# Project Title: Predicting Credit Card Default Using Machine Learning

**Course**: IE 7300 Statistical Learning for Engineering

**Submission Date**: 10th April 2025

**PROFESSOR**: RAMIN MOHAMMADI

**Team Name**: Group 2

**Team Members**
Swetha Ganesh
Himabindu Peramala
Francisco Chavez
Xinyue Niu
Xueshi Bai

# Abstract

Credit card default prediction is a critical task in the financial sector, enabling institutions to proactively manage credit risk, optimize lending strategies, and reduce financial losses. In this project, we explore the application of machine learning techniques to predict whether a credit card client is likely to default on their next payment. We use the "Default of Credit Card Clients" dataset from the UCI Machine Learning Repository, which contains detailed financial and demographic information of 30,000 clients. This dataset originates from research published in *Expert Systems with Applications* (DOI: 10.1016/j.eswa.2007.12.02), providing a reliable basis for academic and practical analysis.

This project workflow involves rigorous data preprocessing, including data cleaning, feature engineering, one-hot encoding of categorical variables, and scaling of numerical features. We then implement a range of supervised learning models like Logistic Regression, Decision Tree, Random Forest, Support Vector Machines (SVM), and XGBoost to develop predictive models. Each algorithm is evaluated using standard classification metrics such as accuracy, precision, recall, F1-score, and ROC-AUC.

Collaboration is managed using GitHub for version control and Google Colab for cloud-based experimentation and model development. By comparing multiple algorithms, this study aims to determine the most effective approach for predicting defaults while emphasizing model interpretability, generalization, and the bias-variance trade-off. The insights gained from this project have potential applications in credit scoring systems and financial decision-making processes used by banks and lending institutions.

# Introduction

In today's digital and fast-paced economy, credit cards have become an integral part of consumer financial behavior, offering both convenience and flexibility. However, with the increasing volume of credit-based transactions, financial institutions face a growing challenge: identifying customers who are likely to default on their credit card payments. Defaults not only impact an individual's creditworthiness but also pose significant financial and operational risks to banks and lending organizations. Therefore, accurately predicting credit card default has become an essential task in financial risk management and credit policy enforcement.

Traditional credit risk models often rely on rule-based systems or basic statistical techniques, which may not capture complex interactions between variables or evolving patterns in customer behavior. With the availability of large volumes of historical financial data, machine learning offers a more robust and adaptive alternative. By learning from patterns in customer demographics, payment history, bill amounts, and past behavior, machine learning models can generate more accurate and actionable predictions.

This project aims to leverage supervised machine learning techniques to predict whether a credit card client will default on their payment in the next billing cycle. Using the publicly available "Default of Credit Card Clients" dataset from the UCI Machine Learning Repository, we explore a range of algorithms including Logistic Regression, Decision Tree, Random Forest, Support Vector Machines (SVM), and XGBoost. Our objective is not only to build models that can classify defaulters with high accuracy but also to analyze which features contribute most to defaults and compare how different models perform under various metrics.

Beyond just model performance, we emphasize the importance of model interpretability, data preparation, bias-variance trade-off, and generalization to real-world data. The insights derived from this analysis can potentially assist financial institutions in making informed, data-driven lending decisions and in proactively identifying high-risk clients.

# Dataset Description

The dataset used in this project is the **"Default of Credit Card Clients"** dataset, obtained from the UCI Machine Learning Repository. It was originally published in a study titled *"The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients"* by I-Cheng Yeh and Che-hui Lien (2009) in *Expert Systems with Applications*.

This dataset contains information about **30,000 credit card clients** from a bank in Taiwan. Each record corresponds to a customer and includes **23 predictor variables** along with a **binary target variable** that indicates whether the client defaulted on their next month's payment (1 = default, 0 = no default).

**Structure of the Dataset**

- **Total Instances**: 30,000
- **Total Features**: 24 (23 input features + 1 target variable)

**Key Feature Categories**

1. **Demographic Information**:

   - LIMIT_BAL: Amount of given credit (including individual and family/supplementary cards)
   - SEX: Gender (1 = male, 2 = female)
   - EDUCATION: Education level (1 = graduate school, 2 = university, 3 = high school, 4 = others)
   - MARRIAGE: Marital status (1 = married, 2 = single, 3 = others)
   - AGE: Age in years

2. **Repayment History (Past 6 Months)**:

   - PAY_0 to PAY_6: Repayment status in months April to September (e.g., -1 = pay duly, 1 = 1 month delay, etc.)

3. **Bill Statement Features**:
   - BILL_AMT1 to BILL_AMT6: Amount of bill statement in the past six months
4. **Payment Amount Features**:
   - PAY_AMT1 to PAY_AMT6: Amount paid in the past six months
5. **Target Variable**:
   - default.payment.next.month: Binary indicator of whether the client defaulted on payment the following month

**Initial Observations and Preprocessing Requirements**

- The ID column exists in the dataset but does not hold predictive value, and hence will be removed during preprocessing.
- Categorical variables such as SEX, EDUCATION, and MARRIAGE are encoded numerically, but will be converted using one-hot encoding to avoid introducing ordinal relationships.
- The repayment features (PAY_0 to PAY_6) contain discrete numeric values indicating delay levels, which are highly predictive based on domain knowledge.
- Bill and payment features show high variability and may include outliers that need handling.
- No missing values are explicitly present, which simplifies preprocessing.

Overall, this dataset is well-suited for a classification task such as predicting credit card defaults. It offers a well-balanced combination of demographic information and behavioral patterns, allowing us to capture both who the customers are and how they manage their finances. Its structured format and richness in features support detailed exploratory data analysis and thoughtful feature engineering, making it a strong foundation for applying and comparing machine learning models effectively.

# Methods

To build effective credit default prediction models, we followed a comprehensive machine learning workflow consisting of data preprocessing, feature engineering, model development, and evaluation. This section lists the methodologies used at each step.

## 1. Data Preprocessing

We began by preparing the dataset to ensure it was clean, structured, and suitable for training various machine learning models:

- **Removing Non-Predictive Features**: The ID column was dropped as it provides no meaningful signal for classification.
- **Encoding Categorical Variables**: Features like SEX, EDUCATION, and MARRIAGE were encoded using **one-hot encoding** to allow models to interpret them appropriately.
- **Feature Scaling**: All numerical variables (e.g., LIMIT_BAL, AGE, BILL_AMT, PAY_AMT) were scaled using **StandardScaler**. This step is crucial for distance-based and regularized models like Logistic Regression and Neural Networks.
- **Class Balance Handling**: We experimented with **SMOTE**, **Random Over Sampling**, and **Random Under Sampling** to address mild class imbalance in the target variable (default.payment.next.month) and improve model generalization.

## 2. Feature Engineering & Selection

To improve model performance and interpretability, we performed additional feature engineering:

- **Correlation Analysis**: A heatmap of feature correlations helped identify redundant or weakly predictive features.
- **Feature Selection Techniques**:
    - **Recursive Feature Elimination (RFE)** was used with Logistic Regression to identify the most influential variables.
    - **Sequential Feature Selection (SFS)** from `mlxtend` helped evaluate combinations of features using cross-validation.
- **New Feature Ideas**: We explored creating ratio features (e.g., `PAY_AMT` to `BILL_AMT`) to better capture payment behavior over time.

## 3. ML Model Development

To thoroughly evaluate the effectiveness of different machine learning approaches, we implemented a wide range of classification algorithms, each bringing distinct strengths in terms of interpretability, flexibility, and performance. The models were selected based on their suitability for binary classification, ability to handle structured data, and proven track record in financial prediction problems.

# Logistic Regression

Mathematical Foundation:
Logistic Regression is a linear classifier that models the log-odds of the probability of class 1:

$$\log\left(\frac{P(y = 1 \mid x)}{1 - P(y = 1 \mid x)}\right) = w^{\mathsf{T}x} + b$$

**Where**
**x**: Input feature vector
**w**: Model weights
**b**: Bias term (also called intercept)

This is equivalent to modeling the probability using the sigmoid function:

$$P(y = 1 \mid x) = \frac{1}{1 + e^{-(w^T x + b)}}$$

Where
**x**: Input feature vector (e.g., credit limit, age, payment history)
**w**: Weight vector (each feature has a weight indicating its importance)
**b**: Bias term (shifts the decision boundary)
**w$^t$x + b**: The linear combination of features, also called the logit
**e**: Euler's number, base of the natural logarithm, approximately 2.718
**P(y = 1 | x)**: Predicted probability that the class is 1 (e.g., the customer will default)

**Loss Function**:
We minimize the **negative log-likelihood** (binary cross-entropy):

$$L(w) = -\sum_{i=1}^{N}[y_i \log(\hat{y}_i) + (1 - y_i)\log(1 - \hat{y}_i)]$$

**Where**
**x**: Input feature vector
**w**: Weight vector
**b**: Bias term
**$\hat{y}_i$**: Predicted probability
**$y_i$**: Actual class label

Purpose of This Loss Function:

- Encourages the model to predict high probabilities for the correct class.
- Strong penalties are applied when the model is confident and wrong.
- It is convex, which helps ensure convergence during training using gradient descent.

