

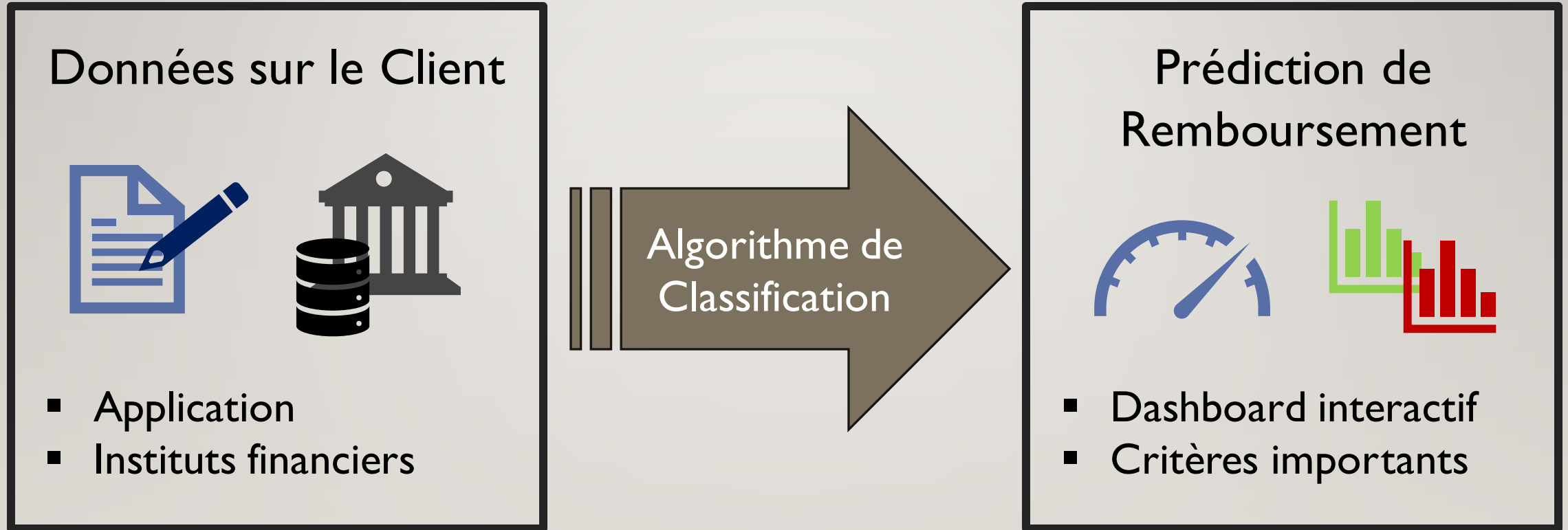
PROJET 7

UN MODÈLE DE SCORING **HOME CREDIT**

EVA BOOKJANS



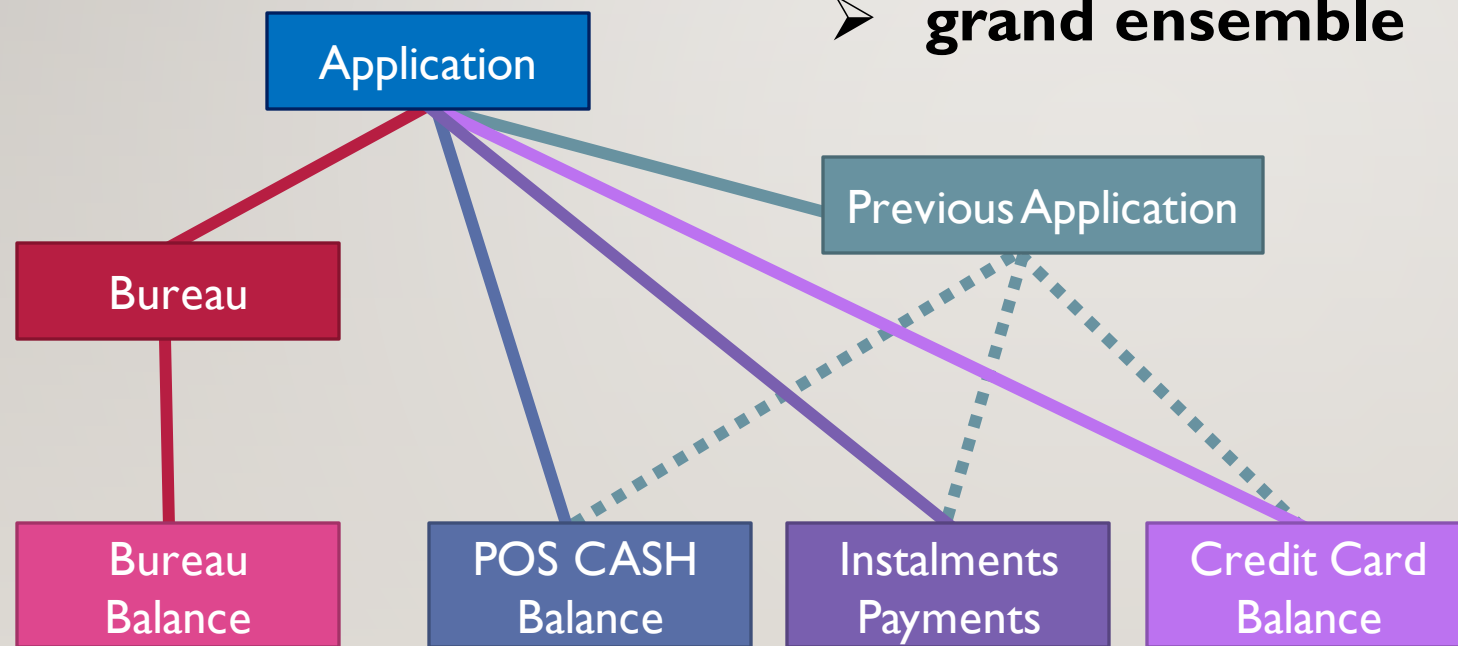
OBJECTIVE – un Outil de Scoring Crédit



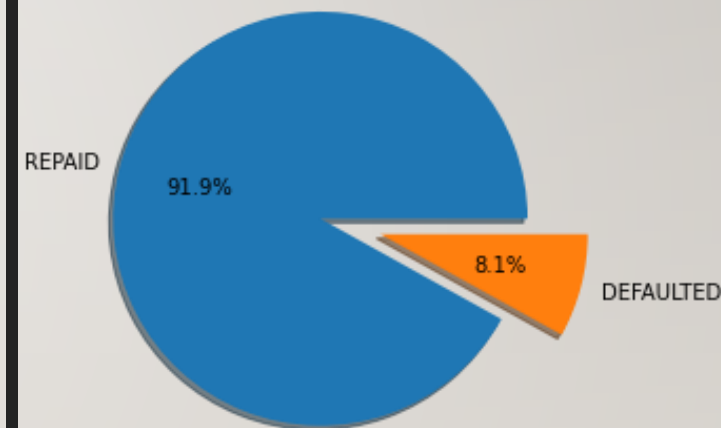
LE JEU DES DONNÉES*

Structure de Données

- 7 fichiers
- 307507 individus
- **grand ensemble**



La Classe Cible



- **fortement déséquilibré**

PRÉTRAITEMENT*

➡ 773 Variables



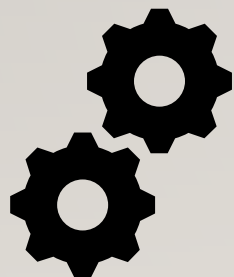
Nettoyage

- remplacement des **valeurs aberrantes** par une valeur manquante (NaN)
- regroupement cohérente des valeurs catégorielles pour des valeurs peu fréquentes

1100
1010
0101

Numérisation des Données

- **transformation des valeurs catégorielles en one-hot variables** (encodage binaire) en **variables ordinales** (e.g. niveau d'éducation, niveau de rendement,...)



