



Banque Haolim



Loan Application Presentation

August 5, 2025

Introduction

Welcome to Bank Haolim's loan application presentation.

The goal of this presentation is to analyze how various factors influence loan approval decisions at the bank, using a loan dataset inspired by Kaggle.

Outline

- ▶ Data understanding
- ▶ EDA (Exploration data analysis)
- ▶ Data preprocessing
- ▶ KPI
- ▶ Modeling
- ▶ Model interpretation and validation
- ▶ Deployment and prediction via API
- ▶ Conclusion

Data understanding

- ▶ Variable descriptions
- ▶ Missing values

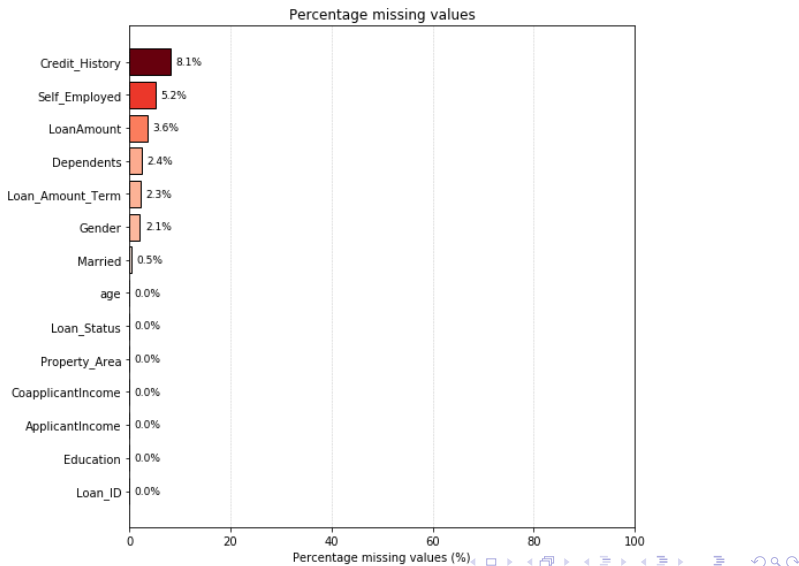


Data Understanding : Variable description

The dataset contains **14 variables** and **614 observations**. The target variable is `Loan_Status`, which indicates whether a loan application was approved (Y) or rejected (N). The dataset includes the following variables:

- ▶ `Loan_ID`: Unique identifier for each loan application
- ▶ `Gender`: Applicant's gender
- ▶ `Married`: Marital status (Yes or No)
- ▶ `Dependents`: Number of dependents supported by the applicant
- ▶ `Education`: Whether the applicant is a graduate
- ▶ `Self_Employed`: Whether the applicant is self-employed
- ▶ `ApplicantIncome`: Monthly income of the applicant
- ▶ `CoapplicantIncome`: Monthly income of the co-applicant
- ▶ `LoanAmount`: Requested loan amount
- ▶ `Loan_Amount_Term`: Loan duration in months
- ▶ `Credit_History`: Presence of credit history (1 = Yes, 0 = No)
- ▶ `Property_Area`: Type of residential area (Urban, Semiurban, Rural)
- ▶ `Age`: Applicant's age (added manually)

Data understanding : Missing values

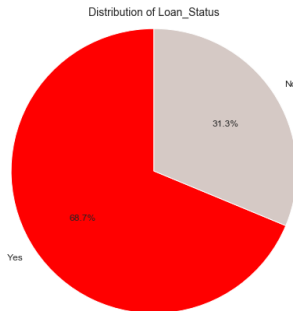


EDA : Exploratory Data Analysis

- ▶ Target variable distribution
- ▶ Categorical variable analysis
- ▶ Numerical feature exploration
- ▶ Bivariate relationships with target
- ▶ Feature correlation analysis

EDA : Target variable distribution

Target variable distribution: The target variable `Loan_Status` indicates whether a loan application was approved (Y : 68.7%) or not (N : 31.3 %). The distribution shows that a majority of applications were approved, revealing a slight class imbalance that may need to be considered during model training.