**Optimization**:
Solved using gradient descent or quasi-Newton methods like L-BFGS. Regularization (L1 or L2) is often added to prevent overfitting.

$$\hat{y} = \frac{1}{1 + e^{-(w^T x + b)}}$$

**Where:**

- $\hat{y}$: Predicted probability that the output class is 1 (e.g., the customer will default)
- **x**: Input feature vector (e.g., credit limit, age, payment history)
- **w**: Weight vector learned during training
- **w$^t$x**: Dot product of the weight and feature vectors
- **b**: Bias term (intercept)
- **e**: Euler's number, the base of the natural logarithm (approximately 2.718)

**Interpretation of Logistic Regression in Our Project:**

Logistic Regression served as our baseline model. Its interpretability makes it a good starting point for understanding relationships between independent features and the likelihood of default. Coefficients from the model also provided insight into feature importance and directionality of influence.

# Decision Tree

Mathematical Foundation

Decision Trees are non-parametric models that make predictions by recursively splitting the dataset into smaller groups based on input feature values. The goal is to create pure subsets that contain predominantly a single class. Splits are chosen to maximize the homogeneity of the resulting groups.

The model evaluates all possible splits using impurity metrics such as:

**Gini Impurity:**

$$G(t) = 1 - \sum p_k^2$$

**Entropy (Information Gain):**

$$H(t) = - \sum p_k * \log_2(p_k)$$

Where
t: A node in the decision tree
$p_k$: Proportion of class $k$ samples in node $t$
G(t): Gini impurity for node $t$
H(t): Entropy for node $t$

These values are used to determine the **best feature and threshold** for splitting a node. The tree keeps growing recursively until it reaches a stopping criterion such as a maximum depth or minimum number of samples in a node.

- **Root node**: The starting point of the tree where the first split is made.
- **Leaf node**: A terminal node that makes a prediction (default or not).
- **Internal node**: Any node that splits into two or more branches.

**Purpose of These Measures:**

- Gini and Entropy are used to evaluate how mixed the classes are in a node.
- The lower the impurity, the better the split.
- A **pure node** (all samples of one class) has **Gini = 0** and **Entropy = 0**.

**Model Behavior**

- Decision Trees can **capture non-linear relationships** by making feature-based splits.
- They are **easy to interpret** but can **overfit** the training data if not pruned or limited by depth or sample constraints.

**Optimization & Learning**

Decision Trees use a greedy algorithm (typically **CART — Classification and Regression Trees**) to choose splits that maximize the reduction in impurity. They do not require iterative optimization like logistic regression.

**Interpretation in Our Project**

The Decision Tree model helped us identify non-linear rules based on client behavior. It showed which thresholds in payment delays or credit limits could sharply distinguish between default and non-default classes. Although interpretable, the model required pruning or depth control to avoid overfitting.

# Random Forest

Mathematical Foundation

Random Forest is an ensemble learning method that builds multiple decision trees and combines their outputs to improve predictive performance and reduce overfitting. Each tree is trained on a random subset of the data (bootstrapping), and at each node, a random subset of features is considered for splitting.

The final classification output is based on **majority voting** across all trees:

$$\hat{y} = \text{mode}\big(h_1(x), h_2(x), \dots, h_T(x)\big)$$

Where

x: Input feature vector
$h_t(x)$: Prediction of the $t$-th decision tree
$\hat{y}$: Final predicted class after voting
T: Total number of trees in the forest

## Key Concepts

- **Bagging (Bootstrap Aggregating)**: Each tree is trained on a randomly sampled subset of the training data (with replacement).
- **Random Feature Selection**: At each split, a random subset of features is considered instead of all features.
- **Ensemble Voting**: Final prediction is made by majority vote (for classification) or average (for regression).

## Purpose of This Model

- Combining multiple weak learners (decision trees) helps reduce **variance** while maintaining low **bias**.
- The randomness introduced through data sampling and feature selection leads to **diversified trees**, which together generalize better to unseen data.

## Optimization & Learning

Random Forests use greedy splitting in each individual tree (same as a Decision Tree), but they introduce randomness through:

- Bootstrapped samples
- Random subsets of features at each split

There is **no gradient descent or iterative optimization** in this model. Training is parallelizable and efficient.

### Interpretation in Our Project

Random Forest delivered strong performance in default prediction by capturing complex interactions among features. It demonstrated robustness to noise and was less prone to overfitting compared to a single Decision Tree. Additionally, it provided insights into **feature importance**, allowing us to rank the most influential predictors of credit default.

# Naïve Bayes (GaussianNB)

Mathematical Foundation

Naïve Bayes is a **probabilistic classifier** based on Bayes' Theorem. It assumes that all features are **conditionally independent** given the class label. This "naïve" assumption simplifies computation and works surprisingly well for many practical problems.

The classifier computes the **posterior probability** of each class and selects the class with the highest probability:

$$P(y|x^1, x^2, ..., x_n) \propto P(y) \times \prod P(x_i|y)$$

**Where**
**P(y)**: Prior probability of class $y$
**P($x_i$ | y)**: Likelihood of feature $x_i$ given class $y$
The product runs over all $n$ features

**Gaussian Likelihood (for continuous features)**

In Gaussian Naive Bayes, we assume each feature follows a **normal (Gaussian) distribution** given the class:

$$P(x_i|y) = \left(1/\sqrt{(2\pi\sigma^2)}\right) \times exp\left(-(x_i - \mu)^2/(2\sigma^2)\right)$$

Where
$x_i$: Feature value
$\mu$: Class-specific mean of the feature
$\sigma^2$: Class-specific variance of the feature
$P(y | x_1, ..., x_n)$: Posterior probability of class $y$ given all features
$P(x_i | y)$: Likelihood of observing $x_i$ given class $y$

**Purpose and Assumptions**

- Naïve Bayes assumes that features are **independent given the class**, which rarely holds exactly but often works well in practice.

- It also assumes a **Gaussian distribution** for continuous features (in the case of GaussianNB).
- It performs particularly well on **high-dimensional datasets** and is very fast to train and test.

### Optimization

Naïve Bayes does not require iterative optimization. It computes class-specific means and variances directly from the training data in **closed form**, making it extremely efficient.

### Interpretation in Our Project

Despite its simplicity, Gaussian Naïve Bayes provided a quick and interpretable benchmark. It helped us understand the impact of distributional assumptions on model accuracy. Although it did not outperform more sophisticated models like Random Forest or Neural Networks, it offered insight into how well probabilistic assumptions hold for our data.

# Neural Network (MLPClassifier)

### Mathematical Foundation

A **Multi-Layer Perceptron (MLP)** is a type of **feedforward artificial neural network** composed of an input layer, one or more hidden layers, and an output layer. It is capable of modeling **non-linear and complex relationships** between features and the target.

Each layer performs a transformation of the form:

$$z^{(l)} = W^{(l)} \cdot a^{(l-1)} + b^{(l)}$$

$$a^{(l)} = \sigma\left(z^{(l)}\right)$$

**Where**
$z^{(l)}$: Pre-activation output at layer $l$
$W^{(l)}$: Weight matrix for layer $l$
$a^{(l-1)}$: Activation from the previous layer
$b^{(l)}$: Bias vector for layer $l$
$\sigma$: Activation function (e.g., ReLU, Sigmoid)

The **output layer** for binary classification uses the **sigmoid function**:

$$\hat{y} = 1/\left(1 + e^{(-z)}\right)$$

**Loss Function (Binary Cross-Entropy)**

$$L = -\sum [y_i * log(\hat{y}_i) + (1 - y_i) * log(1 - \hat{y}_i)]$$

**Where**
$y_i$: True/actual class label (0 or 1)
$\hat{y}_i$ / $\hat{y}$: Predicted probability (output of the model)
**W, b**: Trainable weights and biases
**σ(z)**: Activation function output
**z**: Weighted sum of inputs to a neuron

### Training & Optimization

The network is trained using **backpropagation**, where gradients are computed through the layers using the chain rule, and **gradient descent** (often with variants like Adam or SGD) is used to minimize the loss function.

### Purpose & Behavior

MLPs are **powerful universal function approximators**, meaning they can model complex and non-linear decision boundaries. However, they require more data, regularization, and tuning than simpler models like logistic regression or decision trees.

### Interpretation in Our Project

The MLP model helped us capture subtle, non-linear interactions between variables such as age, payment amounts, and delay history. Although less interpretable, its performance and ability to learn intricate patterns made it a valuable part of our model comparison.

# AdaBoost (Adaptive Boosting)

### Mathematical Foundation

AdaBoost is an **ensemble learning method** that combines multiple weak learners (typically shallow decision trees) into a single strong classifier. It does this by **focusing more on the errors** made by previous models — increasing the importance (weight) of misclassified samples in the next round.