Feature Engineering

- création des **nouvelles variables métier** potentiellement plus pertinentes
- **agrégations statistiques** des données venant des fichiers supplémentaires avec des différents fonctions (e.g. minimum, maximum, moyenne, déviations, ...)

RÉDUCTION DE DIMENSION

→ 377 Variables



Peu Peuplées - plus de 70% des valeurs manquantes



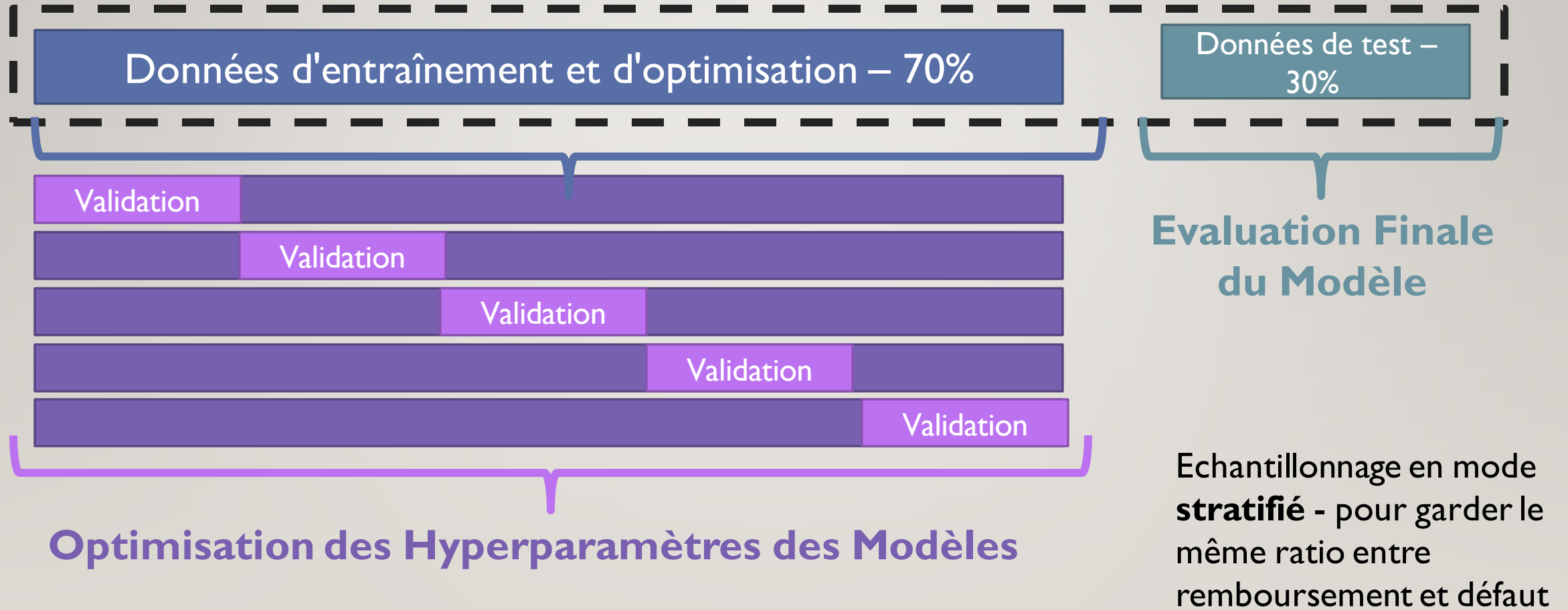
Fortement Corrélés - plus de 90% de corrélation



Presque Constantes – moins de 10% de déviation
(si encodé entre 0 et 1)






MÉTHODE D'ÉVALUATION



LA CLASSE CIBLE DÉSÉQUILIBRÉE

STRATÉGIES:

