

Home Credit Capstone Project

Kleyton Polzonoff

2024-12-07

Table of contents

1. Introduction	2
2. Data set	2
3. Exploratory Data Analysis (EDA)	3
3.1 Target variable: Default	3
3.2 Outliers and Anomalies	3
Numeric Values	3
Income	4
Customer age	4
Days Employed (Years)	5
Character Variables	5
Contract types	6
Income type	7
Education	8
Family Status	9
Housing type	10
Occupation type	10
3.3 Additional Observations	11
External sources	13
3.4 Correlations	16
3.5 Evaluation of Prediction Baselines	17
4. Pre-Modeling	18
4.1 Dummy Encoding	18
4.2 Handling with Missing Data	18
4.3 Sampling and Balancing	20
4.4 Data partition	21
5. Modeling	21
5.1 Logistic Regression	22
5.2 Random Forest	26
5.3 XGBoost Model	29
5.4 LightGBM Model	33
6. Comparison Metrics	37
6.1 Bussiness relevance of each model	38
6.2 Ensemble Model for Better Credit Allocation	38
7. Predicting the Kaggle Test Set	38
7.1 Application Test Transformations	38
7.2 Predicting Application Test	39
8. Submission and Kaggle Score	39

1. Introduction

Access to credit is crucial for development globally, enabling individuals and businesses to leverage financial resources for growth. However, many people face barriers to credit access due to inadequate information, such as risk agency scores, making them susceptible to predatory lending and severe financial consequences. Home Credit seeks to address this issue by providing secure credit access to underserved populations. This project aims to develop a default prediction model using advanced analytics, aligning with Home Credit's goal of fostering financial inclusion while ensuring the company's financial stability.

2. Data set

Data set structure and description

SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR
100002	1	Cash loans	M	N
100003	0	Cash loans	F	N
100004	0	Revolving loans	M	Y
100006	0	Cash loans	F	N
100007	0	Cash loans	M	N
100008	0	Cash loans	M	N

variable_types	Freq
character	16
integer	41
numeric	65

This data set contains information about consumers who have obtained credit from Home Credit. It includes demographic and business details to support the credit granting process.

The application_train table will serve as the starting point for creating predictive models. It comprises 122 columns and 307,511 rows, originally categorized into three types: character, numeric, and integer.

The descriptions of each variable are detailed in a CSV file, which is quite extensive and, therefore, not included in this document. To access these descriptions, please visit: <https://www.kaggle.com/c/home-credit-default-risk/data>.

3. Exploratory Data Analysis (EDA)

To understand the data in the application_train table and to identify any potential errors or anomalies, a series of analyses will be conducted. Additionally, this phase aims to explore the characteristics and relationships between the target variable and other variables, which will help in developing a predictive model aligned with the company's principles, strategies, and objectives.

To make this document more understandable and straightforward, only the key analyses will be presented here. However, some code used for these analyses may be included even if their results are not displayed.

3.1 Target variable: Default

The main objective of this project is to predict the binary outcome of the TARGET variable, which indicates loan default. The data set has a notable imbalance, with only about 8% of the observations corresponding to clients who are in default.

TARGET	customers	avg_loan	min	max	Percent.
0	282686	602648.3	45000	4050000	91.9
1	24825	557778.5	45000	4027680	8.1

3.2 Outliers and Anomalies

Numeric Values

In this analysis, extreme values will be examined to identify and address outliers that may clearly represent errors. Additionally, the logical integrity of the numeric values will be assessed to ensure consistency and accuracy within the data set.

After reviewing the numeric and integer variables (results not shown), it was determined which variables are relevant to show their summary at this time due to anomalies in their values.

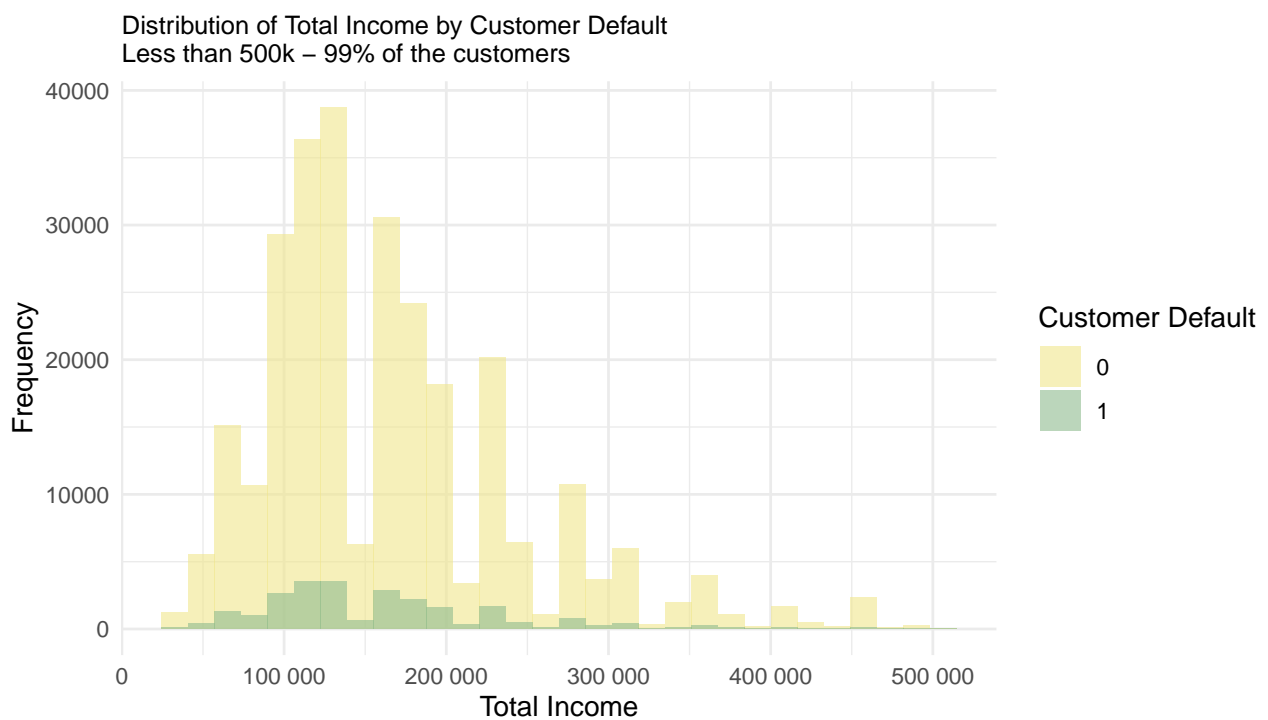
	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_GOODS_PRICE	DAYS_BIRTH	DAYS_EMPLOYED
	Min. : 25650	Min. : 45000	Min. : 40500	Min. :-25229	Min. :-17912
	1st Qu.: 112500	1st Qu.: 270000	1st Qu.: 238500	1st Qu.: -19682	1st Qu.: -2760
	Median : 147150	Median : 513531	Median : 450000	Median :-15750	Median : -1213
	Mean : 168798	Mean : 599026	Mean : 538396	Mean :-16037	Mean : 63815
	3rd Qu.: 202500	3rd Qu.: 808650	3rd Qu.: 679500	3rd Qu.: -12413	3rd Qu.: -289
	Max. :117000000	Max. :4050000	Max. :4050000	Max. : -7489	Max. :365243
	NA	NA	NA's :278	NA	NA

It is possible to observe extreme values and missing data (NAs) in the data set. The date-related variables are counted retroactively from the application date, so to accurately view the applicants' ages, these variables will be transformed for better understanding.

For AMT_ANNUITY, none of the 12 missing values are in the default group (target = 1), so I will impute them using the median of the observed values. Similarly, for AMT_GOODS_PRICE, with 278 missing values (21 in the default group), I will apply the same imputation method.

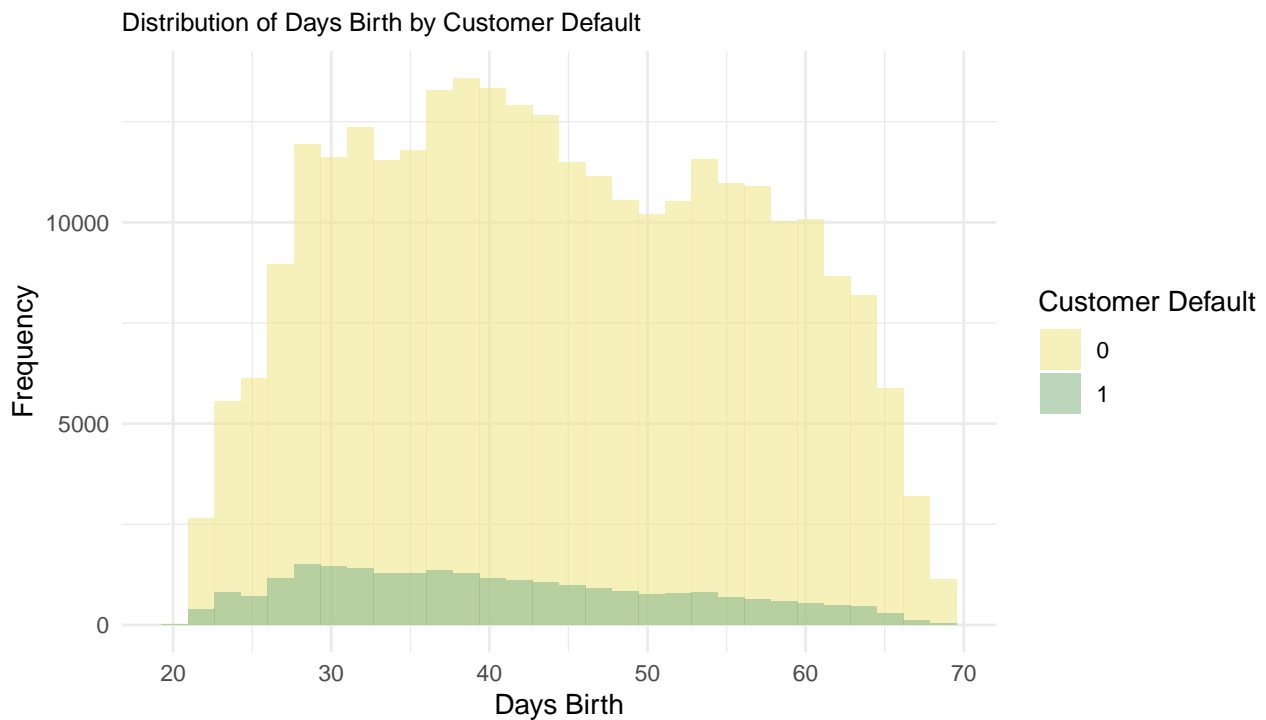
Income

For AMT_INCOME_TOTAL, five individuals reported incomes exceeding 5 million, with one reporting 117 million. These extreme outliers are far from the variable's median of \$147,150 and are not the main focus of the company's business plan. Additionally, 99% of the individuals have incomes lower than \$500,000. Therefore, these five observations will be removed, and the transformation process will also be applied to the application_test data frame later.



Customer age

Below, the DAYS_BIRTH variable will be transformed to provide a clearer understanding of the consumers' ages at the time of the loan application. The data set app_train_clean will be created to store these transformations.



Days Employed (Years)

The DAYS_EMPLOYED variable shows positive values, which should not occur.

TARGET	customers	min	max	Percent.
0	52384	365243	365243	94.6
1	2990	365243	365243	5.4

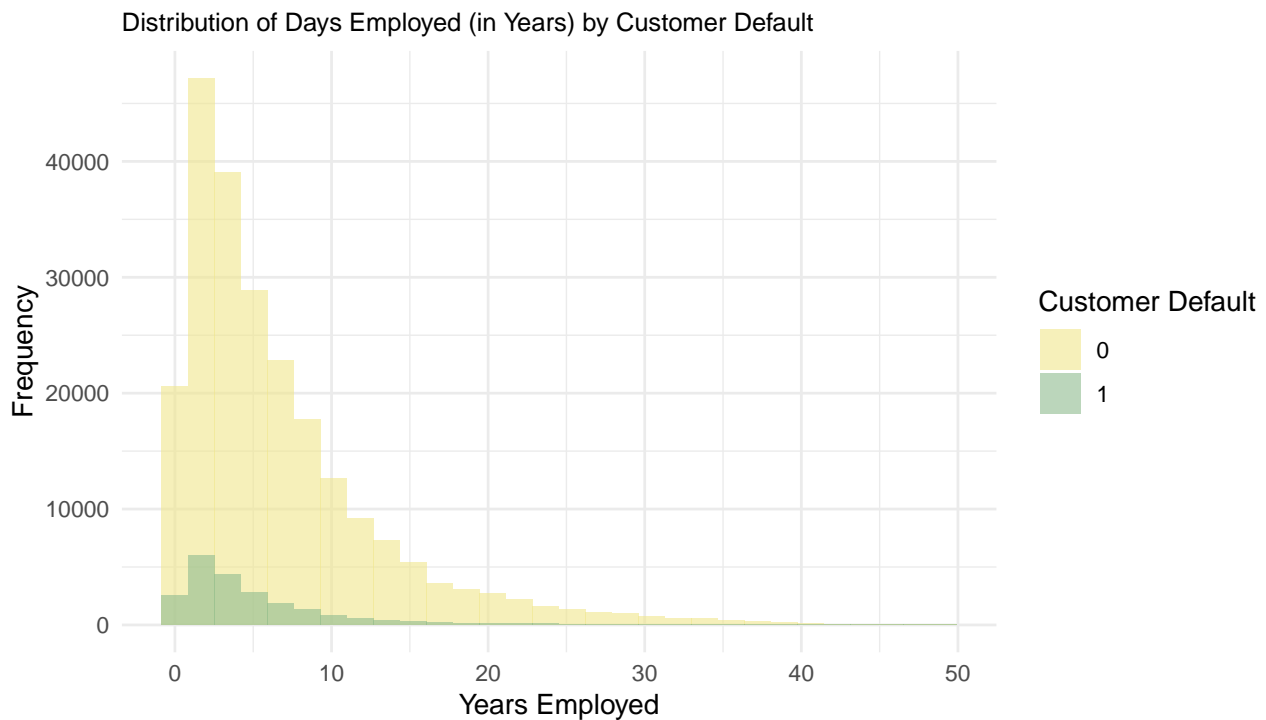
Since anomalous values account for a significant portion of the data set (55,374 entries), the strategy is to mark these entries as non-numeric and create a new column, DAYS_EMPLOYED_ANOM, to flag the presence of these anomalies for each customer.

