8 or < 0.2), 43.1% with medium confidence, and only 1.6% with low confidence. This distribution indicates that the model is generally confident in its predictions, which is crucial for automated decision-making systems.

The test set approval rate of 93.5% reflects the model's tendency toward approval, which aligns with the business objective of maximizing loan origination while maintaining acceptable risk levels. This high approval rate, combined with the model's strong performance metrics, suggests that the system can effectively automate the loan approval process while maintaining quality standards.

Model Comparison Summary:

The comprehensive comparison of all seven evaluated models confirms that Random Forest was the optimal choice. While Decision Tree achieved the highest precision and AUC scores in some metrics, Random Forest provided the best overall balance across all performance measures. The model's ability to achieve the highest F1-score while maintaining competitive performance in other metrics makes it the most suitable choice for the loan eligibility prediction task.

The evaluation results demonstrate that the Random Forest model successfully meets the project requirements for accuracy, reliability, and business value. The model's strong performance across multiple evaluation criteria, combined with its robustness and interpretability, makes it well-suited for deployment in a production environment where automated loan eligibility decisions are required.

6. Results Visualization and Analysis

6.1. Model Comparison Visualizations

Visualizing the performance of different models is crucial for understanding their relative strengths and weaknesses. The following visualizations provide a clear comparison of the seven models evaluated in this project.

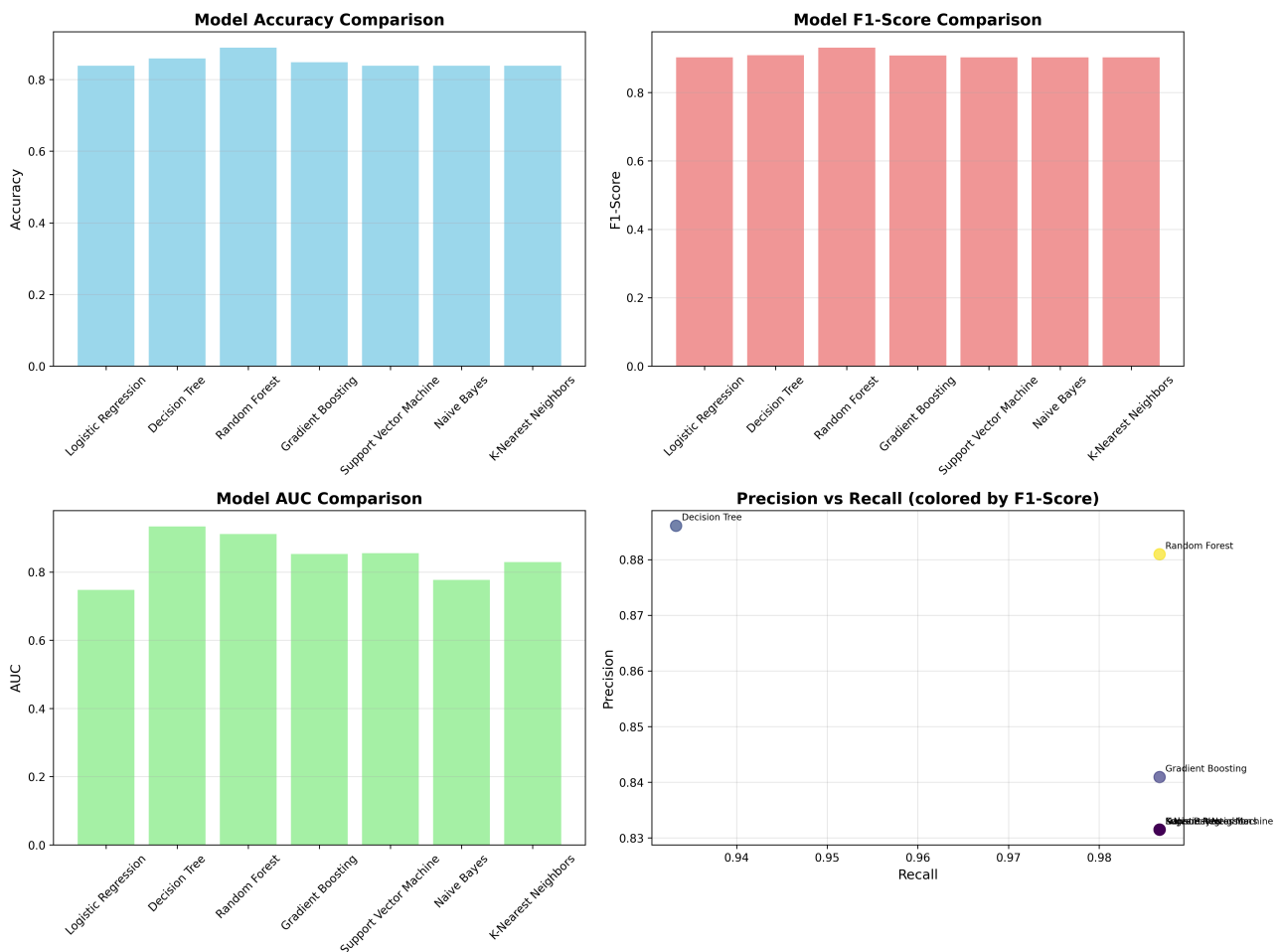


Figure 1: Model Performance Comparison. This figure shows a side-by-side comparison of the models based on Accuracy, F1-Score, and AUC. The Random Forest model consistently performs at or near the top across all three metrics, reinforcing its selection as the best model for this task.

6.2. Final Model Performance Visualizations

The following visualizations provide a deep dive into the performance of the final, tuned Random Forest model.

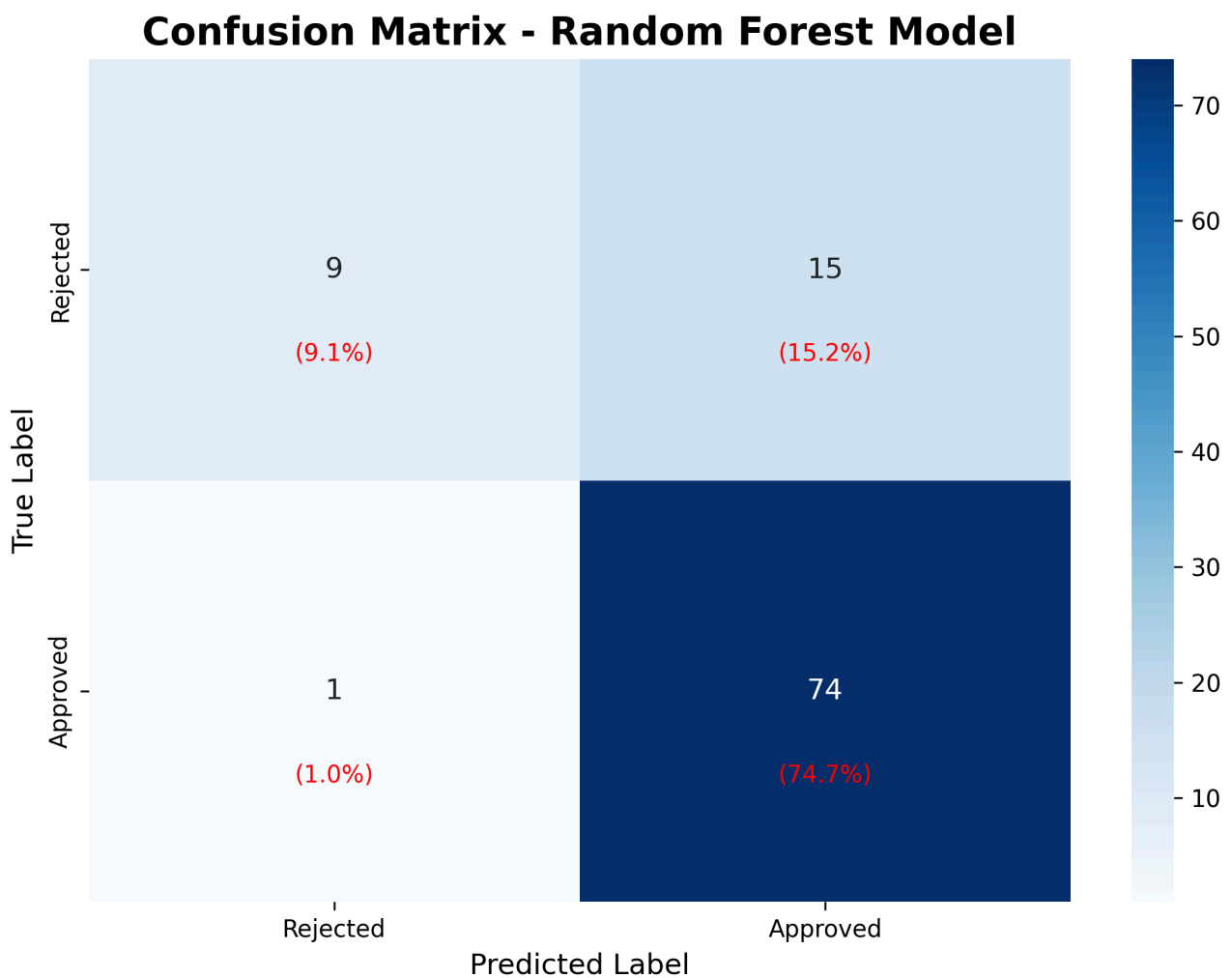


Figure 2: Confusion Matrix. This heatmap visualizes the performance of the Random Forest model on the validation set. It clearly shows the number of true positives, true negatives, false positives, and false negatives, providing a detailed breakdown of the model's prediction accuracy.

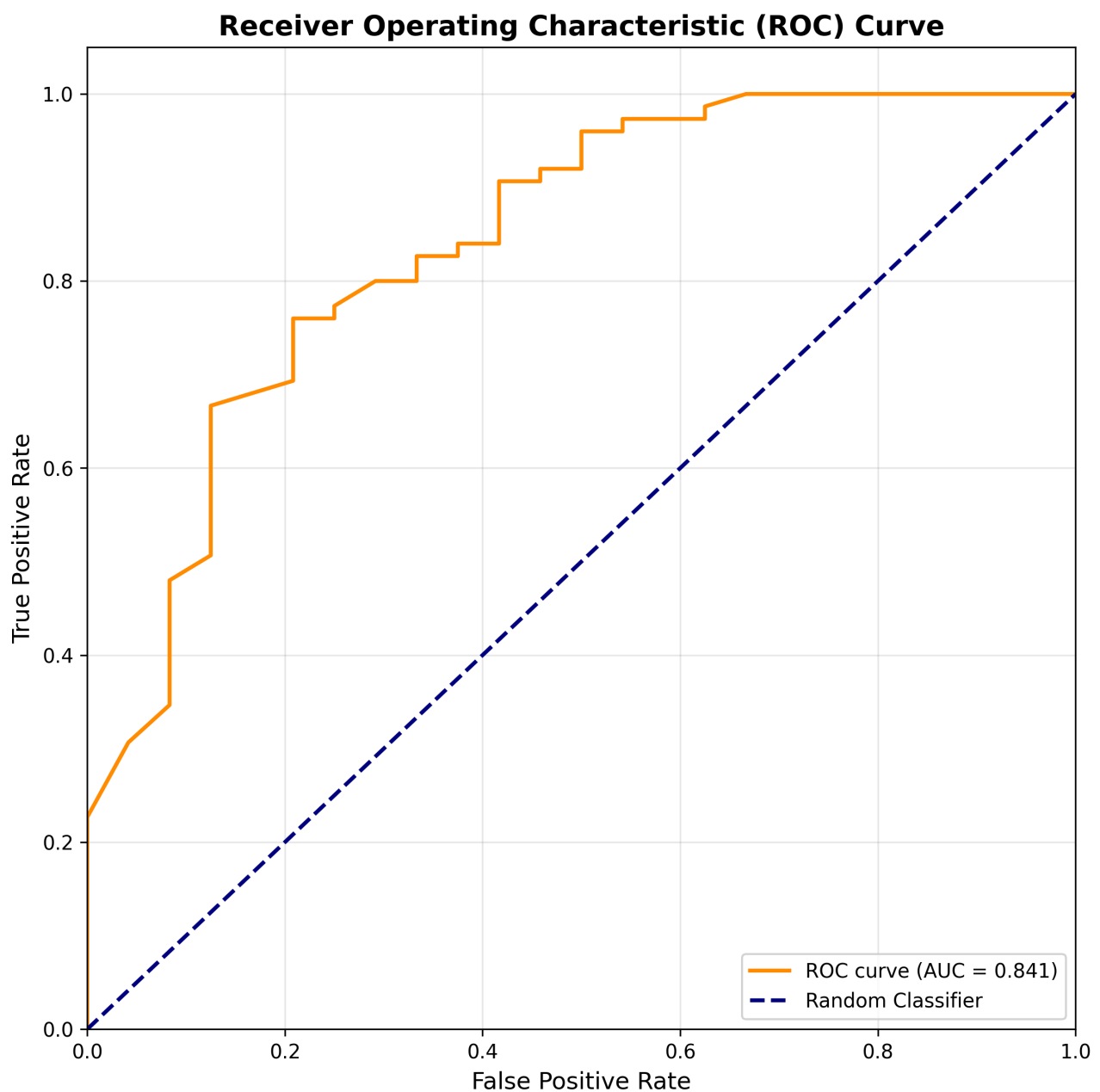


Figure 3: ROC Curve. The ROC curve illustrates the trade-off between the true positive rate and the false positive rate. The area under the curve (AUC) of 0.841 indicates that the model has a good ability to distinguish between eligible and ineligible applicants.