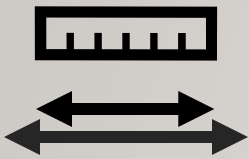
- **sous-échantillonnage** aléatoire de la classe majoritaire 
- **SMOTE** (Synthetic Minority Oversampling Technique)  = suréchantillonnage de la classe minoritaire 
- **pondération** des individus selon la classe cible (`class_weight = 'balanced'`)
- **un métrique d'évaluation adapté** (e.g. ROC-AUC, précision moyenne, ...)
 - éviter les estimations de performances gonflées sur des jeux de données déséquilibrés

MÉTRIQUE D'ÉVALUATION - ROC-AUC



- Sélection des modèles optimaux
 - indépendamment de la distribution de classe
 - avant de spécifier le contexte de coût (mis en place d'un seuil de décision)

LES MODÈLES TESTÉS



Mise à l'Echelle

- `MinMaxScaler()`
- Prétraitement customisé - prend en compte l'asymétrie de la distribution des valeurs (skew) pour éventuellement faire une mise à l'échelle logarithmique



Imputation des valeurs manquants

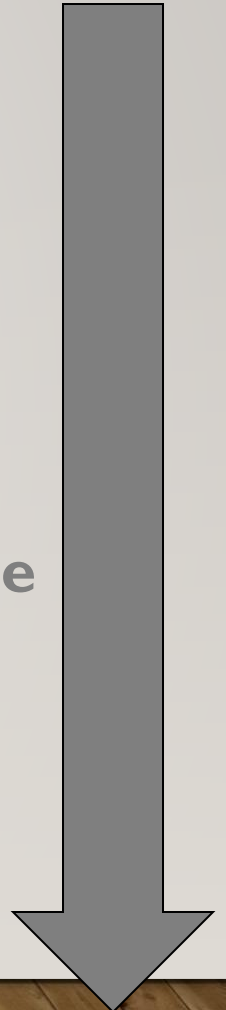
- médiane, moyenne



Modèle de Classification

- **LogisticRegression** (sklearn)
- **BalancedRandomForrest** (imblearn)
- **RUSBoostClassifier** (imblearn)
- **LightGBM** – light gradient boosting machine (lightgbm)

Pipeline



RÉSULTATS

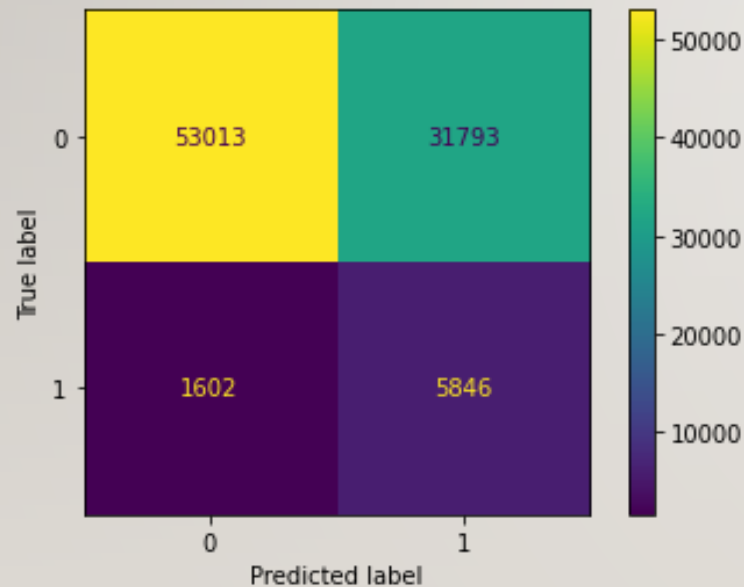
Model	Mise à l'échelle	Imputation	Sous-échant.	AUC-ROC valid.	AUC-ROC test	Précision moyenne test	F1-score test	Temps (s)
LightGBM	--	--	non	0.7807	0.7780	0.1525	0.2935	20.5
LogisticRegr.	Cust.	médiane	non	0.7704	0.7694	0.1442	0.2754	45.2
RUSBoost	--	médiane	non	0.7701	0.7694	0.1452	0.2780	114.6
RandomForest	--	médiane	oui	0.7538	0.7516	0.1369	0.2627	31.7