Below is the histogram of the corrected DAYS_EMPLOYED distribution. In this case, the distribution, transformed into years, is right-skewed. Most defaults are found among consumers with up to 5 years of employment.



Character Variables

As seen earlier, there are 16 character variables in the data set. The values of these variables will be converted to factors at this stage. However, if the machine learning process requires other types of adjustments, these will be made later. The code below was used to inspect these variables, but its results will not be shown as the key points will be discussed in the following sections.



For the variable CODE_GENDER, men represent 34% of the loans and have a default rate of 10%, which is 42% higher compared to women, who have a default rate of 7%.

CODE_GENDER	TARGET	customers	Percent_of_total	Total_by_gender	Percent_by_gender
F	0	188278	61	202448	93
F	1	14170	5	202448	7
M	0	94404	31	105059	90
M	1	10655	3	105059	10
XNA	0	4	0	4	100

The difference in default rates between consumers who own a car (FLAG_OWN_CAR) and those who do not is only 1.3% higher for those without a car. There is no difference in default rates between those who own property (FLAG_OWN_REALTY) and those who do not.

While several variables can be converted to factors to differentiate between groups, they are not necessarily ranked by importance. Many of these variables contain blank values, which will be replaced with 'unknown'.

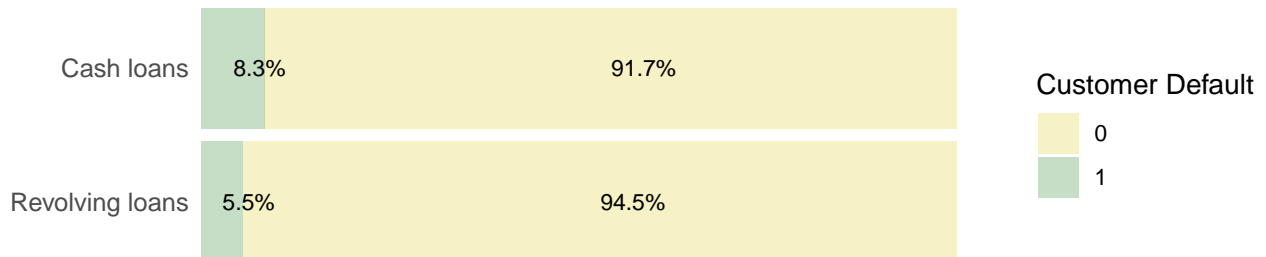
Contract types

The vast majority (93.5%) of defaults are associated with cash loans, compared to just 6.5% for revolving loans. Additionally, the average loan amount for cash loans is more than double at \$578,598, with a maximum limit three times higher at \$4,027,680.

NAME_CONTRACT_TYPE	customers	avg_loan	min	max	Percent.
Cash loans	23221	578598.8	45000	4027680	93.5
Revolving loans	1604	256365.3	135000	1350000	6.5

Relatively, cash loans have a 50% higher default rate compared to revolving loans.

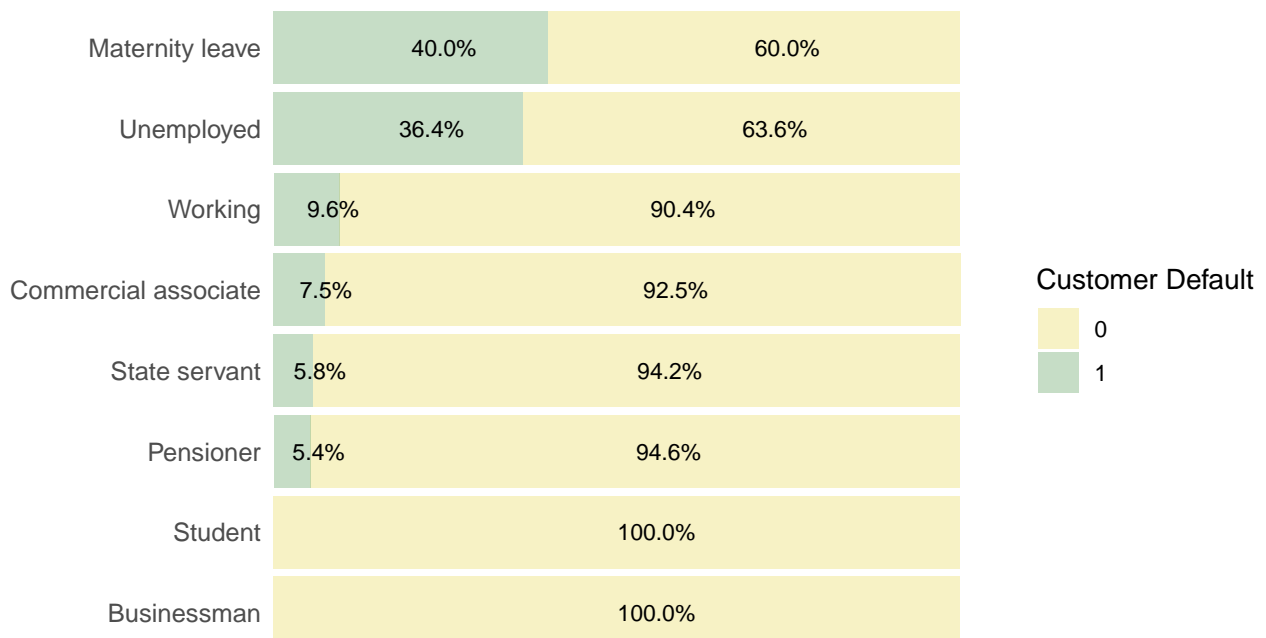
Percentage of Observations by Contract Type and Customer Default



Income type

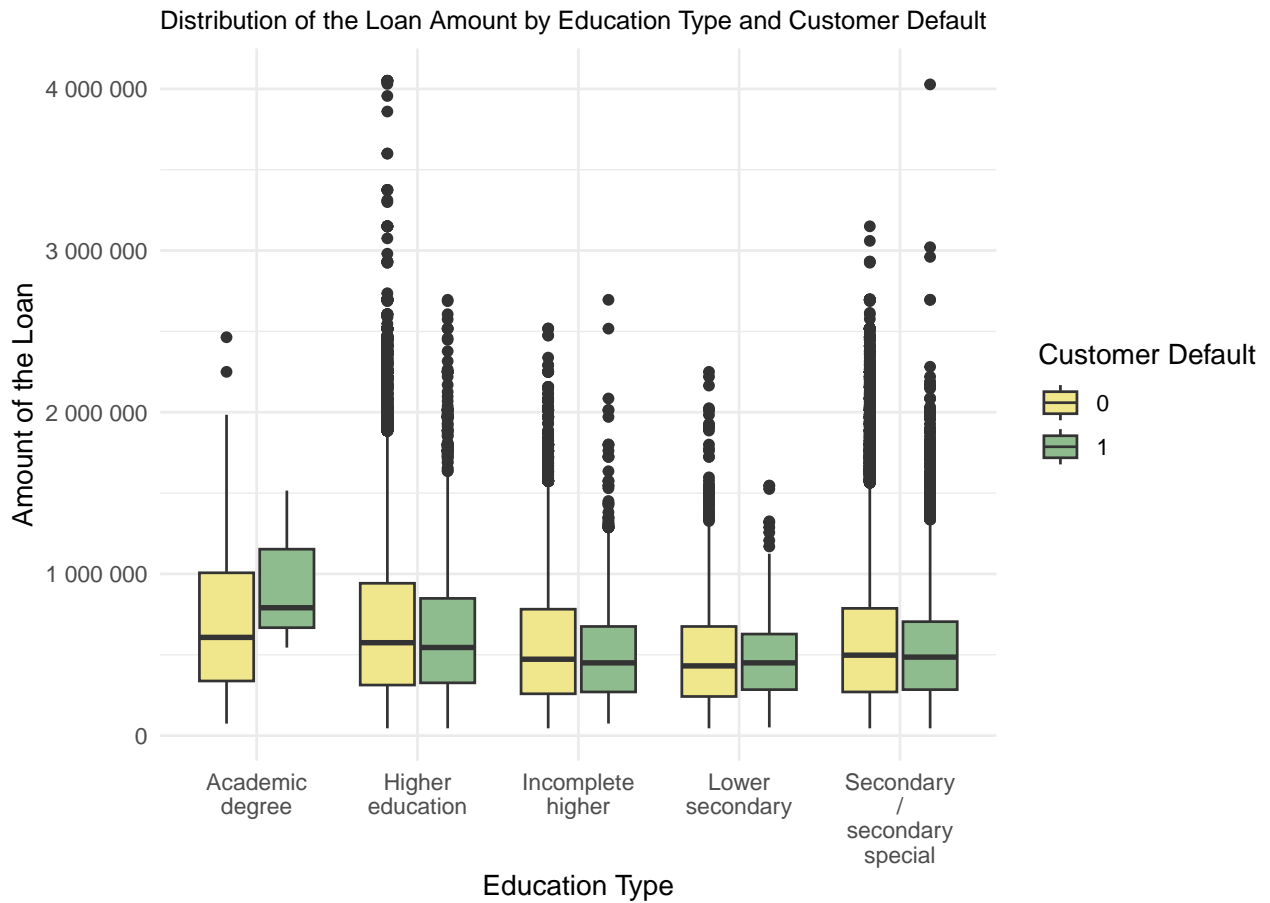
Among the income type groups, those on maternity leave and unemployed have the highest default rates, with more than 35% of loans in default for these groups. The groups of students and business people do not have any individuals with defaults.

Percentage of Observations by Income Type and Customer Default



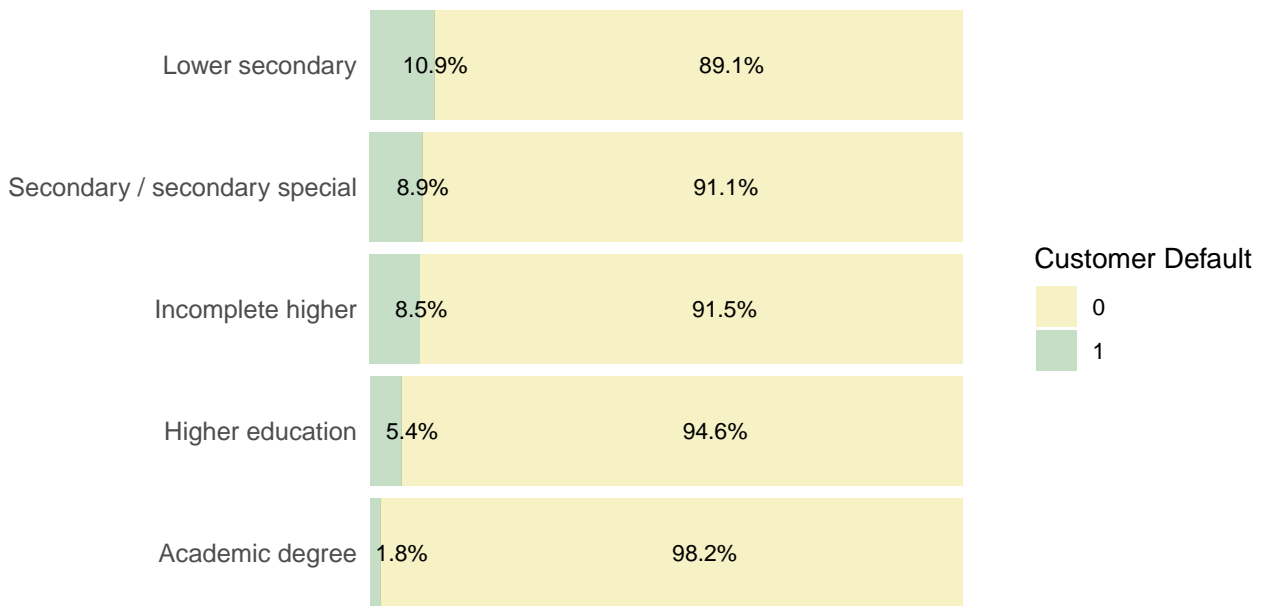
Education

Among the educational levels, it is notable that individuals with an academic degree who are in default have higher loan amounts compared to those who are not in default. This pattern is not significantly observed in other educational levels.



