

# Approve This Loan?

Reducing Defaults & Boosting Revenue with Predictive Models

Presented by  
**Muhammad Rizdky Maulady**



Project Based Internship – Data Scientist Home Credit Indonesia



# TABLE OF CONTENT

**1** Problem Research

**2** Data Preparation

**3** Data Visualization and  
Business Insight

**4** ML Implementation  
and Evaluation

**5** Business  
Recommendation

# Problem Research



Proportion of Good and Bad Clients



Total Clients: 1430153

## 1. Revenue \$653.41 B

Total Good Client		Revenue / Avg Loan
1306813	X	\$500.000

## 2. Total Bad Client Debt \$61.67 B

Total Bad Client		Revenue / Avg Loan
123340	X	\$500.000

## 3. Net Revenue \$591.74 B



### Problem Statement

The current loan approval process lacks effective risk assessment, resulting in a high number of bad clients and significant revenue loss.

### Goal

Reduce default rate through data-driven client assessment using machine learning.

### Objective

Build and deploy a predictive model to identify high-risk clients before loan approval.

### Business Metrics

- Default Rate Reduction (%)
- Increase in Net Revenue (\$)
- Model Accuracy / Recall on Bad Clients (%)
- Cost Saved from Bad Loans (\$)



"High loan default rate (8.6%) significantly reduces net revenue by over \$61.67B annually."

# Data Preparation



## POS\_CASH\_Balance.csv

- No Duplicate
- Handling Missing values use Median
- Agg SK\_DPD by SK\_ID\_PREV (mean)

## credit\_card\_balance.csv

- No Duplicate
- Handling Missing values use Median
- Agg(sum) 5 features by SK\_ID\_Prev

## installments\_payments.csv

- No Duplicate
- Handling Missing values use drop
- Agg(sum) by ID

## bureau\_balance.csv

- No Duplicate
- No Missing values
- Create new features

## previous\_application.csv

## bureau.csv

- No Duplicate
- Handling Missing value >50% drop, 5% s/d 50% median, <5% drop

## application\_train.csv

**Main Dataset**  
1.430.155 rows x 105 feature



# DATA PRE-PROCESSING

## PRE-PROCESSING PIPELINE

01

### Feature Engineering

- **Age Grouping:** Binned DAYS\_BIRTH into age groups
- **Row Filtering:** Dropped NAME\_FAMILY\_STATUS = 'Unknown'
- **Missing Values:** Replaced 'XNA', 'XAP' with NaN
- **Value Correction:** Fixed CNT\_FAM\_MEMBERS anomalies (0.5 → 1, 4.5 → 5)

02

### Data Splitting

Train 75 & Test 25,  
stratify=y

03

### Feature Selection

- Removed high-missing-value features
- Selected variables with strong IV
- Dropped multicollinear features ( $\text{corr} \geq 0.7$ )

