

# Credit default prediction



# Agenda

- 01 Context and objective

---
- 02 Methodology overview

---
- 03 Exploratory Data Analysis

---
- 04 Feature engineering

---
- 05 Modeling

---
- 06 Tuning

---
- 07 To go further

# 01

## Context and objective



# Context and objective

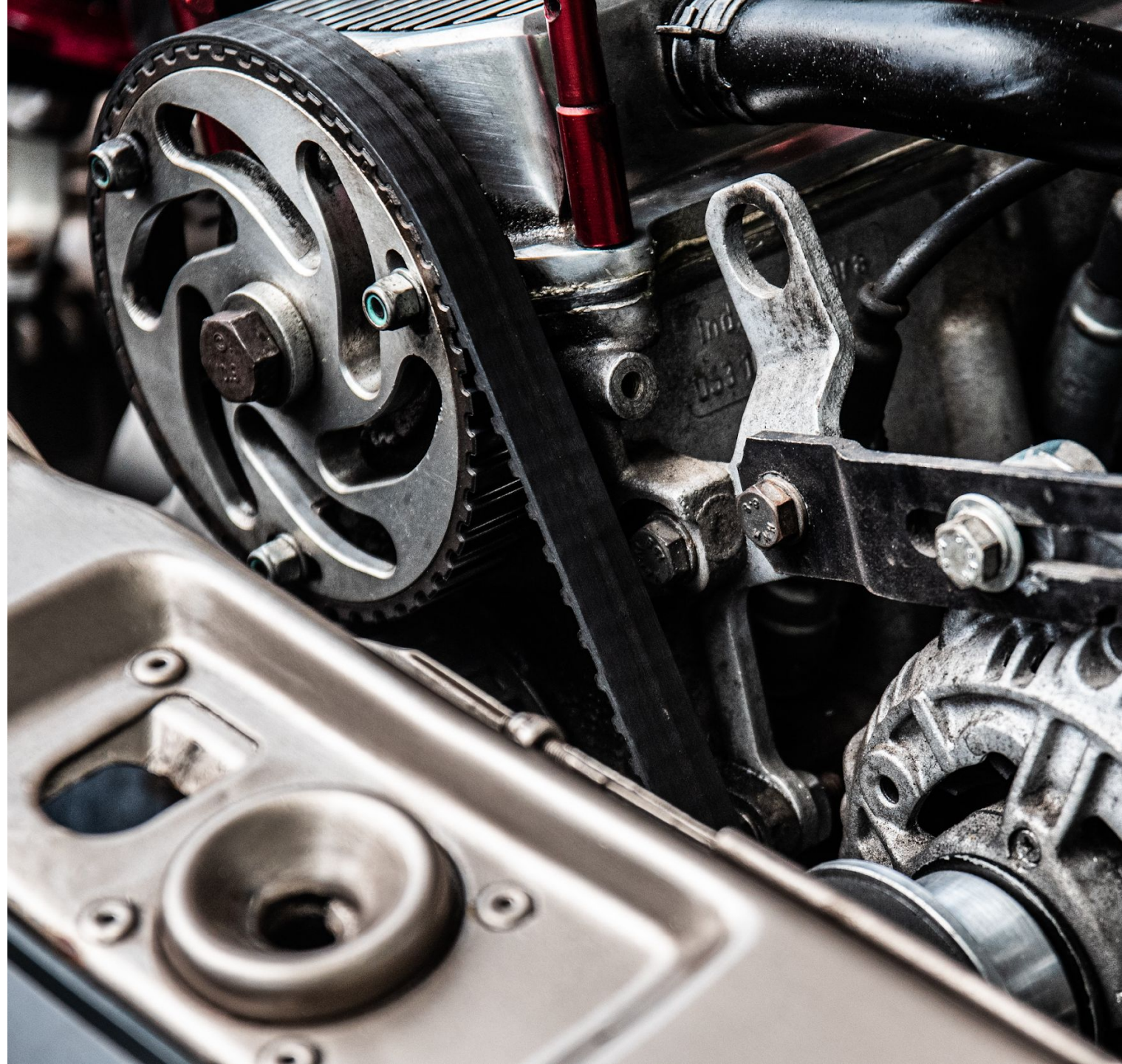
---

- This is a machine learning project that aims to predict whether a customer will default on a loan within 60 days of disbursement.
- The project uses historical customer financial data to train and evaluate several machine learning models and select the best-performing one.
- The target we aim to predict has two values:
  - 0 - Non default
  - 1- Default
- So it is a binary classification problem