Regarding education, all levels have default cases. However, the default rate for the “Lower Secondary” education group is six times higher compared to the “Academic Degree” group.

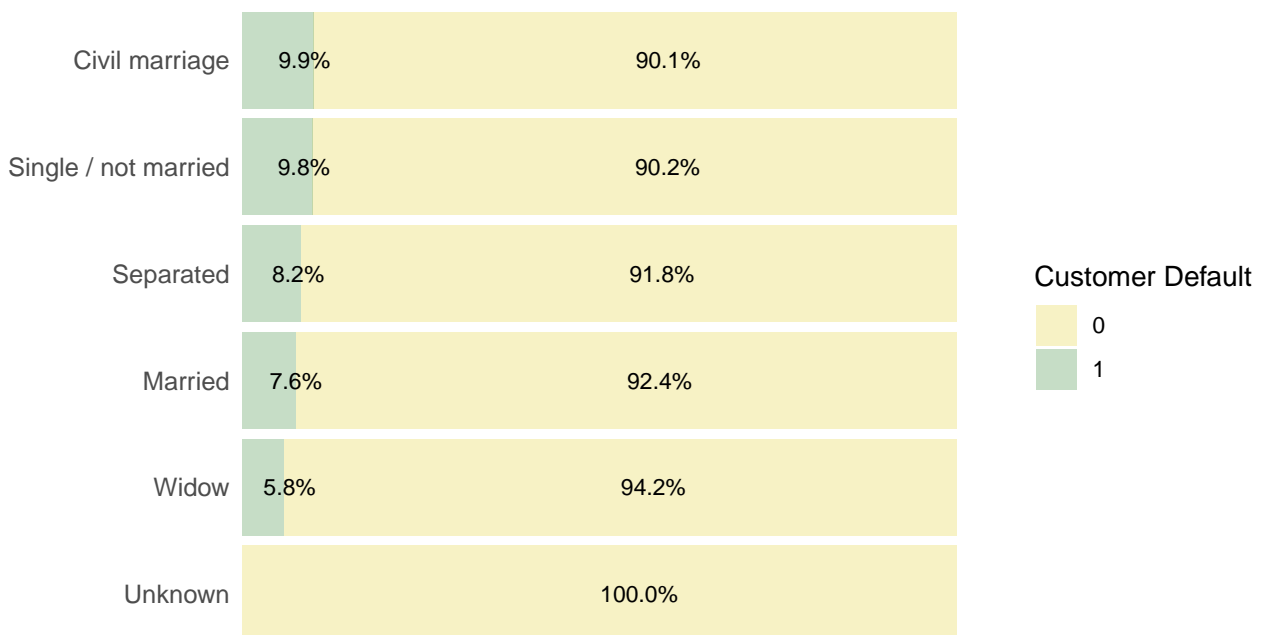
Percentage of Observations by Education level and Customer Default



Family Status

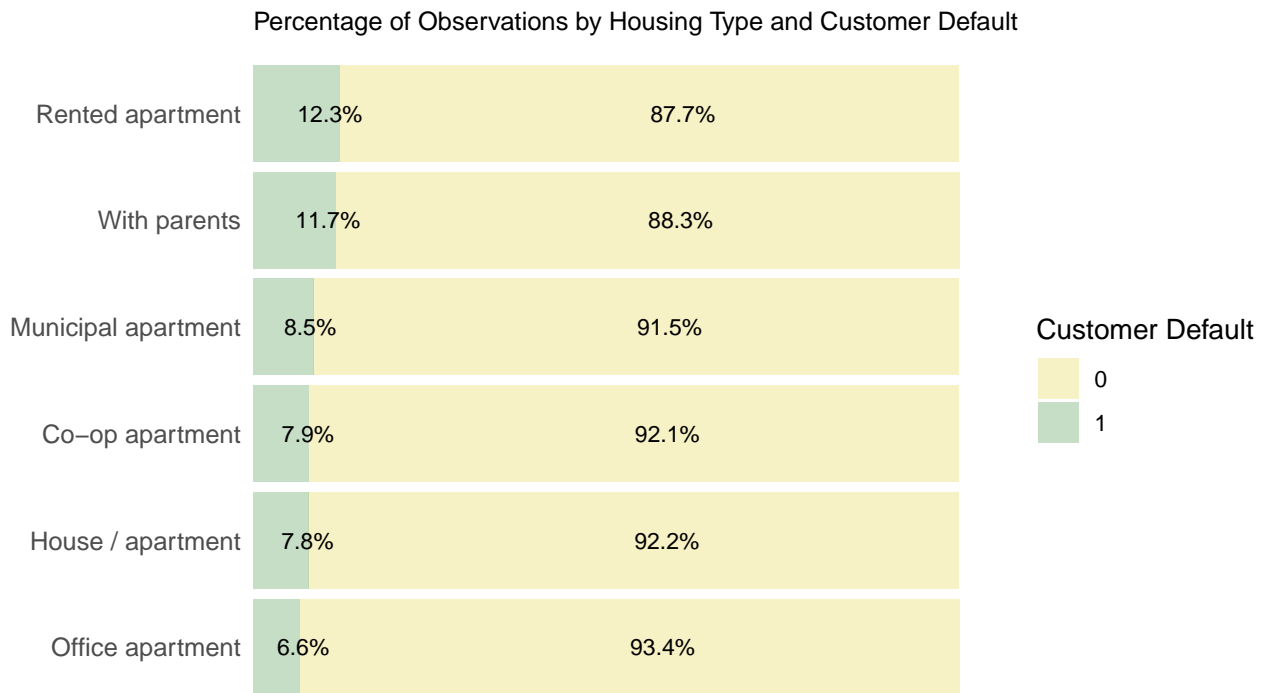
The default rates by family status are quite similar, with a notable exception for the lower percentage of defaults among widows. There are no defaults reported among those with unknown family status.

Percentage of Observations by Family Status and Customer Default



Housing type

Consumers with different housing types exhibit varying default rates. Notably, those living in rented apartments or with their parents have default rates close to 12%. The lowest default rate is observed among those living in office apartments, at 6.6%.



Occupation type

Among the occupation types, 96,391 observations were blank and have been transformed to “Unknown”.

Among the occupation types, low-skill laborers have the highest default rate at 17.2%.

Percentage of Observations by Occupation Type and Customer Default

Low-skill Laborers	17.2%	82.8%
Drivers	11.3%	88.7%
Waiters/barmen staff	11.3%	88.7%
Security staff	10.7%	89.3%
Laborers	10.6%	89.4%
Cooking staff	10.4%	89.6%
Sales staff	9.6%	90.4%
Cleaning staff	9.6%	90.4%
Realty agents	7.9%	92.1%
Secretaries	7.0%	93.0%
Medicine staff	6.7%	93.3%
Private service staff	6.6%	93.4%
Unknown	6.5%	93.5%
IT staff	6.5%	93.5%
HR staff	6.4%	93.6%
Core staff	6.3%	93.7%
Managers	6.2%	93.8%
High skill tech staff	6.2%	93.8%
Accountants	4.8%	95.2%

Customer Default

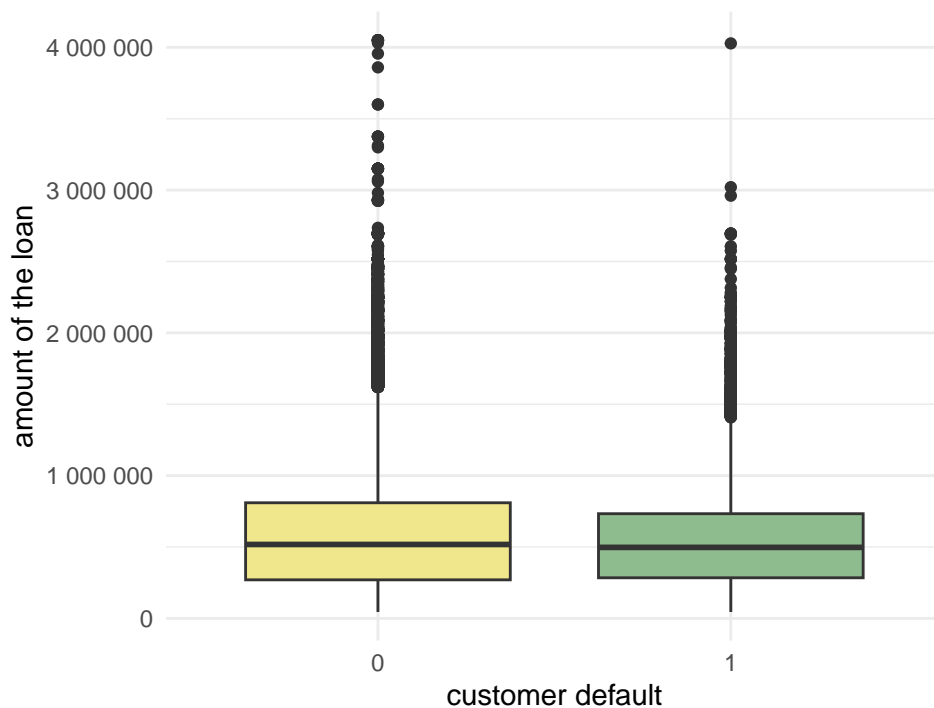


3.3 Additional Observations

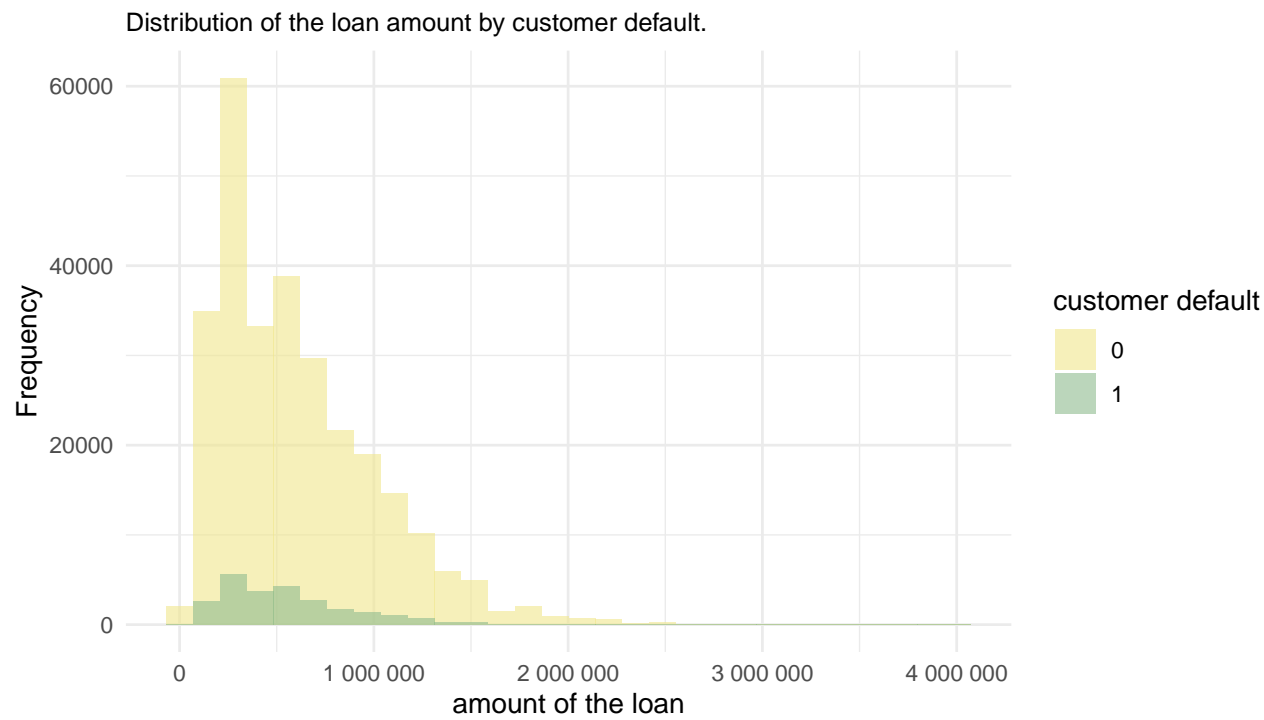
Loan amount and default

The average loan amounts are similar between the groups that default (\$557,778) and those that do not (\$602,648).

Distribution of the loan amount by customer default.

