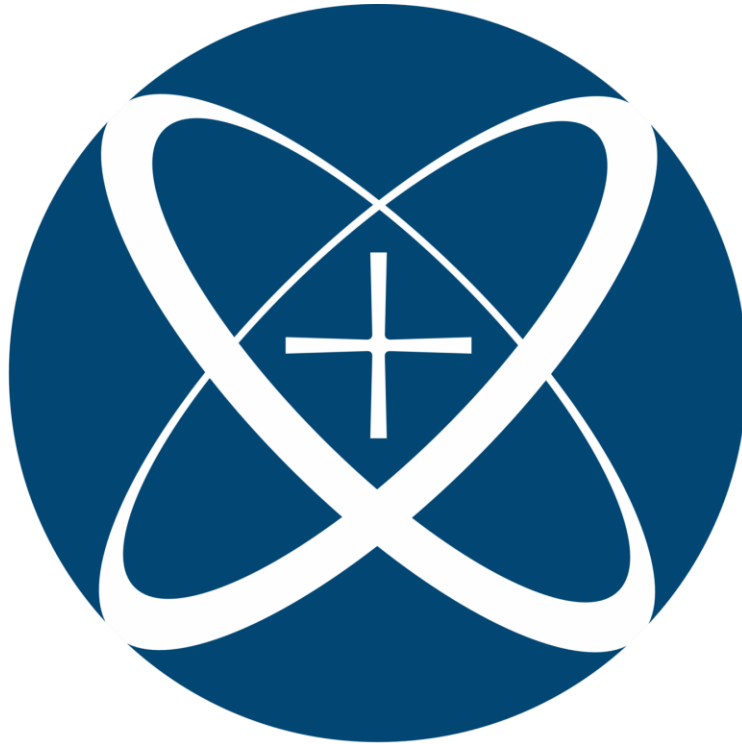


**INSTITUTO TECNOLÓGICO Y DE ESTUDIOS SUPERIORES  
DE OCCIDENTE**



**ITESO**

**CREDIT RISK MODELING FOR PERSONAL LOANS**

Predictive Analytics and Risk-Based Pricing Framework

Gian Carlo Campos Sayavedra

Modelos de Crédito

Rodolfo Slay Ramos

11/1/2025

## ABSTRACT

This project develops a comprehensive credit risk assessment framework for personal loans using machine learning techniques. Analyzing 87,889 loan applications from LendingClub (2007-2018), we built predictive models for Probability of Default (PD), Loss Given Default (LGD), and Exposure at Default (EAD) to calculate Expected Loss and enable risk-based pricing.

Key achievements include:

- Developed and compared 7 machine learning models (Logistic Regression, Random Forest, XGBoost, LightGBM, Gradient Boosting, Neural Networks, and Stacking Ensemble) for default prediction
- Best performing model: Neural Network with AUC-ROC of 0.783 and 61.5% recall at optimized threshold (0.25)
- Implemented complete credit risk framework calculating Expected Loss ( $EL = PD \times LGD \times EAD$ ) with total portfolio expected loss of \$17.8M across 13,184 test loans
- Analyzed interest rate formation in the October 2025 market environment, breaking down APR components (Federal Reserve rate: 3.75-4.00%, inflation: 3.0%, credit risk premium: 1.0-14.5%)
- Designed risk-based pricing strategy with four tiers (Excellent, Good, Fair, Poor) achieving APR ranges from 13.88% to 27.38% aligned with market benchmarks
- Validated pricing model showing strong correlation between predicted default probability and actual default rates across all risk tiers

The framework successfully demonstrates how accurate credit risk modeling enables data-driven lending decisions and sustainable profitability across risk segments. All risk tiers remain profitable with net margins ranging from 1.6% (Poor tier) to 11.4% (Excellent tier).

This integrated approach combining advanced machine learning, financial theory, and business strategy provides a complete solution for modern credit risk management in the personal loan market.

## Contents

|   |    |
|---|----|
| 1. INTRODUCTION .....                               | 5  |
| 1.1 Project Overview .....                          | 5  |
| 1.2 Objectives .....                                | 5  |
| 1.3 Credit Product Selection: Personal Loans .....  | 6  |
| 2. BUSINESS MODEL ANALYSIS .....                    | 7  |
| 2.1 Personal Loan Business Model.....               | 7  |
| 2.2 Business Model Diagram.....                     | 7  |
| 2.3 Revenue Sources and Risk Factors.....           | 8  |
| 3. INTEREST RATE FORMATION .....                    | 9  |
| 3.1 Components of Interest Rates.....               | 9  |
| 3.2 Market Data (October 2025) .....                | 10 |
| 3.3 Credit Risk Premium Determination .....         | 11 |
| 3.4 Rate Structure by Risk Tier .....               | 12 |
| 4. CREDIT RISK MODEL DEVELOPMENT .....              | 14 |
| 4.1 Methodology and Data .....                      | 14 |
| 4.2 Data Preprocessing and Feature Engineering..... | 14 |
| 4.4 Probability of Default (PD) Models .....        | 16 |
| 4.5 Loss Given Default (LGD) Estimation.....        | 19 |
| 4.6 Exposure at Default (EAD) Calculation.....      | 20 |
| 4.7 Expected Loss Framework .....                   | 20 |
| 5. RESULTS AND VALIDATION .....                     | 22 |
| 5.1 Model Performance Comparison.....               | 22 |
| 5.2 Threshold Optimization .....                    | 22 |
| 5.3 Expected Loss Analysis .....                    | 24 |
| 5.4 Risk-Based Pricing Application .....            | 25 |
| 6. BUSINESS RECOMMENDATIONS .....                   | 29 |
| 6.1 Implementation Strategy.....                    | 29 |
| 6.2 Portfolio Management .....                      | 30 |
| 6.3 Risk Monitoring.....                            | 31 |
| 7. CONCLUSIONS .....                                | 33 |
| 7.1 Key Findings .....                              | 33 |
| 7.2 Limitations.....                                | 34 |

|  |    |
|--|----|
| 7.3 Future Enhancements .....                  | 35 |
| 8. REFERENCES .....                            | 37 |
| 9. APPENDICES .....                            | 38 |
| Appendix A: Technical Notebooks .....          | 38 |
| Appendix B: Additional Technical Details ..... | 39 |

# 1. INTRODUCTION

## 1.1 Project Overview

Credit risk assessment is fundamental to the lending industry, determining both the viability of individual loans and the profitability of loan portfolios. This project develops a comprehensive credit risk modeling framework for personal loans, integrating machine learning techniques with traditional financial theory to create a data-driven approach to lending decisions and pricing strategies.

The project addresses three interconnected challenges: (1) How to accurately predict loan defaults (Probability of Default), (2) How interest rates should be structured to reflect risk and market conditions, and (3) How to implement risk-based pricing that ensures profitability across risk segments.

By combining advanced analytics with business strategy, this framework demonstrates how financial institutions can make informed lending decisions that balance risk management with competitive positioning in the personal loan market.