EDA : Categorical variable analysis

We examined all categorical variables, as well as numerical variables with a limited number of unique values (fewer than 10), such as Gender, Married, Education, Self_Employed, Property_Area, Credit_History, and Dependents.

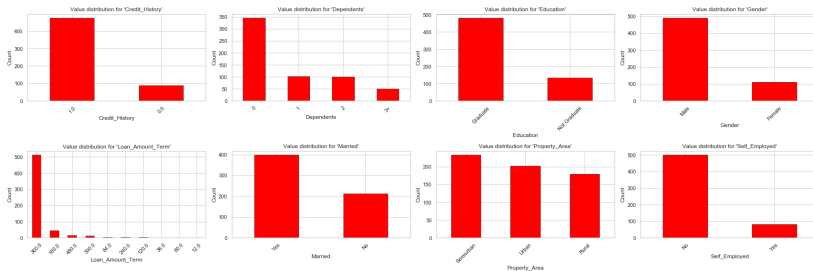


Figure: Histogramm of categorial variables

EDA : Numerical feature exploration

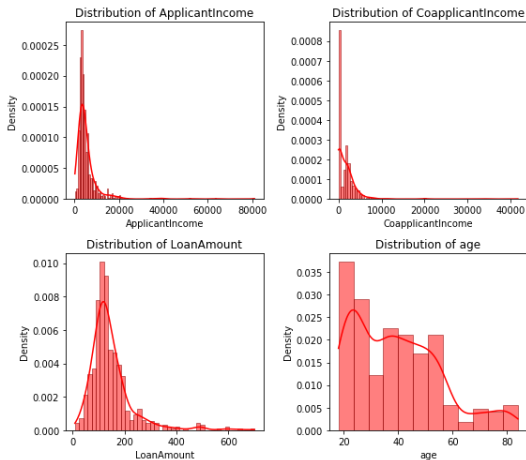


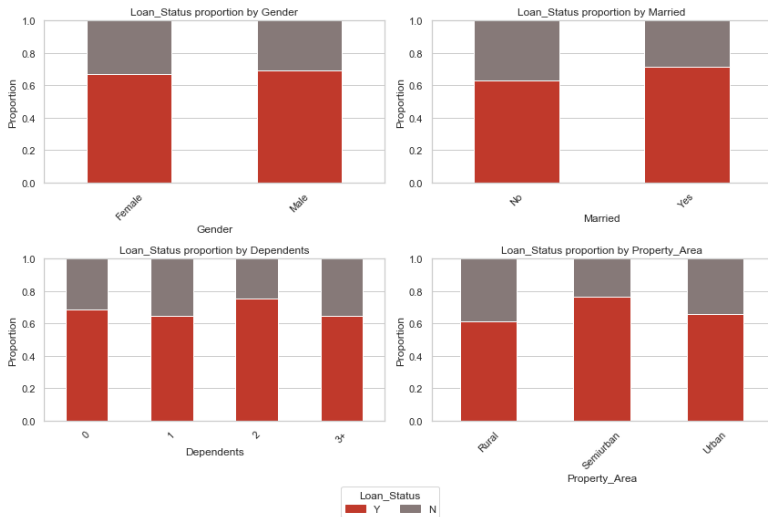
Figure: Density of numerical variables

EDA : Numerical feature exploration

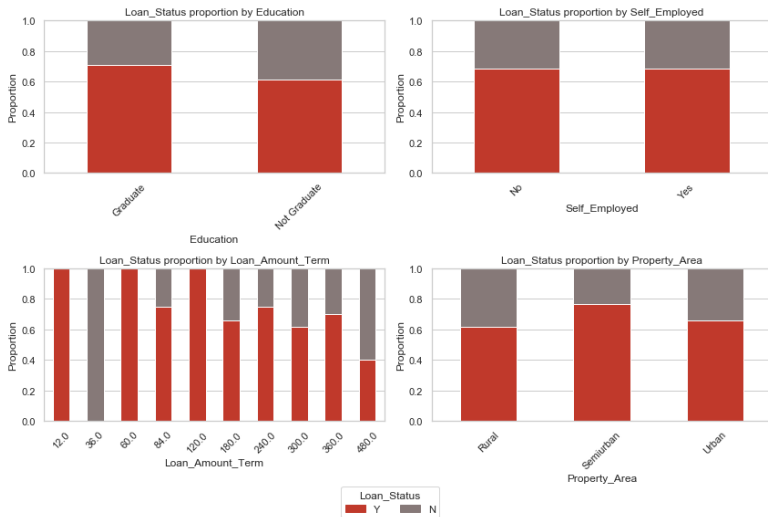
Numerical feature exploration: We explored the distribution of key numerical variables such as `ApplicantIncome`, `CoapplicantIncome`, and `LoanAmount`. These variables showed noticeable skewness and outliers, particularly in income-related features. Understanding their distribution helps assess the need for transformations and potential feature engineering in the modeling stage.