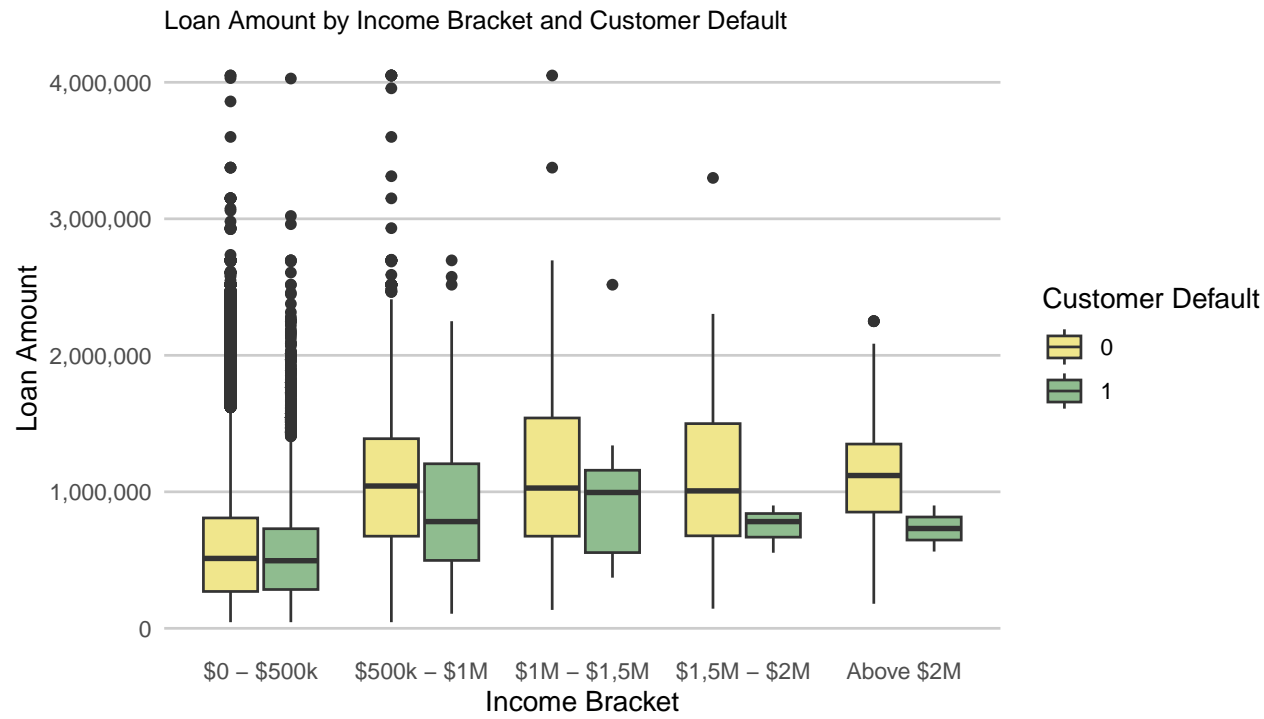


The distribution of borrowers in relation to the loan amounts is left-skewed, with the same pattern of skewness at different intensities across the target variable groups.



Loan amount, income and default

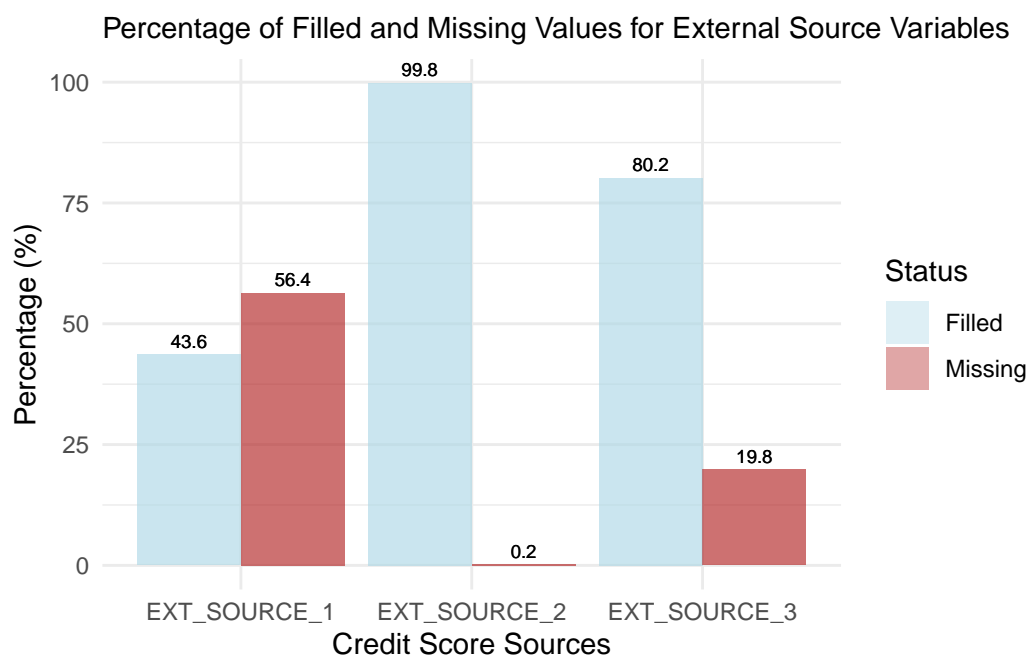
The graph below aims to explore the variation in loan amounts across income groups and default status. The income groups were defined based on the distribution of the majority of the data set.



Among consumers with incomes above \$500,000, the average loan amounts tend to be lower for those who default compared to those who pay on time. The exception is the income range between \$1 million and \$1.5 million, where the average loan amounts are very similar for both groups.

External sources

Credit scores generally serve as a good indicator of whether consumers are likely to delay payments or not. Among the variables used by Home Credit, there are three sources providing such information. These scores are normalized to be comparable, with values ranging from 0 to 1, where lower scores indicate higher risk of default and higher scores represent lower risk of default.

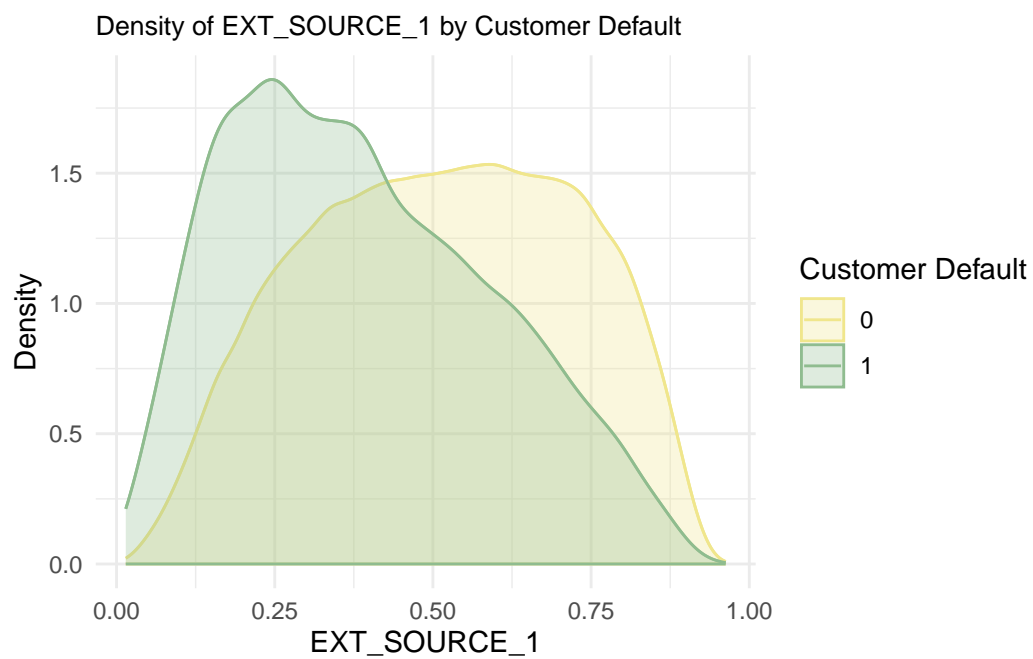


For almost all individuals observed, there is information available from SOURCE 2. For SOURCE 3, more than 70% of individuals have a score, while for SOURCE 1, less than 50% have available data.

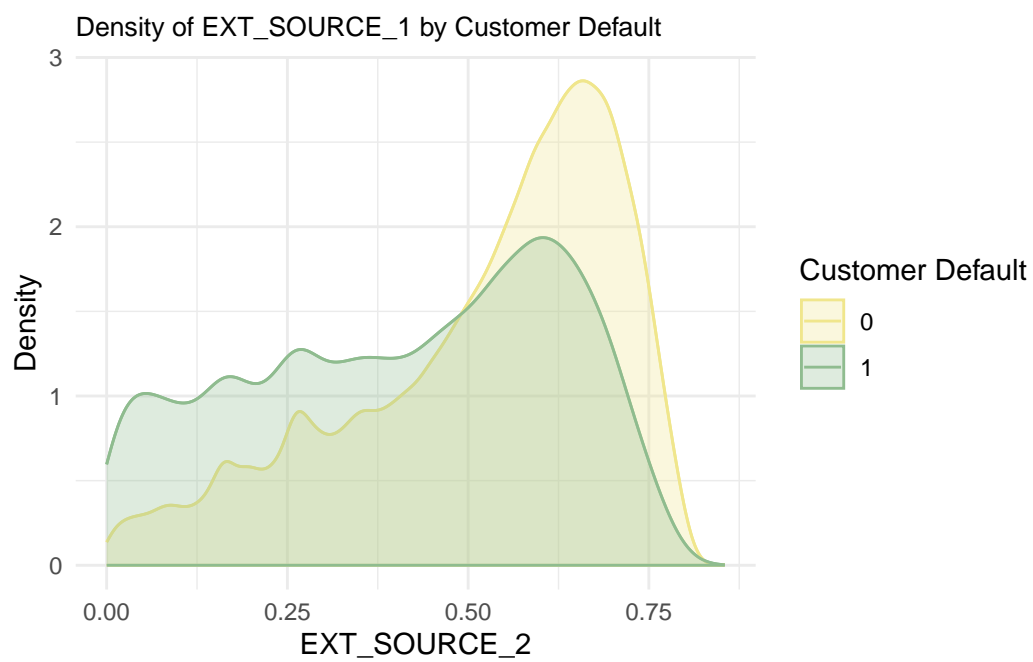
Density plots

The graphs below illustrate the density distributions between consumers who are in default and those who are not.