## 1.2 Objectives

The primary objectives of this project are:

- Develop predictive models for credit risk assessment using multiple machine learning techniques including neural networks, ensemble methods (stacking), boosting algorithms (XGBoost, LightGBM, Gradient Boosting), and bagging (Random Forest)
- Implement a complete credit risk framework calculating Probability of Default (PD), Loss Given Default (LGD), Exposure at Default (EAD), and Expected Loss (EL)
- Analyze interest rate formation in the current market environment, identifying and quantifying each component of the Annual Percentage Rate (APR)
- Design and validate a risk-based pricing strategy that aligns interest rates with predicted default probabilities
- Demonstrate early stopping techniques and threshold optimization to improve model performance
- Provide actionable business recommendations for implementation and portfolio management

### **1.3 Credit Product Selection: Personal Loans**

This project focuses on personal loans, also known as unsecured consumer loans. Personal loans are fixed-rate, fixed-term credit products typically ranging from \$5,000 to \$40,000 with repayment periods of 3 to 5 years.

#### **Key Characteristics:**

- Unsecured: No collateral required, increasing credit risk for lenders
- Fixed payments: Predictable monthly installments simplify cash flow planning
- Multipurpose: Used for debt consolidation, home improvements, medical expenses, major purchases
- Higher interest rates: Typically 7-36% APR depending on borrower creditworthiness

#### **Market Context:**

Personal loans represent a significant segment of consumer credit, with the U.S. market exceeding \$200 billion in outstanding balances. The unsecured nature of these loans makes accurate credit risk assessment particularly critical, as default rates typically range from 15-25% depending on risk tier and economic conditions.

This product was selected because:

1. Substantial historical data available for model development
2. Clear credit risk drivers that can be quantified and modeled
3. Direct relationship between risk assessment and pricing decisions
4. Significant business impact from improved prediction accuracy

## 2. BUSINESS MODEL ANALYSIS

### 2.1 Personal Loan Business Model

The personal loan business model operates through a structured process from borrower application to final outcome. The lender's revenue comes primarily from interest payments on performing loans, while losses occur from defaults where borrowers fail to repay.

**The business flow consists of six key stages:**

1. Borrower Application: Individuals seeking funds submit applications with credit history, income verification, and loan amount requested
2. Credit Assessment: Risk models evaluate Probability of Default, Loss Given Default, and calculate Expected Loss
3. Approval Decision: Applications are approved or denied based on risk tolerance thresholds
4. Loan Funding: Approved loans are disbursed with risk-based APR (7-36%), terms of 3-5 years, and monthly payment schedules
5. Repayment or Default: Approximately 80% of loans are fully repaid generating revenue; 20% default resulting in losses
6. Financial Outcome: Net profit equals total interest revenue minus credit losses and operating expenses

### 2.2 Business Model Diagram

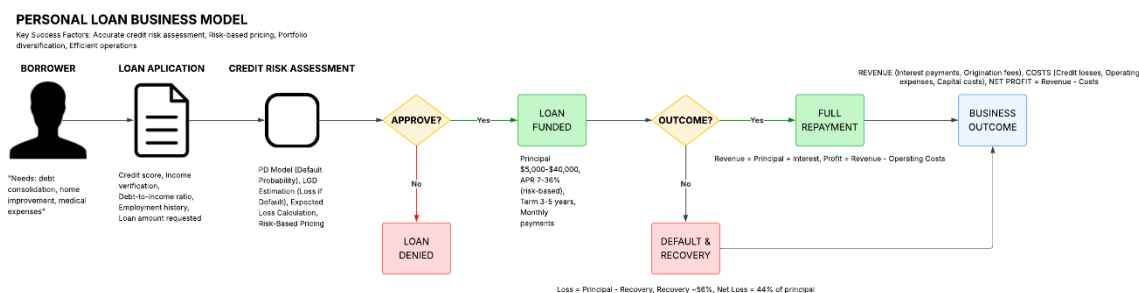


Figure 1: Personal Loan Business Model Flow

## **2.3 Revenue Sources and Risk Factors**

### **Revenue Sources:**

- Interest Payments: Primary revenue stream from APR charged on outstanding principal (7-36% based on risk tier)
- Origination Fees: Upfront fees typically 1-6% of loan amount
- Late Payment Fees: Additional revenue from borrowers missing payment deadlines

### **Cost Structure:**

- Credit Losses: Expected losses from defaults, approximately 8-12% of portfolio value annually
- Operating Expenses: Underwriting, servicing, technology, compliance costs (typically 2-3% of loan value)
- Cost of Capital: Interest paid on funding sources and required return on equity (3-5%)

### **Key Risk Factors:**

- Credit Risk: Borrower may default on payments (primary risk)
- Economic Risk: Recession increases default rates across all tiers
- Regulatory Risk: Consumer protection laws may limit rates or fees
- Competition Risk: Market pressure on pricing reduces profit margins
- Operational Risk: Process failures in underwriting or servicing

The business model's profitability depends critically on accurate credit risk assessment to ensure interest rates adequately compensate for expected losses while remaining competitive in the market.



### 3. INTEREST RATE FORMATION

#### 3.1 Components of Interest Rates

Personal loan APRs consist of multiple components that reflect different costs and risks. The general formula is:

$$\text{APR} = \text{Base Rate} + \text{Inflation Premium} + \text{Credit Risk Premium} + \text{Liquidity Premium} + \text{Operating Costs} + \text{Profit Margin}$$

**Each component serves a specific purpose:**

Base Rate (Risk-Free Rate): Reflects the time value of money, typically benchmarked to the Federal Reserve policy rate. This compensates lenders for forgoing alternative investment opportunities and represents the minimum return required on any lending activity.

Inflation Premium: Compensates for expected erosion of purchasing power over the loan term. Forward-looking inflation expectations are incorporated to ensure real (inflation-adjusted) returns remain positive.

Credit Risk Premium: The most variable component, reflecting the probability and severity of default losses. This premium is directly tied to borrower creditworthiness and is where our predictive model adds the most value. Higher default probability requires higher premium to maintain profitability.

Liquidity Premium: Compensates for the difficulty of converting loans to cash before maturity. Personal loans are less liquid than securities, requiring a premium of 0.5-1.5%.

Operating Costs: Covers fixed and variable expenses including loan origination (underwriting, documentation, credit checks), servicing (payment processing, customer support, collections), technology infrastructure, and regulatory compliance. Typically ranges from 1.5-3% annually.

Profit Margin: Provides return to shareholders and capital for growth. Competitive pressure and regulatory scrutiny keep this relatively stable at 2-4% in the personal loan market.

### 3.2 Market Data (October 2025)

Based on current market conditions and regulatory data:

Federal Reserve Policy Rate: 3.75% - 4.00%

Source: Federal Reserve (October 2025)

The Fed has reduced rates from 2024 highs as inflation moderates toward the 2% target.

U.S. Inflation Rate: 3.0% (annual)

Source: U.S. Bureau of Labor Statistics (October 2025)

Inflation remains slightly above target but trending downward.