EDA : Bivariate relationships with target (categorical variables)



EDA : Bivariate relationships with target (categorical variables)



Data source: Kaggle

EDA : Bivariate relationships with target (continuous variables)

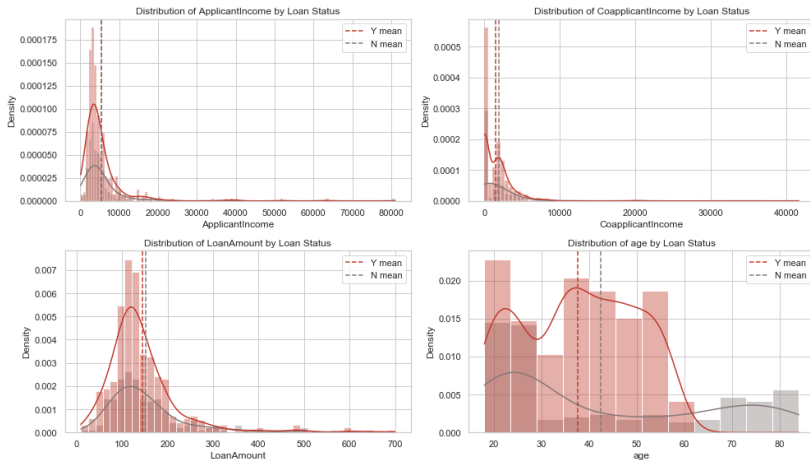


Figure: Density of numerical variables

EDA : Bivariate relationships with target (continuous variables)

▶ **LoanAmount, ApplicantIncome, and CoapplicantIncome:**

- ▶ High dispersion, skewness, and outliers → not well described by the mean alone.
- ▶ Refusals (N): more varied and dispersed profiles.
- ▶ Approvals (Y): more homogeneous profiles, concentrated within a “favorable” range.

▶ **Age:**

- ▶ More stable, centered around 30–40 years old.
- ▶ Slightly more discriminant: approvals tend to cluster in this range.

Feature correlation analysis

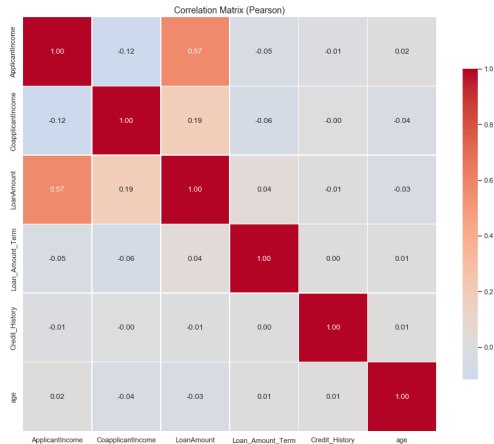


Figure: Correlation matrix

Data preprocessing

- ▶ Data cleaning
- ▶ Encoding Categorical Variables
- ▶ Feature Engineering

Data preprocessing : data cleaning

- ▶ Handling missing values
- ▶ Detecting and treating outliers

Data preprocessing : Handling missing values

- ▶ For **continuous features**, missing values were imputed using the **median**, which is more robust to outliers.
- ▶ For **categorical features**, missing **values** were filled with the **most frequent modalities** .

Data preprocessing : Detecting and treating outliers

To identify and handle outliers, we used a two-step approach:

- ▶ **Visual inspection:** Boxplots were used to visually detect potential outliers in continuous variables.
- ▶ **Statistical detection:** Z-scores were computed to quantify how far each value deviates from the mean. Observations with extreme Z-scores (typically $|Z| > 3$) were considered outliers.