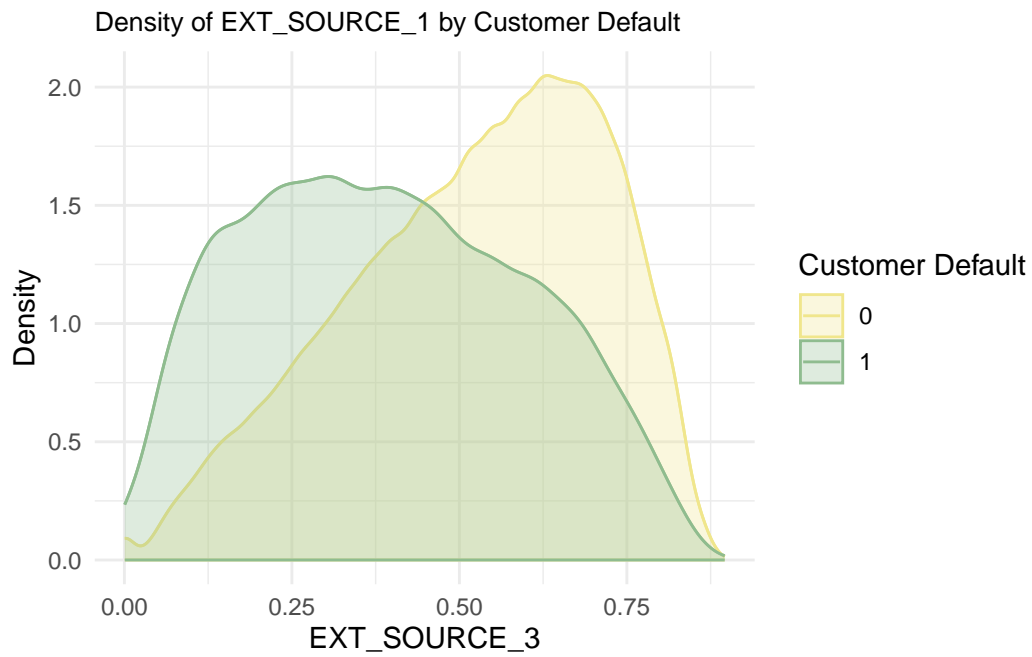
Source 1



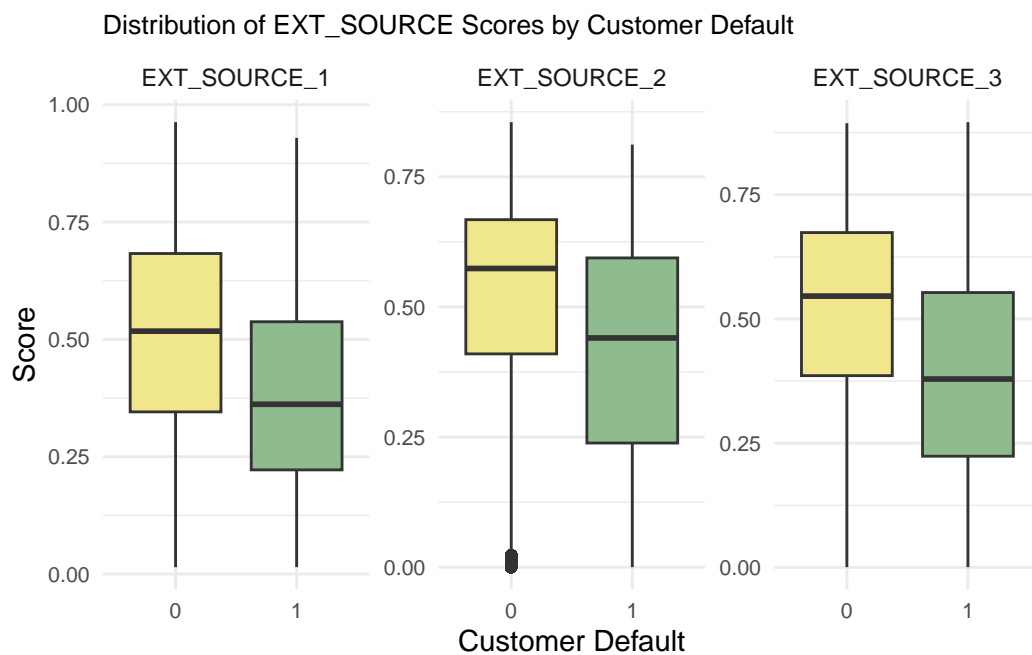
Source 2



Source 3

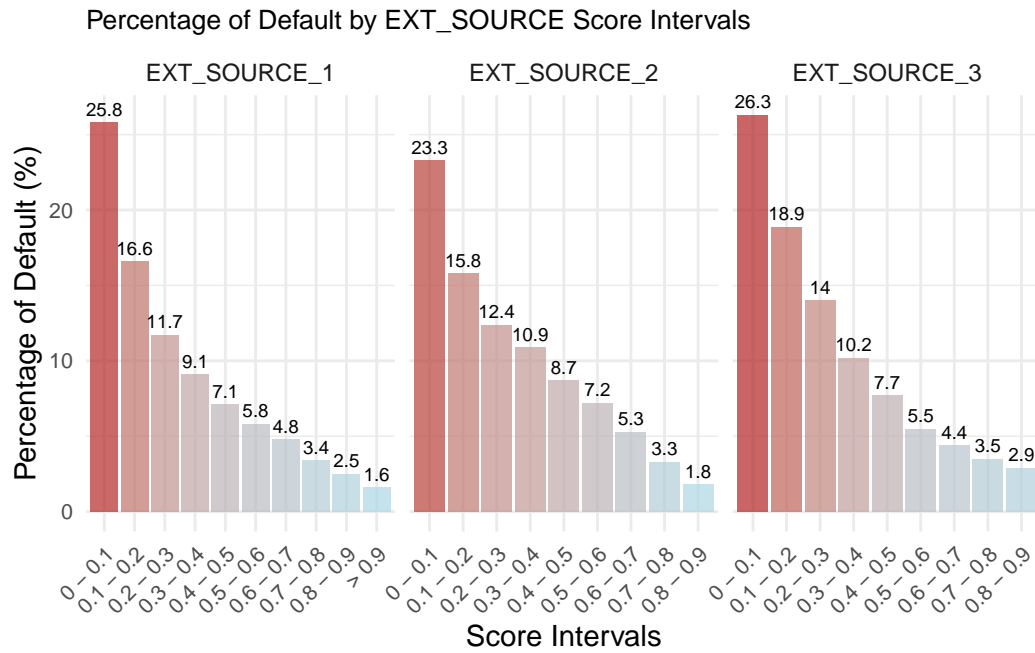


Below is the numeric distribution of scores by default status for each source.



In the box plot above, comparing the three sources, it is clear that, on average, lower scores are observed among consumers who default. However, there are individuals with low scores who make their payments regularly, as well as those with high scores who end up defaulting.

The chart below clearly shows the percentages of defaults based on the scores.



As expected, lower scores tend to show a higher percentage of defaults.

3.4 Correlations

The tables below display the main correlations between numerical variables and default.

	x
REGION_RATING_CLIENT_W_CITY	0.06
REGION_RATING_CLIENT	0.05
FLAG_DOCUMENT_3	0.05
DAYS_BIRTH	0.04
OWN_CAR_AGE	0.04
AMT_REQ_CREDIT_BUREAU_YEAR	0.03
OBS_30_CNT_SOCIAL_CIRCLE	0.03
OBS_60_CNT_SOCIAL_CIRCLE	0.03
DAYS_ID_PUBLISH	0.03
DEF_60_CNT_SOCIAL_CIRCLE	0.03

In addition to the data exploration conducted so far, the table above shows the ten highest positive correlations between numerical variables and default.

It is noteworthy that variables related to consumers' social circles and the regions where they live are present. The submission of document 3, the age of the car, and the time since a new ID issuance are also among the strongest correlations.

	x
EXT_SOURCE_3	-0.16
EXT_SOURCE_2	-0.14
EXT_SOURCE_1	-0.13
DAYS_EMPLOYED_YEARS	-0.06
FLOORSMAX_AVG	-0.05
FLOORSMAX_MODE	-0.05
FLOORSMAX_MEDI	-0.05
AMT_INCOME_TOTAL	-0.05
FLOORSMIN_MEDI	-0.03
FLOORSMIN_AVG	-0.03

Among the negative correlations—where higher values are associated with a lower occurrence of default—external credit information stands out, with values approximately two to three times larger than the fourth variable. An increase in years of employment, age, and income also contributes to a lower likelihood of default. Additionally, variables related to individuals' housing show that the larger these values are, the lower the probability of default.

3.5 Evaluation of Prediction Baselines

Naive baseline and majority classifier

The naive baseline and the Majority Rule Classifier both predict that no customer defaults, as 91.9% of customers do not default (TARGET=0). This results in an accuracy of 91.9% for the model; however, it fails to identify any defaulters, missing the 8.1% who actually do default. Thus, even though the accuracy appears high, this approach is not effective for identifying customers at risk of default. A more advanced model is needed to accurately classify both defaulters and non-defaulters.

Random classifier

The Random Classifier makes predictions randomly based on the distribution of classes in the data set. With 92% of the TARGET being 0 and 8% being 1, the classifier will predict TARGET = 0 with a 92% probability and TARGET = 1 with an 8% probability.

The expected accuracy reflects the class distribution, resulting in an expected accuracy of about 92% for TARGET = 0 and 8% for TARGET = 1, similar to the majority rule classifier.

The simulation below is intended to illustrate the accuracy of using a random classifier for this data set.

```
[1] 0.85
```

Given that a Random Classifier achieves an accuracy of 85%, I expect the developed model to outperform this benchmark, showing improvements in both precision and recall, thereby enhancing its ability to accurately identify true positives and minimize false negatives.

4. Pre-Modeling

4.1 Dummy Encoding

Before proceeding with the dummy encoding process, we will remove instances with missing CODE_Gender values. Additionally, the variable created to indicate individuals with an infinite number of employment days will be converted into a factor.

In addition, all non-numeric variables will be selected to perform dummy encoding, ensuring they are properly transformed for analysis.

For clarity and standardization, the variables will be renamed to remove special characters.

Dummy encoding is a common method for transforming categorical variables into a numerical format, making them usable for machine learning models. This approach creates binary (0 or 1) columns for each category in a categorical variable, allowing the models to interpret these categories without implying any order. Our intention in using dummy encoding is to avoid biases that can arise when treating categorical variables as ordered. Additionally, dummy encoding helps the model better understand the relationships between the input features and the target variable.

The code below combines the variables back into a single data frame after performing dummy encoding. The process will create new columns in the dataset, increasing the total to 253.

4.2 Handling with Missing Data

The process of identifying missing values in the dataset is described. The percentage of missing data for each column is calculated and visualized in the following table. Then, the missing values are replaced with the mean for each variable using the aggregate function.

Show entries

Search:

Percentage of Missing Values for Each Column

	Column	missing_values
1	COMMONAREA_AVG	69.9%
2	COMMONAREA_MODE	69.9%
3	COMMONAREA_MEDI	69.9%
4	NONLIVINGAPARTMENTS_AVG	69.4%
5	NONLIVINGAPARTMENTS_MODE	69.4%
6	NONLIVINGAPARTMENTS_MEDI	69.4%
7	LIVINGAPARTMENTS_AVG	68.4%
8	LIVINGAPARTMENTS_MODE	68.4%
9	LIVINGAPARTMENTS_MEDI	68.4%
10	FLOORSMIN_AVG	67.8%

Showing 1 to 10 of 253 entries

Previous

1

2

3

4

5

...

26

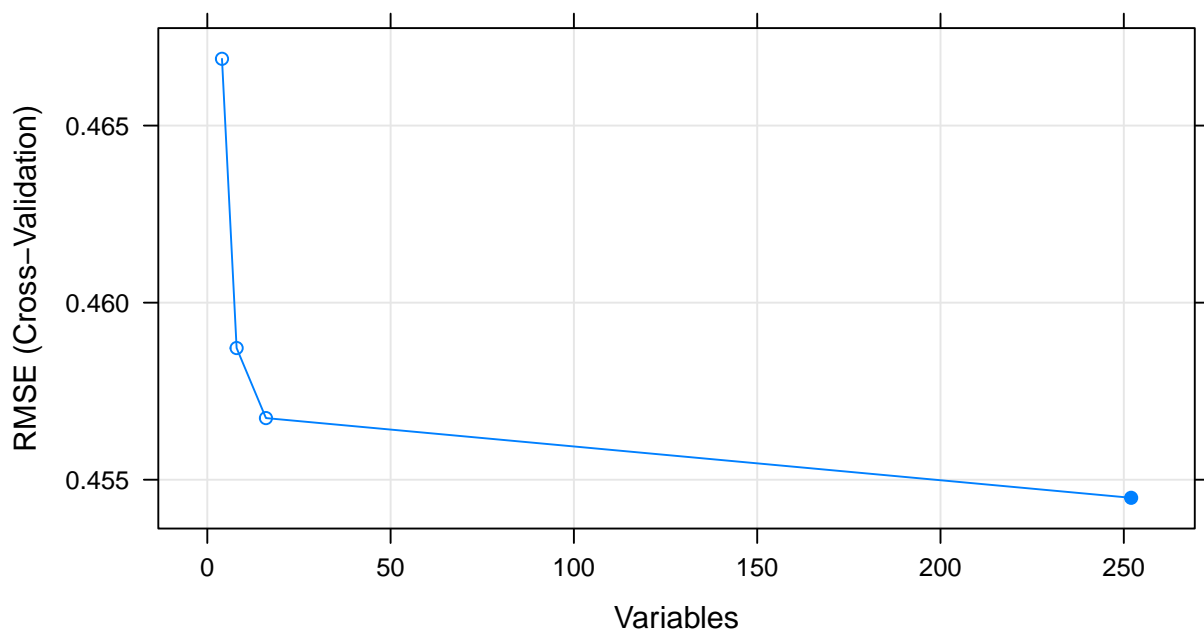
Next

4.3 Sampling and Balancing

Due to the high computational cost of processing the entire dataset, we opted to perform downsampling by using nearly all instances that represent `TARGET == 1`. To balance the dataset and avoid bias in the modeling process, we randomly sampled the same number of instances with `TARGET == 0`. As a result, the dataset is balanced at 50% for each class, using approximately 23% of the total dataset (72,000 instances) for training and testing. Additionally, We removed the variable `SK_ID_CURR` because it is not necessary for the modeling process, as it is just a registration number.

Recursive Feature Elimination for Variable Selection

This code uses the Recursive Feature Elimination (RFE) method to select important variables from the dataset, aiming to reduce computational time for subsequent analyses. This process takes many hours to complete, so I chose to save the vector and load it in its final version.



The results indicate that among the 253 variables in this dataset, fewer than 25 correspond to the lowest root mean squared error. Below is the list of the 10 most important predictors.

Predictors
EXT_SOURCE_3
EXT_SOURCE_2
EXT_SOURCE_1
AMT_CREDIT
DAYS_EMPLOYED_YEARS
AMT_GOODS_PRICE
DAYS_BIRTH
AMT_ANNUITY
DAYS_ID_PUBLISH
NAME_EDUCATION_TYPE_higher_education

To simplify the dataset, we removed the variables NAME_EDUCATION_TYPE_secondary__secondary_special, NAME_EDUCATION_TYPE_higher_secondary_special, NAME_CONTRACT_TYPE_revolving_loans, and CODE_GENDER_f, as they were redundant or exhibited collinearity, ensuring a cleaner and more effective set of predictors. Additionally, to include a broader range of predictors and improve the model's quality, we chose to focus on the 50 most impactful variables, going beyond the standard Recursive Feature Elimination approach to refine the selection further.

4.4 Data partition

Partition

The dataset will be partitioned into training and testing sets. The TARGET variable will be transformed into a binary format to facilitate better model interpretation, and 70% of the data will be allocated for training while the remaining 30% will be reserved for testing, ensuring effective evaluation on unseen data.

The training control is configured to use repeated cross-validation, enhancing the reliability of the model evaluation. This setup incorporates 10 folds with 2 repeats and enables probability calculations for each class.

5. Modeling

Due to the computational cost of the modeling process, we opted to save the final versions of the models, making the knitting process easier.

All the metrics displayed within the model blocks refer to their performance on the test set (30%, 21,600 instances).

5.1 Logistic Regression

Since we have already performed the processes of variable elimination and standardization, I decided not to use techniques like Lasso and Ridge at this moment. (Note: we tested them and still opted to remove it.)