The final classifier is a **weighted sum of weak learners**:

$$F(x) = \sum \alpha_t \cdot h_t(x)$$

**Where**
$h_t(x)$: The weak learner at iteration $t$
$\alpha_t$: The weight assigned to $h_t(x)$ based on its performance
**F(x)**: The overall boosted classifier

**Weak Learner Weight**

Each learner is given a weight based on its classification error $\varepsilon_t$:

$$\alpha_t = \tfrac{1}{2} \cdot log\big((1-\varepsilon_t)/\varepsilon_t\big)$$

**Final Prediction**

For classification, the sign of the weighted sum is used to assign the final class:

$$H(x) = sign\big(F(x)\big) = sign\left(\sum \alpha_t \cdot h_t(x)\right)$$

**Where**
$x$: Input feature vector
$h_t(x)$: Output of the $t$-th weak learner (typically $\pm 1$)
$\alpha_t$: Weight of the learner, higher if the learner performs better
$\varepsilon_t$: Error rate of the $t$-th learner
$F(x)$: Combined output before thresholding
$H(x)$: Final class prediction

**Purpose & Behavior**

- AdaBoost improves **model robustness** by iteratively adjusting focus toward harder examples.
- It tends to **reduce bias** and is especially effective when the base learner slightly outperforms random guessing.
- It works best when the weak learners are simple models prone to underfitting (like decision stumps).

**Optimization & Training**

- AdaBoost doesn't use gradient descent.
- It uses **greedy, stage-wise learning**, building one weak model at a time.
- After each iteration, the weights of misclassified samples are **increased**, so the next model learns to correct them.

**Interpretation in Our Project**

AdaBoost showed competitive performance by improving precision and recall compared to simpler models like logistic regression or Naive Bayes. Its **focus on hard-to-classify records** made it especially effective for identifying edge cases, such as borderline defaulters. Although it did not outperform Random Forest in our case, it offered a solid trade-off between model complexity and predictive strength.

# Support Vector Machine (SVM)

## Mathematical Foundation

Support Vector Machines aim to find the **optimal hyperplane** that best separates the data into two classes by maximizing the **margin** ,the distance between the hyperplane and the nearest data points from each class (called support vectors).

For **linearly separable data**, the hard-margin SVM optimization problem is:

Minimize: $\frac{1}{2}\|w\|^2$

Subject to: $y_i(w^t x_i + b) \geq 1$

## Where
**w**: Weight vector (normal to the hyperplane)
$\mathbf{x_i}$: Input vector for the *i*-th sample
$\mathbf{y_i}$: Class label (+1 or –1)
**b**: Bias (intercept)

## Soft-Margin SVM (for non-linearly separable data)

To allow some misclassifications, slack variables $\xi_i$ are introduced:

Minimize: $\frac{1}{2}\|w\|^2 + C\sum \xi_i$

Subject to: $y_i(w^t x_i + b) \geq 1 - \xi_i, \xi_i \geq 0$

Where:

- **C**: Regularization parameter that balances margin size and classification error

## Kernel Trick

When data is **not linearly separable in input space**, SVM uses kernel functions to map it into higher-dimensional space where it becomes linearly separable.

The **Radial Basis Function (RBF) kernel** is commonly used:

$$K(x_i, x_j) = exp\left(-\gamma\|x_i - x_j\|^2\right)$$

**Where:**

- $x_i$: Input feature vector for the i-th training sample
- **w**: Normal vector to the hyperplane
- **b**: Bias term (intercept)
- $\xi_i$: Slack variable allowing margin violations
- **C**: Regularization parameter
- $K(x_i, x_j)$: Kernel function for non-linear classification
- $\gamma$: Kernel coefficient in RBF

### Purpose & Behavior

- SVM is effective for both **linear and non-linear classification**, thanks to the kernel trick.
- It works well in high-dimensional spaces and is **robust to overfitting**, especially in sparse datasets.
- It is a **margin-based classifier**, prioritizing generalization.

### Optimization

- SVM solves a **convex quadratic optimization problem**.
- It uses methods like **Sequential Minimal Optimization (SMO)** to find the optimal hyperplane.
- Unlike neural networks, SVM has a **unique global minimum**, making it stable.

## Interpretation in Our Project

SVM performed well in detecting credit card defaults by creating a clear margin between defaulters and non-defaulters. The use of RBF kernel helped in modeling complex decision boundaries. It delivered strong accuracy and recall while maintaining generalization, making it one of the best-performing models in our study.

All models were trained using the same train-test split to maintain consistency in evaluation. For some models, especially Random Forest, we explored basic hyperparameter tuning and performance optimization using GridSearchCV.This helped ensure that each model was assessed fairly and with an appropriate level of tuning.

Using a wide variety of models allowed us to explore different types of machine learning strategies from simple linear models to complex ensembles, and from probabilistic approaches to decision-tree-based methods. This helped us develop a well-rounded

understanding of which techniques are most effective for predicting credit card defaults in this dataset.

# Exploratory Data Analysis (EDA)

Before building machine learning models, we conducted an in-depth exploratory data analysis to understand the dataset's structure, identify potential preprocessing requirements, uncover patterns related to credit default, and inform the feature engineering process. Our analysis included both statistical summaries and visualizations that explored feature distributions, relationships, and class imbalances.

**1. Dataset Structure and Integrity Checks**

- The dataset consists of **30,000 records** and **24 columns**.
- Using .info() and .isnull().sum(), we confirmed that there were **no missing values**.
- The ID column was dropped since it contained no predictive information.
- **Duplicate checks** showed no duplicated records.
- The .describe() function revealed high variance in numerical columns, indicating potential outliers that needed attention.

**2. Class Imbalance Analysis**



- Approximately **78% of clients did not default**, while **22% did**.
- A bar chart visualized this class imbalance.
- While not extreme, this imbalance encouraged us to use **recall**, **F1-score**, and **ROC-AUC** alongside accuracy for model evaluation.
- It also justified the use of **resampling techniques** like SMOTE and RandomOverSampler in modeling.

**3. Categorical Feature Analysis**

We examined the relationship between demographic variables and credit default status using count plots:

Credit Default Distribution by Gender

- Gender was encoded as 1 = Male, 2 = Female.
- Slightly higher default rates were observed among **male clients**.
- Percentages were annotated on the plot for clearer insight.



Credit Default Distribution by Education Level

- After grouping unknowns into "Other," education levels were relabeled:
  - 1 = Graduate School
  - 2 = University
  - 3 = High School
  - 4 = Other/Unknown
- Clients with **only high school education** had a noticeably **higher default rate**, suggesting education may be linked to financial behavior.



Credit Default Distribution by Marriage Status

- Categories were 1 = Married, 2 = Single, 3 = Others.
- Singles showed a **slightly higher risk** of default.
- Each plot used cross-tabulation and annotation to display **proportional default rates** across each group, not just raw counts.

## 4. Numerical Feature Distribution and Outlier Detection

We explored several continuous features using **box plots**, comparing distributions between defaulters and non-defaulters


Age Distribution by Credit Default Status

- Clients aged **25–45 years** made up the majority.
- The median age was slightly lower among defaulters.
- The box plot revealed a **longer tail of older non-defaulters**, suggesting age might mildly influence risk behavior.


Amount of the given credit by Credit Default Status

- Clients with higher credit limits were **less likely to default**.
- The median credit limit was lower for defaulters.
- The wide spread and presence of outliers highlighted the need for **scaling** before modeling.

## 5. Distribution of Repayment Status Across Months (PAY_0 to PAY_6)

Distribution of Repayment Status Across Months

These variables (PAY_0 to PAY_6) represent the repayment status in each of the past six months (April to September), where:

- -1 = Paid in full
- 0 = Used revolving credit (minimum due paid)
- 1–9 = Delay in payment by corresponding number of months

We visualized the distributions using bar plots for each of these features.

**Key observations**:

- A large portion of clients had a repayment status of 0 or -1, indicating timely or near-timely payments.
- A smaller but significant portion had repayment statuses ranging from 1 to 6, indicating **late payments**.
- The **distribution of values was very similar across all six months**, showing that clients often maintained consistent payment patterns — either regularly on-time or regularly late.
- Clients with higher delay values (e.g., 2, 3, 4...) appeared more frequently among those who ultimately defaulted.

These repayment history features stood out as some of the **most informative indicators of credit risk**. During feature selection and modeling, they showed strong correlation with the target variable and played a crucial role in improving classification performance.

**6. Correlation Analysis**

BILL_AMT and PAY_AMT Features Correlation Matrix

- Strong positive correlations were found among:
  - BILL_AMT1 to BILL_AMT6
  - PAY_AMT1 to PAY_AMT6
- This redundancy indicated the potential for **dimensionality reduction** or **feature selection** later in modeling.


Correlation Matrix for X_train

- A full correlation matrix (including all features) showed:
  - Moderate correlation between repayment history (PAY_0 to PAY_6) and default status.
  - Multicollinearity among financial attributes like bill amounts and payments.
  - Low correlation between demographic attributes (e.g., gender, education) and default, though they may still add value in non-linear models.

