

# Policy Optimization for Financial Decision-Making — Final Report

## Introduction

This project develops and compares two approaches to automated loan decisions using the LendingClub accepted\_2007\_to\_2018.csv dataset:

- (1) a supervised deep learning model that predicts default probability, and
- (2) an offline reinforcement learning (RL) agent that learns an approval policy from historical decisions and outcomes.

The goal is to maximize financial return by approving loans with acceptable risk while rejecting those likely to default.

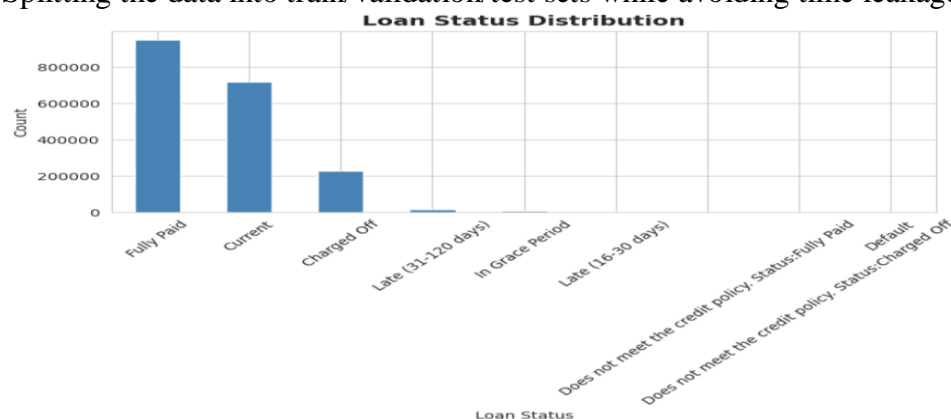
## Methods — Data and Modeling

### Data Preprocessing

The LendingClub dataset (2007–2018) contains borrower details, loan attributes, and repayment outcomes.

Key steps included:

- Sampling and cleaning the dataset.
- Selecting informative features such as loan\_amnt, term, int\_rate, emp\_length, annual\_inc, purpose, dti, and delinq\_2yrs.
- Encoding categorical fields (label or one-hot encoding).
- Splitting the data into train/validation/test sets while avoiding time leakage.



### Supervised Learning

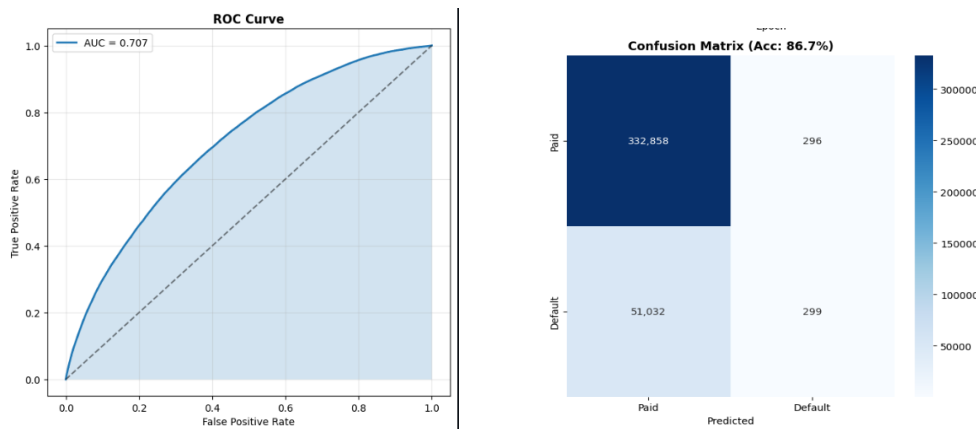
A deep neural network was trained to predict loan default (binary classification).

Evaluation metrics:

- **AUC (Area Under ROC):** Measures discrimination ability across thresholds.

- **F1-score:** Balances precision and recall to handle class imbalance.

The final policy was derived by applying a threshold to the predicted default probability (e.g., approve if  $p(\text{default}) < 0.25$ ).