Call:

NULL

Deviance Residuals:

Min	1Q	Median	3Q	Max
-3.20	-0.84	-0.58	1.04	2.96

Coefficients:

	Estimate	Std. Error	z	value	Pr(> z)
(Intercept)	-0.8202	0.0107	-76.9	<2e-16	***
NAME_CONTRACT_TYPE_cash_loans	0.0923	0.0118	7.8	7e-15	***
CODE_GENDER_m	0.1516	0.0106	14.2	<2e-16	***
AMT_INCOME_TOTAL	-0.1158	0.0175	-6.6	4e-11	***
AMT_CREDIT	0.8704	0.0645	13.5	<2e-16	***
AMT_ANNUITY	0.1518	0.0176	8.6	<2e-16	***
AMT_GOODS_PRICE	-0.9871	0.0661	-14.9	<2e-16	***
DAYS_BIRTH	-0.0217	0.0134	-1.6	0.105	
DAYS_REGISTRATION	0.0291	0.0112	2.6	0.009	**
DAYS_ID_PUBLISH	0.0744	0.0108	6.9	6e-12	***
OWN_CAR_AGE	0.0401	0.0101	4.0	7e-05	***
REGION_RATING_CLIENT	-0.0472	0.0342	-1.4	0.168	
REGION_RATING_CLIENT_W_CITY	0.1240	0.0343	3.6	3e-04	***
REG_CITY_NOT_LIVE_CITY	0.0506	0.0100	5.0	5e-07	***
EXT_SOURCE_1	-0.1998	0.0117	-17.1	<2e-16	***
EXT_SOURCE_2	-0.4460	0.0110	-40.5	<2e-16	***
EXT_SOURCE_3	-0.5050	0.0108	-46.7	<2e-16	***
APARTMENTS_AVG	-0.1179	0.1974	-0.6	0.550	
BASEMENTAREA_AVG	-0.2099	0.1747	-1.2	0.230	
YEARS_BEGINEXPLUATATION_AVG	-0.1244	0.0997	-1.2	0.212	
YEARS_BUILD_AVG	0.5551	0.2218	2.5	0.012	*

COMMONAREA_AVG	-0.0357	0.0703	-0.5	0.612
FLOORSMAX_AVG	0.0817	0.1624	0.5	0.615
LANDAREA_AVG	0.0175	0.0549	0.3	0.750
LIVINGAPARTMENTS_AVG	-0.0761	0.1420	-0.5	0.592
LIVINGAREA_AVG	0.0231	0.1210	0.2	0.849
NONLIVINGAREA_AVG	0.0597	0.0786	0.8	0.448
APARTMENTS_MODE	0.0034	0.0938	0.0	0.971
BASEMENTAREA_MODE	0.1103	0.0797	1.4	0.166
YEARS_BEGINEXPLUATATION_MODE	0.0514	0.0487	1.1	0.291
YEARS_BUILD_MODE	-0.0781	0.0772	-1.0	0.312
COMMONAREA_MODE	0.0260	0.0699	0.4	0.710
FLOORSMAX_MODE	-0.0165	0.0812	-0.2	0.839
LANDAREA_MODE	0.0131	0.0688	0.2	0.849
LIVINGAPARTMENTS_MODE	-0.0617	0.0703	-0.9	0.380
LIVINGAREA_MODE	-0.1091	0.0898	-1.2	0.224
NONLIVINGAREA_MODE	-0.0405	0.0519	-0.8	0.435
APARTMENTS_MEDI	0.1292	0.2074	0.6	0.533
BASEMENTAREA_MEDI	0.0877	0.1742	0.5	0.614
YEARS_BEGINEXPLUATATION_MEDI	0.0528	0.1003	0.5	0.599
YEARS_BUILD_MEDI	-0.5002	0.2148	-2.3	0.020 *
FLOORSMAX_MEDI	-0.1083	0.1764	-0.6	0.539
LANDAREA_MEDI	-0.0352	0.0845	-0.4	0.677
LIVINGAPARTMENTS_MEDI	0.1582	0.1536	1.0	0.303
LIVINGAREA_MEDI	0.0677	0.1460	0.5	0.643
NONLIVINGAREA_MEDI	-0.0117	0.0909	-0.1	0.898
TOTALAREA_MODE	-0.0041	0.0318	-0.1	0.897
DEF_30_CNT_SOCIAL_CIRCLE	0.0744	0.0104	7.1	9e-13 ***
DAYS_LAST_PHONE_CHANGE	0.0412	0.0110	3.7	2e-04 ***
FLAG_DOCUMENT_6	-0.0471	0.0120	-3.9	8e-05 ***
DAYS_EMPLOYED_YEARS	-0.1446	0.0121	-11.9	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 64161 on 50399 degrees of freedom
Residual deviance: 55827 on 50349 degrees of freedom
AIC: 55929

Number of Fisher Scoring iterations: 4

The coefficient values in the logistic regression are expressed as probabilities of occurrence, indicating the likelihood of an event happening for each predictor. However, interpreting these coefficients can be complex, as they are influenced by the scale and type of variables used in the model, with only 20 out of 50 predictors meeting the 5% significance level or lower.

These include: NAME_CONTRACT_TYPE_cash_loans, CODE_GENDER_m, AMT_INCOME_TOTAL, AMT_CREDIT, AMT_ANNUITY, AMT_GOODS_PRICE, DAYS_ID_PUBLISH, OWN_CAR_AGE, REGION_RATING_CLIENT_W_CITY, REG_CITY_NOT_LIVE_CITY, EXT_SOURCE_1, EXT_SOURCE_2, EXT_SOURCE_3, DEF_30_CNT_SOCIAL_CIRCLE, DAYS_LAST_PHONE_CHANGE, FLAG_DOCUMENT_6, DAYS_EMPLOYED_YEARS, DAYS_REGISTRATION, YEARS_BUILD_AVG, and YEARS_BUILD_MEDI.

Among these 20 predictors, 8 are listed among the top 10 most important predictors identified through the RFE process.

Train set Metrics

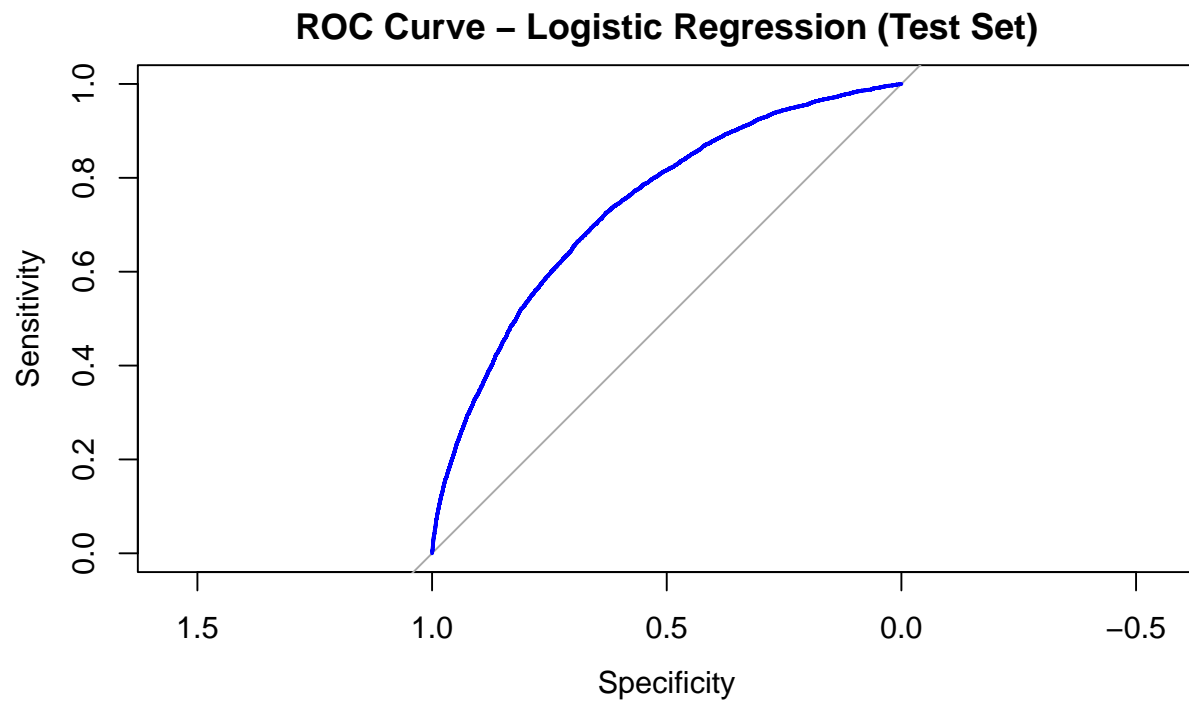
AUC Train set - Logistic: 0.74

Accuracy Train set - Logistic: 0.72

The model showed the above performance on the train set.

AUC - Logistic Model - Test set

AUC Test set - Logistic: 0.74



Confusion Matrix - Logistic Regression Test Set

Confusion Matrix and Statistics

	Reference	
Prediction	Class0	Class1
Class0	12793	4541
Class1	1607	2659

Accuracy : 0.7154

95% CI : (0.7093, 0.7214)

No Information Rate : 0.6667

P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.2869

Mcnemar's Test P-Value : < 2.2e-16

Sensitivity : 0.3693

Specificity : 0.8884

Pos Pred Value : 0.6233
Neg Pred Value : 0.7380
Prevalence : 0.3333
Detection Rate : 0.1231
Detection Prevalence : 0.1975
Balanced Accuracy : 0.6289

'Positive' Class : Class1

Results for the logistic model (test set) include an AUC of 0.74, accuracy of 71%, sensitivity (true positive for Class 1) of 37%, and specificity (true negative for Class 0) of 89%. The logistic regression model provides a reasonable baseline for distinguishing between defaulters and non-defaulters, being better than a classification based on the no-information rate.

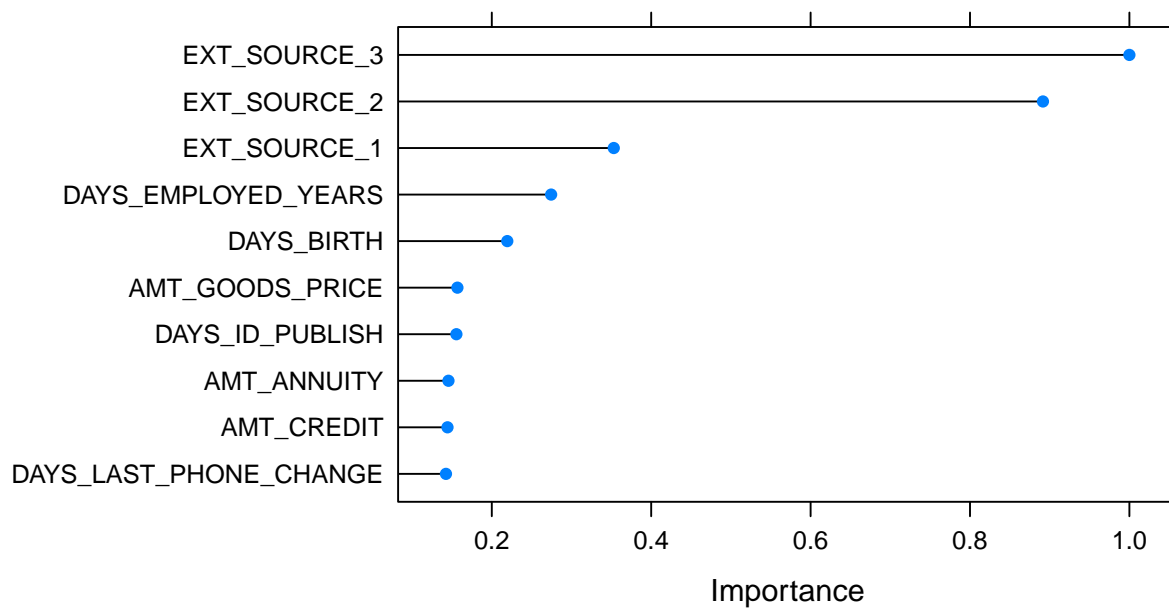
5.2 Random Forest

In the Random Forest model setup, I had to make several adjustments as the model was showing signs of overfitting to the training set. Key hyperparameters included experimenting with different values for the number of features considered at each split to balance model complexity and accuracy. I also explored different split criteria and modified the minimum node size to control tree depth and reduce overfitting. The model was trained using repeated cross-validation with 10 folds and 2 repeats. Class probabilities were calculated for performance evaluation, and parallel processing was enabled to speed up training. Also I used the “ranger” method for Random Forest modeling because it’s efficient with large datasets and faster than traditional implementations, plus it supports parallel processing.

	mtry	splitrule	min.node.size
15	6	gini	250

The best hyperparameters obtained for the model are: mtry = 15, which means that 15 variables were selected for each split; splitrule = “gini”, indicating that the splits were made based on the Gini index; and min.node.size = 250, which sets the minimum number of observations required in each node before it can be split.

Random Forest Top 10 Most Important Variables (Normalized)



These are the top 10 most important variables in the model, with EXT_SOURCE_3 being the most important, followed by EXT_SOURCE_2 and EXT_SOURCE_1. The other variables, such as DAYS_EMPLOYED_YEARS and DAYS_BIRTH, have progressively lower importance in comparison.

Among these 10 predictors, 9 are listed among the top 10 most important predictors identified through the RFE process.