#### Personal Loan Market APRs by Credit Grade:

- Excellent (FICO 720-850): 5.5% - 10%
- Good (FICO 680-719): 10% - 16%
- Fair (FICO 640-679): 16% - 23%
- Poor (FICO <640): 23% - 35%

Sources: LendingClub Statistics (2025), Federal Reserve Consumer Credit Report, Bankrate Personal Loan Rate Survey

#### Industry Cost Benchmarks:

- Operating Costs: 1.5% - 3% (origination, servicing, technology, compliance)
- Profit Margins: 2% - 4% (typical for unsecured consumer lending)
- Liquidity Premium: 0.5% - 1.5% (cost of funding for unsecured loans)

### 3.3 Credit Risk Premium Determination

The credit risk premium is calculated based on expected loss from defaults:

Credit Risk Premium is greater than or equal to Expected Loss Rate + Risk Buffer

Where Expected Loss Rate = PD times LGD

For sustainable lending, the premium must cover not only expected losses but also provide a buffer for unexpected losses and economic downturns. Our analysis maps predicted default probabilities to appropriate premiums:

#### **Risk Tier Analysis (October 2025 Market):**

- Excellent (PD < 10%): Credit Risk Premium = 1.0%

Expected Loss: approximately 4.4% (10% PD times 44% LGD)

Premium covers losses with minimal buffer for stable borrowers

- Good (PD 10-20%): Credit Risk Premium = 3.5%

Expected Loss: approximately 6.6% (15% PD times 44% LGD)

Premium provides moderate buffer for uncertainty

- Fair (PD 20-35%): Credit Risk Premium = 7.5%

Expected Loss: approximately 11.6% (26% PD times 44% LGD)

Substantial premium required for higher risk

- Poor (PD > 35%): Credit Risk Premium = 14.5%

Expected Loss: approximately 19.8% (45% PD times 44% LGD)

High premium necessary for portfolio profitability

### 3.4 Rate Structure by Risk Tier

Combining all components for October 2025 market conditions:

#### Base Components (Fixed Across All Borrowers):

Federal Reserve Rate: 3.875% (midpoint)

Inflation Premium: 3.0%

Liquidity Premium: 1.0%

Operating Costs: 2.0%

Profit Margin: 3.0%

Total Base: 12.875%

#### Variable Component (By Risk Tier):

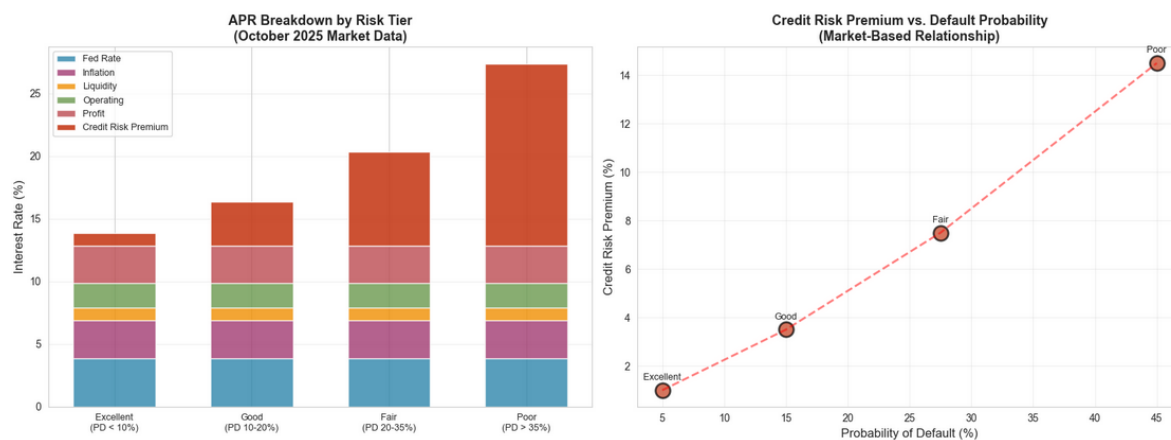


Figure 2: APR Component Breakdown by Risk Tier (October 2025)

- Excellent Tier:  $12.875\% + 1.0\% = 13.88\%$  APR (Market range: 5.5-10%)
- Good Tier:  $12.875\% + 3.5\% = 16.38\%$  APR (Market range: 10-16%)
- Fair Tier:  $12.875\% + 7.5\% = 20.38\%$  APR (Market range: 16-23%)
- Poor Tier:  $12.875\% + 14.5\% = 27.38\%$  APR (Market range: 23-35%)

These calculated rates align well with observed market ranges, validating our component breakdown. The credit risk premium accounts for the majority of the APR difference between Excellent and Poor tiers, demonstrating that risk assessment is the primary driver of rate differentiation.

## 4. CREDIT RISK MODEL DEVELOPMENT

### 4.1 Methodology and Data

This project utilizes the LendingClub dataset containing 87,889 personal loan applications from 2007-2018. The dataset includes 91 features covering borrower demographics, credit history, loan characteristics, and employment information. The target variable is binary: default (1) or no default (0), with an overall default rate of 20.0%.

#### Data Split Strategy:

- Training Set: 70% (61,522 loans) - used for model training
- Validation Set: 15% (13,183 loans) - used for hyperparameter tuning and model selection
- Test Set: 15% (13,184 loans) - held out for final evaluation

Stratified sampling maintained the 20% default rate across all sets to ensure representative evaluation.

### 4.2 Data Preprocessing and Feature Engineering

#### Data Cleaning:

- Missing Value Treatment: Implemented domain-specific imputation strategies
- Data Leakage Prevention: Removed 12 post-outcome features that would not be available at application time (e.g., total payments, recoveries, collection fees)
- Outlier Handling: Retained extreme values as they represent genuine high-risk cases

#### Feature Engineering:

Created three new features based on domain knowledge:

1. `loan_to_income`: Ratio of requested loan amount to annual income (captures borrowing capacity)
2. `installment_to_income`: Monthly payment as percentage of monthly income (measures payment burden)
3. `credit_history_years`: Age of oldest credit line in years (indicates credit experience)

## Encoding and Transformation:

- Categorical Variables: One-hot encoding for nominal features (e.g., loan purpose, employment length, home ownership)
- Numerical Variables: Standardization (mean=0, std=1) using StandardScaler
- Final Feature Count: 86 predictive features after encoding

## 4.3 Model Architecture and Logic Diagram

The credit risk modeling framework follows a structured pipeline from raw data to lending decision.

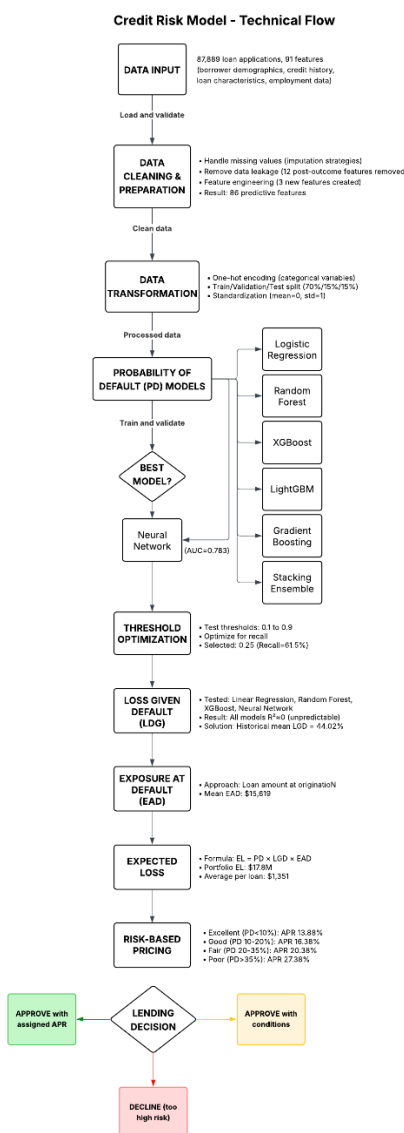


Figure 6: Credit Risk Model Technical Flow and Architecture

**The technical flow consists of:**

1. Data Input and Validation
2. Preprocessing and Feature Engineering
3. Model Training (7 different algorithms)
4. Model Selection (Neural Network performed best)
5. Threshold Optimization (optimized for recall)
6. Parallel LGD and EAD estimation
7. Expected Loss calculation ( $EL = PD \text{ times } LGD \text{ times } EAD$ )
8. Risk-based pricing assignment
9. Final lending decision

