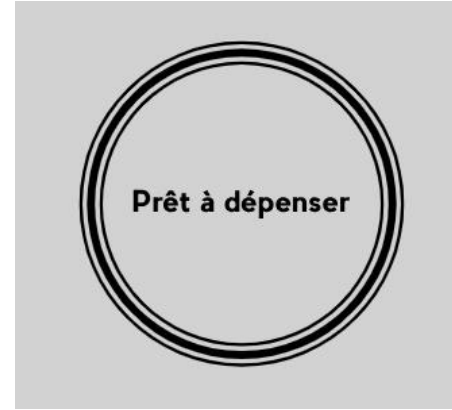


# CREDIT SCORING

# Project



Create a credit scoring  
calculator based on machine  
learning techniques

# Challenges

- Clients having poor credit history
- Highly imbalanced data
- Transparency of the credit score

# MISSION



Implement a credit scoring model



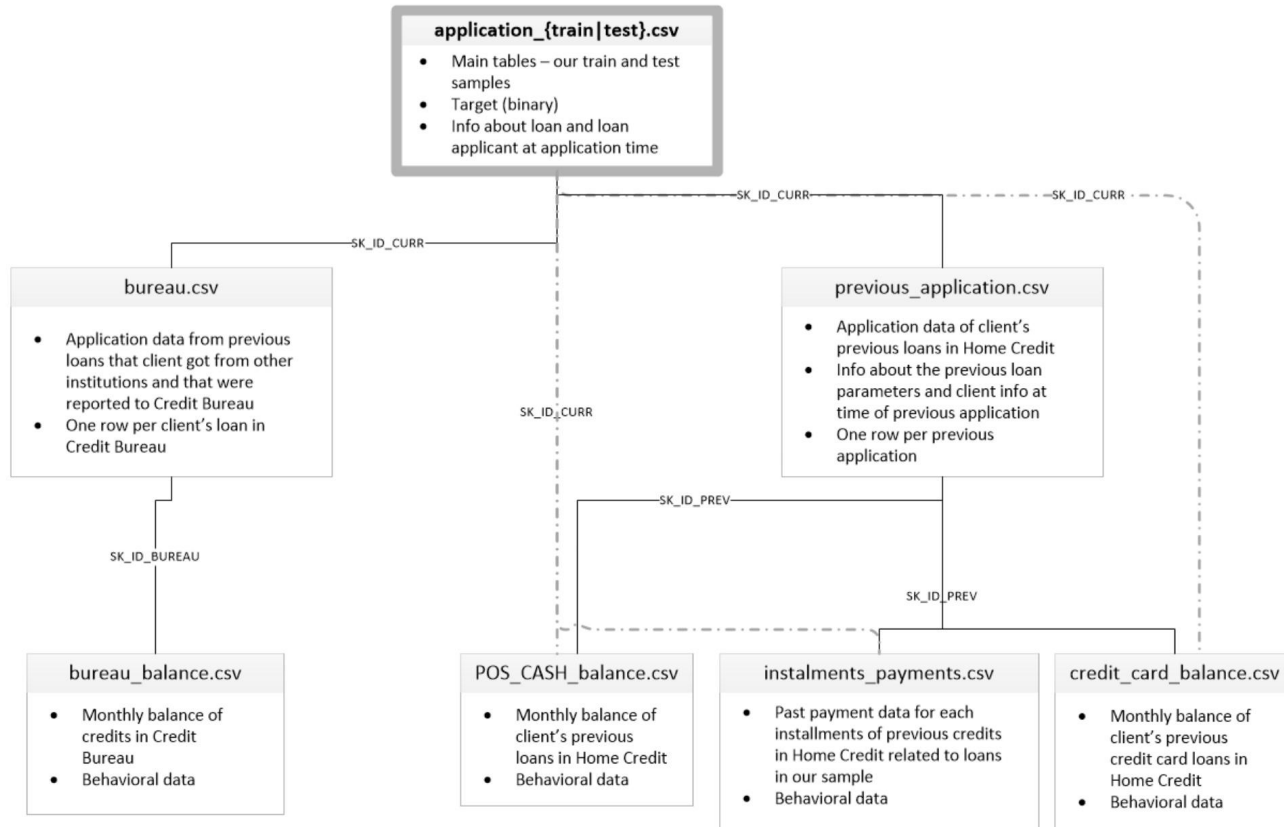
Create a predict API



Create an interactive dashboard

# DATA

# Database schema



# Quick overview

- Data Origin: Home Credit Kaggle Competition
- 7 tables
- 307511 applications in the train set
- 8% applications with payment difficulties
- 92% of no risk applications