Cleaning: Values identified as strongly aberrant were removed to ensure the quality of the analysis.

Data preprocessing : Detecting and treating outliers

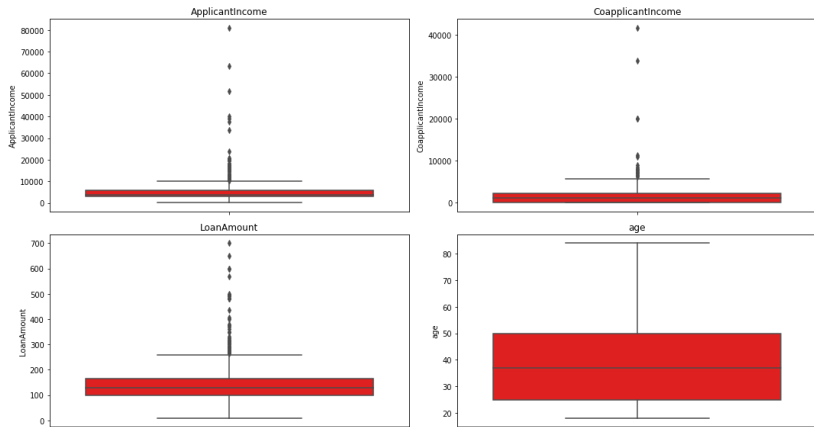


Figure: Boxplot of Numerical Variables

Data preprocessing : Detecting and treating outliers

The **Z-score**, or *standard score*, indicates how many standard deviations a data point is from the mean of a distribution.

Formula:

$$Z = \frac{x - \mu}{\sigma}$$

where:

- ▶ x : individual value
- ▶ μ : mean of the dataset
- ▶ σ : standard deviation

Interpretation:

- ▶ $Z = 0$: value equals the mean
- ▶ $Z > 0$: value is above the mean
- ▶ $Z < 0$: value is below the mean
- ▶ $|Z| > 3$: usually considered an outlier

Data preprocessing : Detecting and treating outliers

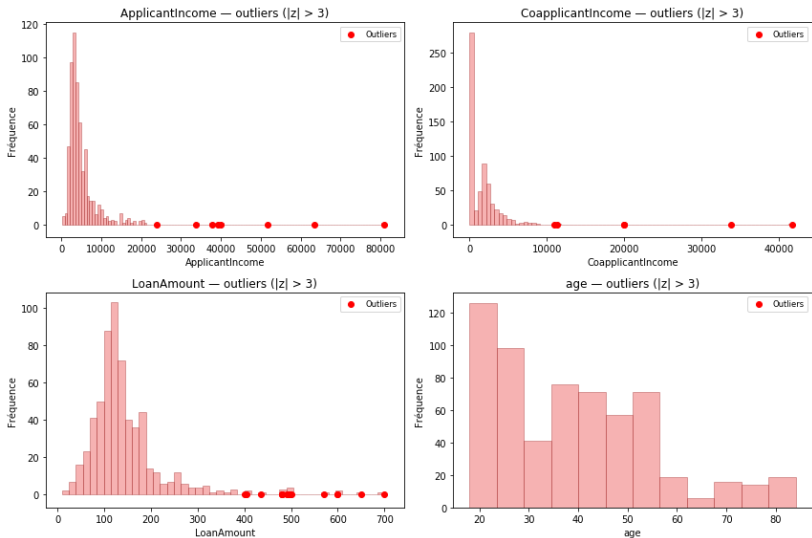


Figure: Boxplot of Numerical Var



Banque Haolim



Data preprocessing : Detecting and treating outliers

Based on both visual analysis using boxplots and statistical validation with the Z-test, we identified several outliers. These observations were removed to prevent them from skewing the results and affecting model accuracy.

After this step, **589 observations** remained in the dataset.

Data preprocessing : Feature Engineering

Feature engineering is the process of creating, transforming, or selecting features to improve the performance of machine learning models.