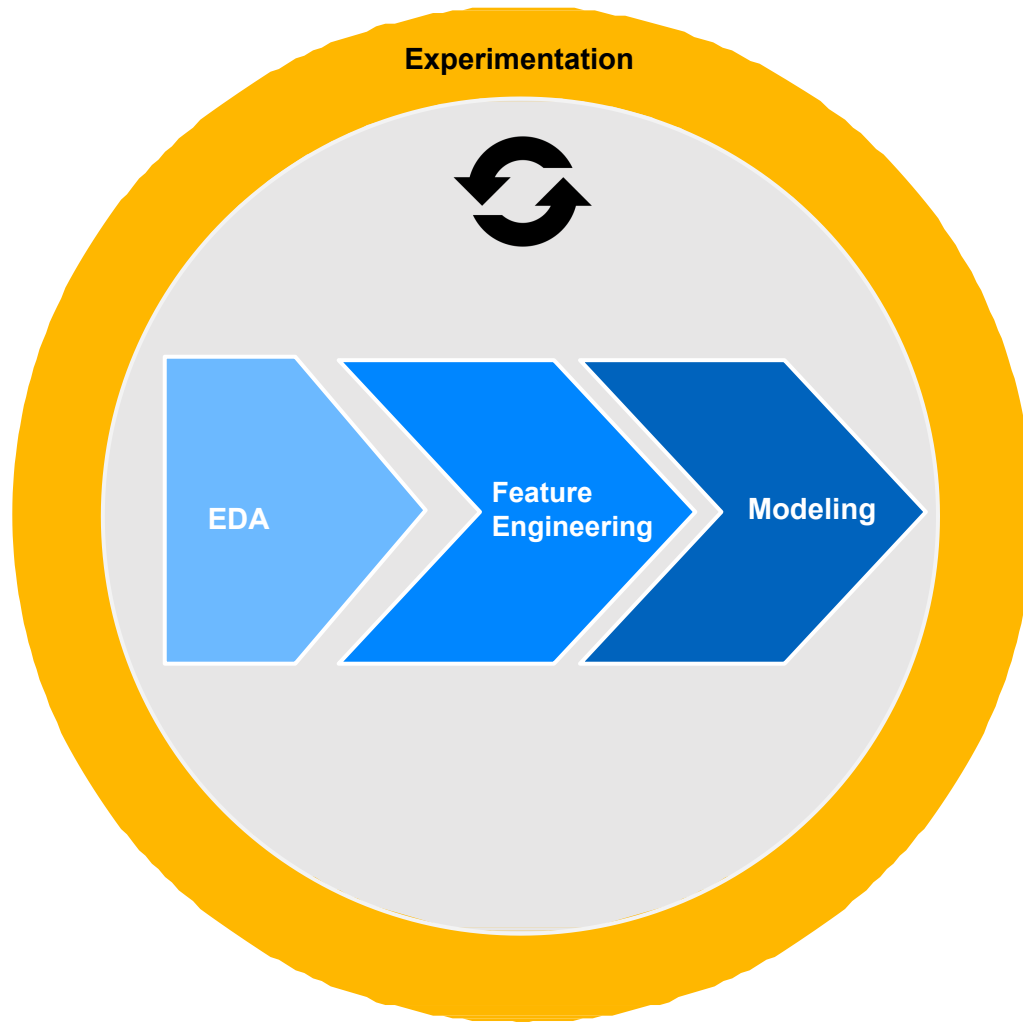


# 02

## Methodology overview



# Methodology overview



## EDA

- Checking for irregularities (null values, outliers)
- Understanding the distribution of the data
- Looking for relationship between variables



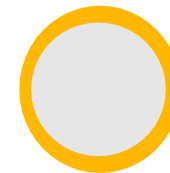
## Feature engineering:

- Creating new features
- Extracting the most important



## Modeling:

- Testing several algorithms
- Tuning the best performing algorithm



## Experimentation loops:

- Experiment different approaches and learn

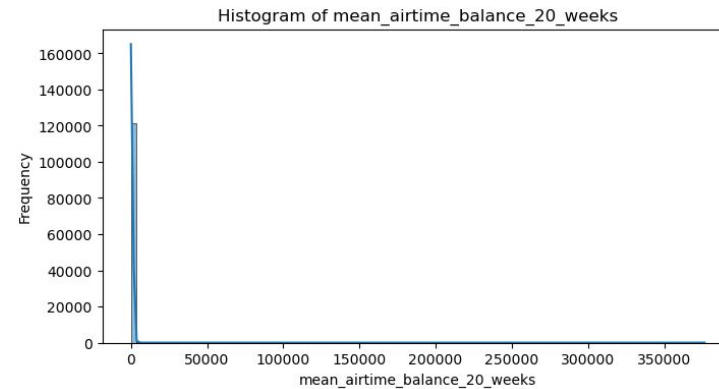
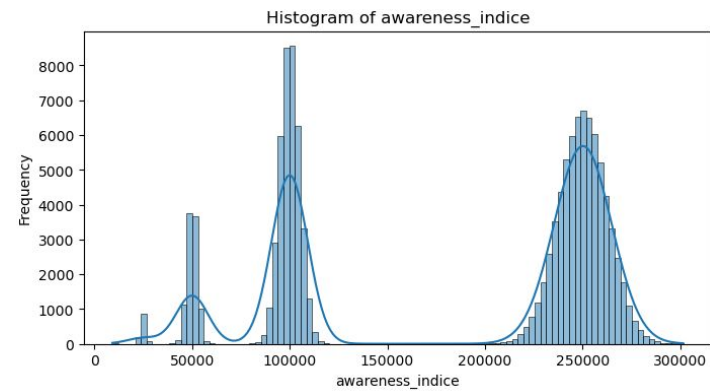
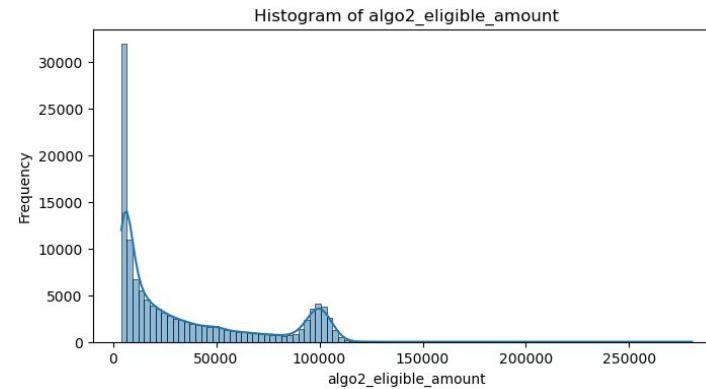
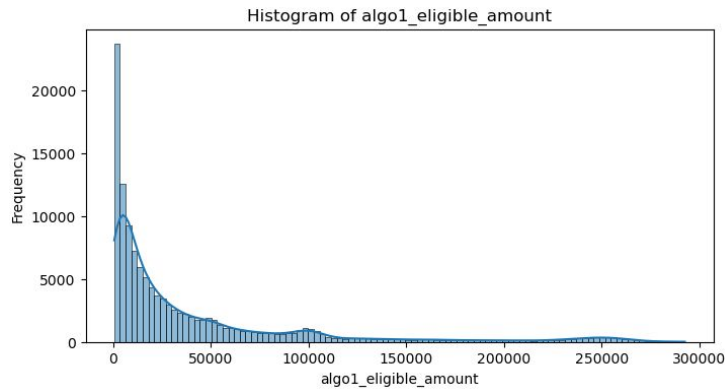