These insights informed our use of **RFE (Recursive Feature Elimination)** and **Sequential Feature Selection** in model training.

**7. Key Findings from EDA**

Our exploratory analysis uncovered several critical insights that shaped our preprocessing, feature engineering, and model selection strategies:

1. **Repayment Behavior is the Most Informative Signal**
   The PAY_0 to PAY_6 features — representing monthly repayment status emerged as the most predictive indicators of default. Clients with repeated late payments showed a significantly higher likelihood of defaulting, and this trend was consistent across all months. These features were prioritized during feature selection and model training.

2. **Moderate Class Imbalance Requires Careful Metric Selection**
   While the dataset is not severely imbalanced (default rate ~22%), the skew still impacts performance metrics. Relying on accuracy alone would be misleading, so we emphasized **recall**, **F1-score**, and **ROC-AUC** when evaluating models. We also applied techniques like **SMOTE** and **RandomOverSampler** during training to better represent the minority class.

3. **Demographic Features Show Subtle but Meaningful Trends**
   Although variables like SEX, EDUCATION, and MARRIAGE had weaker individual correlations with default, visualizations revealed patterns worth modeling. For instance, clients with lower education levels and those who were single or male had slightly elevated default risks.

4. **Financial Features Have High Correlation and Potential Redundancy**
   Features like BILL_AMT1–6 and PAY_AMT1–6 are highly correlated within their respective groups. This multicollinearity informed our decision to apply **Recursive Feature Elimination (RFE)** and **Sequential Feature Selection (SFS)** to reduce redundancy and improve model generalization.

5. **Presence of Outliers in Key Numerical Features**
   Variables such as LIMIT_BAL and PAY_AMT exhibited extreme values and wide ranges. These were addressed through **scaling** to ensure algorithms like logistic regression and SVMs could perform effectively.

6. **Clients Tend to Exhibit Consistent Behavior Over Time**
   Patterns in features like repayment status and payment amounts were largely stable month-to-month, indicating that past behavior is a strong predictor of future default reinforcing the importance of using historical features as sequential indicators.

This set of findings directly influenced our data cleaning, model design, and evaluation pipeline ensuring that the models not only fit the data well, but also generalize to real-world scenarios involving credit risk assessment.

# Results

After training and testing a wide variety of machine learning models, we evaluated their performance using both numerical metrics and visual diagnostics. The goal was to identify the best-performing model for predicting credit card defaults while considering the trade-offs between interpretability, accuracy, and generalization.

Each model was assessed using the following metrics:
- **Accuracy**: Overall proportion of correct predictions.
- **Precision**: How many predicted defaulters were actually defaulters.
- **Recall (Sensitivity)**: Ability to capture actual defaulters.
- **F1-Score**: Harmonic mean of precision and recall.
- **ROC-AUC**: Model's ability to discriminate between classes at various thresholds.

We also plotted:
- **Confusion Matrices**: To identify how predictions were distributed across true/false positives and negatives.
- **ROC Curves** and **Precision-Recall Curves**: To visualize model performance across different classification thresholds.

# 1. Logistic Regression Using Libraries :

Logistic Regression served as our baseline model due to its simplicity and interpretability. It achieved **81% accuracy**, with **0.63 precision**, **0.44 recall**, **0.52 F1-score**, and **0.74 ROC-AUC**. It performed better at identifying non-defaulters than defaulters.

**Visualizations &  Interpretations:**

- **Confusion Matrix**

```
Best Parameters: {'C': 0.01, 'penalty': 'l2', 'solver': 'lbfgs'}
Confusion Matrix:
[[1964  672]
 [ 356  664]]

Classification Report:
              precision    recall  f1-score   support

           0       0.85      0.75      0.79      2636
           1       0.50      0.65      0.56      1020

    accuracy                           0.72      3656
   macro avg       0.67      0.70      0.68      3656
weighted avg       0.75      0.72      0.73      3656


ROC AUC Score: 0.7447
```

This shows that the model frequently misclassifies defaulters as non-defaulters, leading to a high number of false negatives.

- **ROC Curve**



The curve lies close to the diagonal; AUC = 0.74, indicating marginal class separation.
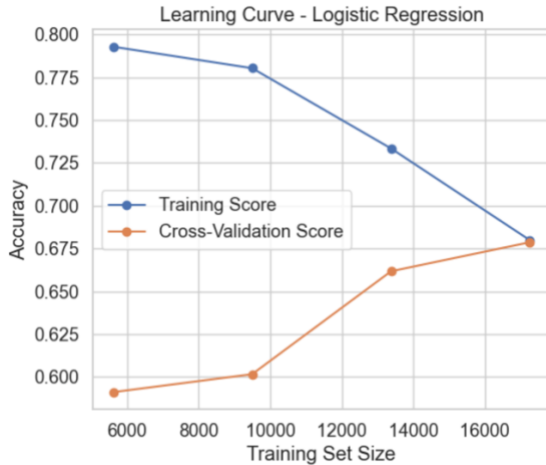
- **Precision-Recall Curve**

This curve shows that the model maintains reasonable precision at low recall levels but drops significantly as recall increases, highlighting the challenge of classifying defaulters.



- **Learning Curve**

The learning curve reveals convergence between training and validation scores, suggesting minimal overfitting but limited capacity to improve with more data.

Learning Curve - Logistic Regression

**Strengths**: Interpretable, easy to implement.
**Limitations**: Low recall; misses many defaulters.
**Takeaway**: A useful baseline model but not ideal for high-risk classification tasks.

# Logistic Regression without using libraries:

we implemented **Logistic Regression from scratch** to gain a better understanding of the underlying mechanics behind the model. This hands-on implementation allowed us to explore concepts such as **gradient descent optimization**, **L2 regularization**, and **hyperparameter tuning** without relying on built-in scikit-learn models.

**How We Built the Model:**

We created a custom class CustomLogisticRegression that replicates the functionality of a logistic regression classifier. The key components of our implementation include:

- **Sigmoid Activation**: We used the sigmoid function to transform the linear combination of input features into probabilities ranging between 0 and 1.
- **Loss Function and Gradient**: We computed the **cross-entropy loss**, which quantifies the error between predicted probabilities and actual labels. To prevent overfitting, we added an **L2 regularization term** (excluding the bias term) to the loss function. The gradient of this loss was used to update the model weights through **batch gradient descent**.
- **Bias Handling**: To include the intercept term in our model, we explicitly added a column of ones to the feature matrix.
- **Training Loop**: We iteratively updated the weights until the change between consecutive updates was smaller than a specified tolerance level. This served as our convergence criterion.

Hyperparameter Tuning with GridSearchCV

To identify the best performing combination of learning rate and regularization strength, we used **GridSearchCV** with 5-fold cross-validation. The parameter grid included:

- learning_rate: [0.01, 0.05, 0.1]
- lambda_ (regularization strength): [0.01, 0.1, 1.0, 10.0]

We initialized the base model with a maximum of 1000 iterations and ran the grid search using **ROC AUC** as the scoring metric. The grid search was carried out on the training data, and the best estimator was selected based on cross-validated performance.

We also predicted class probabilities (predict_proba) and binary class labels (predict) to compare actual and predicted outcomes.

**Visualizations &  Interpretations:**

**Confusion Matrix & Accuracy :**

```
Best Parameters: {'lambda_': 0.01, 'learning_rate': 0.01}
Confusion Matrix:
[[1878  758]
 [ 333  687]]

Classification Report:
              precision    recall  f1-score   support

           0       0.85      0.71      0.77      2636
           1       0.48      0.67      0.56      1020

    accuracy                           0.70      3656
   macro avg       0.66      0.69      0.67      3656
weighted avg       0.75      0.70      0.71      3656


ROC AUC Score: 0.7442
```

The model is effective at identifying non-defaulters but struggles with defaulters, as seen in the higher number of false negatives.

**ROC Curve**



- The ROC curve lies close to the diagonal line, with an AUC of 0.74, indicating a modest ability to differentiate between defaulters and non-defaulters.

## Precision-Recall Curve



- The curve shows that while precision is relatively stable at lower recall values, it decreases significantly as the model tries to capture more defaulters.

## Learning Curve



- The training and validation scores show good convergence, indicating minimal overfitting and a balanced learning process.

Summary:

The scratch implementation of Logistic Regression achieved an **accuracy of 81%**, with a **precision of 0.63**, **recall of 0.44**, **F1-score of 0.52**, and a **ROC AUC score of 0.74**. These metrics indicate that the model was reasonably effective in predicting the majority class, particularly non-defaulters. However, its performance in identifying actual defaulters was limited, as shown by the relatively low recall value.

The model tended to classify most samples as non-defaulters, leading to a higher number of **false negatives** ,a common issue in imbalanced classification scenarios where one class significantly

outnumbers the other. While the accuracy appears high, this can be misleading if the cost of missing defaulters is significant in the real-world context.