In our project, we performed the following:

Created new features:

- ▶ **TotalIncome:** Combined the applicant's and co-applicant's income to reflect total household earnings. We will discretize this continuous variable into five intervals using the `pd.cut` function. This transformation will allow for a more detailed analysis of the relationship between aggregated income levels and loan approval. The resulting variable will be named `TotalIncomeCut`.
- ▶ **MonthlyPayment:** Approximated by dividing `LoanAmount` by `Loan_Amount_Term`.
- ▶ **DebtToIncomeRatio:** Defined as `MonthlyPayment` divided by `TotalIncome`, representing the burden of the loan relative to income.



Data preprocessing : Feature engineering

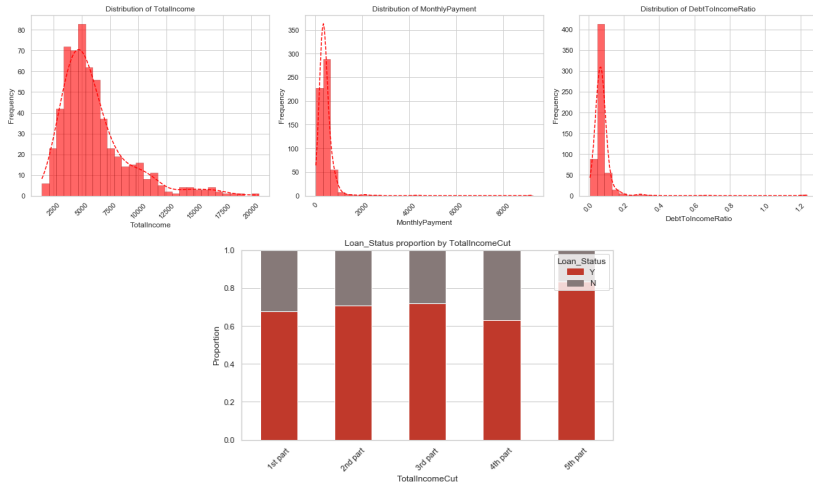


Figure: Density of newly created variables

Data preprocessing : Feature Engineering

- ▶ **Transformed categorical variables:** Converted categorical text fields into numerical codes to enable model processing.

These steps help the model better capture patterns and relationships within the data.

Data preprocessing : Feature engineering

Encoded categorical variables as dummies: Categorical variables were transformed using one-hot encoding, creating one binary column per category. This step is essential for most machine learning models that require numerical input and cannot directly interpret text labels.

Data preprocessing : Feature engineering

For example, the variable Education with categories such as Graduate and Not Graduate was converted into binary columns. One of these, Education_Graduate, takes the value 1 if the applicant is a graduate, and 0 otherwise. This allows the model to quantify the impact of being a graduate on the likelihood of loan approval.

This transformation also avoids introducing artificial ordinality (i.e., false numerical order) into categorical variables and ensures that the model treats each category independently.

Education_Not Graduate
0
0
0
1
0

Education
Graduate
Graduate
Graduate
Not Graduate
Graduate

KPI : Global Indicators

This report analyzes key performance indicators (KPIs) related to loan approval decisions using data from a Kaggle dataset. The objective is to identify which applicant characteristics and financial factors most influence whether a loan is approved.

We begin by examining general statistics related to loan approval:

KPI : Socio-Demographic Analysis

This section explores how demographic features impact loan approval:

- ▶ **Gender:** Comparison of approval rates between male and female applicants.
- ▶ **Education:** Differences between graduates and non-graduates.
- ▶ **Marital status:** Influence of family situation on loan decisions.
- ▶ **Age:** Examination of how applicant age groups impact loan approval rate.



KPI : Socio-Demographic Analysis

The **Chi-squared test of independence** is used to determine whether there is a statistically significant association between two categorical variables.

In our case, it helps answer questions such as: *Is there a relationship between Gender, Education or statut marital and Loan Approval?*

Hypotheses:

- ▶ H_0 : The two variables are independent (no association).
- ▶ H_1 : The two variables are dependent (there is an association).

Decision Rule:

- ▶ If the p-value is < 0.05 , we reject H_0 — the variables are statistically dependent.
- ▶ If the p-value is ≥ 0.05 , we fail to reject H_0 — no evidence of a relationship.