# 03

## Exploratory Data Analysis



# EDA - distributions

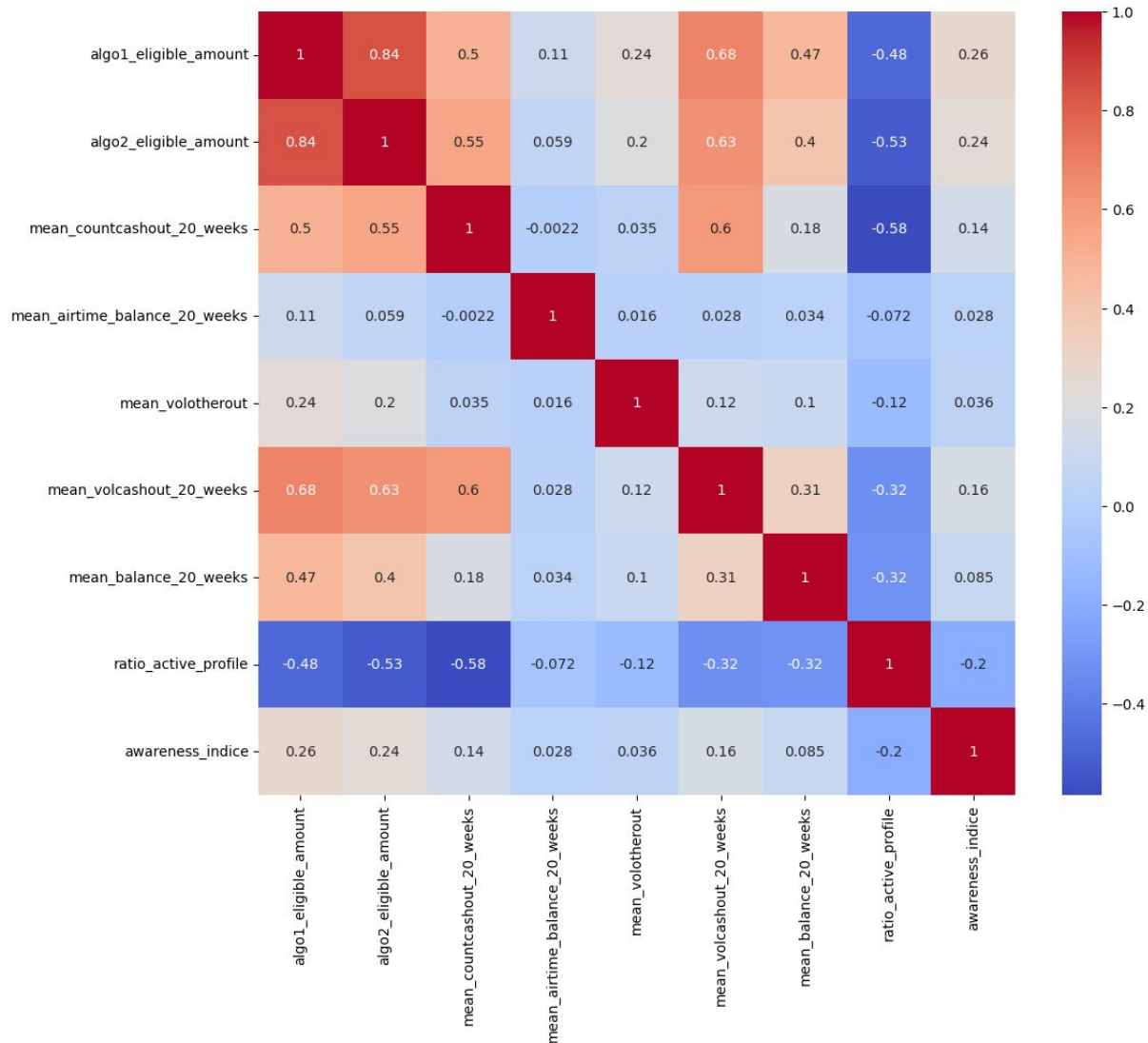


Several features are right skewed => there may be a significant number of outliers towards the higher end of the range.

Regarding the awareness\_indice feature, we can observe that it has multiple peaks centered around 50,000, 100,000 and 250,000. This suggests that there may be different groups of customers with different levels of awareness of the offer.