# METHODOLOGY



# 2 TYPES OF ERRORS

TYPE I

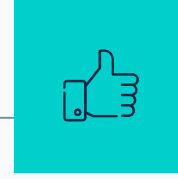


## FALSE POSITIVES

- customers predicted as potential defaulters but did not default
- affect profitability

1:10

TYPE II



## FALSE NEGATIVES

- customers predicted as non risky and who defaulted
- important money loss

# EVALUATION METRICS

1

COST RATIO SCORE

Ratio of  $10 \cdot \text{FN} + \text{FP}$

2

$F\beta$ -SCORE ( $\beta = 10$ )

Give more weight to Recall

3

ROC AUC

Trade-off between True Positive Rate and the False Positive Rate

4

RECALL

True Positive Rate

5

ACCURACY

Proportion of observations that were correctly predicted

# MODELS

## DUMMY: baseline

- Dummy Classifier

## LINEAR

- Logistic Regression

## ENSEMBLE

- Random Forest Classifier

## GRADIENT BOOSTING

- LightGBM
- XGBoost

# TRAINING PROCESS FOR EACH MODEL

Choose the best  
method for the class  
balancing

CLASS IMBALANCE



STEP 1

STEP 2



FEATURE  
ENGINEERING

Feature  
transformations +  
optimization

Selecting the right  
preprocessing  