KPI : Socio-Demographic Analysis

=== Chi-square Test Results ===

Variable	Chi2 Statistic	Degrees of Freedom	p-value	Interpretation
Gender	0.174	1	0.6762	Not significant
Education	4.610	1	0.0318	Significant ($p < 0.05$)
Married	4.907	1	0.0268	Significant ($p < 0.05$)

Based on the Chi-squared test of independence, we find that there is **no statistically significant association** between Gender and Loan_Status. In this case, $p \approx 0.68$, we **fail to reject** the null hypothesis (H_0), meaning the two variables are considered independent.

However, both Education and Marital Status show **statistically significant associations** with Loan_Status (p -value < 0.05). For these variables, we **reject** the null hypothesis (H_0), indicating a dependency with the loan approval outcome.

KPI : Socio-Demographic Analysis

T-test statistic: -2.868

P-value: 0.0045

Result: Significant difference in mean age between loan status groups.

For the variable Age, which is continuous rather than categorical, the Chi-squared test is not appropriate. Instead, since the loan status variable has only two categories (approved or not approved), a **two-sample t-test** is more suitable.

The t-test compares the means of the Age variable between the two loan status groups to determine if the difference in means is statistically significant.

In a t-test, the null hypothesis (H_0) states that the two group means are equal. If the p-value is less than the chosen (usually 0.05), we reject H_0 , suggesting a significant difference in age between the groups. Otherwise, we fail to reject H_0 , (p-value = 0.0045) indicating no significant difference in mean age between approved and non-approved applicants.

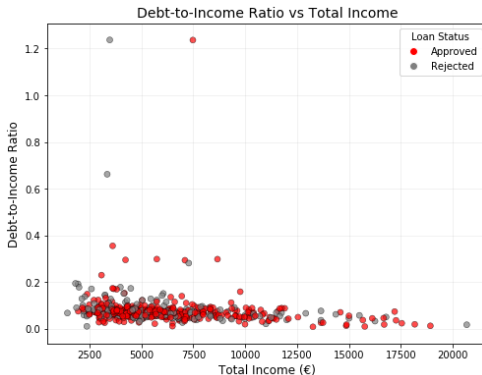


KPI : Financial Situation Analysis

Here, we analyze the role of applicants' financial standing:

- ▶ **Loan-to-Income Ratio:** Ratio between the loan amount and total income.

KPI : Financial Situation Analysis



This graph clearly illustrates the relationship between the **debt-to-income ratio** and the applicants' **total income**. We observe that lower incomes are associated with higher ratios and a higher number of loan rejections, while higher incomes are correlated with a better repayment capacity and thus a higher probability of loan approval.

KPI : Geographical Analysis

We will use ANOVA (Analysis of Variance) to compare the average **Debt-to-Income Ratio** across different **geographic areas** (such as Urban, Semi-Urban, and Rural). This will help us determine whether the financial burden relative to income significantly differs depending on the applicant's location. If the test shows a significant difference, it means that the geographic area has an impact on applicants' debt levels, which could influence loan approval decisions.



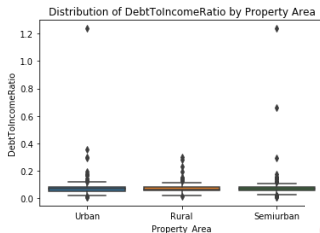
KPI : Geographical Analysis

F-statistic: 0.600

p-value: 0.549

Result: No significant difference in mean DebtToIncomeRatio between property_area groups.

In our case, the null hypothesis (H_0) is that the means of the Urban, Semi-Urban, and Rural groups are equal. The p-value obtained is greater than 0.05, which means we do not reject the null hypothesis. In other words, there is no statistically significant difference in the average loan amounts granted based on the type of area. This suggests that the Urban, Semi-Urban, and Rural zones exhibit similar behavior with respect to this criterion. This conclusion is also visually supported by the graph: the distributions or means are close across the three groups, which reinforces the idea of statistical equivalence.