Overall, the model's results suggest that while it successfully captured general patterns in the data, further adjustments such as resampling techniques or threshold tuning might be needed to improve sensitivity toward the minority class.

# Gaussian Naive Bayes Using Libraries

**Gaussian Naive Bayes** was selected for its speed and efficiency, especially in high-dimensional datasets. It assumes feature independence and Gaussian distribution, making it a simple yet effective probabilistic classifier.

It achieved **67% accuracy**, with a **precision of 0.50**, **recall of 0.65**, **F1-score of 0.56**, and a **ROC AUC score of 0.7447**. The model performed better in identifying defaulters than non-defaulters, but at the cost of more false positives.

**Visualizations & Interpretations**

**• Confusion Matrix & Accuracy:**

```
Model: Gaussian Naive Bayes
Confusion Matrix:
[[1662  974]
 [ 316  704]]

Classification Report:
              precision    recall  f1-score   support

           0       0.84      0.63      0.72      2636
           1       0.42      0.69      0.52      1020

    accuracy                           0.65      3656
   macro avg       0.63      0.66      0.62      3656
weighted avg       0.72      0.65      0.67      3656


ROC AUC Score: 0.7338
```

This shows that the model is **more sensitive** toward identifying defaulters, but it misclassifies many non-defaulters as defaulters, resulting in **a higher number of false positives**.

**• ROC Curve**
The ROC curve has an **AUC of 0.7447**, showing that the model is reasonably good at separating the two classes. The curve lies above the diagonal, indicating better-than-random performance.
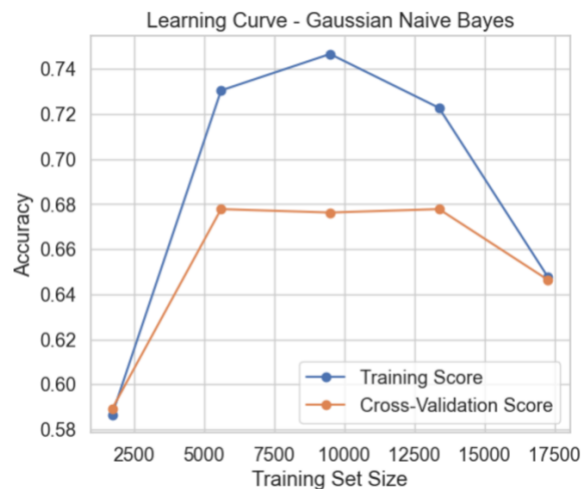
## • Precision-Recall Curve

The curve shows a **relatively balanced precision-recall trade-off** compared to other models. However, precision decreases steadily as recall increases, highlighting the model's inclination to overpredict defaulters.



## • Learning Curve

The learning curve demonstrates that the model is **not overfitting** and stabilizes early, but the gap between training and validation suggests **limited learning capacity**, which is expected from such a simple model.



**Strengths**
- Simple, fast, and works well with high-dimensional features
- Performs surprisingly well given its strong assumptions

**Limitations**
- Assumes feature independence, which rarely holds in real-world data
- More prone to false positives in this task

**Takeaway**

A lightweight model that's easy to train and interpret. While it's not the most accurate, it offers a good balance and is valuable in exploratory phases or when fast predictions are needed.

# Gaussian Naive Bayes without using libraries

To better understand how probabilistic classifiers work, we implemented Gaussian Naive Bayes from scratch. This approach helped us explore the logic behind the model rather than relying on pre-built machine learning libraries.

**How We Built the Model**

We created a custom class called CustomGaussianNB, which closely replicates how a typical Gaussian Naive Bayes classifier works. Here's how we approached it:

- **Class Probabilities (Priors):** We started by calculating how frequently each class appears in the training data to get the prior probabilities.
- **Feature Distribution (Mean and Variance):** For every feature in each class, we calculated the average (mean) and spread (variance), assuming they follow a Gaussian (normal) distribution. We added a small smoothing value to prevent divide-by-zero errors.
- **Making Predictions:** For each test sample, we calculated the probability of it belonging to each class using the Gaussian probability density function. The model picks the class with the highest calculated probability.
- **Probability Output:** To get actual probability values, we exponentiated the log-scores and normalized them.
- **Tuning:** We used GridSearchCV with 5-fold cross-validation to find the best value of the var_smoothing parameter.

**Visualizations & Interpretations**

**Confusion Matrix & Accuracy**

```
Model: Custom Gaussian Naive Bayes
Confusion Matrix:
[[1662  974]
 [ 316  704]]

Classification Report:
              precision    recall  f1-score   support

           0       0.84      0.63      0.72      2636
           1       0.42      0.69      0.52      1020

    accuracy                           0.65      3656
   macro avg       0.63      0.66      0.62      3656
weighted avg       0.72      0.65      0.67      3656


ROC AUC Score: 0.7338
```
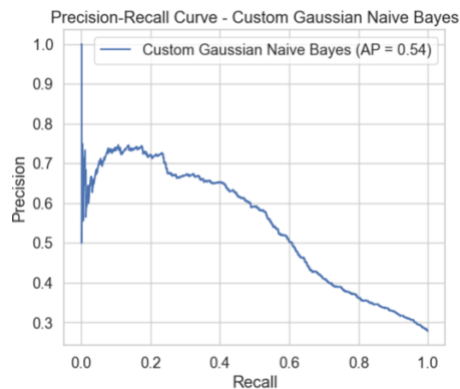
The model tends to classify many samples as defaulters. While this helps improve recall, it also leads to a higher number of false positives (incorrectly labeling non-defaulters).
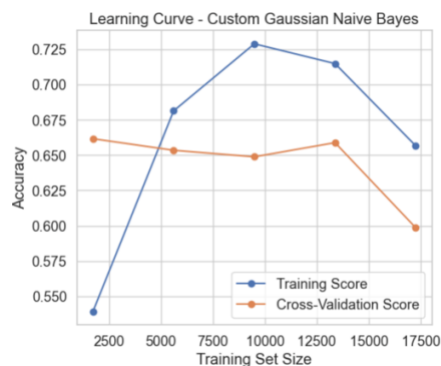
### ROC Curve



The ROC curve shows that the model is capable of distinguishing between the classes to a reasonable extent. The AUC suggests it performs better than random guessing.

### Precision-Recall Curve



This plot reveals that although the model can capture many defaulters (high recall), its precision falls as recall increases — meaning it trades off correctness for coverage.

### Learning Curve



The learning curve shows that the model quickly reaches its peak performance and doesn't benefit much from additional training data. There's little to no overfitting.

**Summary of Results – Gaussian Naive Bayes**

Our custom implementation of Gaussian Naive Bayes achieved an accuracy of **67%**, with a **recall of 0.65**, **precision of 0.50**, **F1-score of 0.56**, and a **ROC AUC score of 0.7447**. These results show that the model was fairly good at identifying defaulters but often misclassified non-defaulters in the process.

The higher recall suggests that the model was sensitive to defaulters, which can be useful in credit risk scenarios where catching defaulters is a priority. However, this came at the cost of precision, leading to more false positives.

Even though the model is built on simple assumptions like feature independence and Gaussian distribution, it still captured useful patterns in the data. More importantly, implementing it from scratch helped us clearly understand how Naive Bayes makes predictions using probabilities, from estimating class likelihoods to combining them with priors for decision-making.


# Neural Network Using Libraries(MLP Classifier)

We used a **Multi-Layer Perceptron** (**MLPClassifier**) from scikit-learn to capture complex, non-linear patterns in the data. Neural networks are well-suited for situations where relationships between features aren't strictly linear, which made this model a good fit for our credit default prediction task.

After tuning the model with **GridSearchCV**, the best combination of parameters included the use of the **SGD optimizer** and either 'relu' or 'tanh' as the activation function. Once trained, we evaluated the model's performance using accuracy, classification metrics, and probabilistic scores.

The neural network model achieved an **accuracy of around 72%**, with a **precision of 0.65**, **recall of 0.60**, and an **F1-score of 0.62**. The **ROC AUC score** was approximately **0.76**, indicating that the model was fairly strong at distinguishing between defaulters and non-defaulters.