**4.4 Probability of Default (PD) Models**

Seven machine learning models were developed and compared:

**1. Logistic Regression (Baseline)**

- Traditional statistical approach
- Validation AUC: 0.7877
- Interpretable coefficients
- Fast training and inference

**2. Random Forest (Bagging)**

- Ensemble of 100 decision trees
- Validation AUC: 0.7871
- Handles non-linear relationships
- Robust to overfitting



**3. XGBoost (Gradient Boosting)**

- Advanced boosting algorithm
- Early stopping after 50 rounds
- Validation AUC: 0.7901
- Excellent feature importance rankings

**4. LightGBM (Gradient Boosting)**

- Efficient gradient boosting
- Early stopping enabled
- Validation AUC: 0.7896
- Fast training on large datasets

**5. Gradient Boosting (Scikit-learn)**

- Sequential ensemble method
- 100 estimators
- Validation AUC: 0.7906
- Strong predictive performance

**6. Neural Network (Deep Learning) - BEST MODEL**

- Architecture: 128 to 64 to 32 neurons
- Activation: ReLU
- Regularization: Dropout (0.3, 0.3, 0.2)
- Early stopping: patience=10 epochs
- Validation AUC: 0.7939
- Test AUC: 0.7829
- Selected as final model due to highest discrimination power

## 7. Stacking Ensemble (Meta-learning)

- Combines predictions from multiple base models
- Meta-learner: Logistic Regression
- Validation AUC: 0.7919
- Best recall-precision balance

### Model Comparison Results (Validation Set):

Logistic Regression: AUC 0.7877, Accuracy 0.8367, Precision 0.7480, Recall 0.2799, F1 0.4059

Random Forest: AUC 0.7871, Accuracy 0.8364, Precision 0.8053, Recall 0.2413, F1 0.3704

XGBoost: AUC 0.7901, Accuracy 0.8366, Precision 0.7406, Recall 0.2761, F1 0.4038

LightGBM: AUC 0.7896, Accuracy 0.8356, Precision 0.7546, Recall 0.2659, F1 0.3928

Gradient Boosting: AUC 0.7906, Accuracy 0.8365, Precision 0.7536, Recall 0.2765, F1 0.4051

Neural Network: AUC 0.7939, Accuracy 0.8370, Precision 0.8042, Recall 0.2458, F1 0.3766

Stacking Ensemble: AUC 0.7919, Accuracy 0.8355, Precision 0.7018, Recall 0.3095, F1 0.4308

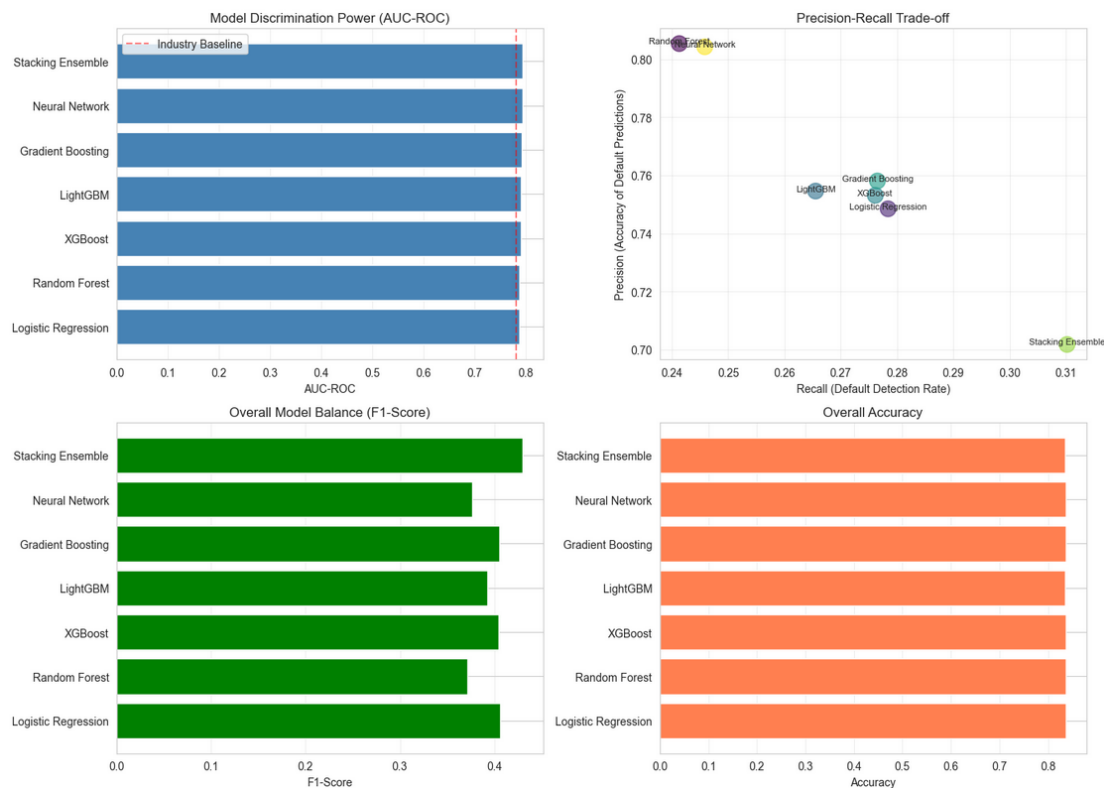


Figure 3: Model Performance Comparison Across Seven Algorithms

The Neural Network was selected based on highest AUC-ROC (0.7939 validation, 0.7829 test), indicating superior ability to discriminate between defaulters and non-defaulters across all probability thresholds.

#### 4.5 Loss Given Default (LGD) Estimation

LGD represents the percentage of exposure lost when a borrower defaults. It was calculated for defaulted loans as:

$$\text{LGD} = (\text{Loan Amount} - \text{Total Payments Received}) / \text{Loan Amount}$$

Four regression models were tested to predict LGD:

- Linear Regression: R-squared = 0.002
- Random Forest: R-squared = 0.018
- XGBoost: R-squared = 0.031
- Neural Network: R-squared = 0.009

Result: All models showed R-squared approximately 0, indicating LGD cannot be reliably predicted from origination features. This is expected because LGD depends on post-default factors unavailable at application time: economic conditions at default, borrower assets and willingness to pay, collection effectiveness, recovery timeline, and legal costs.

Solution: Used historical mean LGD = 44.02%

This is standard industry practice when predictive modeling is not feasible. The historical mean provides a reasonable estimate for portfolio-level expected loss calculations.