# EDA - correlations between features



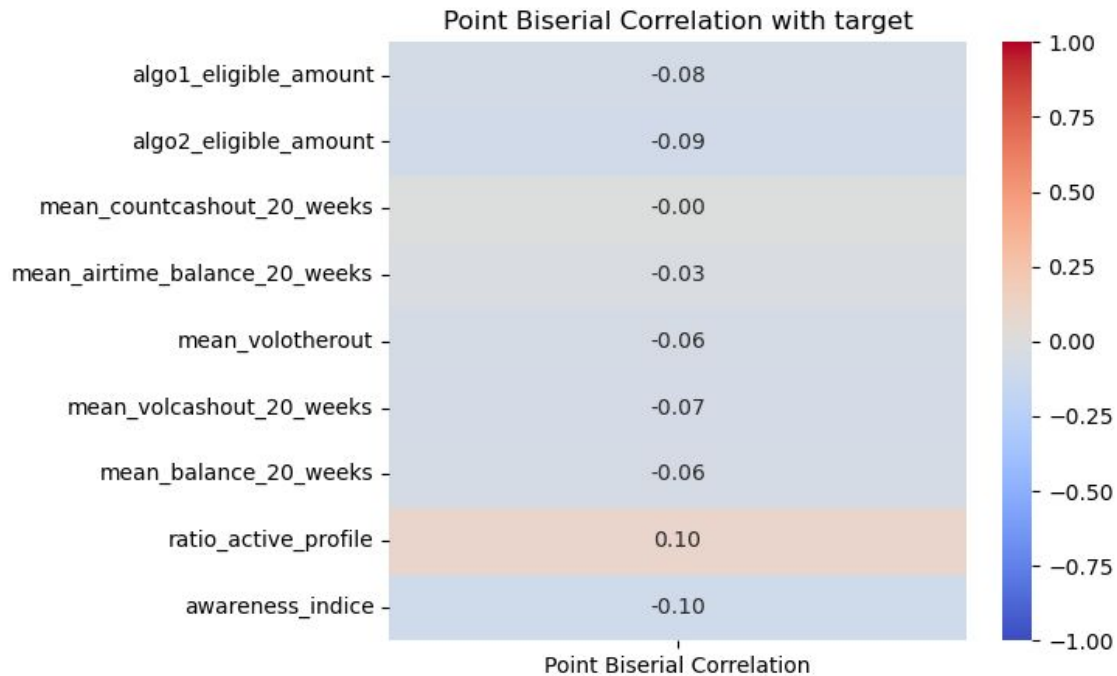
➤ algo1\_eligible\_amount and algo2\_eligible\_amount have a strong positive correlation of 0.84. This is indicating that the two algorithms generally agree on the eligible amount for a customer.

➤ mean\_countcashout\_20\_weeks, mean\_voicashout\_20\_weeks, and mean\_balance\_20\_weeks are positively correlated with both algo1\_eligible\_amount and algo2\_eligible\_amount. This suggests that customers with higher average cashout transactions, cashout volumes, and balances in the last 20 weeks are likely to be eligible for larger loan amounts.

➤ ratio\_active\_profile has a negative correlation with both algo1\_eligible\_amount (-0.47) and algo2\_eligible\_amount (-0.53). This suggests that customers with more volatile transaction patterns may be considered less eligible for larger loan amounts.

