

Gaussian Naive Bayes for Personal Loan Classification

Sarthak Sablania

April 2025

1 Problem Statement

The objective is to identify customers who are likely to accept a personal loan. By doing so, the bank can selectively target these individuals in marketing campaigns, increasing conversion while reducing cost.

2 Dataset Overview

The dataset contains 11 input features and 1 target variable (**Personal Loan**).

Feature Types

- **Continuous Features (5/11):**
 - Age
 - Experience
 - Income
 - CCAvg
 - Mortgage
- **Binary Categorical Features (4/11):**
 - Securities Account
 - CD Account
 - Online
 - CreditCard
- **Ordinal Features (2/11):**
 - Family
 - Education

Target

Personal Loan: 0 or 1

Analysis

I performed exploratory analysis to understand feature distributions, class imbalance, and feature correlations.

- **Data Size:** 5000 Rows
- **Feature Independence and Correlations:** Since **Experience** is highly correlated with **Age**, I removed the **Experience** column and checked the performance of the classifier, but it was the same with and without the column.

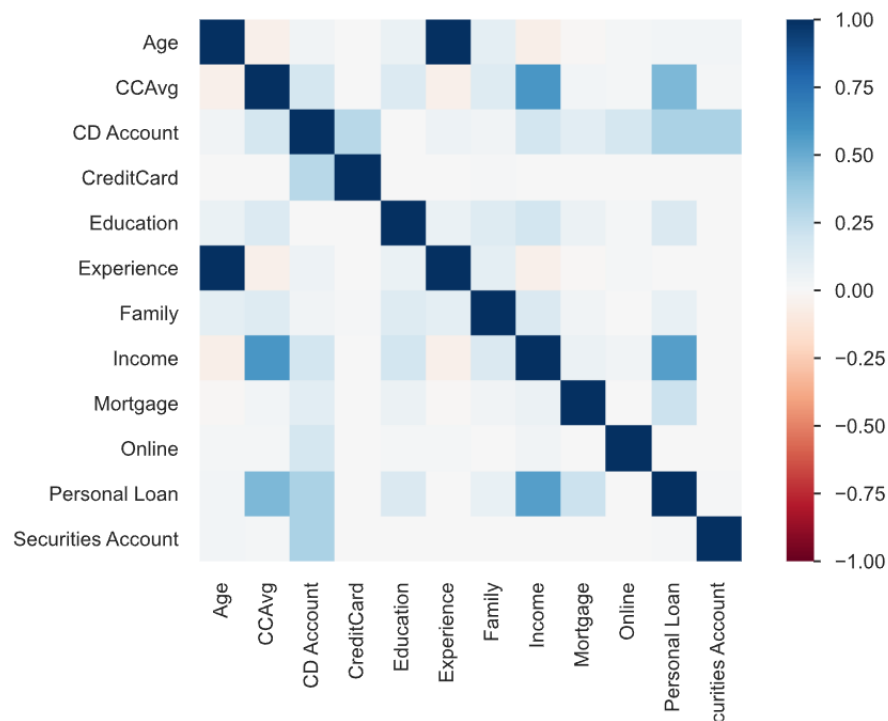


Figure 1: Correlation Matrix of Features and Target.

- **Class Imbalance:** 9.6% of the data is labelled as 1, rest labelled as 0.

3 Results and Threshold Tuning for Business Context

Motivation

The bank prefers **high recall** (capturing as many potential loan-takers as possible) over precision. Missing a potential customer is more costly than targeting an uninterested one.

Strategy

I lowered the classification threshold from the default of 0.5 to 0.02 to increase recall. Therefore, we classify a person as a potential personal loan customer if their predicted probability exceeds 0.02 instead of 0.5.

Results at Different Thresholds

Threshold = 0.02		Threshold = 0.05		Threshold = 0.16	
Metric	Value	Metric	Value	Metric	Value
Recall (Class 1)	0.94	Recall (Class 1)	0.89	Recall (Class 1)	0.80
Precision (Class 1)	0.43	Precision (Class 1)	0.45	Precision (Class 1)	0.50
F1 Score	0.59	F1 Score	0.60	F1 Score	0.62
Customers Targeted	23%	Customers Targeted	20%	Customers Targeted	17%

Threshold = 0.50	
Metric	Value
Recall (Class 1)	0.61
Precision (Class 1)	0.50
F1 Score	0.55
Customers Targeted	13%

Threshold = 0.93	
Metric	Value
Recall (Class 1)	0.46
Precision (Class 1)	0.59
F1 Score	0.52
Customers Targeted	8%

- As threshold decreases, **recall increases** but **precision decreases**.
- At threshold ≈ 0.02 , recall reaches ~ 0.94 , while only $\sim 23\%$ of the population is targeted.
- At threshold ≈ 0.93 , recall ≈ 0.46 , while only $\sim 8\%$ of the population is targeted.

If we choose the threshold of 0.02, we successfully identify 94% of actual positive cases while only targeting 23% of the population, making it a highly focused and cost-effective strategy.

Alternatively, a higher threshold (e.g., 0.93) could be chosen to precisely target the most probable potential customers, depending on business goals.

Results from Various Implementations

Results from all three types of implementations — from scratch, using `scikit-learn`, and AI-generated — were consistent.