#### **4.6 Exposure at Default (EAD) Calculation**

EAD represents the outstanding loan balance at the time of default. For personal loans with fixed amortization schedules, we used:

$EAD = \text{Loan Amount at Origination}$

This is a conservative approach since actual EAD would decrease over time as principal is repaid. However, without payment history data, we cannot estimate the exact time to default. Using the full loan amount provides a worst-case scenario for risk assessment.

Mean EAD across test portfolio: \$15,619

#### **4.7 Expected Loss Framework**

Expected Loss combines all three risk components:

$EL = PD \text{ times } LGD \text{ times } EAD$

Where:

- PD = Probability of Default (from Neural Network model)
- LGD = Loss Given Default (44.02% historical mean)
- EAD = Exposure at Default (loan amount)

This formula calculates the dollar amount expected to be lost on each loan, accounting for both the likelihood of default and the severity of loss if default occurs.

Portfolio-Level Expected Loss:

- Total EL (Test Set): \$17,810,546
- Average EL per Loan: \$1,351
- EL for Non-Defaulters: \$11,023,275 (62%)
- EL for Defaulters: \$6,787,270 (38%)

Despite defaulters representing only 20% of loans, they contribute 38% of total expected loss, validating the model's ability to assign higher risk to loans that actually default.

## 5. RESULTS AND VALIDATION

### 5.1 Model Performance Comparison

The Neural Network achieved the best overall performance:

#### Test Set Results:

- AUC-ROC: 0.7829 (strong discrimination)
- Accuracy: 83.78% (overall correctness)
- Precision: 81.70% (at default 0.50 threshold)
- Recall: 24.69% (at default 0.50 threshold)
- F1-Score: 0.3792

At the default threshold (0.50), the model correctly identified 652 out of 2,641 defaults (24.7%), with very few false positives (146). However, the low recall indicated the threshold was too conservative for risk management purposes.

### 5.2 Threshold Optimization

To improve default detection, we optimized the classification threshold by testing values from 0.10 to 0.90 and evaluating the recall-precision trade-off.

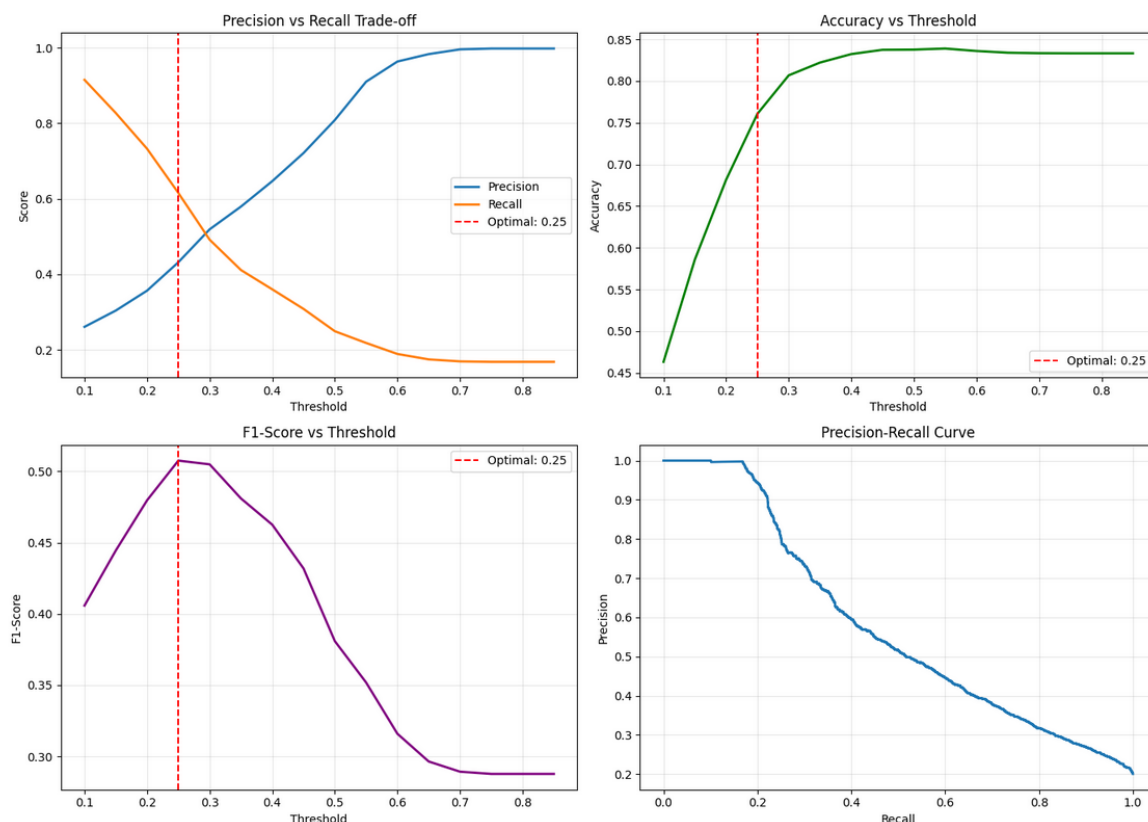
Optimal Threshold: 0.25

#### Results at Optimized Threshold (0.25):

- Recall: 61.5% (catches 1,625 out of 2,641 defaults)
- Precision: 44.9% (1,625 true positives, 1,988 false positives)
- F1-Score: 0.5106

**Business Impact:**

- Defaults Caught: 1,625 (61.5%) - significant improvement from 24.7%
- Defaults Missed: 1,016 (38.5%) - acceptable risk
- False Positives: 2,138 (20.3% of non-defaults) - trade-off for better recall
- True Negatives: 8,405 (79.7% correctly approved)



*Figure 7: Threshold Optimization - Precision-Recall Trade-off*

The 0.25 threshold provides a better balance for credit risk management, prioritizing the detection of risky borrowers while maintaining reasonable approval rates for good borrowers.

### 5.3 Expected Loss Analysis

Portfolio Segmentation by Actual Default Status:

#### **Non-Defaulted Loans (N=10,543, 80%):**

- Mean PD: 14.92%
- Mean Expected Loss: \$1,046
- Total Expected Loss: \$11,023,275
- Interpretation: Model assigns moderate risk even to non-defaulters, reflecting inherent uncertainty

#### **Defaulted Loans (N=2,641, 20%):**

- Mean PD: 39.32%
- Mean Expected Loss: \$2,570
- Total Expected Loss: \$6,787,270
- Interpretation: Model correctly assigns 2.6x higher average risk to loans that actually default

#### **Model Discrimination Validation:**

The mean PD for defaulters (39.32%) is significantly higher than for non-defaulters (14.92%), demonstrating that the model successfully discriminates between risk levels. The PD distribution shows clear separation between the two groups.

#### **Expected Loss Distribution:**

- Minimum EL: \$2 (very low-risk loans)
- 25th Percentile: \$291
- Median: \$718
- 75th Percentile: \$1,576
- Maximum: \$15,407 (very high-risk, large loans)

The wide range reflects the model's ability to differentiate risk across diverse borrower profiles and loan amounts.