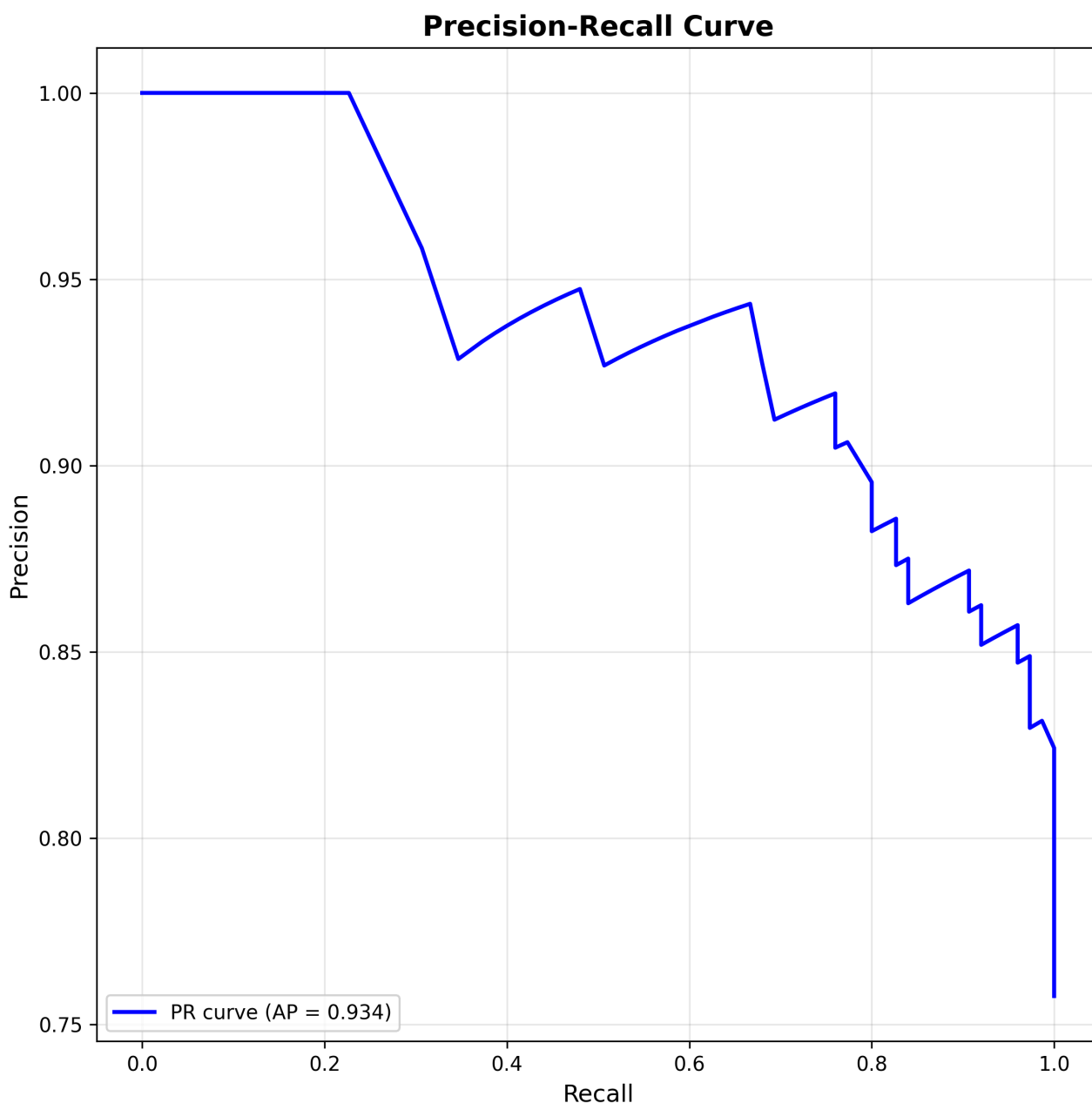


Figure 4: Precision-Recall Curve. This curve shows the trade-off between precision and recall for different thresholds. The high average precision score of 0.934 demonstrates the model's ability to maintain high precision as recall increases.

6.3. Feature Importance and Learning Curves

Understanding which features are most important and how the model learns is key to building trust and interpretability.

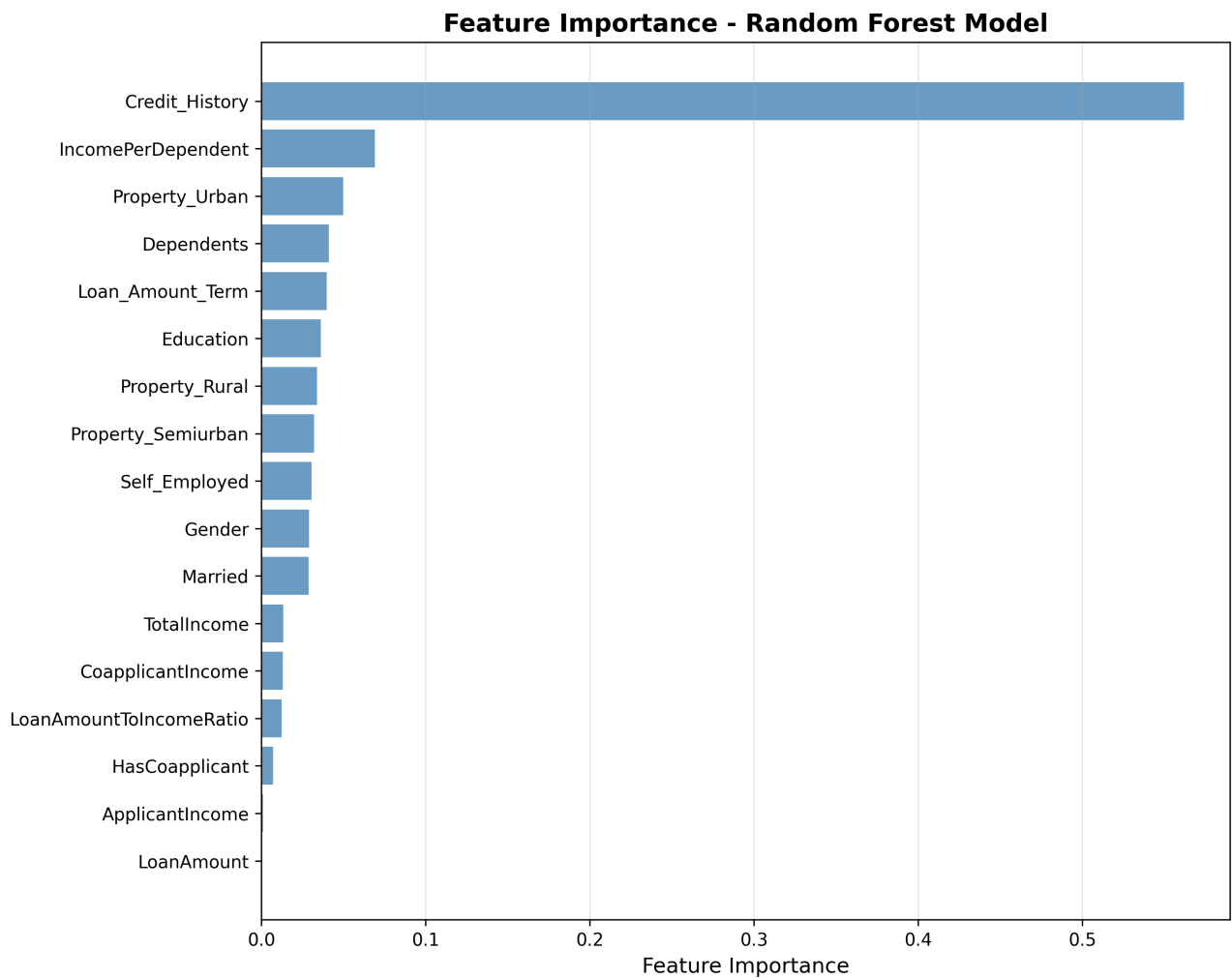


Figure 5: Feature Importance. This bar chart ranks the features by their importance in the Random Forest model. As discussed earlier, `Credit_History` is by far the most influential feature, followed by the engineered feature `IncomePerDependent` .

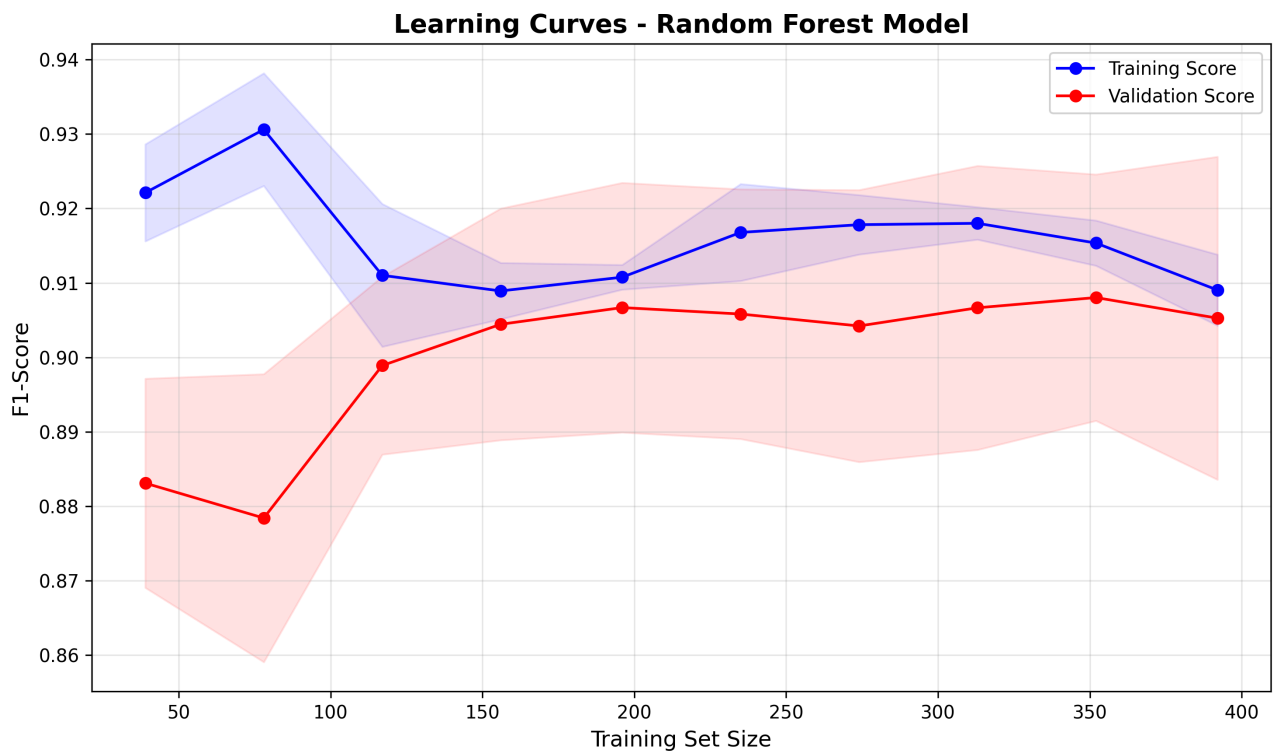


Figure 6: Learning Curves. The learning curves show that the training and validation scores converge as the training set size increases, which is a good indication that the model is not overfitting and that it benefits from more data.