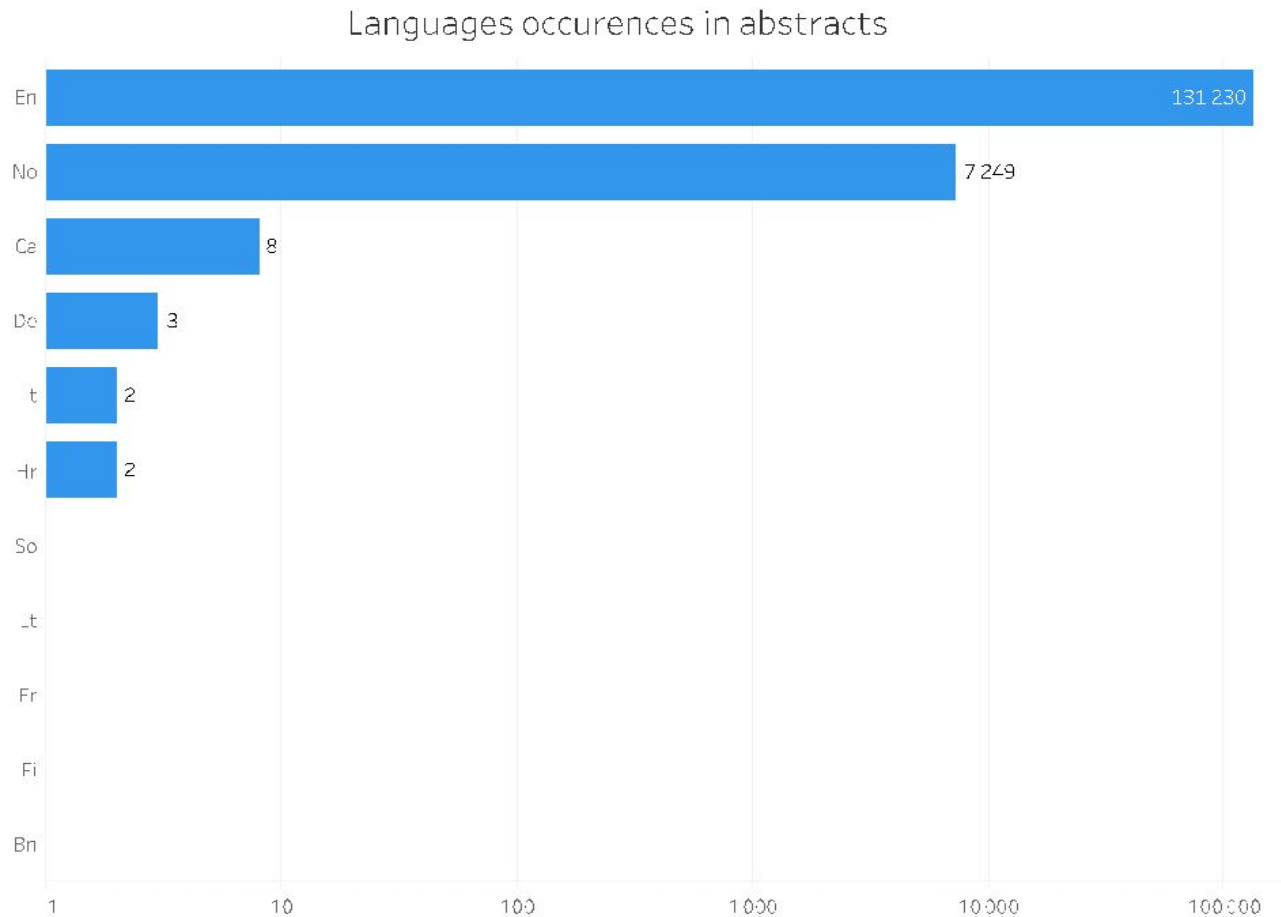
➤ awareness\_indice has a weak positive correlation with both algo1\_eligible\_amount (0.26) and algo2\_eligible\_amount (0.24). This indicates that customers with a better understanding of the offer might be slightly more eligible for larger loan amounts, but this relationship is not very strong.

# EDA - correlation between features and target



- The majority of the features have a weak correlation with the target variable (default\_60days).
- The low correlation values suggest that the individual features alone may not be sufficient for accurate prediction, and it is likely that a combination of features will be necessary.
- This can also indicate that the relationship between the features and the target is non-linear, thus sophisticated models that can capture such relationships may be needed for better performance.

# Data exploration and processing



➤ 11 detected languages in the abstracts. English is the most represented



➤ Language-specific processing technics



# 04

## Features engineering



# Feature engineering

## Dates



- request\_year
- request\_month
- request\_day
- request\_hour
- request\_dayofweek
- reimbursement\_year
- reimbursement\_month
- reimbursement\_day
- reimbursement\_dayofweek
- date\_diff

## Amounts



- mean\_cashout\_to\_balance\_ratio
- mean\_cashout\_to\_airtime\_ratio
- mean\_volootherout\_to\_balance\_ratio
- algo\_diff
- total\_eligible\_amount
- min\_algo\_eligible\_amount
- max\_algo\_eligible\_amount
- mean\_algo\_eligible\_amount