 **LightGBM**

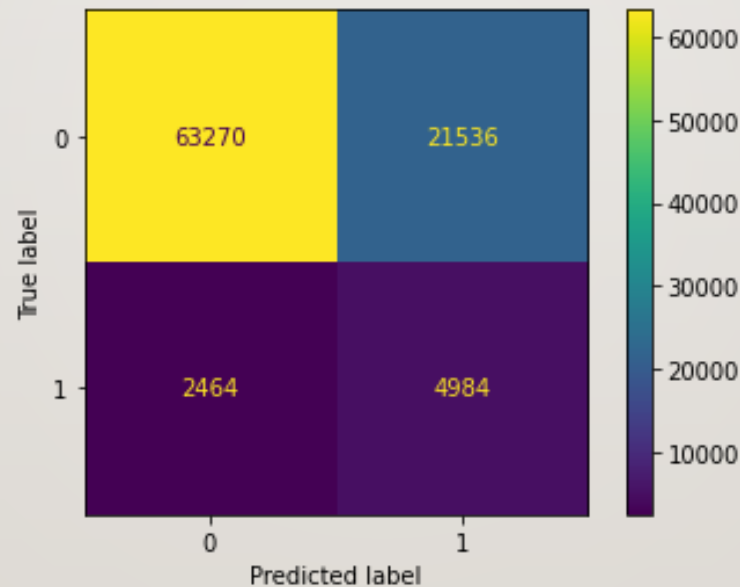
LE MODÈLE LIGHT-GBM

Matrix de Confusion avec différents seuils de décision