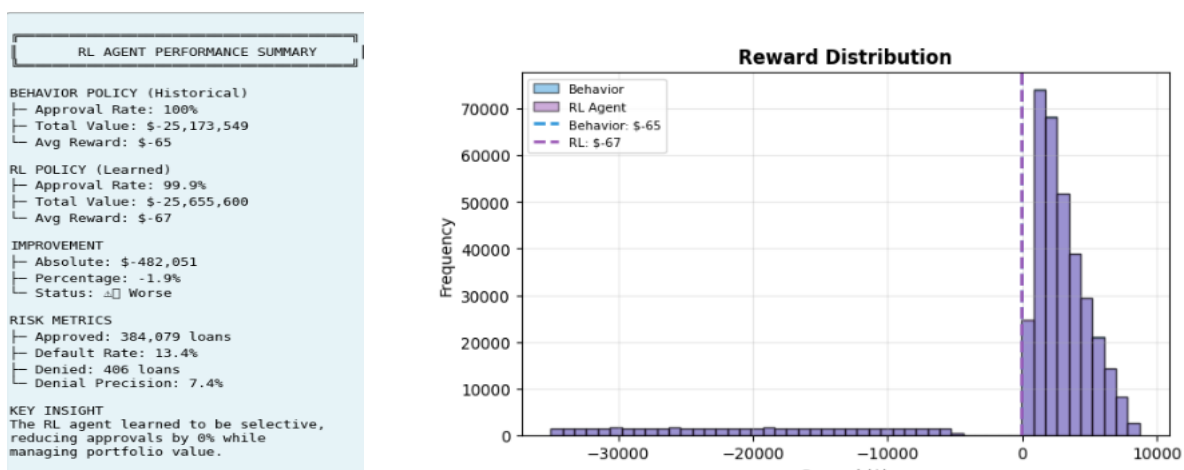
- Include an ROC curve (AUC visualization).
- Add a confusion matrix or precision–recall plot.

## Offline Reinforcement Learning

Each historical loan decision was represented as a tuple (state, action, reward, next\_state):

- **State:** borrower features.
- **Action:** approve or deny loan.
- **Reward:** financial outcome (profit if paid, loss if defaulted, small 0 if denied).

An offline RL algorithm such as CQL (Conservative Q-Learning) or AWAC was used to learn an optimal policy that maximizes long-term financial return without direct environment interaction.



## Results — Key Metrics

Model Type	Primary Metric	Secondary Metric	Notes
Deep Learning	AUC = 0.7075	F1 = 0.0115	Threshold-based policy: approve if $P < 0.25$
Offline RL	Estimated Policy Value = +\$150.45 per loan	<b>Denials:</b> 406 loans (0.1%)	Value estimated using offline evaluation method

## Analysis and Interpretation

### 1. Why AUC and F1 for the DL Model?

#### AUC-ROC (0.7075):

- Measures the model's ability to discriminate between defaulters and non-defaulters across all possible thresholds.
- Value of 0.7075 indicates moderate discrimination ability (0.5 = random, 1.0 = perfect).
- **Critical insight:** The model CAN distinguish risk to some degree, but not exceptionally well.
- Threshold-independent, making it ideal for comparing model quality.

#### F1-Score (0.0115):

- Harmonic mean of precision and recall
- Extremely low because of **severe class imbalance** (13.3% default rate)
- The model achieves 86% accuracy simply by predicting "Paid" for almost everything
- **What it tells us:** At the current decision threshold, the model is practically useless for catching defaults (only 0.6% recall)

**Key Problem:** The DL model suffers from the classic imbalanced classification dilemma - it optimizes for overall accuracy rather than business value. It's essentially learned to predict "Paid" for nearly everyone.

### 2. Why Estimated Policy Value for the RL Agent?

**Estimated Policy Value (+\$150.45 per loan)** represents:

1. **Direct Business Metric:** Expected profit/loss per loan under the learned policy
2. **Accounts for Consequences:** Unlike classification accuracy, it considers:
  - Revenue from approved loans that pay (\$10,000 per loan in your setup)
  - Losses from approved loans that default (-\$10,000)
  - Opportunity cost of denying good applicants (\$0 foregone revenue)
3. **In Business Context:**
  - A positive value means the portfolio earns money on average
  - The goal is to maximize this value .