04

### Handling Outlier

**IQR Winsorization:**  
Capped outliers at IQR boundaries

05

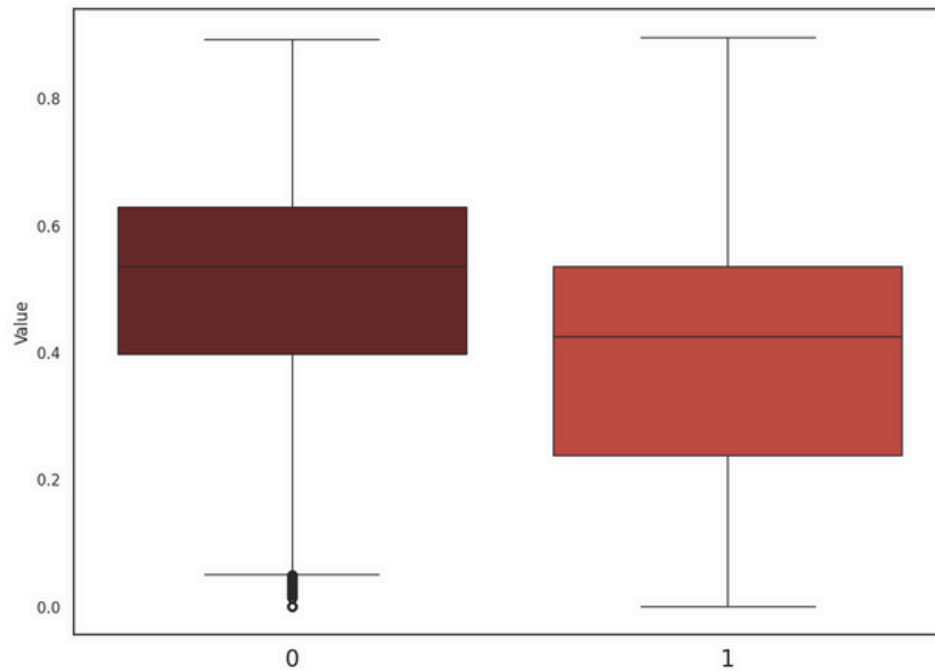
### Feature Encoding & Scalling

Feature **Binning** & **WoE** Transformation

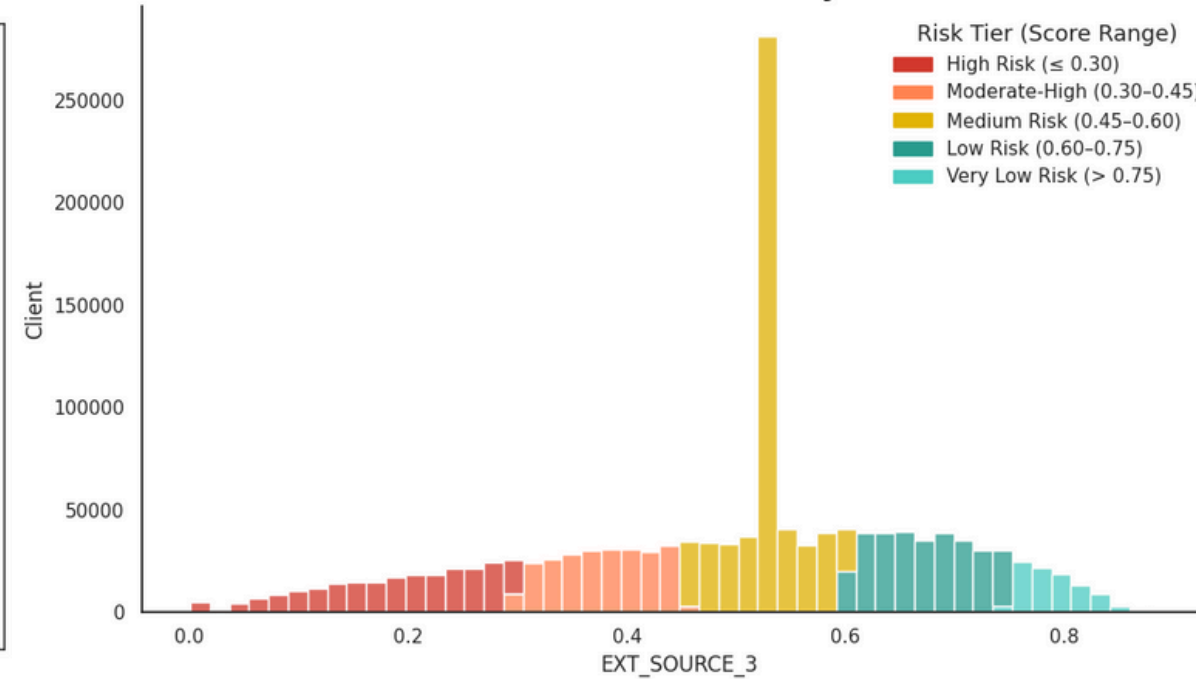


# BUSINESS INSIGHT

Default Rate by EXT\_SOURCE\_3



Distribution EXT SOURCE 3 by Risk Tier



## Insight:

- Clients with higher EXT\_SOURCE\_3 scores (above 0.45) tend to have lower default rates.
- The majority of clients fall into the medium to low risk category (score 0.45 – 0.75).
- This indicates a strong positive correlation between EXT\_SOURCE\_3 and client creditworthiness.