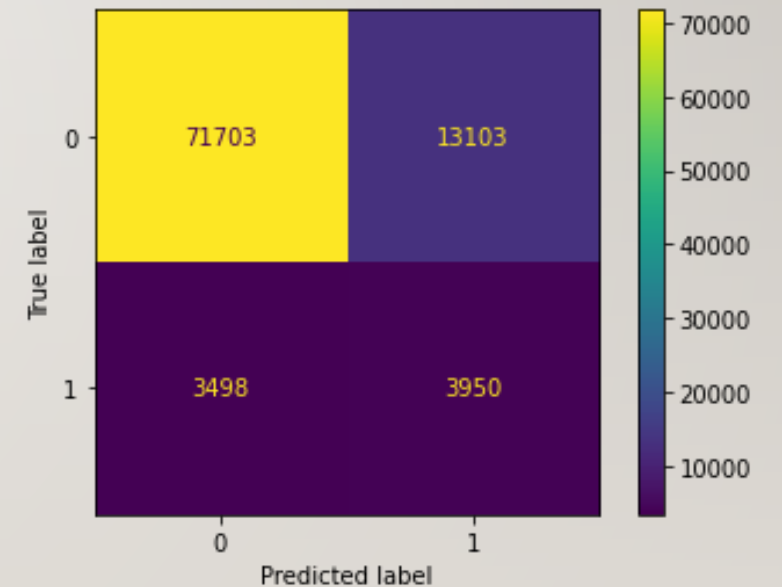
$S = 0.4$



$S = 0.5$

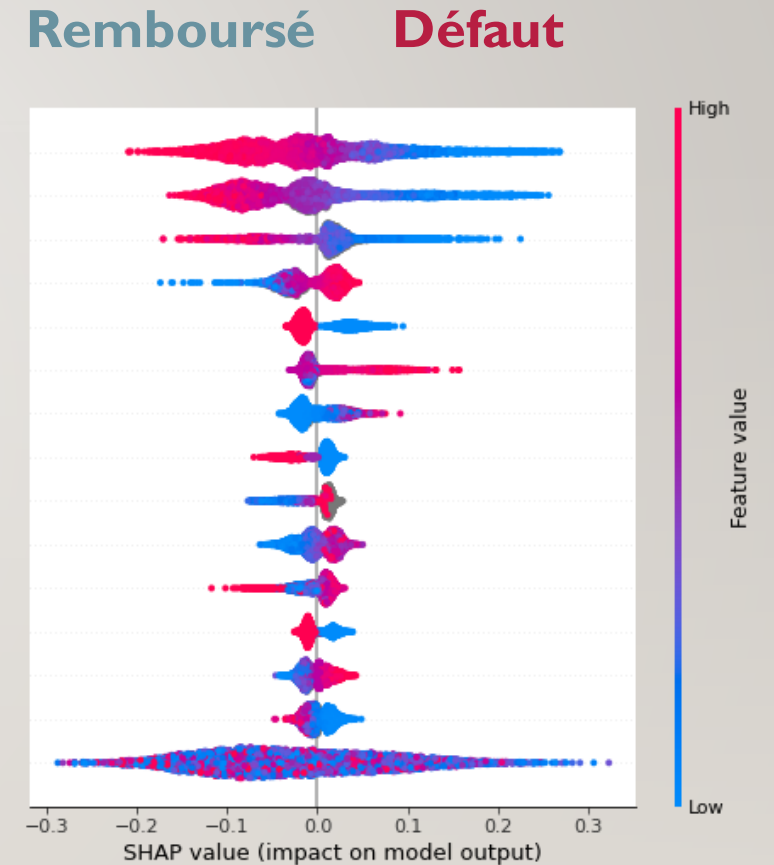
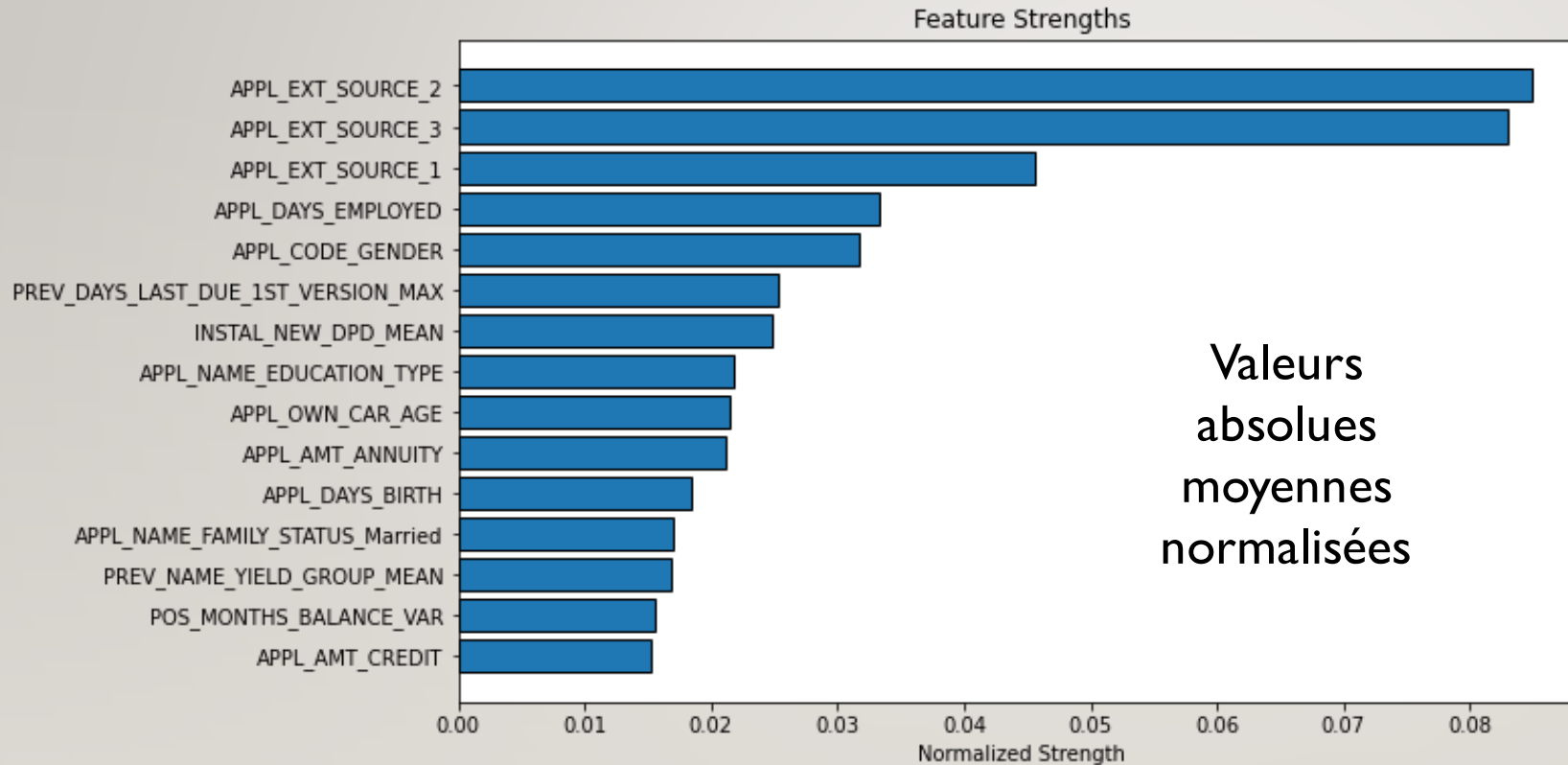