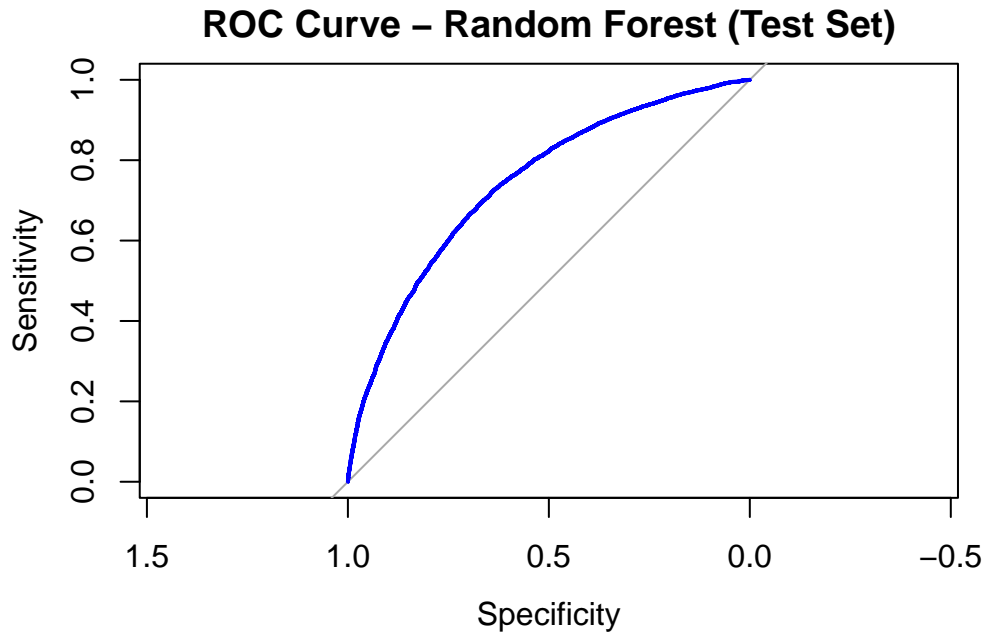
Train set Metrics

AUC Train set - Random Forest: 0.84

Accuracy Train set - Random Forest: 0.76

Initially with maximum overfitting to the train set, the model now has better balance, aiming for greater generalization.

AUC - Random Forest - Test Set



AUC Test set - Random Forest: 0.74

Confusion Matrix - Random Forest - Test Set

Confusion Matrix and Statistics

	Reference	
Prediction	Class0	Class1
Class0	13197	4912
Class1	1203	2288

Accuracy : 0.7169

95% CI : (0.7108, 0.7229)

No Information Rate : 0.6667

P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.2689

McNemar's Test P-Value : < 2.2e-16

Sensitivity : 0.3178

```
Specificity : 0.9165
Pos Pred Value : 0.6554
Neg Pred Value : 0.7288
Prevalence : 0.3333
Detection Rate : 0.1059
Detection Prevalence : 0.1616
Balanced Accuracy : 0.6171

'Positive' Class : Class1
```

Results for the Random Forest model (test set) include an AUC of 0.74, accuracy of 72%, sensitivity (true positive for Class 1) of 32%, and specificity (true negative for Class 0) of 92%. The Random Forest model offers a slight improvement in specificity compared to logistic regression, making it more effective at identifying non-defaulters. However, its lower sensitivity means it may miss more defaulters. While less interpretable than logistic regression, Random Forest can capture non-linear relationships and interactions among variables, which is useful for complex datasets, but it does not show a significant performance gain in this case.

5.3 XGBoost Model

The XGBoost model was optimized in multiple stages to improve computational performance. The dataset was split into training and test sets, then transformed into the xgb.DMatrix format. Hyperparameter tuning was performed in rounds, adjusting key parameters like max_depth, min_child_weight, and eta to maximize the AUC score. The second round fine-tuned sampling parameters (subsample, colsample_bytree) to reduce overfitting, followed by regularization parameters (gamma, lambda, alpha) in the final round.

Finally, the number of boosting iterations (nrounds) was optimized, with cross-validation used throughout to identify the best parameter combination. This process improved performance while minimizing overfitting.

Best hyperparameters and Predictors

Show entries

Search:

XGBoost Feature Importance

	Feature	Gain	Cover	Frequency
1	EXT_SOURCE_3	25.00%	11.00%	8.00%
2	EXT_SOURCE_2	23.00%	11.00%	7.00%
3	EXT_SOURCE_1	9.00%	9.00%	7.00%
4	DAYS_EMPLOYED_YEARS	5.00%	5.00%	5.00%
5	AMT_CREDIT	4.00%	6.00%	6.00%
6	DAYS_BIRTH	4.00%	7.00%	7.00%
7	AMT_GOODS_PRICE	4.00%	7.00%	5.00%
8	AMT_ANNUITY	3.00%	5.00%	6.00%
9	OWN_CAR_AGE	2.00%	5.00%	3.00%
10	DAYS_ID_PUBLISH	2.00%	4.00%	6.00%

Showing 1 to 10 of 50 entries

Previous

1

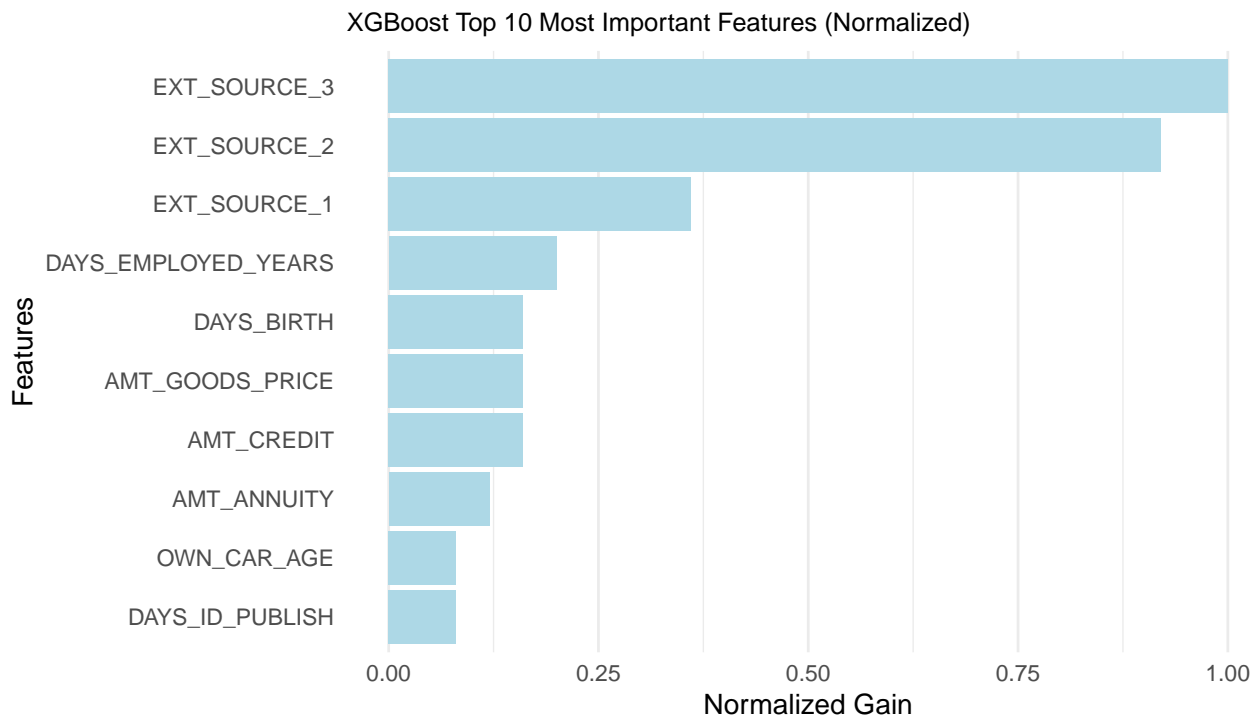
2

3

4

5

Next



The final model was trained with 250 trees, and the AUC metric was used to evaluate its performance. The feature importance analysis revealed that EXT_SOURCE_3 and EXT_SOURCE_2 are the most relevant variables, with contributions of 24% and 23%, respectively.

Additionally, “Cover” and “Frequency” metrics provide insight into how these features contribute to the model. “Cover” measures the proportion of samples that each feature influences, indicating how broadly the feature affects the data. “Frequency” shows how often each feature is used to split the data in the decision trees. EXT_SOURCE_3 and EXT_SOURCE_2 not only have high contributions to the model but also exhibit higher cover and frequency, meaning they are both widely and frequently utilized in the model’s decision-making process.

Other important variables include EXT_SOURCE_1, DAYS_EMPLOYED_YEARS, and DAYS_BIRTH, which also play significant roles in the model’s predictions. These variables help model the relationship between customer characteristics and the likelihood of success in classifying the TARGET variable.

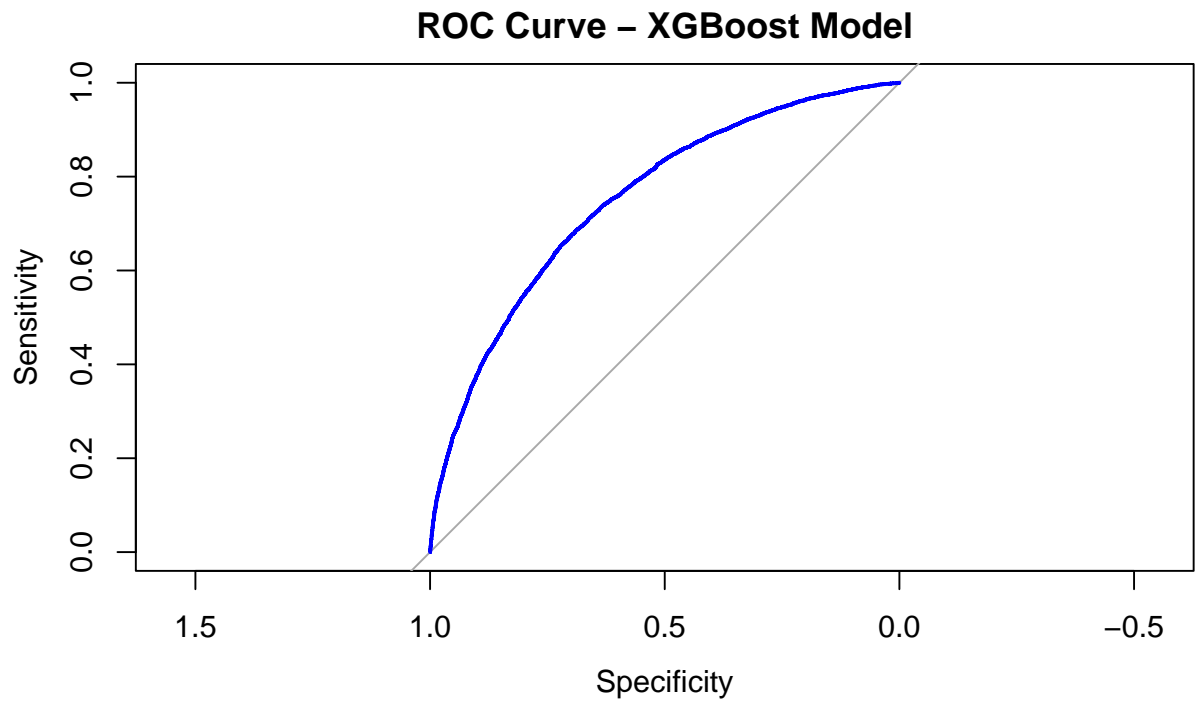
Only OWN_CAR_AGE and NAME_EDUCATION_TYPE_higher_education are not among the top 10 most important variables from the RFE.

Train Set Metrics

AUC Train set - XGBoost: 0.8

Accuracy Train set - XGBoost: 0.75

AUC - XGBoost Model



AUC Test set - XGBoost: 0.7518127

Confusion Matrix - XGBoost Model

Confusion Matrix and Statistics

Reference		
Prediction	Class0	Class1
Class0	12724	4239
Class1	1676	2961

Accuracy : 0.7262

95% CI : (0.7202, 0.7321)

No Information Rate : 0.6667

P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.3237

Mcnemar's Test P-Value : < 2.2e-16


```
Sensitivity : 0.4113
Specificity : 0.8836
Pos Pred Value : 0.6386
Neg Pred Value : 0.7501
Prevalence : 0.3333
Detection Rate : 0.1371
Detection Prevalence : 0.2147
Balanced Accuracy : 0.6474

'Positive' Class : Class1
```

Results for the XGBoost model (test set) include an AUC of 0.75, accuracy of 73%, sensitivity of 42%, and specificity of 88%. The XGBoost model outperforms both logistic regression and Random Forest, offering the highest AUC and accuracy with a balanced performance in sensitivity and specificity. XGBoost's ability to capture complex patterns and interactions likely contributed to this improved performance. Additionally, the modeling process was much faster than with Random Forest, given the way it was implemented.

5.4 LightGBM Model

The LightGBM model was optimized by tuning key parameters in multiple stages. First, the dataset was split into training and test sets and converted to the `lgb.Dataset` format. The model was then fine-tuned by adjusting tree parameters (`max_depth`, `min_data_in_leaf`, `learning_rate`), followed by sampling and regularization parameters to minimize overfitting. Cross-validation was used to determine the optimal number of boosting iterations. The training process was faster than Random Forest, showcasing LightGBM's efficiency while maintaining strong performance.