Banque Haolim



Modeling

Before building our predictive models, we divided the dataset into two distinct subsets: a **training set**, used to fit the models, and a **test set**, used to evaluate their performance on previously unseen data. This separation ensures **reliable validation** and reduces the risk of **overfitting**.

It is also essential that the **target variable** is properly distributed between the training and test sets to preserve the representativeness of the class distribution and avoid any evaluation bias.

We will then evaluate two models: **logistic regression** and **random forest**, and select the one that achieves the best performance. For logistic regression, we will **standardize the predictor variables** (i.e., **center and scale** them to have mean zero and unit variance) to **facilitate model convergence** and **improve the numerical stability of the coefficients**. In Python, this standardization is performed using the `StandardScaler` from `scikit-learn`, integrated within a `Pipeline` to ensure that the scaling is applied consistently during both training and testing.



Modeling

- ▶ Logistic regression
- ▶ Random forest

Logistic regression

Logistic regression is a statistical method used for binary classification problems. It predicts the probability that a given input belongs to a certain category, typically 0 or 1.

Model Formulation

The logistic regression model estimates the probability using the sigmoid function:

$$P(Y = 1 | X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_n X_n)}}$$

Why Logistic Regression?

- ▶ Suitable for binary outcomes (e.g., accept or reject a loan)
- ▶ Interpretable model: coefficients indicate influence on the log-odds
- ▶ Can handle multiple features

Logistic regression

The logistic regression model was trained on 589 observations and evaluated on a test set of 177 samples. The confusion matrix (Table 3) shows that the model correctly classified 38 out of 55 negative cases (class 0) and 104 out of 122 positive cases (class 1).

	Predicted 0	Predicted 1
Actual 0	38	17
Actual 1	18	104

Table: Confusion matrix

Logistic regression

Performance metrics (precision, recall, and F1-score) for each class are reported in Table 4.

Precision for class 0 is 0.68 and for class 1 is 0.86, indicating that predictions for class 1 are more accurate. Recall is 0.69 for class 0 and 0.85 for class 1, showing better identification of class 1. The F1-score balances precision and recall with values of 0.68 and 0.86 respectively.

Class	Precision	Recall	F1-score	Support
0	0.68	0.69	0.68	55
1	0.86	0.85	0.86	122
Accuracy	0.80			
Macro avg	0.77	0.77	0.77	177
Weighted avg	0.80	0.80	0.80	177

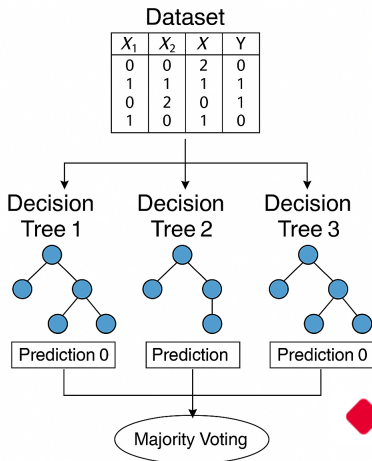
Table: Classification report

Logistic regression performance

The overall accuracy of the model is 80%, with a macro-average F1-score of 0.77 and a weighted average of 0.80. The AUC score of 0.804 indicates good discrimination between the two classes.

Random forest

Random forest is an ensemble learning method that builds many decision trees and combines their results to make more accurate and stable predictions. Each tree is trained on a random subset of the data and variables, which helps reduce overfitting and improve generalization.



Random forest

The model was trained on a dataset and evaluated on a test set of 177 samples. The confusion matrix shows that the model correctly classified 37 out of 55 negative cases (class 0) and 120 out of 122 positive cases (class 1).

	Predicted 0	Predicted 1
Actual 0	37	18
Actual 1	2	120

Table: Confusion matrix

Random forest

The table below reports the key classification metrics. For class 0, the model achieved a precision of 0.95 and a recall of 0.67, indicating high precision but moderate sensitivity in detecting negative cases. For class 1, both precision (0.87) and recall (0.98) are high, showing excellent performance in identifying positive cases.