## Ids



- customer\_id X
- simulation\_id X
- loan\_id X

05

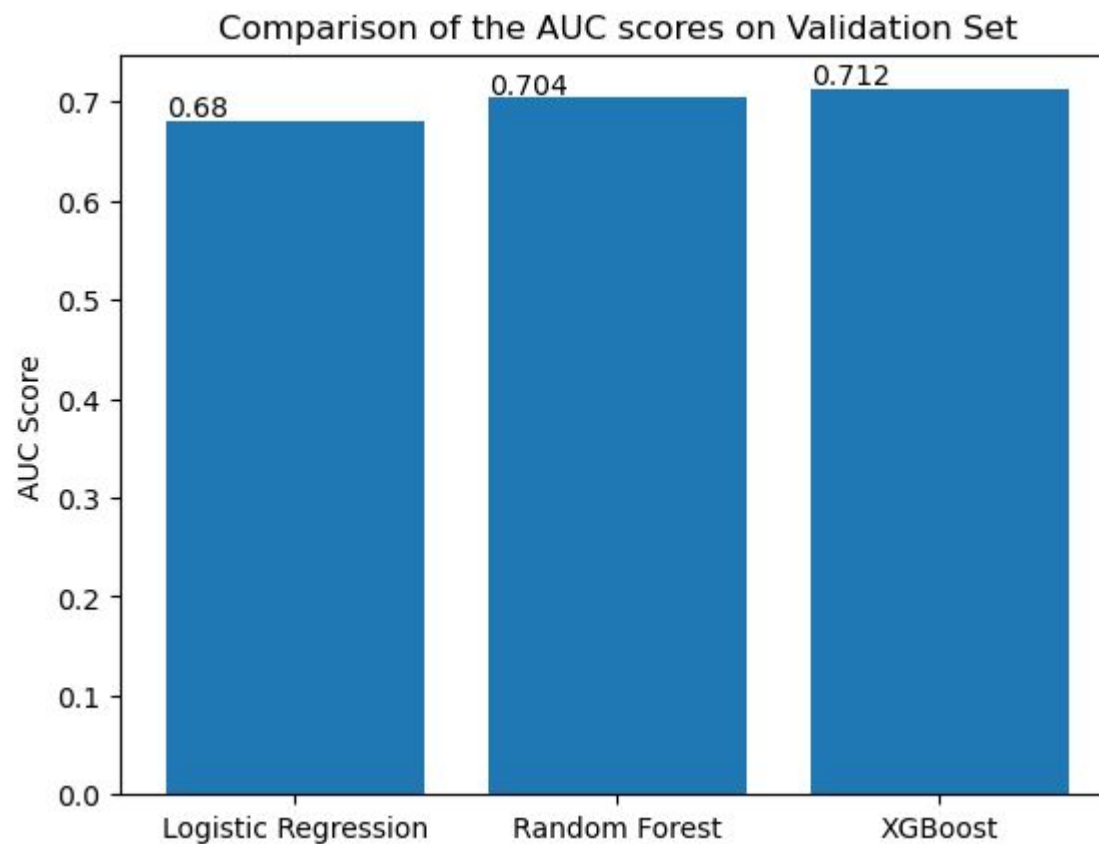
Modeling





# Trying different models

---



06

Tuning



# Tuning XGBoost

## Validation set

Best parameters: {'learning\_rate': 0.1, 'max\_depth': 6, 'min\_child\_weight': 3, 'n\_estimators': 100}

**Classification Report:**

	precision	recall	f1-score	support
0.0	0.71	0.81	0.76	7474
1.0	0.60	0.47	0.53	4669
accuracy			0.68	12143
macro avg	0.66	0.64	0.64	12143
weighted avg	0.67	0.68	0.67	12143

**AUC Score:**

0.7201627883638364

## Test set

**Classification Report:**

	precision	recall	f1-score	support
0.0	0.70	0.81	0.75	7357
1.0	0.61	0.48	0.54	4787
accuracy			0.68	12144
macro avg	0.66	0.64	0.64	12144
weighted avg	0.67	0.68	0.67	12144

**AUC Score:**

0.7214909018435736



# Handling class imbalance with SMOTE

**Classification Report:**

	precision	recall	f1-score	support
0.0	0.73	0.73	0.73	7357
1.0	0.58	0.58	0.58	4787
accuracy			0.67	12144
macro avg	0.65	0.65	0.65	12144
weighted avg	0.67	0.67	0.67	12144

**AUC Score:**  
0.7168164543550068

07

To go further





# To go further

---

- Retrain the XGBoost model using only the top 10 most important features identified from feature importance analysis. This may help to simplify the model and reduce overfitting.
- Explore the use of SMOTE (Synthetic Minority Over-sampling Technique) as a hyperparameter in the XGBoost tuning. This technique may help to balance the class distribution and improve the model's ability to predict defaults.
- Gather additional information about the loan amount and use it as a feature in the model. This could potentially improve the model's performance by capturing the relationship between loan amount and default risk.