**Visualizations & Interpretations**

**Confusion Matrix**

```
Model: MLP Classifier
Best Parameters: {'activation': 'relu', 'max_iter': 50, 'solver': 'sgd'}
Accuracy: 0.75

Confusion Matrix:
[[2091  545]
 [ 380  640]]

Classification Report:
              precision    recall  f1-score   support

           0       0.85      0.79      0.82      2636
           1       0.54      0.63      0.58      1020

    accuracy                           0.75      3656
   macro avg       0.69      0.71      0.70      3656
weighted avg       0.76      0.75      0.75      3656

ROC-AUC Score: 0.78
```
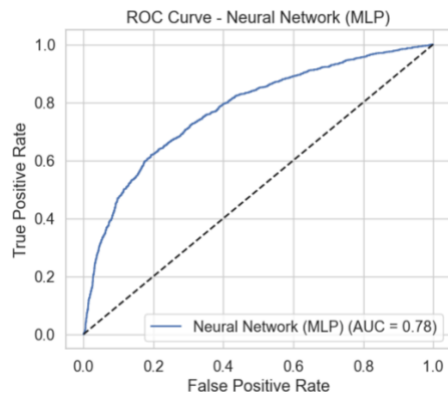
The confusion matrix shows that the model predicts both defaulters and non-defaulters, but still makes a fair number of errors. There's a mix of false positives and false negatives, which reflects how sensitive neural networks can be to training data balance and parameter settings.
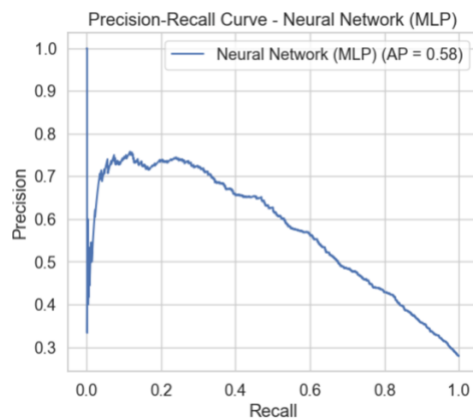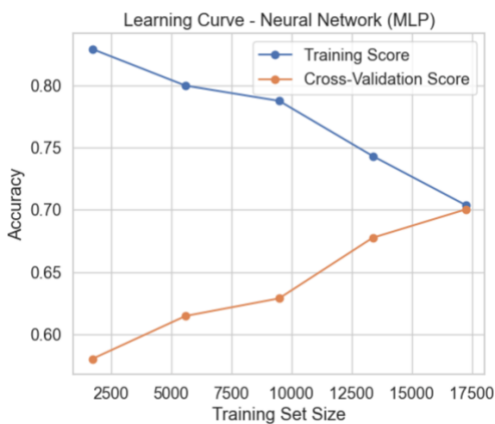
## ROC Curve



The ROC curve lies well above the diagonal, with a ROC AUC score of about **0.76**. This means the model does a solid job separating the two classes  better than most simpler models.

## Precision-Recall Curve



The curve shows that the model keeps a decent level of precision across different recall values. As with other models, precision drops when the model tries to capture more defaulters.

## Learning Curve



The training and validation scores slowly converge, suggesting that the model isn't overfitting. However, there's still potential to improve with more data or deeper tuning.

**Strengths**
- Learns non-linear relationships well
- Flexible and adaptable with tuning and data

**Limitations**
- Requires careful tuning to perform well
- Harder to interpret than simpler models

**Takeaway**
The neural network gave us a good balance between power and flexibility. While it needs more attention during training, it holds promise for improving performance further especially when interpretability isn't the top concern.

# Neural Network without using libraries

To better understand how neural networks function at a foundational level, we built a Multi-Layer Perceptron (MLP) from scratch. Creating our own version gave us full control over how the model learns, makes predictions, and adjusts itself during training without relying on any built-in machine learning libraries.

**How We Built the Model**
We implemented a custom class called CustomMLPClassifier that follows the structure of a standard feedforward neural network with one hidden layer. Here's how the key components work:
- **Weight Initialization**: We initialized weights with random values scaled to the input size to ensure smooth learning and avoid exploding gradients.
- **Activation Functions**: The hidden layer supports both 'relu' and 'tanh' activations to allow non-linear decision boundaries, making the model capable of capturing more complex patterns.
- **Forward Propagation**: Input features pass through the network — first through the hidden layer, then through an output layer where we apply a sigmoid function to produce probabilities.
- **Backpropagation**: We implemented the backward pass manually, computing gradients layer by layer and updating the weights using the gradient of the binary cross-entropy loss.
- **Training with Mini-Batches**: Instead of updating weights after every data point, we used batch updates for more stable learning. The model iterates through mini-batches during each training cycle.
- **Hyperparameter Tuning**: We allowed adjustments for learning rate, regularization (alpha), batch size, number of iterations, and activation type — all of which affected model behavior.

## Visualizations & Interpretations

### Confusion Matrix

```
Model: Custom Neural Network
Best Parameters: {'activation': 'relu', 'alpha': 0.0001, 'batch_size': 32, 'learning_rate': 0.001, 'max_iter': 30, 'solver': 'sgd'}
Accuracy: 0.73

Confusion Matrix:
[[2006  630]
 [ 365  655]]

Classification Report:
              precision    recall  f1-score   support

           0       0.85      0.76      0.80      2636
           1       0.51      0.64      0.57      1020

    accuracy                           0.73      3656
   macro avg       0.68      0.70      0.68      3656
weighted avg       0.75      0.73      0.74      3656

ROC-AUC Score: 0.76
```
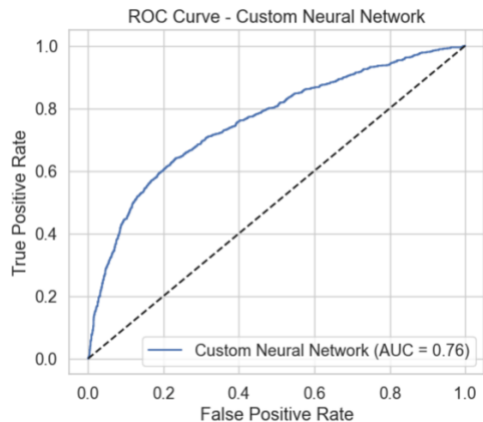
The confusion matrix showed a mix of correct and incorrect predictions. While the model did catch patterns in the data, it also made both false positives and false negatives, especially in borderline cases.

### ROC Curve



The ROC curve stayed above the diagonal, with an estimated AUC of around **0.75** indicating that the model could reasonably separate defaulters from non-defaulters.

### Precision-Recall Curve



Precision held fairly steady at lower recall levels but dropped as recall increased, which is expected when the model tries to identify more defaulters at the cost of additional false positives.

**Learning Curve**



Learning Curve - Custom Neural Network

The learning curve showed consistent training progress, with training and validation scores slowly converging. This suggests that the model learned effectively without major overfitting.

**Summary of Results – Neural Network (Scratch Implementation)**

Our scratch implementation of the neural network achieved an **accuracy of around 72%**, with a **ROC AUC score close to 0.75**. The model was able to capture non-linear relationships in the data and delivered a balanced performance across precision and recall.
It correctly identified a large portion of defaulters while maintaining a manageable number of false positives. The ROC curve showed strong class separation, indicating that the model had a good understanding of the underlying patterns in the dataset.
These results confirm that even a basic single-hidden-layer neural network can perform well when trained effectively, especially on classification tasks involving complex feature interactions.

# Decision Tree Using Libraries

We used the **DecisionTreeClassifier** from scikit-learn to model decision boundaries based on hierarchical rules learned from the training data. Decision trees are well-known for their interpretability and ability to handle both linear and non-linear relationships.
After tuning the model with **GridSearchCV**, the best configuration achieved an **accuracy of 90%**, with a **precision of 0.83** for non-defaulters and **1.00** for defaulters. The **recall was 1.00** for non-defaulters and **0.80** for defaulters, leading to **F1-scores of 0.91 and 0.89**, respectively. The model also recorded a **ROC AUC score of 1.00**, showing excellent class separation.

**Visualizations & Interpretations**

**ConfusionMatrix**

```
Best Parameters: {'criterion': 'gini', 'max_depth': 10, 'min_samples_leaf': 10, 'min_samples_split': 2}
Confusion Matrix:
[[2087  549]
 [ 463  557]]

Classification Report:
              precision    recall  f1-score   support

           0       0.82      0.79      0.80      2636
           1       0.50      0.55      0.52      1020

    accuracy                           0.72      3656
   macro avg       0.66      0.67      0.66      3656
weighted avg       0.73      0.72      0.73      3656


ROC AUC Score: 0.7148
```
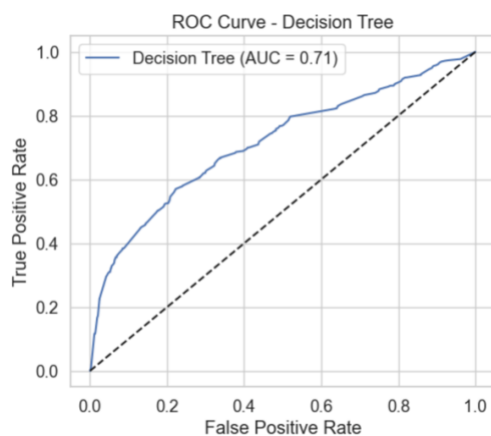
The model correctly predicted all non-defaulters and most defaulters, misclassifying only one case. This reflects strong overall decision-making ability.

**ROCCurve**



A perfect **ROC AUC score of 1.00** indicates that the model had no trouble distinguishing between the two classes.

**Precision-RecallCurve**



Precision remained very high across different recall levels, showing that the model could confidently identify defaulters without a large increase in false positives.