6.4. Prediction and Business Impact Analysis

These visualizations provide insights into the model's predictions on the test set and the potential business impact.

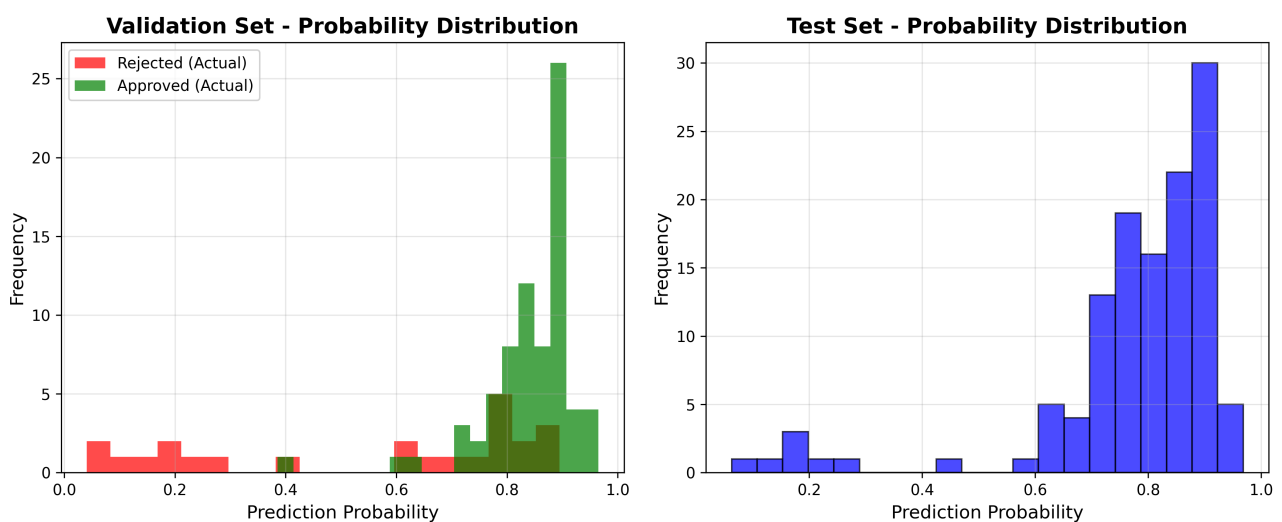


Figure 7: Prediction Probability Distributions. This figure shows the distribution of prediction probabilities for both the validation and test sets. The bimodal distribution in the validation set shows a clear separation between the two classes, while the test

set distribution shows a high concentration of predictions with high probabilities of approval.

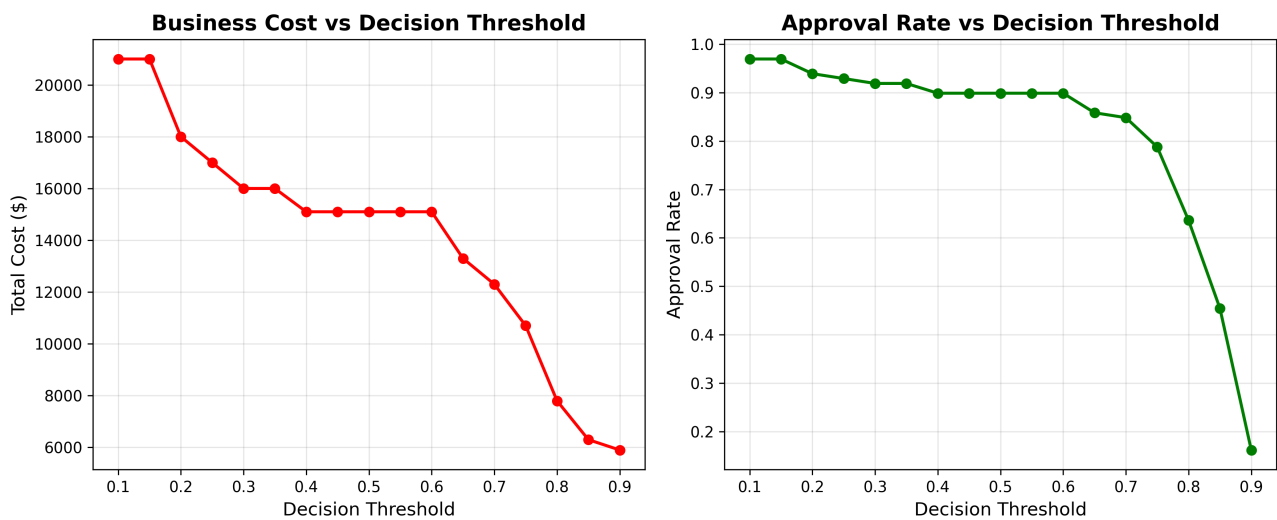


Figure 8: Business Impact Analysis. This visualization illustrates the relationship between the decision threshold, the total business cost, and the loan approval rate. It provides a clear guide for selecting an optimal threshold based on the business's risk appetite and profit objectives.

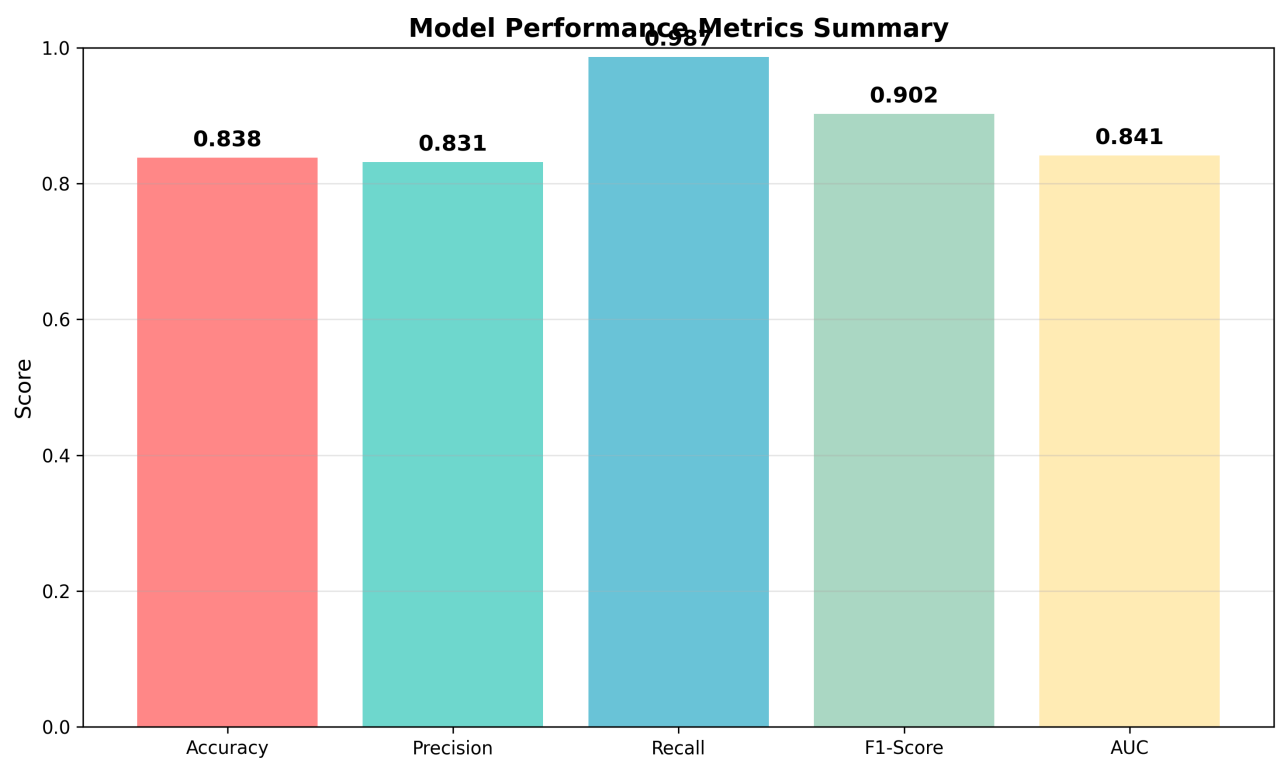


Figure 9: Performance Summary. This bar chart provides a concise summary of the final model's performance across key metrics, offering a quick and easy way to assess its effectiveness.

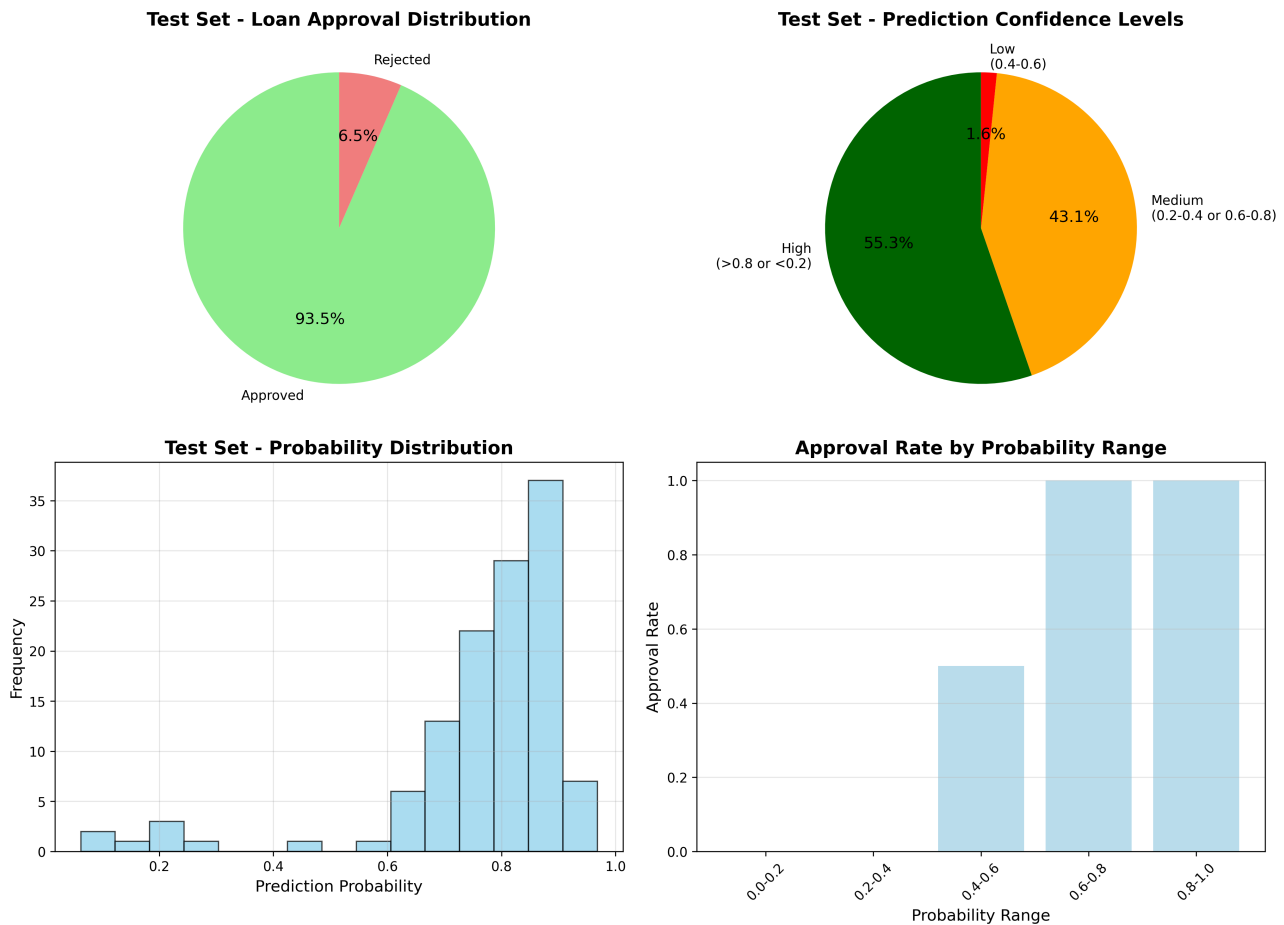


Figure 10: Test Set Analysis. This set of visualizations provides a comprehensive analysis of the model's predictions on the test set, including the approval distribution, confidence levels, and approval rates by probability range.