Best hyperparameters and Predictors

Show entries

Search:

LightGBM Feature Importance

	Feature	Gain	Cover	Frequency
1	EXT_SOURCE_2	20.00%	7.00%	7.00%
2	EXT_SOURCE_3	19.00%	10.00%	6.00%
3	EXT_SOURCE_1	7.00%	7.00%	6.00%
4	DAYS_EMPLOYED_YEARS	5.00%	5.00%	6.00%
5	DAYS_BIRTH	5.00%	7.00%	7.00%
6	AMT_CREDIT	4.00%	7.00%	6.00%
7	AMT_ANNUITY	4.00%	5.00%	6.00%
8	DAYS_LAST_PHONE_CHANGE	3.00%	3.00%	6.00%
9	DAYS_ID_PUBLISH	3.00%	3.00%	6.00%
10	AMT_GOODS_PRICE	3.00%	8.00%	4.00%

Showing 1 to 10 of 50 entries

Previous

1

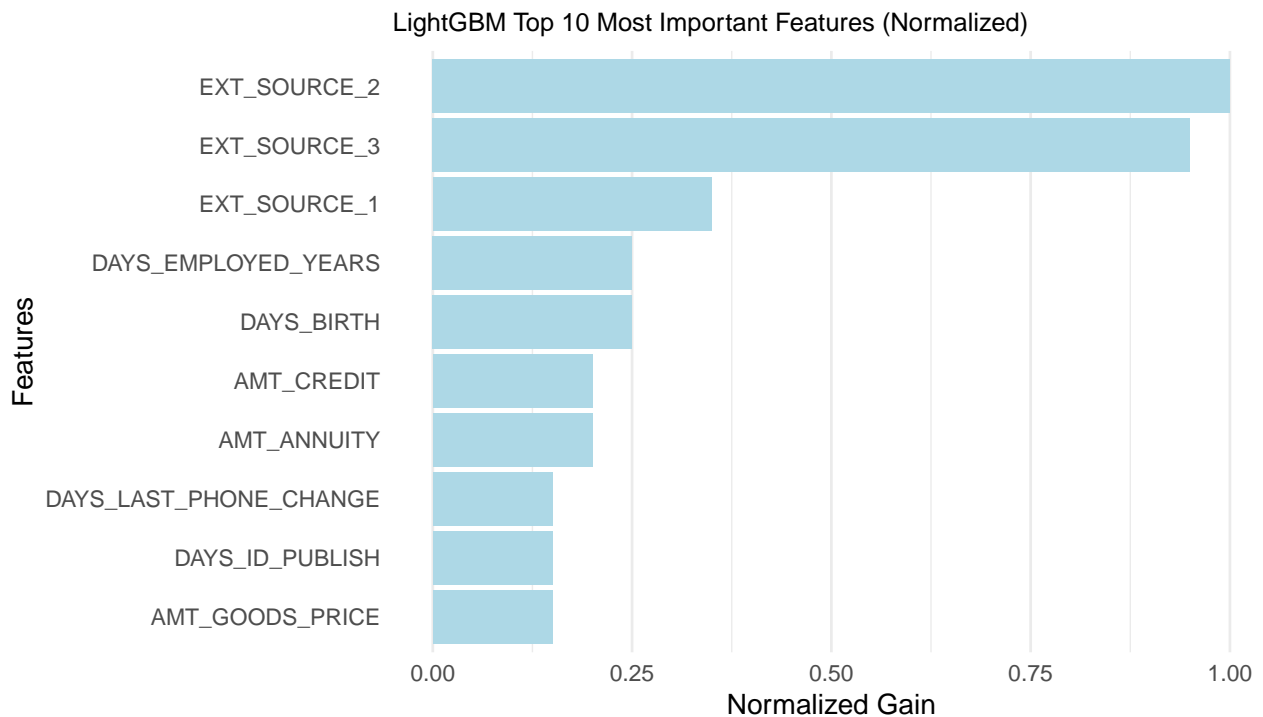
2

3

4

5

Next



The LightGBM model was trained with 2000 boosting iterations, with binary classification as the objective and binary error as the evaluation metric. The best hyperparameters include boosting type “gbdt” and 4 threads for parallelization.

Feature importance analysis shows that EXT_SOURCE_2 and EXT_SOURCE_3 are the most influential variables, with contributions of 20% and 19%, respectively. These features also have the highest cover, indicating that they affect a significant portion of the data, and are frequently used in decision tree splits. Other important features, such as EXT_SOURCE_1, DAYS_EMPLOYED_YEARS, and DAYS_BIRTH, also contribute to the model’s predictions but to a lesser extent.

For this model, only AMT_ANNUITY, DAYS_LAST_PHONE_CHANGE, and NAME_EDUCATION_TYPE_higher_education are not among the top 10 most important features in the RFE.

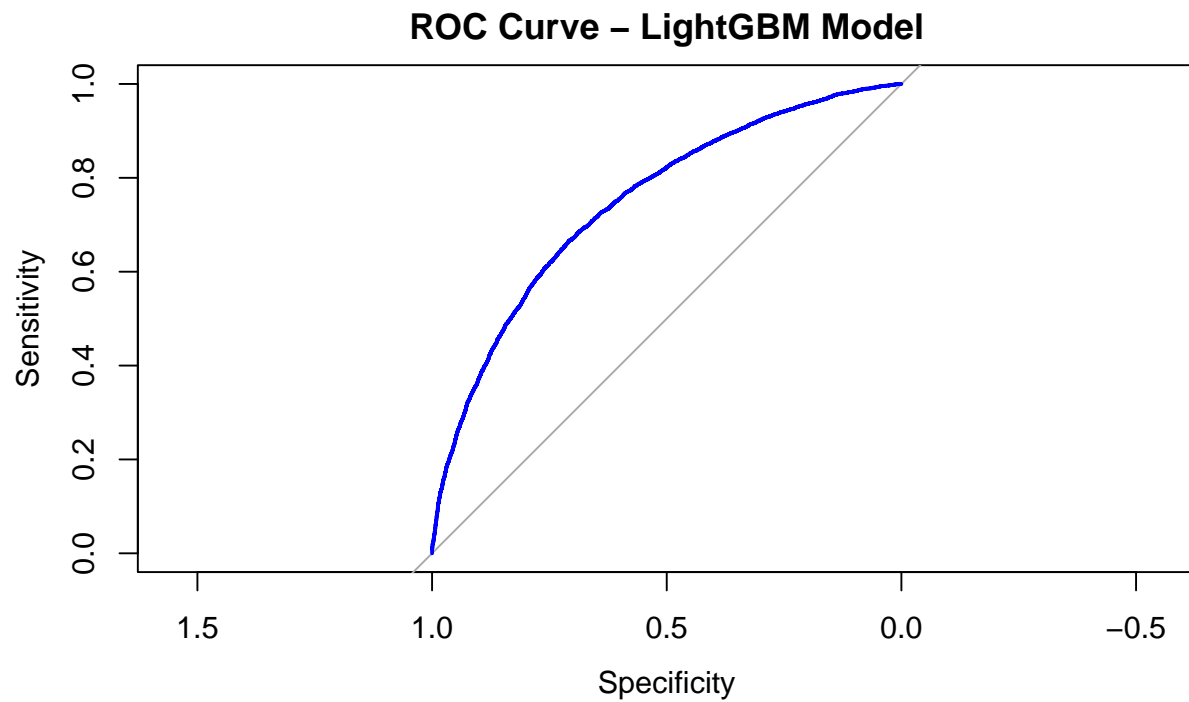
Train Set Metrics

AUC Train set - LightGBM: 0.87

Accuracy Train set - LightGBM: 0.8

AUC - LightGBM Model

AUC Test set - LightGBM: 0.75



Confusion Matrix - LightGBM Model

Confusion Matrix and Statistics

	Reference	
Prediction	0	1
0	12636	4035
1	1843	3086

Accuracy : 0.7279

95% CI : (0.7219, 0.7338)

No Information Rate : 0.6703

P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.332

Mcnemar's Test P-Value : < 2.2e-16

Sensitivity : 0.4334

Specificity : 0.8727

Pos Pred Value : 0.6261
 Neg Pred Value : 0.7580
 Prevalence : 0.3297
 Detection Rate : 0.1429
 Detection Prevalence : 0.2282
 Balanced Accuracy : 0.6530

'Positive' Class : 1

Results for the LightGBM model (test set) include an AUC of 0.75, accuracy of 73%, sensitivity (true positive for Class 1) of 43%, and specificity (true negative for Class 0) of 88%. The LightGBM model performs almost identically to the XGBoost model, with both models achieving the same AUC and similar accuracy and specificity. While the sensitivity of LightGBM is slightly lower than that of XGBoost, both models effectively balance the identification of defaulters and non-defaulters.

6. Comparison Metrics

Final Metrics

The metrics below are simply to summarize what has already been interpreted and provide a unique comparative visualization of the models' performance on both the training and test sets.

Table 1: Model Performance Metrics

Model	AUC	Accuracy	Sensitivity TPR_1	Specificity TNR_0
Logistic Regression (Train)	0.74	0.72	0.38	0.89
Logistic Regression (Test)	0.74	0.72	0.37	0.89
Random Forest (Train)	0.84	0.76	0.40	0.94
Random Forest (Test)	0.74	0.72	0.32	0.92
XGBoost (Train)	0.80	0.75	0.46	0.90
XGBoost (Test)	0.75	0.73	0.41	0.88
LightGBM (Train)	0.87	0.80	0.56	0.92
LightGBM (Test)	0.75	0.73	0.43	0.87

The AUC represents the overall performance of the model in distinguishing between classes, measuring its ability to rank positive instances higher than negative ones across all possible thresholds. Accuracy indicates the proportion of correctly predicted instances among all instances.

Sensitivity (TPR_1) represents the percentage of correct predictions among all actual class 1 instances, which refers to individuals who defaulted on their loans ($TARGET=1$). Specificity (TNR_0) represents the percentage of correct predictions among all actual class 0 instances, which refers to individuals who paid their loans on time ($TARGET=0$).

6.1 Bussiness relevance of each model

Logistic Regression

It was straightforward but performed poorly in predicting defaulters, capturing only 37% of them, which was insufficient for credit prediction.

Random Forest

It was effective at identifying non-defaulters (92% specificity), reducing the risk of lending to them. However, it missed 32% of defaulters, leading to more false negatives. Key predictors, such as EXT_SOURCE_1, EXT_SOURCE_2, and EXT_SOURCE_3, were significant, but its inability to handle complex relationships limited its effectiveness compared to other models.

XGBoost

It offered the best performance with an AUC of 0.75. It excelled at capturing complex relationships, particularly between EXT_SOURCE_3 and EXT_SOURCE_2, which were crucial for predicting defaults. XGBoost effectively balanced sensitivity and specificity, improving the accuracy of credit decisions.

LightGBM

It showed similar results to XGBoost, with key predictors like EXT_SOURCE_2 and EXT_SOURCE_3 playing a significant role in the decision process.

6.2 Ensemble Model for Better Credit Allocation

Given this scenario, we propose adopting a combined model (Ensemble) that integrates predictions from all four models. Averaging the probabilities from each model has shown slightly better performance compared to individual models. The ensemble will balance the strengths of each technique—for example, the interpretability of logistic regression, the ability of Random Forest to capture non-linear patterns, and the accuracy of XGBoost and LightGBM. This approach provides a more robust and reliable solution for credit allocation, crucial for achieving financial inclusion without compromising Home Credit's financial stability.

7. Predicting the Kaggle Test Set

7.1 Application Test Transformations

The code chunks below aim to apply the same feature engineering to the target prediction file (application_test) as was performed on the application_train file. This ensures that the models created can make accurate predictions.

7.2 Predicting Application Test

Below, we will make probability predictions for each model for submission.

Logistic Regression

Random Forest

XGBoost Model

LightGBM Model

8. Submission and Kaggle Score

Average of the models

This submission combines the predictions using the ensemble method, calculating the average of all predictions for submission to Kaggle leveraging the strengths of each individual model to enhance overall performance.

Kaggle Score

The score achieved by averaging the results of all the models was 0.739 only 8% lower than the competition winner. The individual scores for each model were as follows:

Logistic Regression: 0.730

Random Forest: 0.729

XGBoost: 0.733

LightGBM: 0.733

The screenshot shows the Kaggle interface for the 'Home Credit Default Risk' competition. On the left is a sidebar with navigation links: Home, Competitions (selected), Datasets, Models, Code, Discussions, Learn, More, Your Work, and a 'VIEWED' section showing 'Home Credit Default R...'. The main content area has a search bar and a 'Late Submission' button. The competition title 'Home Credit Default Risk' is displayed with the subtitle 'Can you predict how capable each applicant is of repaying a loan?'. Below this is a navigation bar with tabs: Overview, Data, Code, Models, Discussion, Leaderboard (selected), Rules, Team, and Submissions. The 'Leaderboard' section shows a 'YOUR RECENT SUBMISSION' card for 'submission_average.csv' submitted by 'Kleyton Polzonoff Silveira' 8 minutes ago. The submission score is 0.73902, and the private score is 0.73267. There are buttons for 'Raw Data' and 'Refresh'.

Conclusion

Predicting defaults is inherently challenging, as all models rely heavily on external credit scoring data. However, integrating additional variables such as `AMT_CREDIT`, `AMT_GOODS_PRICE`, `DAYS_EMPLOYED_YEARS`, `DAYS_BIRTH`, `DAYS_ID_PUBLISH`, `OWN_CAR_AGE`, `DAYS_LAST_PHONE_CHANGE`, and `NAME_EDUCATION_TYPE_higher_education` can significantly enhance model performance. These insights enable the company to focus on the most relevant factors, reducing the need for unnecessary customer inputs.

By adopting an ensemble model that combines predictions from all four approaches, Home Credit can leverage the strengths of each technique to create a more balanced and accurate system. This approach ensures better credit allocation decisions, reduces the risk of defaults, and fosters financial inclusion while safeguarding the company's financial stability.