## Recommendation:

- Target segments with EXT\_SOURCE\_3 ≥ 0.45.
- Use AI-powered digital ads & lookalike audiences.
- Partner with relevant digital platforms (Socmed, E-commerce, etc.)



## Insight:

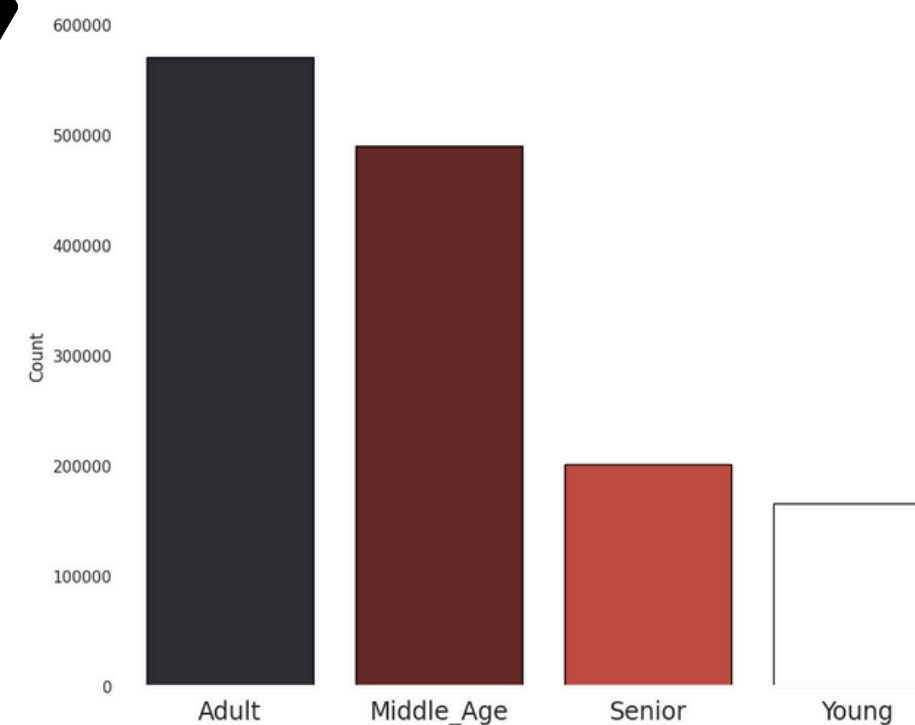
- Most clients are in the Adult and Middle Age groups.
- Seniors show the lowest default rate (5%).
- Young clients have the highest credit risk (12% default rate).

## Recommendation:

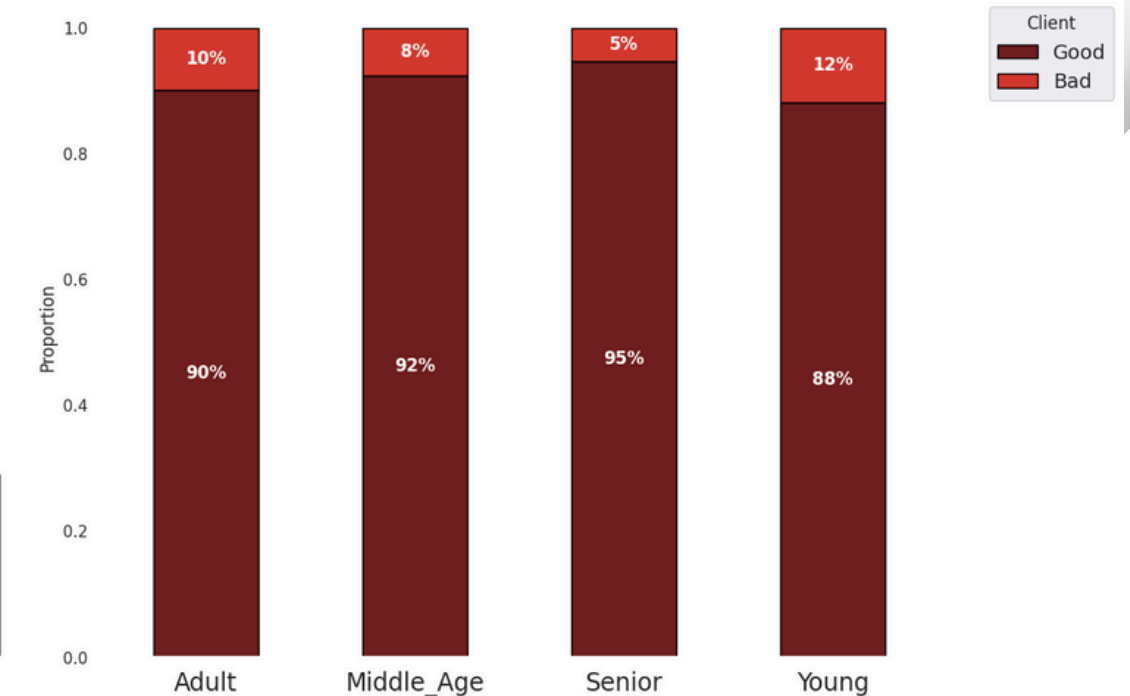
- Focus on retaining Adult and Middle Age clients.
- Use financial education & low-risk products for Young clients.
- Collaborate with platforms popular among younger demographics (social media, fintech apps, etc).