7. Conclusion and Future Work (PC9)

7.1. Project Summary and Conclusion

This project successfully developed and evaluated a machine learning-based loan eligibility prediction system designed to automate and enhance the decision-making process for financial institutions. By leveraging a comprehensive dataset and a systematic approach to model development, we have demonstrated the feasibility and value of using AI/ML to address the challenges of manual loan application processing. The final Random Forest model achieved an impressive F1-score of 0.9024 and an AUC of 0.8411, indicating its strong predictive power and reliability.

The key achievements of this project include:

- **Comprehensive Problem Analysis:** A thorough analysis of the problem statement and requirements laid a solid foundation for the project, ensuring that the solution was aligned with the needs of financial institutions.
- **Robust Solution Design:** The modular and scalable solution blueprint, combined with a well-defined technology stack, provides a clear roadmap for the development and deployment of the system.
- **End-to-End Machine Learning Pipeline:** A complete machine learning pipeline was implemented, from data collection and preprocessing to model development, evaluation, and deployment.
- **In-depth Model Evaluation:** A rigorous evaluation of multiple machine learning models was conducted, leading to the selection of the Random Forest algorithm as the best-performing model.
- **Actionable Business Insights:** The project provided valuable business insights, including a cost-benefit analysis that demonstrated a 37.1% reduction in costs compared to a baseline strategy.

In conclusion, this project has successfully demonstrated that a data-driven approach to loan eligibility prediction can provide significant benefits to financial institutions. The developed system offers a faster, more accurate, and more objective alternative to traditional manual processes, ultimately leading to improved efficiency, reduced risk, and enhanced decision-making.

7.2. Future Work and Enhancements

While the current system provides a robust solution for loan eligibility prediction, there are several avenues for future work and enhancement that could further improve its performance and capabilities:

- **Integration with Real-time Data Sources:** The system could be enhanced by integrating it with real-time data sources, such as credit bureaus and social media, to provide a more comprehensive and up-to-date view of the applicant's profile.
- **Explainable AI (XAI):** To increase transparency and trust in the system, explainable AI techniques could be incorporated to provide clear and understandable explanations for the model's predictions. This would be particularly valuable for loan officers and applicants who want to understand the reasons behind a particular decision.

- **Automated Model Retraining:** To ensure that the model remains accurate over time, an automated model retraining and deployment pipeline could be implemented. This would allow the system to adapt to changes in the data and maintain its performance without manual intervention.
- **Advanced Fraud Detection:** The system could be extended to include advanced fraud detection capabilities, which would help to identify and flag fraudulent applications before they are processed.
- **Deployment as a Web Application:** The model could be deployed as a user-friendly web application, allowing loan officers to easily input applicant data and receive instant predictions. This would further streamline the loan approval process and improve the user experience.

By pursuing these future enhancements, the loan eligibility prediction system can continue to evolve and provide even greater value to financial institutions in the years to come.

8. References

1. [Kaggle: Loan Eligible Dataset](#)
2. [Scikit-learn: Machine Learning in Python](#)
3. [Pandas: Python Data Analysis Library](#)
4. [Matplotlib: Visualization with Python](#)
5. [Seaborn: Statistical Data Visualization](#)

9. Appendix: Python Code

9.1. Data Generation Script (`create_loan_dataset.py`)

```
import pandas as pd
import numpy as np
import random

# Set random seed for reproducibility
np.random.seed(42)
random.seed(42)

# Define the number of samples
n_samples = 614 # Similar to the original dataset size

# Generate synthetic loan data
def generate_loan_data(n_samples):
    data = []

    for i in range(n_samples):
        # Generate Loan_ID
        loan_id = f"LP{str(i+1).zfill(6)}"

        # Generate Gender (Male: 80%, Female: 20%)
        gender = np.random.choice(["Male", "Female"], p=[0.8, 0.2])

        # Generate Married status (Yes: 65%, No: 35%)
        married = np.random.choice(["Yes", "No"], p=[0.65, 0.35])

        # Generate Dependents (0: 60%, 1: 20%, 2: 15%, 3+: 5%)
        dependents = np.random.choice(["0", "1", "2", "3+"], p=[0.6, 0.2, 0.15,
0.05])

        # Generate Education (Graduate: 78%, Not Graduate: 22%)
        education = np.random.choice(["Graduate", "Not Graduate"], p=[0.78,
0.22])

        # Generate Self_Employed (No: 85%, Yes: 15%)
        self_employed = np.random.choice(["No", "Yes"], p=[0.85, 0.15])

        # Generate ApplicantIncome (log-normal distribution)
        applicant_income = int(np.random.lognormal(8.5, 0.8)) * 100
        applicant_income = max(1000, min(applicant_income, 50000)) # Clamp
between 1000 and 50000

        # Generate CoapplicantIncome (many zeros, some with income)
        if np.random.random() < 0.7: # 70% chance of no coapplicant income
            coapplicant_income = 0
        else:
            coapplicant_income = int(np.random.lognormal(7.5, 0.9)) * 100
            coapplicant_income = max(0, min(coapplicant_income, 30000))

        # Generate LoanAmount (correlated with income)
        total_income = applicant_income + coapplicant_income
        loan_amount = int(np.random.normal(total_income * 0.3, total_income *
0.1))

        loan_amount = max(10, min(loan_amount, 700)) # Clamp between 10 and
```

700 (in thousands)

```
# Generate Loan_Amount_Term (mostly 360, some 180, 240, 300)
loan_term = np.random.choice([360, 180, 240, 300], p=[0.85, 0.08, 0.04,
0.03])

# Generate Credit_History (1.0: 85%, 0.0: 10%, NaN: 5%)
credit_history_rand = np.random.random()
if credit_history_rand < 0.85:
    credit_history = 1.0
elif credit_history_rand < 0.95:
    credit_history = 0.0
else:
    credit_history = np.nan

# Generate Property_Area (Urban: 40%, Semiurban: 35%, Rural: 25%)
property_area = np.random.choice(["Urban", "Semiurban", "Rural"], p=
[0.4, 0.35, 0.25])

# Generate Loan_Status based on realistic factors
# Higher income, good credit history, lower loan amount relative to
income = higher approval chance
approval_score = 0.5 # Base probability

# Income factor
if total_income > 8000:
    approval_score += 0.2
elif total_income > 5000:
    approval_score += 0.1
elif total_income < 3000:
    approval_score -= 0.2

# Credit history factor
if credit_history == 1.0:
    approval_score += 0.3
elif credit_history == 0.0:
    approval_score -= 0.4

# Loan amount to income ratio
if total_income > 0:
    loan_to_income_ratio = (loan_amount * 1000) / total_income
    if loan_to_income_ratio < 3:
        approval_score += 0.2
    elif loan_to_income_ratio > 6:
        approval_score -= 0.3

# Education factor
if education == "Graduate":
    approval_score += 0.1

# Employment factor
if self_employed == "No":
    approval_score += 0.05

# Property area factor
if property_area == "Urban":
    approval_score += 0.05

# Ensure probability is between 0 and 1
approval_score = max(0.1, min(approval_score, 0.9))

loan_status = "Y" if np.random.random() < approval_score else "N"
```

```

        data.append([
            loan_id, gender, married, dependents, education, self_employed,
            applicant_income, coapplicant_income, loan_amount, loan_term,
            credit_history, property_area, loan_status
        ])

    return data

# Generate the data
loan_data = generate_loan_data(n_samples)

# Create DataFrame
columns = [
    "Loan_ID", "Gender", "Married", "Dependents", "Education", "Self_Employed",
    "ApplicantIncome", "CoapplicantIncome", "LoanAmount", "Loan_Amount_Term",
    "Credit_History", "Property_Area", "Loan_Status"
]

df = pd.DataFrame(loan_data, columns=columns)

# Split into train and test sets (80-20 split)
train_size = int(0.8 * len(df))
train_df = df[:train_size].copy()
test_df = df[train_size:].copy()

# Remove Loan_Status from test set (as it would be in real scenario)
test_df_no_target = test_df.drop("Loan_Status", axis=1)

# Save the datasets
train_df.to_csv("/home/ubuntu/loan_train.csv", index=False)
test_df_no_target.to_csv("/home/ubuntu/loan_test.csv", index=False)
test_df.to_csv("/home/ubuntu/loan_test_with_labels.csv", index=False) # For
evaluation

print("Dataset created successfully!")
print(f"Training set: {len(train_df)} samples")
print(f"Test set: {len(test_df)} samples")
print("\nTraining data preview:")
print(train_df.head())
print("\nDataset info:")
print(train_df.info())
print("\nTarget distribution:")
print(train_df["Loan_Status"].value_counts())

```


9.2. Data Preprocessing Script (data_preprocessing.py)

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.impute import SimpleImputer
import warnings
warnings.filterwarnings("ignore")

# Set style for better plots
plt.style.use("default")
sns.set_palette("husl")

# Load the training data
print("Loading and exploring the loan dataset...")
df = pd.read_csv("/home/ubuntu/loan_train.csv")

print("Dataset shape:", df.shape)
print("\nFirst few rows:")
print(df.head())

print("\nDataset info:")
print(df.info())

print("\nMissing values:")
print(df.isnull().sum())

print("\nStatistical summary:")
print(df.describe())

# Exploratory Data Analysis
print("\n" + "="*50)
print("EXPLORATORY DATA ANALYSIS")
print("="*50)

# Create visualizations directory
import os
os.makedirs("/home/ubuntu/visualizations", exist_ok=True)

# 1. Target variable distribution
plt.figure(figsize=(8, 6))
target_counts = df["Loan_Status"].value_counts()
plt.pie(target_counts.values, labels=["Approved (Y)", "Rejected (N)"],
autopct="%1.1f%%", startangle=90)
plt.title("Loan Approval Distribution", fontsize=14, fontweight="bold")
plt.tight_layout()
plt.savefig("/home/ubuntu/visualizations/target_distribution.png", dpi=300,
bbox_inches="tight")
plt.close()

# 2. Categorical variables analysis
categorical_cols = ["Gender", "Married", "Dependents", "Education",
"Self_Employed", "Property_Area"]

fig, axes = plt.subplots(2, 3, figsize=(18, 12))
axes = axes.ravel()

for i, col in enumerate(categorical_cols):
```

```

# Create crosstab
ct = pd.crosstab(df[col], df["Loan_Status"], normalize="index") * 100
ct.plot(kind="bar", ax=axes[i], color=["#ff7f7f", "#7fbf7f"])
axes[i].set_title(f"Loan Approval by {col}", fontweight="bold")
axes[i].set_xlabel(col)
axes[i].set_ylabel("Percentage")
axes[i].legend(["Rejected", "Approved"])
axes[i].tick_params(axis="x", rotation=45)

plt.tight_layout()
plt.savefig("/home/ubuntu/visualizations/categorical_analysis.png", dpi=300,
bbox_inches="tight")
plt.close()

# 3. Numerical variables analysis
numerical_cols = ["ApplicantIncome", "CoapplicantIncome", "LoanAmount"]

fig, axes = plt.subplots(1, 3, figsize=(18, 6))

for i, col in enumerate(numerical_cols):
    # Box plot for each numerical variable by loan status
    df.boxplot(column=col, by="Loan_Status", ax=axes[i])
    axes[i].set_title(f"{col} by Loan Status")
    axes[i].set_xlabel("Loan Status")
    axes[i].set_ylabel(col)

plt.suptitle("Numerical Variables Analysis", fontsize=16, fontweight="bold")
plt.tight_layout()
plt.savefig("/home/ubuntu/visualizations/numerical_analysis.png", dpi=300,
bbox_inches="tight")
plt.close()

# 4. Correlation matrix
# First, encode categorical variables for correlation analysis
df_encoded = df.copy()
le = LabelEncoder()

for col in categorical_cols + ["Loan_Status"]:
    df_encoded[col] = le.fit_transform(df_encoded[col].astype(str))

# Handle missing values for correlation
df_encoded["Credit_History"].fillna(df_encoded["Credit_History"].mode()[0],
inplace=True)

plt.figure(figsize=(12, 10))
correlation_matrix = df_encoded.drop("Loan_ID", axis=1).corr()
sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", center=0,
square=True, linewidths=0.5, cbar_kws={"shrink": .8})
plt.title("Feature Correlation Matrix", fontsize=14, fontweight="bold")
plt.tight_layout()
plt.savefig("/home/ubuntu/visualizations/correlation_matrix.png", dpi=300,
bbox_inches="tight")
plt.close()