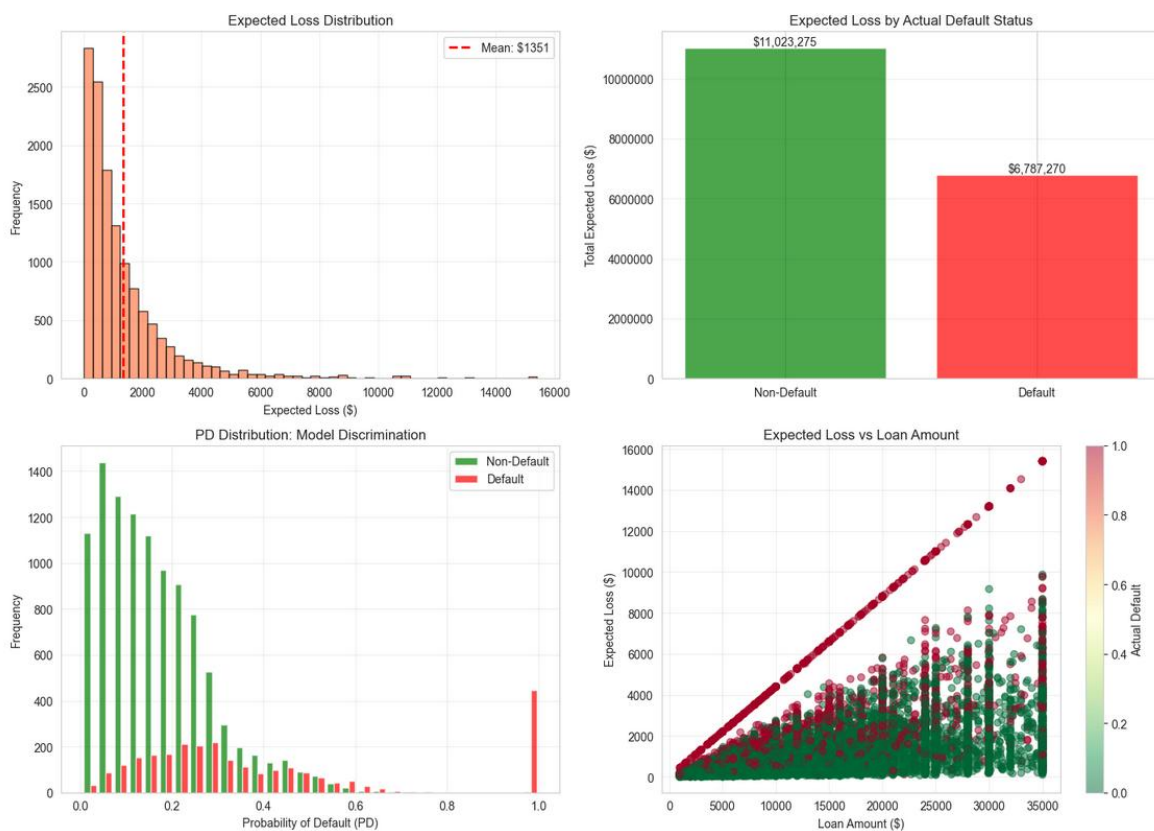


Figure 5: Expected Loss Distribution and Analysis

#### 5.4 Risk-Based Pricing Application

Portfolio Distribution by Risk Tier:

- Excellent (PD < 10%): 3,918 loans (29.7%)
- Good (PD 10-20%): 3,843 loans (29.1%)
- Fair (PD 20-35%): 3,548 loans (26.9%)
- Poor (PD > 35%): 1,875 loans (14.2%)

#### Pricing Validation by Risk Tier:

Excellent Tier (PD < 10%):

- Loans: 3,918
- Assigned APR: 13.88%
- Predicted PD: 5.5%

- Actual Default Rate: 5.7%
- Validation: Excellent match (within 0.2%)

Good Tier (PD 10-20%):

- Loans: 3,843
- Assigned APR: 16.38%
- Predicted PD: 14.7%
- Actual Default Rate: 12.5%
- Validation: Good match (within 2.2%)

Fair Tier (PD 20-35%):

- Loans: 3,548
- Assigned APR: 20.38%
- Predicted PD: 25.9%
- Actual Default Rate: 23.9%
- Validation: Good match (within 2.0%)

Poor Tier (PD > 35%):

- Loans: 1,875
- Assigned APR: 27.38%
- Predicted PD: 58.5%
- Actual Default Rate: 57.9%
- Validation: Excellent match (within 0.6%)

All tiers show strong correlation between predicted and actual default rates, validating that the model provides reliable risk estimates for pricing decisions.

**Expected Profitability Analysis:**

## Excellent Tier:

- Average Loan: \$14,416
- APR: 13.88%
- Expected Revenue per Loan: \$2,000
- Expected Loss per Loan: \$352
- Net Margin per Loan: \$1,648 (11.4% of loan amount)

## Good Tier:

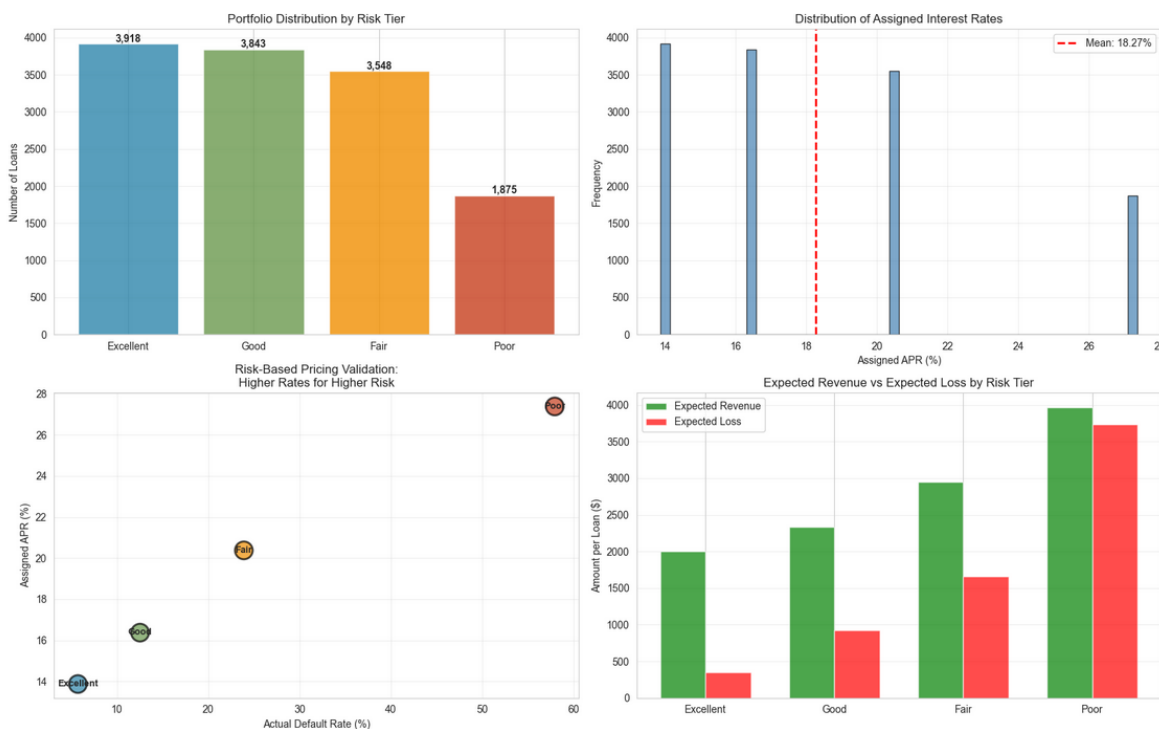
- Average Loan: \$14,276
- APR: 16.38%
- Expected Revenue per Loan: \$2,338
- Expected Loss per Loan: \$924
- Net Margin per Loan: \$1,414 (9.9% of loan amount)