Clients by AGE\_GROUP



Default Rate by AGE\_GROUP



Client  
Good  
Bad

# IMPLEMENTATION & EVALUATION



## Metrics Evaluation : Recall

Recall measures the model's ability to detect high-risk clients (actual positives). A high recall score indicates better performance in capturing most default cases.



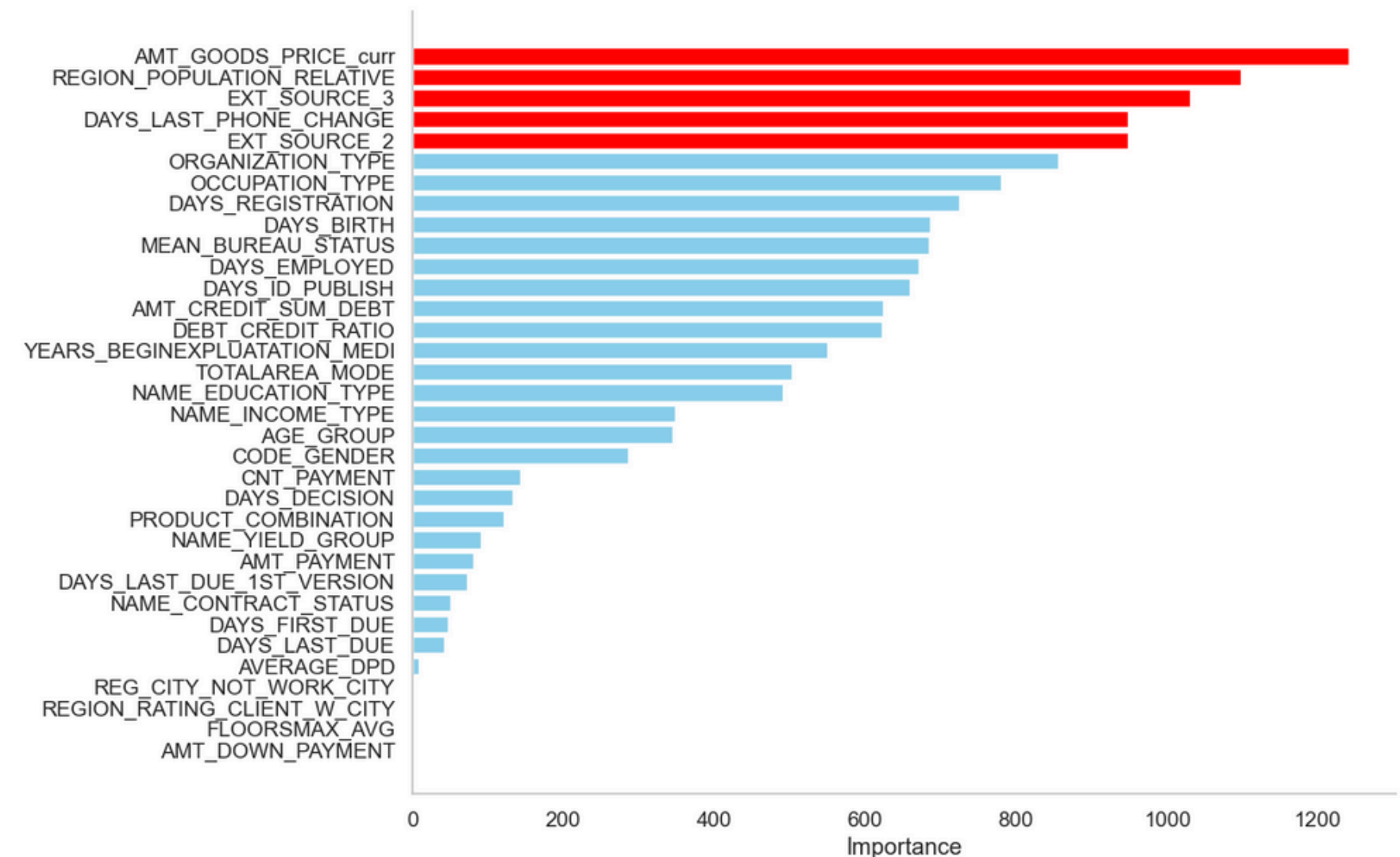
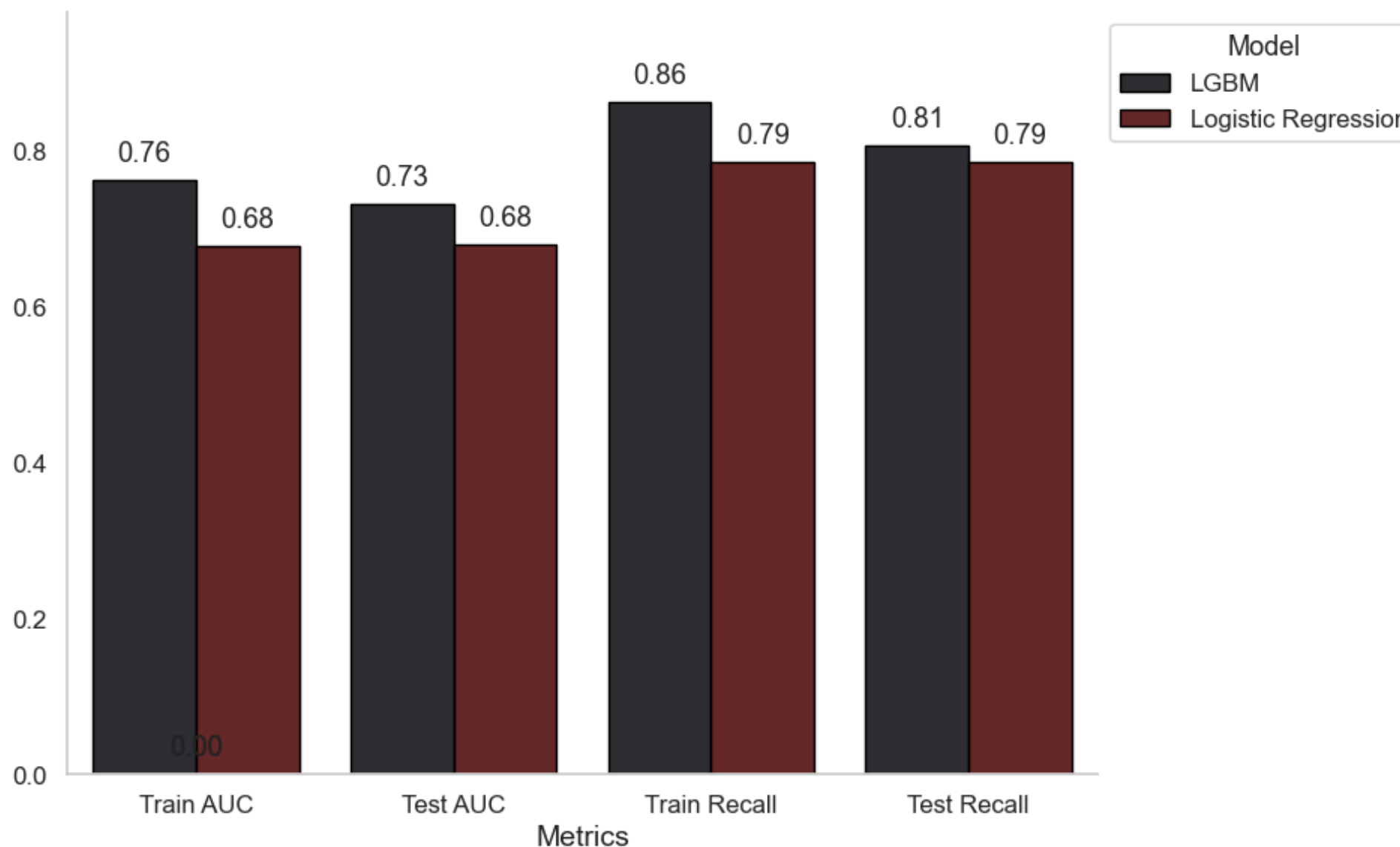
## Feature Importances : LightGBM