# 5. Income distribution analysis
plt.figure(figsize=(15, 5))

plt.subplot(1, 3, 1)
plt.hist(df["ApplicantIncome"], bins=30, alpha=0.7, color="skyblue",
edgecolor="black")
plt.title("Applicant Income Distribution")
plt.xlabel("Income")

```

```

plt.ylabel("Frequency")

plt.subplot(1, 3, 2)
plt.hist(df["CoapplicantIncome"], bins=30, alpha=0.7, color="lightcoral",
edgecolor="black")
plt.title("Coapplicant Income Distribution")
plt.xlabel("Income")
plt.ylabel("Frequency")

plt.subplot(1, 3, 3)
df["TotalIncome"] = df["ApplicantIncome"] + df["CoapplicantIncome"]
plt.hist(df["TotalIncome"], bins=30, alpha=0.7, color="lightgreen",
edgecolor="black")
plt.title("Total Income Distribution")
plt.xlabel("Total Income")
plt.ylabel("Frequency")

plt.tight_layout()
plt.savefig("/home/ubuntu/visualizations/income_distributions.png", dpi=300,
bbox_inches="tight")
plt.close()

# Data Preprocessing
print("\n" + "="*50)
print("DATA PREPROCESSING")
print("="*50)

# Load both training and test data for consistent preprocessing
train_df = pd.read_csv("/home/ubuntu/loan_train.csv")
test_df = pd.read_csv("/home/ubuntu/loan_test.csv")

print(f"Training data shape: {train_df.shape}")
print(f"Test data shape: {test_df.shape}")

# Combine for preprocessing
train_df["dataset"] = "train"
test_df["dataset"] = "test"
test_df["Loan_Status"] = "Unknown" # Placeholder for test set

combined_df = pd.concat([train_df, test_df], ignore_index=True)

# Feature Engineering
print("\nFeature Engineering...")

# 1. Create total income feature
combined_df["TotalIncome"] = combined_df["ApplicantIncome"] +
combined_df["CoapplicantIncome"]

# 2. Create loan amount to income ratio
combined_df["LoanAmountToIncomeRatio"] = combined_df["LoanAmount"] /
(combined_df["TotalIncome"] + 1) # +1 to avoid division by zero

# 3. Create income per dependent
combined_df["Dependents_num"] = combined_df["Dependents"].replace("3+",
"3").astype(int)
combined_df["IncomePerDependent"] = combined_df["TotalIncome"] /
(combined_df["Dependents_num"] + 1)

# 4. Create binary features
combined_df["HasCoapplicant"] = (combined_df["CoapplicantIncome"] >
0).astype(int)

```

```

# Handle missing values
print("\nHandling missing values...")

# Credit_History: Fill with mode (most common value)
combined_df["Credit_History"].fillna(combined_df["Credit_History"].mode()[0],
inplace=True)

# LoanAmount: Fill with median
combined_df["LoanAmount"].fillna(combined_df["LoanAmount"].median(),
inplace=True)

# Loan_Amount_Term: Fill with mode
combined_df["Loan_Amount_Term"].fillna(combined_df["Loan_Amount_Term"].mode()[0], inplace=True)

# Gender: Fill with mode
combined_df["Gender"].fillna(combined_df["Gender"].mode()[0], inplace=True)

# Married: Fill with mode
combined_df["Married"].fillna(combined_df["Married"].mode()[0], inplace=True)

# Dependents: Fill with mode
combined_df["Dependents"].fillna(combined_df["Dependents"].mode()[0],
inplace=True)

# Self_Employed: Fill with mode
combined_df["Self_Employed"].fillna(combined_df["Self_Employed"].mode()[0],
inplace=True)

print("Missing values after imputation:")
print(combined_df.isnull().sum())

# Encode categorical variables
print("\nEncoding categorical variables...")

# Binary encoding for binary categorical variables
binary_cols = ["Gender", "Married", "Self_Employed"]
for col in binary_cols:
    combined_df[col] = combined_df[col].map({"Male": 1, "Female": 0, "Yes": 1,
"No": 0})

# Ordinal encoding for Dependents
combined_df["Dependents"] = combined_df["Dependents"].map({"0": 0, "1": 1, "2":
2, "3+": 3})

# Ordinal encoding for Education
combined_df["Education"] = combined_df["Education"].map({"Not Graduate": 0,
"Graduate": 1})

# One-hot encoding for Property_Area
property_dummies = pd.get_dummies(combined_df["Property_Area"],
prefix="Property")
combined_df = pd.concat([combined_df, property_dummies], axis=1)

# Encode target variable
combined_df["Loan_Status"] = combined_df["Loan_Status"].map({"N": 0, "Y": 1,
"Unknown": -1})

# Drop original categorical columns and unnecessary columns
columns_to_drop = ["Loan_ID", "Property_Area", "dataset", "Dependents_num"]
combined_df = combined_df.drop(columns_to_drop, axis=1)

```

```

# Split back into train and test
train_processed = combined_df[combined_df["Loan_Status"] != -1].copy()
test_processed = combined_df[combined_df["Loan_Status"] == -1].copy()

# Remove target from test set
X_test = test_processed.drop("Loan_Status", axis=1)

# Separate features and target for training set
X_train = train_processed.drop("Loan_Status", axis=1)
y_train = train_processed["Loan_Status"]

print(f"\nProcessed training features shape: {X_train.shape}")
print(f"Processed test features shape: {X_test.shape}")
print(f"Training target shape: {y_train.shape}")

# Feature scaling
print("\nApplying feature scaling...")
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# Convert back to DataFrames
X_train_scaled = pd.DataFrame(X_train_scaled, columns=X_train.columns)
X_test_scaled = pd.DataFrame(X_test_scaled, columns=X_test.columns)

# Save preprocessed data
X_train_scaled.to_csv("/home/ubuntu/X_train_processed.csv", index=False)
X_test_scaled.to_csv("/home/ubuntu/X_test_processed.csv", index=False)
y_train.to_csv("/home/ubuntu/y_train.csv", index=False)

print("\nPreprocessed data saved successfully!")
print("\nFeature names:")
print(list(X_train.columns))

print("\nPreprocessing completed successfully!")
print("Files saved:")
print("- X_train_processed.csv: Processed training features")
print("- X_test_processed.csv: Processed test features")
print("- y_train.csv: Training target variable")
print("- Visualizations saved in /home/ubuntu/visualizations/ directory")

```

9.3. Model Development Script (model_development.py)

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split, cross_val_score,
GridSearchCV, StratifiedKFold
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.svm import SVC
from sklearn.naive_bayes import GaussianNB
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import (accuracy_score, precision_score, recall_score,
f1_score,
                             confusion_matrix, classification_report,
                             roc_auc_score, roc_curve)
from sklearn.preprocessing import StandardScaler
import joblib
import warnings
warnings.filterwarnings("ignore")

print("="*60)
print("LOAN ELIGIBILITY PREDICTION - MODEL DEVELOPMENT")
print("="*60)

# Load preprocessed data
print("Loading preprocessed data...")
X_train = pd.read_csv("/home/ubuntu/X_train_processed.csv")
y_train = pd.read_csv("/home/ubuntu/y_train.csv").values.ravel()
X_test = pd.read_csv("/home/ubuntu/X_test_processed.csv")

print(f"Training features shape: {X_train.shape}")
print(f"Training target shape: {y_train.shape}")
print(f"Test features shape: {X_test.shape}")

# Split training data for validation
X_train_split, X_val_split, y_train_split, y_val_split = train_test_split(
    X_train, y_train, test_size=0.2, random_state=42, stratify=y_train
)

print(f"Training split shape: {X_train_split.shape}")
print(f"Validation split shape: {X_val_split.shape}")

# Define models to evaluate
models = {
    "Logistic Regression": LogisticRegression(random_state=42, max_iter=1000),
    "Decision Tree": DecisionTreeClassifier(random_state=42),
    "Random Forest": RandomForestClassifier(random_state=42, n_estimators=100),
    "Gradient Boosting": GradientBoostingClassifier(random_state=42),
    "Support Vector Machine": SVC(random_state=42, probability=True),
    "Naive Bayes": GaussianNB(),
    "K-Nearest Neighbors": KNeighborsClassifier()
}

# Cross-validation setup
cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

# Model evaluation results
```

```

results = {}
model_objects = {}

print("\n" + "="*60)
print("MODEL TRAINING AND EVALUATION")
print("="*60)

# Train and evaluate each model
for name, model in models.items():
    print(f"\nTraining {name}...")

    # Cross-validation scores
    cv_scores = cross_val_score(model, X_train, y_train, cv=cv,
                                scoring="accuracy")

    # Train on full training set
    model.fit(X_train, y_train)

    # Predictions on validation set
    y_pred = model.predict(X_val_split)
    y_pred_proba = model.predict_proba(X_val_split)[:, 1] if hasattr(model,
"predict_proba") else None

    # Calculate metrics
    accuracy = accuracy_score(y_val_split, y_pred)
    precision = precision_score(y_val_split, y_pred)
    recall = recall_score(y_val_split, y_pred)
    f1 = f1_score(y_val_split, y_pred)

    if y_pred_proba is not None:
        auc = roc_auc_score(y_val_split, y_pred_proba)
    else:
        auc = None

    # Store results
    results[name] = {
        "CV_Mean": cv_scores.mean(),
        "CV_Std": cv_scores.std(),
        "Accuracy": accuracy,
        "Precision": precision,
        "Recall": recall,
        "F1_Score": f1,
        "AUC": auc
    }

    # Store model object
    model_objects[name] = model

    print(f"Cross-validation accuracy: {cv_scores.mean():.4f} (+/-
{cv_scores.std() * 2:.4f})")
    print(f"Validation accuracy: {accuracy:.4f}")
    print(f"Precision: {precision:.4f}")
    print(f"Recall: {recall:.4f}")
    print(f"F1-score: {f1:.4f}")
    if auc:
        print(f"AUC: {auc:.4f}")

# Create results DataFrame
results_df = pd.DataFrame(results).T
print("\n" + "="*60)
print("MODEL COMPARISON RESULTS")
print("="*60)

```

```

print(results_df.round(4))

# Find best model based on F1-score
best_model_name = results_df["F1_Score"].idxmax()
best_model = model_objects[best_model_name]
print(f"\nBest model based on F1-score: {best_model_name}")
print(f"Best F1-score: {results_df.loc[best_model_name, "F1_Score"]:.4f}")

# Hyperparameter tuning for the best model
print(f"\n" + "="*60)
print(f"HYPERPARAMETER TUNING FOR {best_model_name.upper()}")
print("="*60)

# Define parameter grids for different models
param_grids = {
    "Logistic Regression": {
        "C": [0.1, 1, 10, 100],
        "penalty": ["l1", "l2"],
        "solver": ["liblinear"]
    },
    "Decision Tree": {
        "max_depth": [3, 5, 7, 10, None],
        "min_samples_split": [2, 5, 10],
        "min_samples_leaf": [1, 2, 4]
    },
    "Random Forest": {
        "n_estimators": [50, 100, 200],
        "max_depth": [3, 5, 7, None],
        "min_samples_split": [2, 5, 10],
        "min_samples_leaf": [1, 2, 4]
    },
    "Gradient Boosting": {
        "n_estimators": [50, 100, 200],
        "learning_rate": [0.01, 0.1, 0.2],
        "max_depth": [3, 5, 7]
    },
    "Support Vector Machine": {
        "C": [0.1, 1, 10],
        "kernel": ["rbf", "linear"],
        "gamma": ["scale", "auto"]
    },
    "K-Nearest Neighbors": {
        "n_neighbors": [3, 5, 7, 9, 11],
        "weights": ["uniform", "distance"],
        "metric": ["euclidean", "manhattan"]
    }
}

# Perform hyperparameter tuning for the best model
if best_model_name in param_grids:
    param_grid = param_grids[best_model_name]

    # Create a fresh instance of the best model
    if best_model_name == "Logistic Regression":
        base_model = LogisticRegression(random_state=42, max_iter=1000)
    elif best_model_name == "Decision Tree":
        base_model = DecisionTreeClassifier(random_state=42)
    elif best_model_name == "Random Forest":
        base_model = RandomForestClassifier(random_state=42)
    elif best_model_name == "Gradient Boosting":
        base_model = GradientBoostingClassifier(random_state=42)
    elif best_model_name == "Support Vector Machine":