Class	Precision	Recall	F1-score	Support
0	0.95	0.67	0.79	55
1	0.87	0.98	0.92	122
Accuracy	0.89			
Macro avg	0.91	0.83	0.86	177
Weighted avg	0.89	0.89	0.88	177

Table: Classification report

Random forest Performance

The overall accuracy of the model is 89%, with a macro-average F1-score of 0.86 and a weighted average of 0.88. The AUC score of 0.828 further confirms the model's strong ability to distinguish between the two classes.

Random forest Optimization

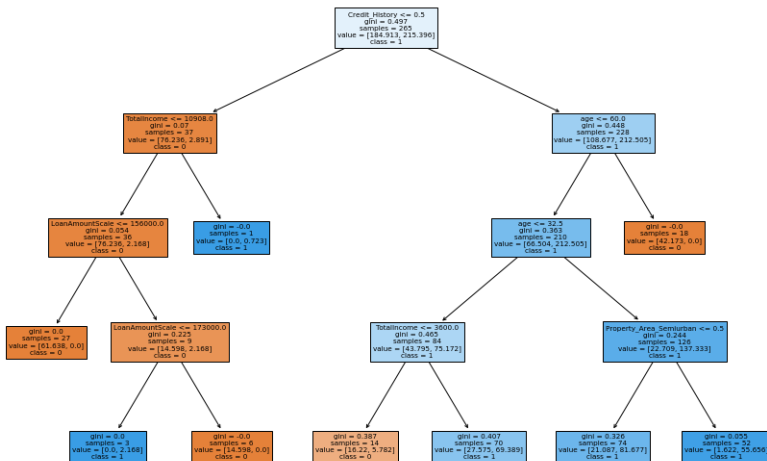
To optimize this model, we will perform a **randomized search**, which randomly samples a set of hyperparameter combinations to identify the configuration that provides the best performance. Unlike a **grid search**, which exhaustively tests all possible combinations, this approach reduces computation time while still efficiently exploring the hyperparameter space. In Python, this procedure is implemented using `RandomizedSearchCV` from `scikit-learn`.



Model interpretation and validation

With the optimized model, we achieved an AUC score of 0.892, indicating a strong ability to distinguish between accepted and rejected loan applications. This optimized model will therefore be used for our future predictions and deployment.

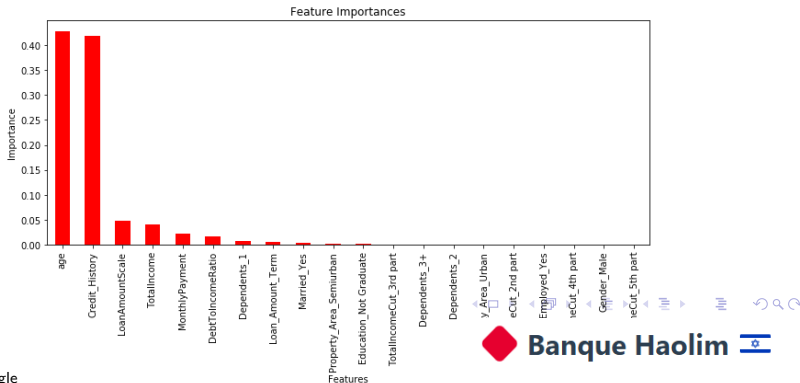
Model interpretation and validation



Deployment and prediction via API

As part of a project for a banking institution, we are developing an API that allows clients to input key information to determine whether their loan application would be approved. To ensure efficiency, interpretability, and regulatory alignment, we rely on our model's feature importance to retain only the most relevant variables for the prediction.

We will keep age, credit history, loan amount, and total income.



Conclusion

In this project, we chose a **Random Forest** model for its ability to effectively handle complex and diverse data. After selecting relevant variables such as *age*, *credit history*, *loan amount*, and *total income*, we optimized the model using a *Grid Search*. This method allowed us to finely tune the hyperparameters to maximize predictive performance. The calibrated model offers a strong ability to identify loan applications likely to be approved.

We intentionally used only the training dataset `Train.csv` to develop and test our model, without exploiting the Kaggle test set `Test.csv`. This approach allowed us to truly experiment with the data, test different hypotheses, and be more creative in building the model. To go further, we could have explored other modeling techniques such as boosting (e.g., XGBoost, LightGBM), or deepened data preprocessing to further improve prediction quality.