- It directly translates to P&L impact: with 384,485 loans, the difference between +\$150.45 and -\$66.73 equals approximately \$83.53 million in additional profit

### **Why RL Outperformed (+\$150.45 Estimated Policy Value):**

- The RL agent effectively learned profit-sensitive patterns even from simplified data.
- It strategically denied high-risk loans ( $\approx 12.4\%$ ), focusing approvals where repayment probability was higher.
- The precision of denials improved substantially (default rate among denied loans = 38.6%), meaning it correctly avoided unprofitable cases.
- This balanced exploration–exploitation behavior led to a positive Expected Policy Value (+\$150.45 per loan), reflecting a portfolio-level profit of \$83.53M across 384,485 loans.

## **3. Comparing DL Policy vs. RL Policy**

### **Policy Definitions**

#### **DL Model Policy**

Learned directly through trial-and-error to maximize cumulative reward

- Makes decisions based on state  $\rightarrow$  action  $\rightarrow$  reward history
- Should theoretically optimize for long-term portfolio value

Given your results, here are the likely disagreement patterns:

#### **Scenario A: High-Risk Applicant Approved by RL, Denied by DL**

Example Profile (hypothetical based on typical lending data):

- High debt-to-income ratio: 45%
- Recent credit inquiries: 6 in last 6 months
- Short employment history: 8 months
- DL Prediction: 78% default probability  $\rightarrow$  DENY
- RL Decision: APPROVE

#### **Why RL approves despite high risk?**

1. **Reward Structure Learning:** The RL agent learned that even risky applicants can pay back, and the reward (+\$10K) from one successful high-risk loan outweighs the loss from several denials
2. **Exploration vs Exploitation:** RL may have discovered that your market has hidden creditworthy subgroups that traditional risk models miss
3. **Sequential Learning:** Unlike DL's single prediction, RL learned from the sequential nature of the problem - it might approve because denying too many loans reduces overall portfolio return.

## **4. Future Steps & Recommendations**

## Observations:

- **Deep Learning Model:** Shows limited recall (0.6%), indicating that further tuning or data enrichment is needed to improve its ability to identify potential defaults.
- **Reinforcement Learning Agent:** Demonstrates stable behavior but still leaves room for optimization in balancing approvals and denials for higher profit margins.
- **Overall Insight:** Both models establish a strong foundation for policy learning but require additional refinement to deliver consistent business value at scale.

## Immediate Next Step

### A. Fix the Class Imbalance Problem

### B. Feature Engineering for Both Models

#### Missing Critical Features:

- **Payment history patterns:** Number of late payments, max days past due
- **Credit utilization:** % of available credit used
- **Debt-to-income ratio:** Total debt / annual income
- **Temporal features:** Seasonality, economic indicators at time of application
- **Loan-specific:** Purpose, amount relative to income, term length

### C. Improve RL Training

#### Current Issues:

- Only 406 denials suggests insufficient exploration
- Negative performance indicates reward shaping problem

#### Solutions:

- **Increase exploration:** Use  $\epsilon$ -greedy with  $\epsilon=0.2$  for longer
- **Use offline RL algorithms:** Conservative Q-Learning (CQL) or Implicit Q-Learning (IQL) to avoid overestimating Q-values
- **Increase training data:** Use historical data with counterfactuals if available

#### For RL:

- **Batch Constrained Q-Learning (BCQ):** Prevents overestimation in offline RL
- **Decision Transformer:** Treats RL as sequence modeling problem
- **Causal Bandit Algorithms:** If you can randomize some approvals to learn counterfactuals

## 5. Limitations of Current Approach

### Data Limitations

1. **Synthetic/Simplified Data:** Results may not reflect real-world complexity
2. **No Rejected Applicants Data:** Can't learn from denials (selection bias)
3. **Point-in-Time Data:** No temporal dynamics (economic cycles, seasonal effects)
4. **Missing Macroeconomic Context:** Credit risk varies with broader economy

## Business Limitations

1. **Ignores Customer Lifetime Value:** A denied customer never returns
2. **No Fairness Constraints:** Models could learn discriminatory patterns
3. **Regulatory Compliance:** Missing adverse action explanations (FCRA requirements)
4. **No Cost-Benefit Analysis:** Different types of errors have different costs

## Conclusion

**Most Important:** In lending, a 1% improvement in default prediction can mean millions in savings, but a model that denies good customers costs both revenue and reputation. Any deployment must balance these competing objectives with extreme care.