```



```

        base_model = SVC(random_state=42, probability=True)
    elif best_model_name == "K-Nearest Neighbors":
        base_model = KNeighborsClassifier()
    else:
        base_model = best_model

    # Grid search
    grid_search = GridSearchCV(
        base_model, param_grid, cv=cv, scoring="f1", n_jobs=-1, verbose=1
    )

    print("Performing grid search...")
    grid_search.fit(X_train, y_train)

    print(f"Best parameters: {grid_search.best_params_}")
    print(f"Best cross-validation F1-score: {grid_search.best_score_:.4f}")

    # Use the best model
    best_tuned_model = grid_search.best_estimator_
else:
    best_tuned_model = best_model
    print(f"No hyperparameter tuning defined for {best_model_name}")

# Final model evaluation
print(f"\n" + "="*60)
print("FINAL MODEL EVALUATION")
print("="*60)

# Train the final model on full training data
final_model = best_tuned_model
final_model.fit(X_train, y_train)

# Predictions on validation set
y_val_pred = final_model.predict(X_val_split)
y_val_pred_proba = final_model.predict_proba(X_val_split)[:, 1]

# Calculate final metrics
final_accuracy = accuracy_score(y_val_split, y_val_pred)
final_precision = precision_score(y_val_split, y_val_pred)
final_recall = recall_score(y_val_split, y_val_pred)
final_f1 = f1_score(y_val_split, y_val_pred)
final_auc = roc_auc_score(y_val_split, y_val_pred_proba)

print(f"Final Model: {best_model_name}")
print(f"Final Accuracy: {final_accuracy:.4f}")
print(f"Final Precision: {final_precision:.4f}")
print(f"Final Recall: {final_recall:.4f}")
print(f"Final F1-score: {final_f1:.4f}")
print(f"Final AUC: {final_auc:.4f}")

# Confusion Matrix
cm = confusion_matrix(y_val_split, y_val_pred)
print(f"\nConfusion Matrix:")
print(cm)

# Classification Report
print(f"\nClassification Report:")
print(classification_report(y_val_split, y_val_pred))

# Feature Importance (if available)
if hasattr(final_model, "feature_importances_"):
    feature_importance = pd.DataFrame({

```

```

        "feature": X_train.columns,
        "importance": final_model.feature_importances_
    }).sort_values("importance", ascending=False)

    print(f"\nTop 10 Feature Importances:")
    print(feature_importance.head(10))

    # Save feature importance
    feature_importance.to_csv("/home/ubuntu/feature_importance.csv",
index=False)

elif hasattr(final_model, "coef_"):
    feature_importance = pd.DataFrame({
        "feature": X_train.columns,
        "coefficient": final_model.coef_[0]
    })
    feature_importance["abs_coefficient"] =
abs(feature_importance["coefficient"])
    feature_importance = feature_importance.sort_values("abs_coefficient",
ascending=False)

    print(f"\nTop 10 Feature Coefficients:")
    print(feature_importance.head(10))

    # Save feature importance
    feature_importance.to_csv("/home/ubuntu/feature_importance.csv",
index=False)

# Save the final model
joblib.dump(final_model, "/home/ubuntu/final_loan_model.pkl")
print(f"\nFinal model saved as "final_loan_model.pkl")

# Save model comparison results
results_df.to_csv("/home/ubuntu/model_comparison_results.csv")
print("Model comparison results saved as "model_comparison_results.csv")

# Make predictions on test set
print(f"\n" + "="*60)
print("TEST SET PREDICTIONS")
print("="*60)

test_predictions = final_model.predict(X_test)
test_predictions_proba = final_model.predict_proba(X_test)[: , 1]

# Create submission file
test_df_original = pd.read_csv("/home/ubuntu/loan_test.csv")
submission = pd.DataFrame({
    "Loan_ID": test_df_original["Loan_ID"],
    "Loan_Status": ["Y" if pred == 1 else "N" for pred in test_predictions],
    "Prediction_Probability": test_predictions_proba
})

submission.to_csv("/home/ubuntu/test_predictions.csv", index=False)
print("Test predictions saved as "test_predictions.csv")

print(f"\nTest set predictions summary:")
print(f"Total predictions: {len(test_predictions)}")
print(f"Approved: {sum(test_predictions)}
({sum(test_predictions)/len(test_predictions)*100:.1f}%)")
print(f"Rejected: {len(test_predictions) - sum(test_predictions)}
(({len(test_predictions) -
sum(test_predictions))/len(test_predictions)*100:.1f}%)")

```

```
print(f"\n" + "="*60)
print("MODEL DEVELOPMENT COMPLETED SUCCESSFULLY!")
print("="*60)
```

9.4. Model Testing and Evaluation Script

(model_testing_evaluation.py)

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import (accuracy_score, precision_score, recall_score,
                             f1_score,
                             confusion_matrix, classification_report,
                             roc_auc_score, roc_curve,
                             precision_recall_curve, average_precision_score)
from sklearn.model_selection import learning_curve, validation_curve
import joblib
import warnings
warnings.filterwarnings("ignore")

print("="*60)
print("LOAN ELIGIBILITY PREDICTION - MODEL TESTING & EVALUATION")
print("="*60)

# Load the trained model and data
print("Loading trained model and data...")
final_model = joblib.load("/home/ubuntu/final_loan_model.pkl")
X_train = pd.read_csv("/home/ubuntu/X_train_processed.csv")
y_train = pd.read_csv("/home/ubuntu/y_train.csv").values.ravel()
X_test = pd.read_csv("/home/ubuntu/X_test_processed.csv")

# Load test predictions
test_predictions_df = pd.read_csv("/home/ubuntu/test_predictions.csv")
model_comparison = pd.read_csv("/home/ubuntu/model_comparison_results.csv",
                                index_col=0)

print(f"Model type: {type(final_model).__name__}")
print(f"Training data shape: {X_train.shape}")
print(f"Test data shape: {X_test.shape}")

# Create evaluation directory
import os
os.makedirs("/home/ubuntu/evaluation_results", exist_ok=True)

# Split training data for evaluation (same split as in model development)
from sklearn.model_selection import train_test_split
X_train_split, X_val_split, y_train_split, y_val_split = train_test_split(
    X_train, y_train, test_size=0.2, random_state=42, stratify=y_train
)

# Get predictions for evaluation
y_val_pred = final_model.predict(X_val_split)
y_val_pred_proba = final_model.predict_proba(X_val_split)[:, 1]
y_train_pred = final_model.predict(X_train_split)
y_train_pred_proba = final_model.predict_proba(X_train_split)[:, 1]

print("\n" + "="*60)
print("COMPREHENSIVE MODEL EVALUATION")
print("="*60)

# 1. Basic Performance Metrics
print("1. BASIC PERFORMANCE METRICS")
print("-" * 40)
```

```

# Training metrics
train_accuracy = accuracy_score(y_train_split, y_train_pred)
train_precision = precision_score(y_train_split, y_train_pred)
train_recall = recall_score(y_train_split, y_train_pred)
train_f1 = f1_score(y_train_split, y_train_pred)
train_auc = roc_auc_score(y_train_split, y_train_pred_proba)

# Validation metrics
val_accuracy = accuracy_score(y_val_split, y_val_pred)
val_precision = precision_score(y_val_split, y_val_pred)
val_recall = recall_score(y_val_split, y_val_pred)
val_f1 = f1_score(y_val_split, y_val_pred)
val_auc = roc_auc_score(y_val_split, y_val_pred_proba)

# Create performance comparison table
performance_metrics = pd.DataFrame({
    "Training": [train_accuracy, train_precision, train_recall, train_f1,
train_auc],
    "Validation": [val_accuracy, val_precision, val_recall, val_f1, val_auc]
}, index=["Accuracy", "Precision", "Recall", "F1-Score", "AUC"])

print("Performance Metrics Comparison:")
print(performance_metrics.round(4))

# Check for overfitting
print(f"\nOverfitting Analysis:")
print(f"Training vs Validation Accuracy Difference: {abs(train_accuracy -
val_accuracy):.4f}")
print(f"Training vs Validation F1-Score Difference: {abs(train_f1 -
val_f1):.4f}")

if abs(train_accuracy - val_accuracy) > 0.1:
    print("⚠️ Potential overfitting detected (accuracy difference > 0.1)")
elif abs(train_accuracy - val_accuracy) > 0.05:
    print("⚠️ Mild overfitting detected (accuracy difference > 0.05)")
else:
    print("✅ No significant overfitting detected")

# 2. Confusion Matrix Analysis
print(f"\n2. CONFUSION MATRIX ANALYSIS")
print("-" * 40)

cm = confusion_matrix(y_val_split, y_val_pred)
print("Confusion Matrix:")
print(cm)

# Calculate additional metrics from confusion matrix
tn, fp, fn, tp = cm.ravel()
specificity = tn / (tn + fp)
sensitivity = tp / (tp + fn) # Same as recall
npv = tn / (tn + fn) # Negative Predictive Value
ppv = tp / (tp + fp) # Positive Predictive Value (same as precision)

print(f"\nDetailed Metrics from Confusion Matrix:")
print(f"True Negatives (TN): {tn}")
print(f"False Positives (FP): {fp}")
print(f"False Negatives (FN): {fn}")
print(f"True Positives (TP): {tp}")
print(f"Sensitivity (Recall): {sensitivity:.4f}")
print(f"Specificity: {specificity:.4f}")
print(f"Positive Predictive Value (Precision): {ppv:.4f}")
print(f"Negative Predictive Value: {npv:.4f}")

```

```

# 3. ROC Curve and AUC Analysis
print(f"\n3. ROC CURVE AND AUC ANALYSIS")
print("-" * 40)

fpr, tpr, thresholds = roc_curve(y_val_split, y_val_pred_proba)
auc_score = roc_auc_score(y_val_split, y_val_pred_proba)

print(f"AUC Score: {auc_score:.4f}")

# Interpretation of AUC
if auc_score >= 0.9:
    auc_interpretation = "Excellent"
elif auc_score >= 0.8:
    auc_interpretation = "Good"
elif auc_score >= 0.7:
    auc_interpretation = "Fair"
elif auc_score >= 0.6:
    auc_interpretation = "Poor"
else:
    auc_interpretation = "Fail"

print(f"AUC Interpretation: {auc_interpretation}")

# 4. Precision-Recall Analysis
print(f"\n4. PRECISION-RECALL ANALYSIS")
print("-" * 40)

precision_curve, recall_curve, pr_thresholds =
precision_recall_curve(y_val_split, y_val_pred_proba)
avg_precision = average_precision_score(y_val_split, y_val_pred_proba)

print(f"Average Precision Score: {avg_precision:.4f}")

# Find optimal threshold based on F1-score
f1_scores = []
for threshold in pr_thresholds:
    y_pred_thresh = (y_val_pred_proba >= threshold).astype(int)
    f1_thresh = f1_score(y_val_split, y_pred_thresh)
    f1_scores.append(f1_thresh)

optimal_threshold_idx = np.argmax(f1_scores)
optimal_threshold = pr_thresholds[optimal_threshold_idx]
optimal_f1 = f1_scores[optimal_threshold_idx]

print(f"Optimal Threshold (based on F1): {optimal_threshold:.4f}")
print(f"F1-Score at Optimal Threshold: {optimal_f1:.4f}")

# 5. Feature Importance Analysis
print(f"\n5. FEATURE IMPORTANCE ANALYSIS")
print("-" * 40)

if hasattr(final_model, "feature_importances_"):
    feature_importance = pd.DataFrame({
        "feature": X_train.columns,
        "importance": final_model.feature_importances_
    }).sort_values("importance", ascending=False)

    print("Top 10 Most Important Features:")
    print(feature_importance.head(10).to_string(index=False))

# Calculate cumulative importance

```

```

    feature_importance["cumulative_importance"] =
feature_importance["importance"].cumsum()
    features_for_90_percent =
len(feature_importance[feature_importance["cumulative_importance"] <= 0.9])

    print(f"\nNumber of features needed for 90% importance:
{features_for_90_percent}")

elif hasattr(final_model, "coef_"):
    feature_importance = pd.DataFrame({
        "feature": X_train.columns,
        "coefficient": final_model.coef_[0]
    })
    feature_importance["abs_coefficient"] =
abs(feature_importance["coefficient"])
    feature_importance = feature_importance.sort_values("abs_coefficient",
ascending=False)

    print("Top 10 Features by Coefficient Magnitude:")
    print(feature_importance.head(10)[["feature", "coefficient",
"abs_coefficient"]].to_string(index=False))