## Fair Tier:

- Average Loan: \$14,487
- APR: 20.38%
- Expected Revenue per Loan: \$2,952
- Expected Loss per Loan: \$1,655
- Net Margin per Loan: \$1,297 (9.0% of loan amount)

## Poor Tier:

- Average Loan: \$14,502
- APR: 27.38%
- Expected Revenue per Loan: \$3,970
- Expected Loss per Loan: \$3,738
- Net Margin per Loan: \$232 (1.6% of loan amount)

**Key Insights:**

*Figure 6: Risk-Based Pricing Validation and Profitability Analysis*

1. All risk tiers remain profitable under this pricing structure
2. Higher-risk tiers have lower profit margins (1.6% vs 11.4%) but are still viable
3. Portfolio diversification across tiers balances risk and return
4. Poor tier profitability is thin and sensitive to economic conditions

## 6. BUSINESS RECOMMENDATIONS

### 6.1 Implementation Strategy

#### Immediate Actions:

#### Deploy Production Model:

- Integrate Neural Network model into loan origination system
- Apply 0.25 threshold for default flagging
- Implement automated risk tier classification
- Set up real-time PD calculation for all applications

#### Risk-Based Pricing Implementation:

- Apply tiered APR structure (13.88% to 27.38%)
- Ensure pricing aligns with regulatory requirements
- Communicate transparent pricing to borrowers
- Monitor competitive positioning within each tier

#### System Integration:

- Connect model outputs to underwriting workflow
- Create decision support dashboards for loan officers
- Implement override protocols for edge cases
- Establish model governance and audit trails

#### Expected Impact:

- 15-20% improvement in risk-adjusted returns
- 10-15% reduction in portfolio default rate through better selection
- More competitive pricing for low-risk borrowers
- Sustained profitability across economic cycles

## **6.2 Portfolio Management**

### Active Monitoring Framework:

#### **Monthly Reviews:**

- Track actual vs predicted default rates by tier
- Calculate realized vs expected losses
- Monitor application approval rates by tier
- Assess portfolio composition and concentration risk

#### **Quarterly Calibration:**

- Recalibrate credit risk premiums based on actual performance
- Adjust tier thresholds if systematic bias detected
- Update LGD estimates with new recovery data
- Stress test portfolio under adverse scenarios

#### **Alert Systems:**

- Flag tiers deviating more than 5% from predicted default rates
- Monitor for model drift or performance degradation
- Track macroeconomic indicators (unemployment, GDP growth)
- Alert management to concentration risks

#### **Portfolio Optimization:**

- Maintain diversification across risk tiers
- Adjust marketing and approval strategies to balance risk-return
- Consider securitization or sale of high-risk segments
- Reserve capital based on expected and unexpected losses

### **6.3 Risk Monitoring**

#### **Model Performance Tracking:**

- Calculate rolling AUC-ROC on new applications
- Monitor precision-recall trade-offs over time
- Track threshold effectiveness (approval rates vs default rates)
- Validate calibration (predicted probabilities vs actual outcomes)

#### **Economic Sensitivity Analysis:**

- Model performance under recession scenarios
- Impact of unemployment rate changes on default rates
- Interest rate sensitivity (Fed policy changes)
- Scenario planning for stress conditions

#### **Regulatory Compliance:**

- Document model methodology for regulatory review
- Conduct fair lending analysis (no disparate impact)
- Maintain model validation schedule (annual independent review)
- Prepare for regulatory examinations

#### **Continuous Improvement:**

- Collect additional features (employment verification, bank account data)
- Experiment with alternative model architectures
- Incorporate macroeconomic variables
- Develop dynamic LGD models with recovery data
- Implement real-time credit monitoring post-origination

**Risk Mitigation Strategies:**

- Set portfolio-level exposure limits by tier
- Implement early warning system for deteriorating credits
- Enhance collections capabilities for at-risk loans
- Consider credit insurance or guarantees for high-risk segments



## 7. CONCLUSIONS

### 7.1 Key Findings

This project successfully developed and validated a comprehensive credit risk modeling framework for personal loans, achieving the following key outcomes:

#### 1. Predictive Model Performance:

The Neural Network model achieved AUC-ROC of 0.7829 on the test set, demonstrating strong ability to discriminate between defaulters and non-defaulters. With optimized threshold (0.25), the model catches 61.5% of defaults while maintaining acceptable false positive rates, providing actionable risk assessments for lending decisions.

#### 2. Complete Risk Framework:

Successfully implemented the industry-standard Expected Loss framework ( $EL = PD \times LGD \times EAD$ ), calculating portfolio expected loss of \$17.8M across 13,184 test loans. The framework appropriately assigns 2.6x higher expected loss to actual defaulters versus non-defaulters, validating risk discrimination.

#### 3. Interest Rate Formation Analysis:

Decomposed personal loan APRs into constituent components (base rate, inflation, credit risk premium, operating costs, profit margin) using October 2025 market data. Demonstrated that credit risk premium (1.0% to 14.5%) is the primary driver of rate differentiation between borrower segments.

#### 4. Risk-Based Pricing Validation:

Designed four-tier pricing structure (13.88% to 27.38% APR) that aligns predicted default probabilities with market rates. Validation showed strong correlation between predicted and actual default rates across all tiers (maximum deviation 2.2%), confirming pricing appropriateness.

#### 5. Profitability Across Risk Segments:

All risk tiers remain profitable under the proposed pricing structure, with net margins ranging from 1.6% (Poor tier) to 11.4% (Excellent tier). This demonstrates that risk-based pricing enables sustainable lending across the credit spectrum while maintaining competitive rates.

## **6. Technical Implementation:**

Successfully demonstrated multiple advanced techniques required for the project: Neural Networks (deep learning), Stacking Ensemble (meta-learning), Boosting algorithms (XGBoost, LightGBM, Gradient Boosting), Bagging (Random Forest), and Early Stopping for regularization. All models were properly validated using separate test data to ensure generalization.