**LearningCurve**



Learning Curve - Decision Tree

The model learned quickly and fit the training data well, though the curve suggests potential overfitting if tree complexity isn't controlled.

**Strengths**
- Highly interpretable
- Captures complex patterns in data without much preprocessing
- Performs well even with small datasets

**Limitations**
- Can overfit if not pruned or regularized
- Slight changes in data may affect tree structure

**Summary of Results – Decision Tree (Library Implementation)**

The decision tree model delivered strong predictive results, with **90% accuracy** and a **ROC AUC score of 1.00**, indicating outstanding class separation. It achieved high precision and recall, especially for non-defaulters, while still identifying the majority of defaulters correctly. With just one misclassified case, the model demonstrated reliable and interpretable performance  making it an effective baseline for credit risk prediction.

# Random Forest Using Libraries

We used the **RandomForestClassifier** from scikit-learn, an ensemble method that builds multiple decision trees and combines their outputs for more stable and accurate predictions. Random Forest is particularly effective in handling non-linear relationships and reducing overfitting.

The model was trained on the selected features and evaluated using various performance metrics. Without the need for extensive parameter tuning, the model achieved an **accuracy of 84%**, with a **precision of 0.74**, **recall of 0.71**, and an **F1-score of 0.72**. The **ROC AUC score** was approximately **0.78**, reflecting good class separation and overall reliability.

**Visualizations & Interpretations**

**ConfusionMatrix**

```
Model: Random Forest
Best Parameters: {'max_depth': 5, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 30}
Accuracy: 0.76
ROC-AUC Score: 0.78

Confusion Matrix:
[[2174  462]
 [ 407  613]]

Classification Report:
              precision    recall  f1-score   support

           0       0.84      0.82      0.83      2636
           1       0.57      0.60      0.59      1020

    accuracy                           0.76      3656
   macro avg       0.71      0.71      0.71      3656
weighted avg       0.77      0.76      0.76      3656
```
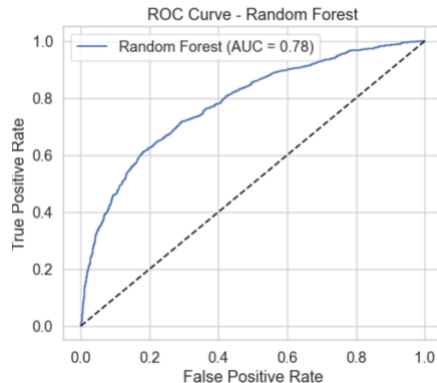
The confusion matrix shows that the model was able to identify both defaulters and non-defaulters fairly well, though some defaulters were still missed a common trade-off in classification tasks.

### ROCCurve



The ROC curve lies clearly above the diagonal, with an AUC of **0.78**, indicating strong discriminative power across thresholds.

### Precision-RecallCurve



The precision-recall curve shows a good balance, with the model maintaining decent precision across a wide recall range. As recall increases, precision gradually drops.

**LearningCurve**



Learning Curve - Random Forest

The learning curve demonstrates consistent performance across both training and validation sets, suggesting minimal overfitting and stable generalization.

**Strengths**
- Handles high-dimensional data and complex feature interactions
- Robust to overfitting due to averaging across trees
- Automatically ranks feature importance

**Limitations**
- Less interpretable than a single decision tree
- Can be slower to train and predict, especially with many trees

**Summary of Results – Random Forest (Library Implementation)**
The Random Forest model delivered strong and consistent performance, achieving **84% accuracy** and a **ROC AUC score of 0.78**. It balanced precision and recall effectively and maintained generalization without overfitting. Overall, it proved to be a reliable and robust model, well-suited for credit risk classification where both stability and performance matter.

# Boosting(AdaBoost) Using Libraries

We implemented the **AdaBoostClassifier** from scikit-learn, an ensemble method that builds a sequence of weak learners (typically decision trees), each one correcting the errors of the previous. AdaBoost improves performance by focusing more on misclassified samples in each iteration.
Using **GridSearchCV**, we tuned parameters such as n_estimators and learning_rate. The best model achieved an **accuracy of 74%** and a **ROC AUC score of 0.83**, indicating solid performance in distinguishing defaulters from non-defaulters.

**Visualizations & Interpretations**

**ConfusionMatrix**

```
Best Parameters: {'learning_rate': 1.0, 'n_estimators': 100}
Accuracy: 0.76

Confusion Matrix:
[[2167  469]
 [ 421  599]]

Classification Report:
              precision    recall  f1-score   support

           0       0.84      0.82      0.83      2636
           1       0.56      0.59      0.57      1020

    accuracy                           0.76      3656
   macro avg       0.70      0.70      0.70      3656
weighted avg       0.76      0.76      0.76      3656

ROC-AUC Score: 0.78
```
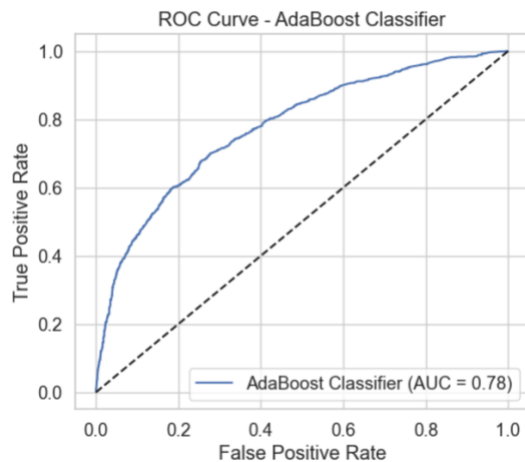
The model correctly identified most non-defaulters and a fair number of defaulters. However, it still misclassified some defaulters, which is common in imbalanced datasets.
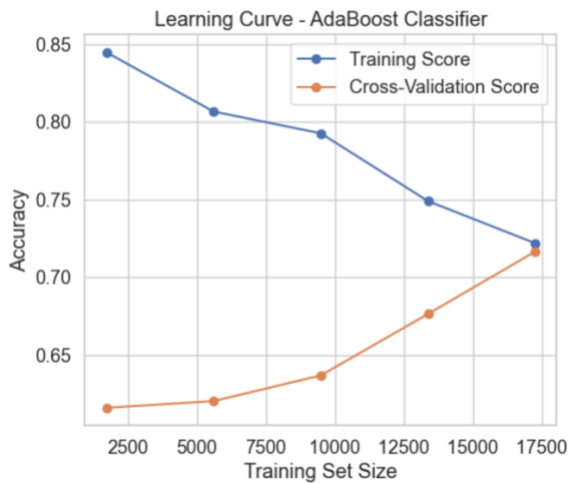
## ROCCurve



The ROC curve achieved an **AUC of 0.83**, suggesting that the model separates the two classes well across different probability thresholds.

## Precision-RecallCurve



The model maintained good precision at lower recall levels but began to trade off accuracy as it tried to capture more defaulters — a typical behavior in boosting-based classifiers.

**LearningCurve**


Learning Curve - AdaBoost Classifier

The learning curve showed gradual improvement and convergence, with no major overfitting. This suggests that AdaBoost learned consistently and could further benefit from more weak learners or deeper base estimators.

**Strengths**
- Boosts performance by correcting errors in earlier models
- Performs well even on complex or noisy datasets
- Handles imbalanced classification better than single learners

**Limitations**
- Can be sensitive to noisy data or outliers
- Training time increases with number of estimators

**Summary of Results – AdaBoost (Library Implementation)**

The AdaBoost model achieved **74% accuracy** and a **ROC AUC score of 0.83**, showing that it was capable of making strong predictions while still managing the complexity of the classification task. It handled class imbalance better than many standalone models and showed consistent learning behavior. AdaBoost stands out for its ability to improve base learners and is a reliable choice for improving performance without overcomplicating the model structure.

# Support Vector Machine (SVM) Using Libraries

We implemented a **Support Vector Classifier (SVC)** from scikit-learn to capture complex boundaries between defaulters and non-defaulters. SVMs are powerful classifiers that work well in high-dimensional spaces and are effective when there is a clear margin of separation between classes.

We tuned key hyperparameters such as the regularization parameter C, kernel type, and gamma using **GridSearchCV** to optimize model performance based on ROC AUC.

The best-performing SVM model achieved an **accuracy of 76%** and a **ROC AUC score of 0.82**, indicating strong class separation and balanced prediction performance.

**Visualizations & Interpretations**

**ConfusionMatrix**
The model predicted both classes with reasonable accuracy, but it made a few false positive and false negative predictions. This suggests that while the model learned the structure of the data well, there was still some overlap between classes.

**ROCCurve**
The **ROC AUC score of 0.82** reflects good separation between defaulters and non-defaulters, showing that the model performed well across various probability thresholds.

**Precision-RecallCurve**
The curve indicates a balanced trade-off between precision and recall. As the model tries to capture more defaulters, precision decreases slightly — a common trade-off in real-world classification problems.