**Best Model LightGBM**

**class\_weight='balanced'{'num\_leaves': 150, 'max\_depth': 30}**



## Model Comparison - Hyperparameter Tuning



# IMPLEMENTATION & EVALUATION



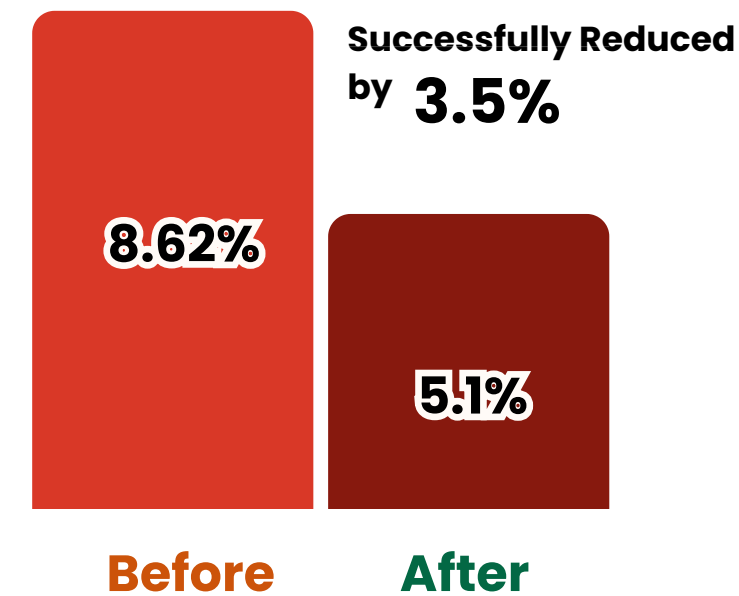
## Impact Simulation



Confusion Matrix LightGBM (Threshold = 0.42)

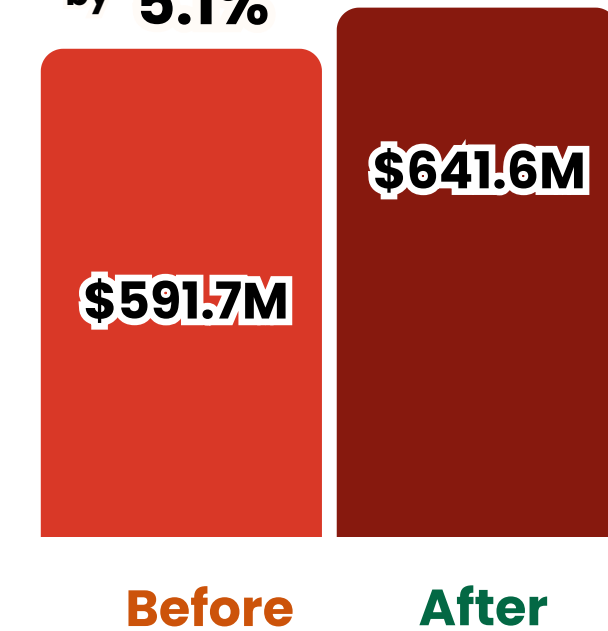
Actual \ Predicted	Good Clients	Bad Clients
Good Clients	<b>TN</b> 212867 (59.54%)	<b>FP</b> 113837 (31.84%)
Bad Clients	<b>FN</b> 5903 (1.65%)	<b>TP</b> 24932 (6.97%)