techniques

PREPROCESSING



STEP 3

STEP 4



FINE-TUNING

Choosing optimal  
hyperparameters

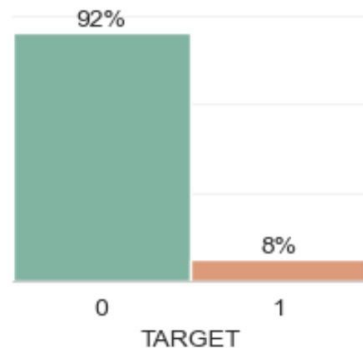
# CLASS IMBALANCE

## ■ Oversampling

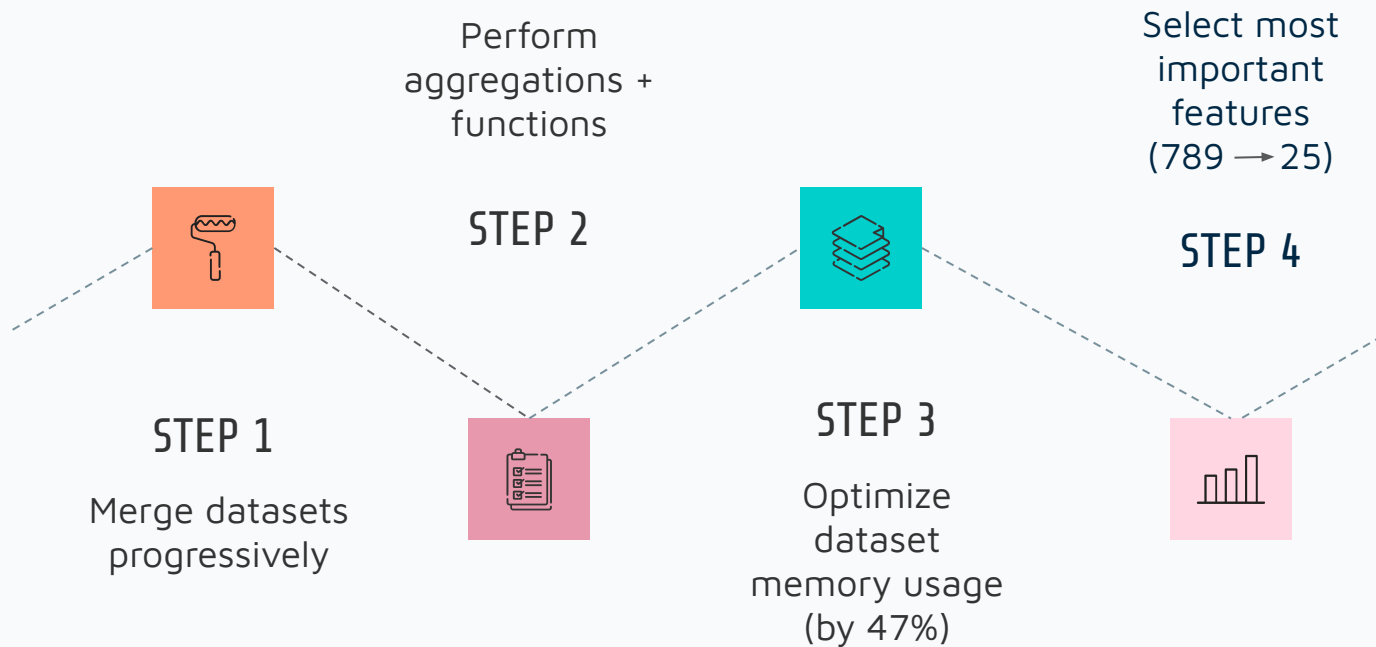
- SMOTE
- SMOTE-NC

## ■ Class weights

- `'class_weights': 'balanced'` (Logistic Regression, Random Forest Classifier, LightGBM)
- `'scale_pos_weight': 11` (XGBoost)



# FEATURE ENGINEERING



# FEATURE OPTIMIZATION

## ■ Drop features:

- with low variance (0.1 variance threshold)
- correlated ( $\geq 0.7$ )

## ■ Select 25 final features:

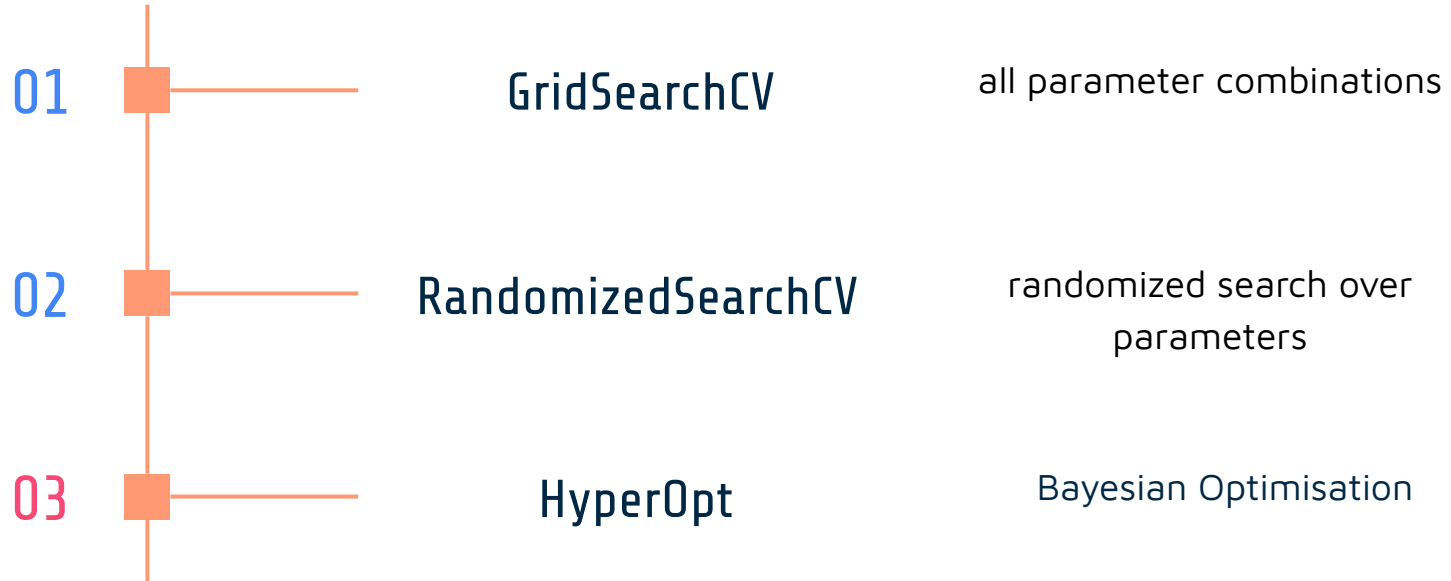
- Logistic Regression *coef\_*
- LightGBM *feature\_importances\_*
- XGBoost *feature\_importances\_*

# PREPROCESSING

OUTLIERS	MISSING VALUES	ENCODING	SCALING	NORMALIZATION
IQR factor 1.5 IQR factor 2.5	Drop columns with > 60% of NaN  Simple Imputer: <ul style="list-style-type: none"><li>■ median</li><li>■ most frequent</li></ul>	One Hot Encoding (categorical)	StandardScaler MinMaxScaler RobustScaler	Power Transformer