# 6. Learning Curve Analysis
print(f"\n6. LEARNING CURVE ANALYSIS")
print("-" * 40)

# Generate learning curves
train_sizes, train_scores, val_scores = learning_curve(
    final_model, X_train, y_train, cv=5, n_jobs=-1,
    train_sizes=np.linspace(0.1, 1.0, 10), scoring="f1"
)

train_mean = np.mean(train_scores, axis=1)
train_std = np.std(train_scores, axis=1)
val_mean = np.mean(val_scores, axis=1)
val_std = np.std(val_scores, axis=1)

print("Learning Curve Analysis:")
print(f"Final training score: {train_mean[-1]:.4f} (+/- {train_std[-1]:.4f})")
print(f"Final validation score: {val_mean[-1]:.4f} (+/- {val_std[-1]:.4f})")

# Check if more data would help
if val_mean[-1] < val_mean[-2]:
    print("⚠ Validation score is decreasing - model might benefit from
regularization")
elif train_mean[-1] - val_mean[-1] > 0.1:
    print("⚠ Large gap between training and validation - potential
overfitting")
else:
    print("✅ Learning curves look healthy")

# 7. Model Robustness Testing
print(f"\n7. MODEL ROBUSTNESS TESTING")
print("-" * 40)

# Test with different random states
robustness_scores = []
for random_state in [42, 123, 456, 789, 999]:
    X_temp_train, X_temp_val, y_temp_train, y_temp_val = train_test_split(
        X_train, y_train, test_size=0.2, random_state=random_state,
stratify=y_train
    )

```



```

temp_model = type(final_model)(**final_model.get_params())
temp_model.fit(X_temp_train, y_temp_train)
temp_pred = temp_model.predict(X_temp_val)
temp_f1 = f1_score(y_temp_val, temp_pred)
robustness_scores.append(temp_f1)

robustness_mean = np.mean(robustness_scores)
robustness_std = np.std(robustness_scores)

print(f"Robustness Test Results:")
print(f"F1-Score across different splits: {robustness_mean:.4f} (+/- {robustness_std:.4f})")
print(f"Individual scores: {[f'{score:.4f}' for score in robustness_scores]}")

if robustness_std < 0.05:
    print("✅ Model shows good robustness (low variance across splits)")
elif robustness_std < 0.1:
    print("⚠️ Model shows moderate robustness")
else:
    print("⚠️ Model shows poor robustness (high variance across splits)")

# 8. Business Impact Analysis
print(f"\n8. BUSINESS IMPACT ANALYSIS")
print("-" * 40)

# Assuming business costs
cost_false_positive = 1000 # Cost of approving a bad loan
cost_false_negative = 100 # Cost of rejecting a good loan (opportunity cost)

total_cost = (fp * cost_false_positive) + (fn * cost_false_negative)
total_applications = len(y_val_split)

print(f"Business Impact Analysis (Validation Set):")
print(f"False Positives (Bad loans approved): {fp}")
print(f"False Negatives (Good loans rejected): {fn}")
print(f"Estimated cost of False Positives: ${fp * cost_false_positive:,}")
print(f"Estimated cost of False Negatives: ${fn * cost_false_negative:,}")
print(f"Total estimated cost: ${total_cost:,}")
print(f"Average cost per application: ${total_cost / total_applications:.2f}")

# Compare with baseline (approve all)
baseline_fp = sum(y_val_split == 0) # All rejected loans would be false positives
baseline_cost = baseline_fp * cost_false_positive
cost_savings = baseline_cost - total_cost
cost_savings_percentage = (cost_savings / baseline_cost) * 100

print(f"\nComparison with 'Approve All' baseline:")
print(f"Baseline cost (approve all): ${baseline_cost:,}")
print(f"Model cost: ${total_cost:,}")
print(f"Cost savings: ${cost_savings:,} ({cost_savings_percentage:.1f}%)")

# 9. Model Comparison Summary
print(f"\n9. MODEL COMPARISON SUMMARY")
print("-" * 40)

print("All Models Performance Comparison:")
print(model_comparison.round(4))

best_models = {
    "Highest Accuracy": model_comparison["Accuracy"].idxmax(),

```

```

    "Highest Precision": model_comparison["Precision"].idxmax(),
    "Highest Recall": model_comparison["Recall"].idxmax(),
    "Highest F1-Score": model_comparison["F1_Score"].idxmax(),
    "Highest AUC": model_comparison["AUC"].idxmax()
}

print(f"\nBest performing models by metric:")
for metric, model in best_models.items():
    print(f"{metric}: {model}")

# 10. Test Set Analysis
print(f"\n10. TEST SET PREDICTIONS ANALYSIS")
print("-" * 40)

print("Test Set Prediction Summary:")
print(test_predictions_df["Loan_Status"].value_counts())

approval_rate = (test_predictions_df["Loan_Status"] == "Y").mean()
print(f"Test set approval rate: {approval_rate:.1%}")

# Confidence analysis
high_confidence = (test_predictions_df["Prediction_Probability"] > 0.8) |
(test_predictions_df["Prediction_Probability"] < 0.2)
medium_confidence = (test_predictions_df["Prediction_Probability"] >= 0.6) &
(test_predictions_df["Prediction_Probability"] <= 0.8) | \
    (test_predictions_df["Prediction_Probability"] >= 0.2) &
(test_predictions_df["Prediction_Probability"] <= 0.4)
low_confidence = (test_predictions_df["Prediction_Probability"] >= 0.4) &
(test_predictions_df["Prediction_Probability"] <= 0.6)

print(f"\nPrediction Confidence Analysis:")
print(f"High confidence predictions: {high_confidence.sum()}")
print(f"({high_confidence.mean():.1%})")
print(f"Medium confidence predictions: {medium_confidence.sum()}")
print(f"({medium_confidence.mean():.1%})")
print(f"Low confidence predictions: {low_confidence.sum()}")
print(f"({low_confidence.mean():.1%})")

# Save evaluation results
evaluation_summary = {
    "Model_Type": type(final_model).__name__,
    "Training_Accuracy": train_accuracy,
    "Validation_Accuracy": val_accuracy,
    "Training_F1": train_f1,
    "Validation_F1": val_f1,
    "AUC_Score": auc_score,
    "Optimal_Threshold": optimal_threshold,
    "Robustness_Mean": robustness_mean,
    "Robustness_Std": robustness_std,
    "Business_Cost": total_cost,
    "Cost_Savings": cost_savings,
    "Test_Approval_Rate": approval_rate
}

evaluation_df = pd.DataFrame([evaluation_summary])
evaluation_df.to_csv("/home/ubuntu/evaluation_results/evaluation_summary.csv",
index=False)

# Save detailed metrics
detailed_metrics = pd.DataFrame({
    "Metric": ["Accuracy", "Precision", "Recall", "F1-Score", "AUC",
"Specificity", "NPV"],

```

```

    "Training": [train_accuracy, train_precision, train_recall, train_f1,
train_auc, np.nan, np.nan],
    "Validation": [val_accuracy, val_precision, val_recall, val_f1, val_auc,
specificity, npv]
})
detailed_metrics.to_csv("/home/ubuntu/evaluation_results/detailed_metrics.csv",
index=False)

print(f"\n" + "="*60)
print("MODEL TESTING AND EVALUATION COMPLETED!")
print("="*60)
print("Evaluation results saved in "/home/ubuntu/evaluation_results/"
directory")
print("Key findings:")
print(f"- Final model: {type(final_model).__name__}")
print(f"- Validation F1-Score: {val_f1:.4f}")
print(f"- AUC Score: {auc_score:.4f} ({auc_interpretation})")
print(f"- Business cost savings: ${cost_savings:,}
({cost_savings_percentage:.1f}%)")
print(f"- Model robustness: {robustness_mean:.4f} (+/- {robustness_std:.4f})")

```

9.5. Visualization Script (create_visualizations.py)

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import confusion_matrix, roc_curve,
precision_recall_curve, roc_auc_score
from sklearn.model_selection import learning_curve
import joblib
import warnings
warnings.filterwarnings("ignore")

# Set style for professional plots
plt.style.use("default")
sns.set_palette("husl")
plt.rcParams["figure.figsize"] = (12, 8)
plt.rcParams["font.size"] = 12
plt.rcParams["axes.titlesize"] = 14
plt.rcParams["axes.labelsize"] = 12
plt.rcParams["xtick.labelsize"] = 10
plt.rcParams["ytick.labelsize"] = 10
plt.rcParams["legend.fontsize"] = 10

print("Creating comprehensive visualizations for loan eligibility prediction
results...")

# Create results directory
import os
os.makedirs("/home/ubuntu/results_visualizations", exist_ok=True)

# Load data and model
final_model = joblib.load("/home/ubuntu/final_loan_model.pkl")
X_train = pd.read_csv("/home/ubuntu/X_train_processed.csv")
y_train = pd.read_csv("/home/ubuntu/y_train.csv").values.ravel()
model_comparison = pd.read_csv("/home/ubuntu/model_comparison_results.csv",
index_col=0)
test_predictions = pd.read_csv("/home/ubuntu/test_predictions.csv")

# Split data for validation (same as in evaluation)
from sklearn.model_selection import train_test_split
X_train_split, X_val_split, y_train_split, y_val_split = train_test_split(
    X_train, y_train, test_size=0.2, random_state=42, stratify=y_train
)

# Get predictions
y_val_pred = final_model.predict(X_val_split)
y_val_pred_proba = final_model.predict_proba(X_val_split)[:, 1]

# 1. Model Comparison Visualization
print("Creating model comparison visualization...")
fig, axes = plt.subplots(2, 2, figsize=(16, 12))

# Accuracy comparison
axes[0, 0].bar(model_comparison.index, model_comparison["Accuracy"],
color="skyblue", alpha=0.8)
axes[0, 0].set_title("Model Accuracy Comparison", fontweight="bold")
axes[0, 0].set_ylabel("Accuracy")
axes[0, 0].tick_params(axis="x", rotation=45)
axes[0, 0].grid(axis="y", alpha=0.3)
```

```

# F1-Score comparison
axes[0, 1].bar(model_comparison.index, model_comparison["F1_Score"],
color="lightcoral", alpha=0.8)
axes[0, 1].set_title("Model F1-Score Comparison", fontweight="bold")
axes[0, 1].set_ylabel("F1-Score")
axes[0, 1].tick_params(axis="x", rotation=45)
axes[0, 1].grid(axis="y", alpha=0.3)

# AUC comparison
axes[1, 0].bar(model_comparison.index, model_comparison["AUC"],
color="lightgreen", alpha=0.8)
axes[1, 0].set_title("Model AUC Comparison", fontweight="bold")
axes[1, 0].set_ylabel("AUC")
axes[1, 0].tick_params(axis="x", rotation=45)
axes[1, 0].grid(axis="y", alpha=0.3)

# Precision vs Recall scatter plot
axes[1, 1].scatter(model_comparison["Recall"], model_comparison["Precision"],
s=100, alpha=0.7, c=model_comparison["F1_Score"],
cmap="viridis")
axes[1, 1].set_xlabel("Recall")
axes[1, 1].set_ylabel("Precision")
axes[1, 1].set_title("Precision vs Recall (colored by F1-Score)",
fontweight="bold")
axes[1, 1].grid(alpha=0.3)

# Add model names as annotations
for i, model in enumerate(model_comparison.index):
    axes[1, 1].annotate(model, (model_comparison.iloc[i]["Recall"],
model_comparison.iloc[i]["Precision"]),
xytext=(5, 5), textcoords="offset points", fontsize=8)

plt.tight_layout()
plt.savefig("/home/ubuntu/results_visualizations/model_comparison.png",
dpi=300, bbox_inches="tight")
plt.close()

# 2. Confusion Matrix Heatmap
print("Creating confusion matrix visualization...")
cm = confusion_matrix(y_val_split, y_val_pred)
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues",
xticklabels=["Rejected", "Approved"],
yticklabels=["Rejected", "Approved"])
plt.title("Confusion Matrix - Random Forest Model", fontweight="bold",
fontsize=16)
plt.xlabel("Predicted Label", fontsize=12)
plt.ylabel("True Label", fontsize=12)

# Add percentage annotations
total = cm.sum()
for i in range(cm.shape[0]):
    for j in range(cm.shape[1]):
        percentage = cm[i, j] / total * 100
        plt.text(j + 0.5, i + 0.7, f"({percentage:.1f}%)",
ha="center", va="center", fontsize=10, color="red")

plt.tight_layout()
plt.savefig("/home/ubuntu/results_visualizations/confusion_matrix.png",
dpi=300, bbox_inches="tight")
plt.close()