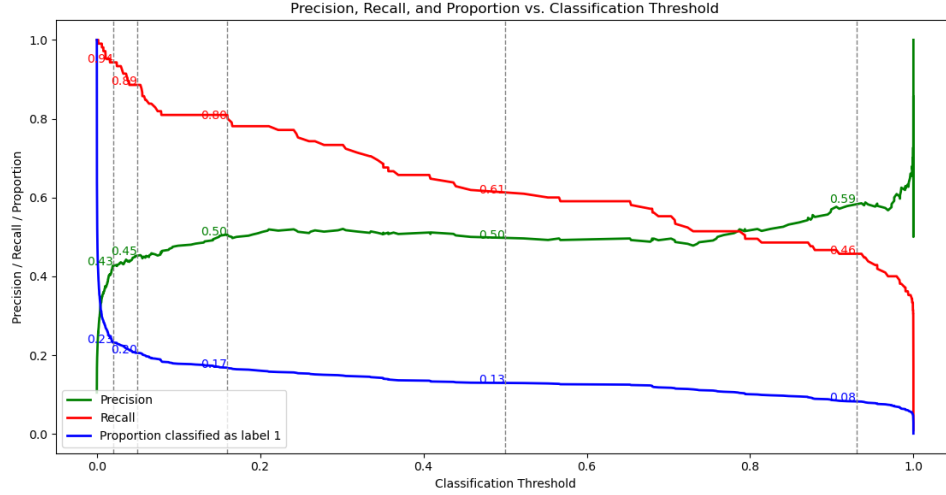


Figure 2: Precision, Recall, and Proportion of Label 1 Predictions vs. Threshold

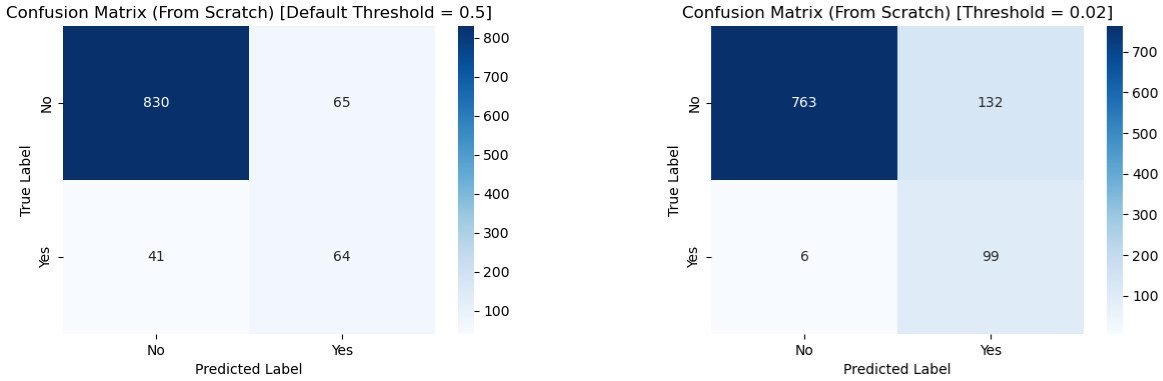


Figure 3: Confusion Matrices for From-Scratch Implementation: (a) Threshold = 0.5 and (b) Threshold = 0.02. We correctly identify 99 out of 105 of the potential PL customers if we choose 0.02 as threshold.

4 Conclusion

- Gaussian Naive Bayes is effective and fast for this task, even though there are 4 binary categorical features and features with a distribution highly deviating from Normal distribution.
- Threshold tuning enabled us to **prioritize recall**, aligning with the business requirement of minimizing missed opportunities.
- The campaign is efficient, targeting a small fraction of customers (23%) while capturing most of the positives (94%).

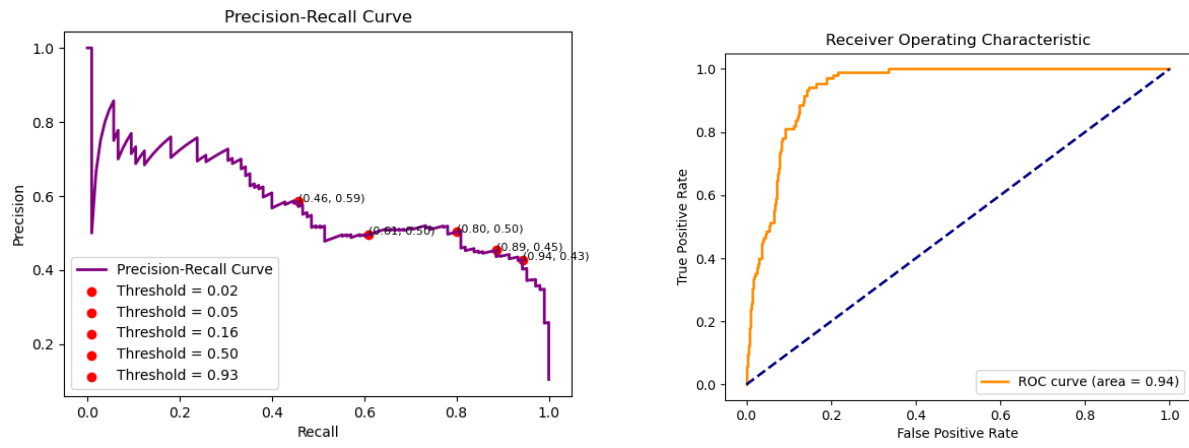


Figure 4: Model Evaluation Plots: (a) Precision-Recall Curve and (b) ROC Curve