# FINE TUNING



# TRAINING RESULTS

	Dummy	LightGBM	XGBoost	Random Forest	Logistic Regression
Cost Ratio	0.79	<b>0.56</b>	0.59	0.6	0.61
F $\beta$ -score	0	0.64	<b>0.65</b>	0.61	0.61
AUC	0.5	<b>0.73</b>	0.71	0.7	0.69
Recall	0	0.66	<b>0.68</b>	0.64	0.63
Accuracy	0.92	<b>0.68</b>	0.64	0.66	0.65
Duration	7.8s	1.9min	4.5min	2.4min	<b>1.7min</b>

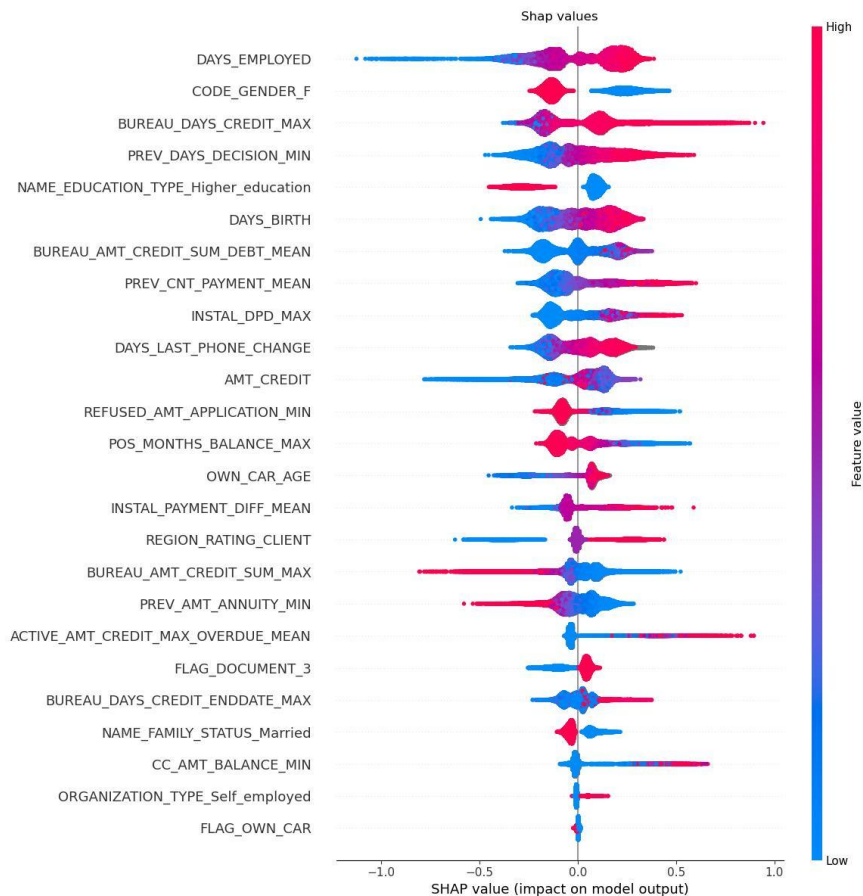
# LIGHTGBM CONFIGURATION

Class balance setting: class\_weights: 'balanced'