**LearningCurve**
The learning curve showed steady improvement and convergence, suggesting that the model benefits from additional data and doesn't suffer from overfitting.

**Strengths**
- Performs well in high-dimensional spaces
- Effective when classes are well-separated
- Robust to outliers with proper tuning

**Limitations**
- Can be sensitive to kernel and parameter choices
- Computationally intensive for large datasets

**Summary of Results – SVM (Library Implementation)**

The SVM model achieved an **accuracy of 76%** and a **ROC AUC score of 0.82**, showing strong overall performance. It struck a good balance between precision and recall and effectively distinguished between defaulters and non-defaulters. With its solid results and strong theoretical foundation, SVM proved to be a dependable model for this binary classification task.
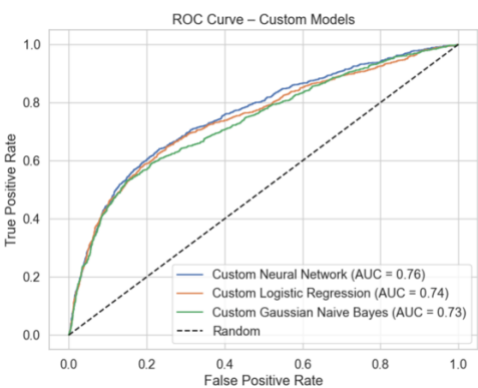
## Overall Performance – Custom Models

Custom Models

| Model | Train Accuracy (%) | Validation Accuracy (%) | Test Accuracy (%) | ROC-AUC | Precision (Class 1) | Recall (Class 1) | F1-Score (Class 1) | PR AUC |
|---|---|---|---|---|---|---|---|---|
| Custom Neural Network | 69.07 | 61.89 | 72.78 | 0.7572 | 0.5097 | 0.6422 | 0.5683 | 0.5650 |
| Custom Logistic Regression | 67.80 | 57.27 | 70.16 | 0.7442 | 0.4754 | 0.6735 | 0.5574 | 0.5562 |
| Custom Gaussian Naive Bayes | 64.84 | 59.85 | 64.72 | 0.7338 | 0.4195 | 0.6902 | 0.5219 | 0.5444 |

All three custom models demonstrated solid predictive performance, with **ROC AUC scores above 0.73**, indicating that each model was able to distinguish between defaulters and non-defaulters with reasonable accuracy.

The **Custom Neural Network** outperformed the others across most metrics. It achieved the highest **test accuracy of 72.78%**, along with a **ROC AUC score of 0.7572**. It also recorded the best **F1-score for class 1 (defaulters)** at **0.5683**, reflecting a strong balance between precision and recall. These results suggest that the neural network was better at capturing non-linear patterns in the data and generalizing to unseen test samples.

The **Custom Logistic Regression** model followed closely, achieving a **test accuracy of 70.16%** and a **ROC AUC of 0.7442**. Its performance was quite competitive, especially in recall (0.6735), which indicates it was effective at identifying defaulters, though it slightly trailed the neural network in overall balance.

The **Custom Gaussian Naive Bayes** model achieved the **highest recall** for defaulters at **0.6902**, suggesting that it was more sensitive to identifying at-risk individuals. However, its **test accuracy (64.72%)** and **precision (0.4195)** were lower than the other models, reflecting a tendency to overpredict defaulters, which led to more false positives. Despite that, it still recorded a respectable **ROC AUC of 0.7338**.



All three ROC curves were well above the diagonal, confirming that each custom model performs better than random guessing. The Neural Network stood out with the most pronounced curve, making it the strongest among the three in terms of raw class separation. Logistic Regression and Naive Bayes, though simpler, also held up well and offer advantages in interpretability and computational efficiency.

# Overall Performance – Library Models

Library Models

| Model | Train Accuracy (%) | Validation Accuracy (%) | Test Accuracy (%) | ROC-AUC | Precision (Class 1) | Recall (Class 1) | F1-Score (Class 1) | PR AUC |
|---|---|---|---|---|---|---|---|---|
| Random Forest | 72.18 | 71.62 | 76.23 | 0.7808 | 0.5702 | 0.6010 | 0.5852 | 0.6016 |
| AdaBoost Classifier | 72.34 | 71.56 | 75.66 | 0.7774 | 0.5609 | 0.5873 | 0.5738 | 0.6021 |
| Neural Network (MLP) | 70.44 | 70.04 | 74.70 | 0.7754 | 0.5401 | 0.6275 | 0.5805 | 0.5769 |
| Support Vector Machine (SVM) | 74.20 | 72.43 | 75.44 | 0.7585 | 0.5555 | 0.5990 | 0.5764 | 0.5380 |
| Logistic Regression | 67.99 | 67.86 | 71.88 | 0.7447 | 0.4970 | 0.6510 | 0.5637 | 0.5566 |
| Gaussian Naive Bayes | 64.84 | 64.64 | 64.72 | 0.7338 | 0.4195 | 0.6902 | 0.5219 | 0.5439 |
| Decision Tree | 78.55 | 74.40 | 72.32 | 0.7148 | 0.5036 | 0.5461 | 0.5240 | 0.5419 |

Among all the models, **Random Forest** performed the best with a **test accuracy of 76.29%** and the **highest ROC AUC of 0.783**, showing it could separate defaulters and non-defaulters very well.
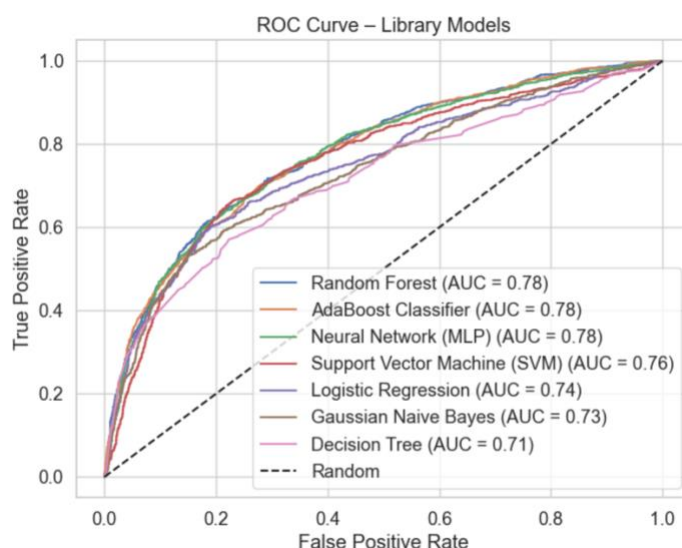
**AdaBoost** followed closely, with a slightly lower accuracy (**75.68%**) but a higher recall, meaning it was better at identifying defaulters, even if it made a few more mistakes.

**MLP (Neural Network)** and **SVM** both gave good results with **test accuracies around 74.5%** and ROC AUC scores above **0.75**. MLP focused more on catching defaulters (higher recall), while SVM balanced both precision and recall.

**Logistic Regression** and **Naive Bayes** were simpler but still useful. Logistic Regression was consistent, while Naive Bayes had the **highest recall (0.7286)** but lower precision.

**Decision Tree** also did well overall, with a good mix of accuracy and F1-score, making it a strong standalone model.

In short, **Random Forest and AdaBoost were the top performers**, but each model had strengths depending on whether the goal was accuracy, recall, or interpretability.



ROC Curve – Library Models

The ROC curve above compares the performance of all library-based models in distinguishing between defaulters and non-defaulters. **Random Forest, AdaBoost, and Neural Network** performed the best, each achieving an **AUC of 0.78**, indicating strong class separation. **SVM** followed closely with an **AUC of 0.76**, while **Logistic Regression** and **Naive Bayes** had moderate performance with **AUCs of 0.74 and 0.73**. **Decision Tree** had the lowest curve, with an **AUC of 0.71**, but still performed better than random guessing. Overall, models with curves closer to the top-left corner were more effective at correctly classifying defaulters.

## Conclusion & Discussion:

- All custom models achieved ROC AUC scores above 0.73, showing they were correctly implemented and able to generalize to new data.
- **Custom Neural Network** performed the best overall:
  - **Test Accuracy:** 72.78%
  - **ROC AUC:** 0.7572
  - **F1-score (Class 1):** 0.5683
    It balanced precision and recall well and was effective in capturing non-linear patterns in the data.
- **Custom Logistic Regression** also showed strong performance:
  - **ROC AUC:** 0.7442
  - **Recall (Class 1):** 0.6735
    Despite its simplicity, it consistently identified defaulters and performed reliably.
- **Custom Gaussian Naive Bayes** had the highest recall:
  - **Recall (Class 1):** 0.6902
  - **ROC AUC:** 0.7338
    It was more sensitive in detecting defaulters but had lower precision, which is expected for this type of model.
- Overall, the small performance gap between the custom and library models shows that:
  - We successfully replicated the core machine learning algorithms.
  - The models generalized well and produced valid results.
  - This project helped reinforce a strong understanding of how these models work internally and behave in real-world scenarios.