### Default Rate

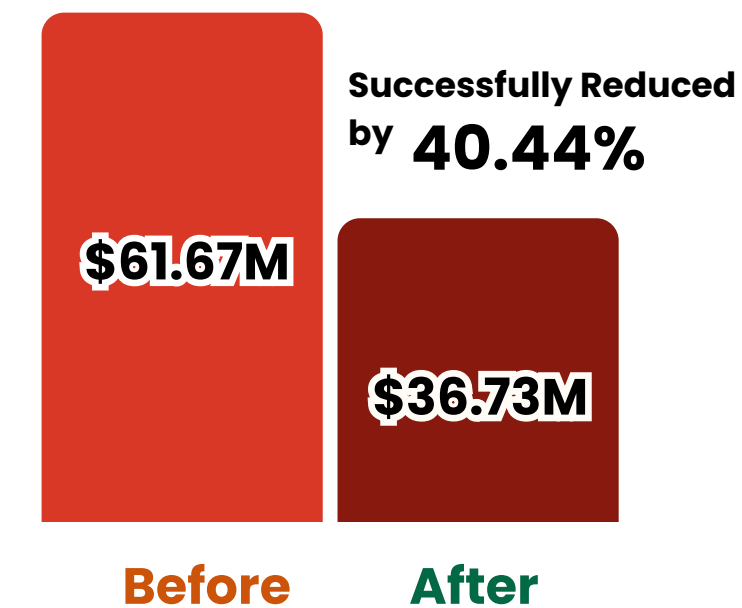


### Increase in Net Revenue (\$)

Successfully Increased by **5.1%**

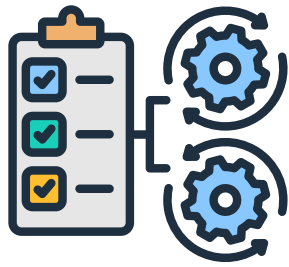


### Cost Saved from Bad Loans (\$)





## Offer Credit Schemes Based on Risk Score & Product Value



**Insight:** Clients with `EXT_SOURCE_3`  $\geq 0.45$  show significantly lower default rates.

Strategy:

- Leverage ML-driven credit scoring (`EXT_SOURCE_3` & `EXT_SOURCE_2`) to personalize loan limits, interest rates, and tenors.
- Automate installment offers based on product price (`AMT_GOODS_PRICE`) and customer risk profile.
- Implement real-time loan approval with transparent interest rates tailored to individual risk scores.
- quality.

## Target High-Population Areas with Localized Campaigns



**Insight:** Customers from densely populated regions (`REGION_POPULATION_RELATIVE`) tend to be more digitally active and lower-risk.

Strategy:

- Use geo-AI and population heatmaps to prioritize marketing and digital expansion areas.
- Launch hyper-local campaigns with customized content (e.g., regional languages, local influencers).
- Deploy digital financial literacy programs in high-potential but underbanked regions.

## Monitor Digital Behavior Changes as Risk Signals



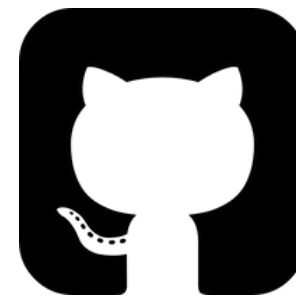
**Insight:** Customers from densely populated regions (`REGION_POPULATION_RELATIVE`) tend to be more digitally active and lower-risk.

Strategy:

- Implement early warning systems based on changes in digital identity (phone, email, device).
- Trigger verification or alerts when sudden digital behavior shifts are detected.
- Combine digital footprint data to strengthen fraud detection and identity validation models.

# Thank You

"Insights are clear, today's decisions shape tomorrow's growth."



**Github**



**Linkedin**



**Website**