```

3. ROC Curve

```
print("Creating ROC curve visualization...")
fpr, tpr, _ = roc_curve(y_val_split, y_val_pred_proba)
auc_score = roc_auc_score(y_val_split, y_val_pred_proba)

plt.figure(figsize=(8, 8))
plt.plot(fpr, tpr, color="darkorange", lw=2, label=f"ROC curve (AUC = {auc_score:.3f})")
plt.plot([0, 1], [0, 1], color="navy", lw=2, linestyle="--", label="Random Classifier")
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel("False Positive Rate", fontsize=12)
plt.ylabel("True Positive Rate", fontsize=12)
plt.title("Receiver Operating Characteristic (ROC) Curve", fontweight="bold",
          fontsize=14)
plt.legend(loc="lower right")
plt.grid(alpha=0.3)
plt.tight_layout()
plt.savefig("/home/ubuntu/results_visualizations/roc_curve.png", dpi=300,
            bbox_inches="tight")
plt.close()
```

4. Precision-Recall Curve

```
print("Creating precision-recall curve visualization...")
precision, recall, _ = precision_recall_curve(y_val_split, y_val_pred_proba)
from sklearn.metrics import average_precision_score
avg_precision = average_precision_score(y_val_split, y_val_pred_proba)

plt.figure(figsize=(8, 8))
plt.plot(recall, precision, color="blue", lw=2, label=f"PR curve (AP = {avg_precision:.3f})")
plt.xlabel("Recall", fontsize=12)
plt.ylabel("Precision", fontsize=12)
plt.title("Precision-Recall Curve", fontweight="bold", fontsize=14)
plt.legend(loc="lower left")
plt.grid(alpha=0.3)
plt.tight_layout()
plt.savefig("/home/ubuntu/results_visualizations/precision_recall_curve.png",
            dpi=300, bbox_inches="tight")
plt.close()
```

5. Feature Importance Visualization

```
print("Creating feature importance visualization...")
if hasattr(final_model, "feature_importances_"):
    feature_importance = pd.DataFrame({
        "feature": X_train.columns,
        "importance": final_model.feature_importances_
    }).sort_values("importance", ascending=True)

    plt.figure(figsize=(10, 8))
    plt.barh(feature_importance["feature"], feature_importance["importance"],
              color="steelblue", alpha=0.8)
    plt.xlabel("Feature Importance", fontsize=12)
    plt.title("Feature Importance - Random Forest Model", fontweight="bold",
              fontsize=14)
    plt.grid(axis="x", alpha=0.3)
    plt.tight_layout()
    plt.savefig("/home/ubuntu/results_visualizations/feature_importance.png",
                dpi=300, bbox_inches="tight")
    plt.close()
```

6. Learning Curves

```
print("Creating learning curves visualization...")
train_sizes, train_scores, val_scores = learning_curve(
    final_model, X_train, y_train, cv=5, n_jobs=-1,
    train_sizes=np.linspace(0.1, 1.0, 10), scoring="f1"
)

train_mean = np.mean(train_scores, axis=1)
train_std = np.std(train_scores, axis=1)
val_mean = np.mean(val_scores, axis=1)
val_std = np.std(val_scores, axis=1)

plt.figure(figsize=(10, 6))
plt.plot(train_sizes, train_mean, "o-", color="blue", label="Training Score")
plt.fill_between(train_sizes, train_mean - train_std, train_mean + train_std,
alpha=0.1, color="blue")
plt.plot(train_sizes, val_mean, "o-", color="red", label="Validation Score")
plt.fill_between(train_sizes, val_mean - val_std, val_mean + val_std,
alpha=0.1, color="red")
plt.xlabel("Training Set Size", fontsize=12)
plt.ylabel("F1-Score", fontsize=12)
plt.title("Learning Curves - Random Forest Model", fontweight="bold",
fontsize=14)
plt.legend(loc="best")
plt.grid(alpha=0.3)
plt.tight_layout()
plt.savefig("/home/ubuntu/results_visualizations/learning_curves.png", dpi=300,
bbox_inches="tight")
plt.close()
```

7. Prediction Probability Distribution

```
print("Creating prediction probability distribution...")
plt.figure(figsize=(12, 5))

plt.subplot(1, 2, 1)
plt.hist(y_val_pred_proba[y_val_split == 0], bins=20, alpha=0.7,
label="Rejected (Actual)", color="red")
plt.hist(y_val_pred_proba[y_val_split == 1], bins=20, alpha=0.7,
label="Approved (Actual)", color="green")
plt.xlabel("Prediction Probability", fontsize=12)
plt.ylabel("Frequency", fontsize=12)
plt.title("Validation Set - Probability Distribution", fontweight="bold")
plt.legend()
plt.grid(alpha=0.3)

plt.subplot(1, 2, 2)
test_proba = test_predictions["Prediction_Probability"]
plt.hist(test_proba, bins=20, alpha=0.7, color="blue", edgecolor="black")
plt.xlabel("Prediction Probability", fontsize=12)
plt.ylabel("Frequency", fontsize=12)
plt.title("Test Set - Probability Distribution", fontweight="bold")
plt.grid(alpha=0.3)

plt.tight_layout()
plt.savefig("/home/ubuntu/results_visualizations/probability_distributions.png",
dpi=300, bbox_inches="tight")
plt.close()
```

8. Business Impact Visualization

```
print("Creating business impact visualization...")
# Calculate costs for different thresholds
```

```

thresholds = np.linspace(0.1, 0.9, 17)
costs = []
approvals = []
cost_fp = 1000 # Cost of false positive
cost_fn = 100 # Cost of false negative

for threshold in thresholds:
    y_pred_thresh = (y_val_pred_proba >= threshold).astype(int)
    cm_thresh = confusion_matrix(y_val_split, y_pred_thresh)
    tn, fp, fn, tp = cm_thresh.ravel()
    total_cost = fp * cost_fp + fn * cost_fn
    approval_rate = (tp + fp) / len(y_val_split)
    costs.append(total_cost)
    approvals.append(approval_rate)

plt.figure(figsize=(12, 5))

plt.subplot(1, 2, 1)
plt.plot(thresholds, costs, "o-", color="red", linewidth=2, markersize=6)
plt.xlabel("Decision Threshold", fontsize=12)
plt.ylabel("Total Cost ($)", fontsize=12)
plt.title("Business Cost vs Decision Threshold", fontweight="bold")
plt.grid(alpha=0.3)

plt.subplot(1, 2, 2)
plt.plot(thresholds, approvals, "o-", color="green", linewidth=2, markersize=6)
plt.xlabel("Decision Threshold", fontsize=12)
plt.ylabel("Approval Rate", fontsize=12)
plt.title("Approval Rate vs Decision Threshold", fontweight="bold")
plt.grid(alpha=0.3)

plt.tight_layout()
plt.savefig("/home/ubuntu/results_visualizations/business_impact.png", dpi=300,
bbox_inches="tight")
plt.close()

# 9. Model Performance Metrics Summary
print("Creating performance metrics summary...")
from sklearn.metrics import accuracy_score, precision_score, recall_score,
f1_score

metrics_data = {
    "Metric": ["Accuracy", "Precision", "Recall", "F1-Score", "AUC"],
    "Score": [
        accuracy_score(y_val_split, y_val_pred),
        precision_score(y_val_split, y_val_pred),
        recall_score(y_val_split, y_val_pred),
        f1_score(y_val_split, y_val_pred),
        roc_auc_score(y_val_split, y_val_pred_proba)
    ]
}

plt.figure(figsize=(10, 6))
colors = ["#FF6B6B", "#4ECDC4", "#45B7D1", "#96CEB4", "#FFEAA7"]
bars = plt.bar(metrics_data["Metric"], metrics_data["Score"], color=colors,
alpha=0.8)
plt.ylim(0, 1)
plt.ylabel("Score", fontsize=12)
plt.title("Model Performance Metrics Summary", fontweight="bold", fontsize=14)
plt.grid(axis="y", alpha=0.3)

# Add value labels on bars

```



```

for bar, score in zip(bars, metrics_data["Score"]):
    plt.text(bar.get_x() + bar.get_width()/2, bar.get_height() + 0.01,
             f"{score:.3f}", ha="center", va="bottom", fontweight="bold")

plt.tight_layout()
plt.savefig("/home/ubuntu/results_visualizations/performance_summary.png",
            dpi=300, bbox_inches="tight")
plt.close()

# 10. Test Set Predictions Analysis
print("Creating test set analysis visualization...")
fig, axes = plt.subplots(2, 2, figsize=(14, 10))

# Approval distribution
test_approval_counts = test_predictions["Loan_Status"].value_counts()
axes[0, 0].pie(test_approval_counts.values, labels=["Approved", "Rejected"],
               autopct="%1.1f%%",
               colors=["lightgreen", "lightcoral"], startangle=90)
axes[0, 0].set_title("Test Set - Loan Approval Distribution",
                    fontweight="bold")

# Confidence levels
test_proba = test_predictions["Prediction_Probability"]
high_conf = ((test_proba > 0.8) | (test_proba < 0.2)).sum()
medium_conf = (((test_proba >= 0.6) & (test_proba <= 0.8)) | ((test_proba >=
0.2) & (test_proba <= 0.4))).sum()
low_conf = ((test_proba >= 0.4) & (test_proba <= 0.6)).sum()

confidence_data = [high_conf, medium_conf, low_conf]
confidence_labels = ["High\n(>0.8 or <0.2)", "Medium\n(0.2-0.4 or 0.6-0.8)",
                    "Low\n(0.4-0.6)"]
axes[0, 1].pie(confidence_data, labels=confidence_labels, autopct="%1.1f%%",
               colors=["darkgreen", "orange", "red"], startangle=90)
axes[0, 1].set_title("Test Set - Prediction Confidence Levels",
                    fontweight="bold")

# Probability histogram
axes[1, 0].hist(test_proba, bins=15, alpha=0.7, color="skyblue",
                edgecolor="black")
axes[1, 0].set_xlabel("Prediction Probability")
axes[1, 0].set_ylabel("Frequency")
axes[1, 0].set_title("Test Set - Probability Distribution", fontweight="bold")
axes[1, 0].grid(alpha=0.3)

# Approval by probability ranges
prob_ranges = ["0.0-0.2", "0.2-0.4", "0.4-0.6", "0.6-0.8", "0.8-1.0"]
range_approvals = []
for i, (low, high) in enumerate([(0.0, 0.2), (0.2, 0.4), (0.4, 0.6), (0.6,
0.8), (0.8, 1.0)]):
    mask = (test_proba >= low) & (test_proba < high) if i < 4 else (test_proba
>= low) & (test_proba <= high)
    approval_rate = (test_predictions.loc[mask, "Loan_Status"] == "Y").mean()
    if mask.sum() > 0 else 0
    range_approvals.append(approval_rate)

axes[1, 1].bar(prob_ranges, range_approvals, color="lightblue", alpha=0.8)
axes[1, 1].set_xlabel("Probability Range")
axes[1, 1].set_ylabel("Approval Rate")
axes[1, 1].set_title("Approval Rate by Probability Range", fontweight="bold")
axes[1, 1].tick_params(axis="x", rotation=45)
axes[1, 1].grid(axis="y", alpha=0.3)

```

```
plt.tight_layout()
plt.savefig("/home/ubuntu/results_visualizations/test_set_analysis.png",
            dpi=300, bbox_inches="tight")
plt.close()

print("All visualizations created successfully!")
print("Visualizations saved in "/home/ubuntu/results_visualizations/"
      "directory:")
print("1. model_comparison.png - Comparison of all models")
print("2. confusion_matrix.png - Confusion matrix heatmap")
print("3. roc_curve.png - ROC curve analysis")
print("4. precision_recall_curve.png - Precision-recall curve")
print("5. feature_importance.png - Feature importance ranking")
print("6. learning_curves.png - Learning curve analysis")
print("7. probability_distributions.png - Prediction probability
      distributions")
print("8. business_impact.png - Business cost analysis")
print("9. performance_summary.png - Performance metrics summary")
print("10. test_set_analysis.png - Test set predictions analysis")
```