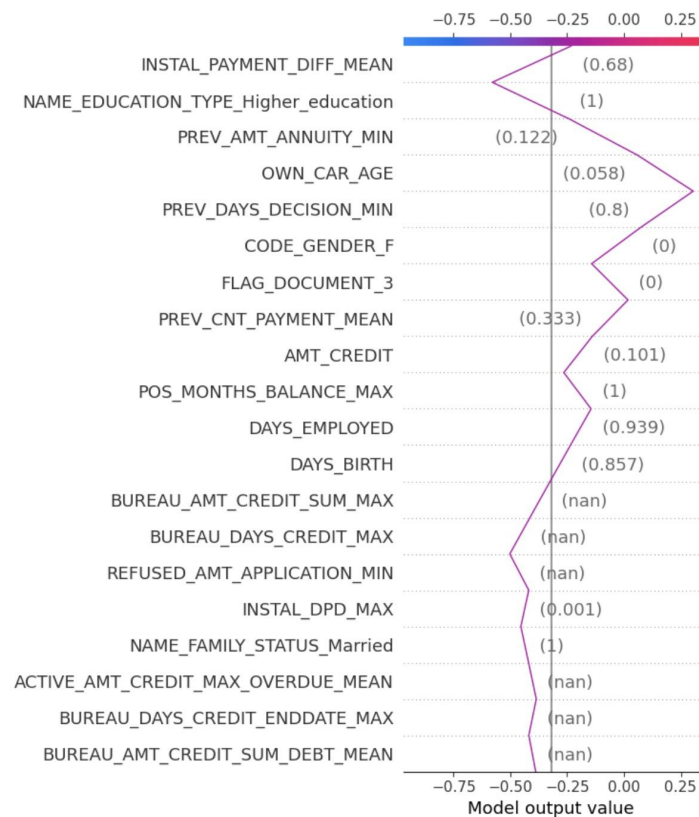
Best threshold: 0.5

PREPROCESSING	HYPERPARAMETERS
<ul style="list-style-type: none"><li>■ Filling missing values: ✗</li><li>■ Removing outliers: ✗</li><li>■ OneHotEncoding: ✓</li><li>■ Scaling (MinMaxScaler): ✓</li><li>■ Normalization: ✗</li></ul>	<ul style="list-style-type: none"><li>■ n_estimators: 158</li><li>■ num_leaves: 111</li><li>■ max_depth: 4</li><li>■ learning_rate: 0.09</li><li>■ min_data_in_leaf: 135</li><li>■ subsample_for_bin: 20000</li><li>■ colsample_by_tree: 0.59</li><li>■ lambda_l1: 0.71</li><li>■ lambda_l2: 0.19</li></ul>

## FEATURE IMPORTANCE: GLOBAL



## FEATURE IMPORTANCE: LOCAL

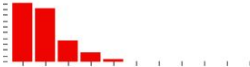
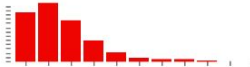




# DATA DRIFT ANALYSIS

Tool: Evidently

Detected drift: 8.0% of columns (2 out of 25)

Detection threshold: 0.5

Column	Type	Reference Distribution	Current Distribution	Data Drift	Stat Test	Drift Score
AMT_CREDIT	num			Detected	Wasserstein distance (normed)	0.207334
DAYS_LAST_PHONE_CHANGE	num			Detected	Wasserstein distance (normed)	0.158644

# CONTINUOUS INTEGRATION

## MODEL TRACKING

- MLFlow

## VERSION CONTROL

- Git
- Github ([repository link](#))

## JOBS

- Tests with PyTest
- Github actions

## API HOSTING

- Render platform

## DASHBOARD

- Streamlit ([dashboard link](#))

## COMMITTS

### feat(ops): generate a data drift report

 NatChe committed last week · ✓ 1 / 1

---

### refacto(dashboard): move dashboard.py to the root

 NatChe committed last week · ✓ 1 / 1

---

### fix(dashboard): run shap initjs

 NatChe committed last week · ✓ 1 / 1

---

### fix(dashboard): use abs path for feature importance image

 NatChe committed last week · ✓ 1 / 1

---

### fix(dashboard): use absolute path for css

 NatChe committed last week · ✓ 1 / 1

## GITHUB ACTIONS

### ✓ feat(ops): generate a data drift report

PyTest #21: Commit [10e73c1](#) pushed by NatChe

---

### ✓ refacto(dashboard): move dashboard.py to the root

PyTest #20: Commit [a64cedd](#) pushed by NatChe

---

### ✓ fix(dashboard): run shap initjs

PyTest #19: Commit [116775f](#) pushed by NatChe

---

### ✓ fix(dashboard): use abs path for feature importance image

PyTest #18: Commit [e8686ab](#) pushed by NatChe

---

### ✓ fix(dashboard): use absolute path for css

PyTest #17: Commit [209fc45](#) pushed by NatChe

## TO GO FURTHER

- Study credit risk literature and research papers (feature engineering, score).
- Fine-tune Cost Ratio between Error Type I and Error Type II.
- Exclude features that can lead to discrimination.
- More advanced optimization of the LightGBM model.



# THANKS!

Powered by



## Links:

1. Github repository: <https://github.com/NatChe/credit-scoring-dashboard>
2. Dashbord on Streamlit:  
<https://credit-scoring-dashboard-vo4t4kjrntvwp3fabtzsss.streamlit.app/>
3. Render platform: <https://render.com/>
4. Kaggle notebook used for Feature Engineering:  
<https://www.kaggle.com/code/jsaguiar/lightgbm-with-simple-features>
5. Main page of the Home Credit Kaggle Competition for the data explanation:  
<https://www.kaggle.com/competitions/home-credit-default-risk/data>