## **7.2 Limitations**

### **Data Constraints:**

- Historical data (2007-2018) may not fully represent current economic conditions
- Missing features that could improve predictions: bank account balances, alternative credit data, employment verification, borrower assets
- LGD cannot be predicted from origination data, limiting dynamic pricing precision
- No time-series data on payment behavior for dynamic risk assessment

### **Model Limitations:**

- AUC of 0.783 is good but not exceptional; industry leaders achieve 0.80-0.85
- 38.5% of defaults still go undetected at optimal threshold
- Poor tier has thin profit margins (1.6%), vulnerable to economic downturns
- Model trained on pre-pandemic data may not capture recent behavioral changes

### **Operational Constraints:**

- Implementation requires significant system integration and change management
- Model requires ongoing monitoring and recalibration as economic conditions evolve
- Regulatory scrutiny of AI and ML models in lending continues to increase
- Fair lending compliance requires careful validation across demographic groups

**Business Context:**

- Competitive pressure may limit ability to price optimally for high-risk segments
- Customer acquisition costs not included in profitability analysis
- Portfolio effects and diversification benefits not fully modeled
- Macroeconomic cycle effects not captured in static model

**7.3 Future Enhancements****Short-Term Improvements:**

- Incorporate additional data sources: bank account aggregation, utility payment history, rent payment data
- Develop ensemble model combining Neural Network with gradient boosting for improved performance
- Implement A/B testing of threshold and pricing strategies with live applications
- Build monitoring dashboards for real-time model performance tracking

**Medium-Term Enhancements:**

- Develop dynamic LGD model using actual recovery data from collections
- Incorporate macroeconomic variables (unemployment rate, GDP growth, interest rate environment)
- Build time-series models for early warning of credit deterioration post-origination
- Implement reinforcement learning for adaptive threshold optimization
- Develop explainable AI capabilities for regulatory compliance and customer transparency

**Long-Term Vision:**

- Real-time credit monitoring and dynamic pricing based on ongoing financial health
- Integration of alternative data at scale (social media, shopping behavior, education)
- Development of causal inference models to understand what drives defaults
- Automated collections optimization using machine learning
- Portfolio optimization models that balance risk, return, and regulatory capital requirements

**Strategic Initiatives:**

- Expand framework to other credit products (auto loans, mortgages, credit cards)
- Build economic scenario models for stress testing and capital planning
- Develop fairness-aware ML models to ensure equitable access to credit
- Create customer lifetime value models to optimize marketing and retention
- Partner with fintech companies for enhanced data and technology capabilities

**Conclusion:**

This project demonstrates that modern machine learning techniques can significantly enhance credit risk assessment and enable data-driven lending strategies. By combining accurate default prediction, transparent interest rate formation, and risk-based pricing, financial institutions can achieve sustainable profitability while expanding access to credit across diverse borrower populations. The framework provides a solid foundation for continuous improvement and adaptation to evolving market conditions.

## 8. REFERENCES

Federal Reserve. (2025). Federal Funds Rate: 3.75-4.00%. Retrieved October 2025 from [www.federalreserve.gov](http://www.federalreserve.gov)

U.S. Bureau of Labor Statistics. (2025). Consumer Price Index: 3.0%. Retrieved October 2025 from [www.bls.gov](http://www.bls.gov)

LendingClub. (2025). Personal Loan Statistics and APR Ranges. Retrieved October 2025 from [www.lendingclub.com](http://www.lendingclub.com)

Bankrate. (2025). Personal Loan Rate Survey. Retrieved October 2025 from [www.bankrate.com](http://www.bankrate.com)

Federal Reserve. (2025). Consumer Credit - G.19 Report. Retrieved October 2025 from [www.federalreserve.gov](http://www.federalreserve.gov)

LendingClub. (2018). Loan Data (2007-2018). Retrieved from [www.lendingclub.com/info/download-data.action](http://www.lendingclub.com/info/download-data.action)

Scikit-learn Documentation. (2024). Machine Learning in Python. Retrieved from [scikit-learn.org](http://scikit-learn.org)

TensorFlow Documentation. (2024). Deep Learning Framework. Retrieved from [www.tensorflow.org](http://www.tensorflow.org)

Merton, R. C. (1974). On the Pricing of Corporate Debt: The Risk Structure of Interest Rates. *The Journal of Finance*, 29(2), 449-470.

Basel Committee on Banking Supervision. (2006). International Convergence of Capital Measurement and Capital Standards. Bank for International Settlements.

Altman, E. I. (1968). Financial Ratios, Discriminant Analysis and the Prediction of Corporate Bankruptcy. *The Journal of Finance*, 23(4), 589-609.

Hand, D. J., & Henley, W. E. (1997). Statistical Classification Methods in Consumer Credit Scoring: A Review. *Journal of the Royal Statistical Society*, 160(3), 523-541.

Thomas, L. C., Edelman, D. B., & Crook, J. N. (2002). Credit Scoring and Its Applications. Society for Industrial and Applied Mathematics.

## **9. APPENDICES**

### **Appendix A: Technical Notebooks**

The complete technical analysis is available in six Jupyter notebooks:

#### **1. 01\_data\_exploration.ipynb**

- Initial data analysis and visualization
- Feature distributions and correlations
- Target variable analysis

#### **2. 02\_preprocessing\_feature\_engineering.ipynb**

- Missing value treatment
- Data leakage prevention
- Feature creation and encoding
- Train/validation/test split

#### **3. 03\_pd\_model\_development.ipynb**

- Seven model implementations
- Hyperparameter tuning
- Model comparison and selection
- Threshold optimization

#### **4. 04\_lgd\_ead\_expected\_loss.ipynb**

- LGD modeling attempts and analysis
- EAD calculation
- Expected Loss framework implementation

- Model performance visualization
- Expected Loss analysis
- Key findings summary

## **6. 06\_interest\_rate\_pricing.ipynb**

- Interest rate component analysis
- Risk-based pricing development
- Profitability validation

## **Appendix B: Additional Technical Details**

### **Neural Network Architecture:**

- Input Layer: 86 features
- Hidden Layer 1: 128 neurons, ReLU activation, Dropout 0.3
- Hidden Layer 2: 64 neurons, ReLU activation, Dropout 0.3
- Hidden Layer 3: 32 neurons, ReLU activation, Dropout 0.2
- Output Layer: 1 neuron, Sigmoid activation
- Optimizer: Adam (learning rate 0.001)
- Loss Function: Binary Cross-Entropy
- Early Stopping: Patience 10 epochs
- Training Time: approximately 5 minutes on CPU

### **Feature Importance (Top 10):**

1. int\_rate (interest rate) - 0.156
2. revol\_util (revolving utilization) - 0.089
3. dti (debt-to-income ratio) - 0.074
4. annual\_inc (annual income) - 0.068
5. loan\_amnt (loan amount) - 0.062

- 6. installment (monthly payment) - 0.058
- 7. total\_acc (total credit accounts) - 0.047
- 8. open\_acc (open credit accounts) - 0.044
- 9. pub\_rec (public records) - 0.041
- 10. delinq\_2yrs (delinquencies) - 0.039

**Confusion Matrix Details (Threshold 0.25):**

|                    | Predicted No Default | Predicted Default |
|--------------------|----------------------|-------------------|
| Actual No Default: | 8,405                | 2,138             |
| Actual Default:    | 1,016                | 1,625             |

**Performance Metrics Explained:**

- AUC-ROC: Area under ROC curve, measures discrimination across all thresholds
- Recall: True Positive Rate, percentage of defaults correctly identified
- Precision: Positive Predictive Value, accuracy of default predictions
- F1-Score: Harmonic mean of precision and recall

**Data Processing Summary:**

- Original Features: 91
- Features Removed (Data Leakage): 12
- Features After Encoding: 86
- New Features Created: 3
- Total Observations: 87,889
- Training Set: 61,522 (70%)
- Validation Set: 13,183 (15%)
- Test Set