$S = 0.6$



INTERPRÉTATION DU MODÈLE

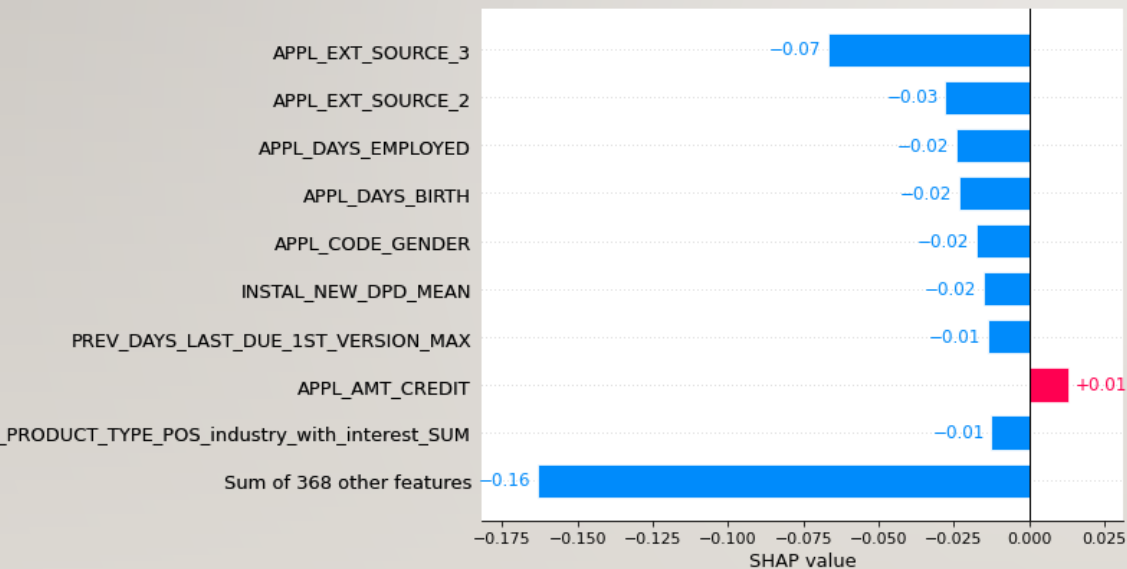
Valeurs de Shapley (SHAP module)



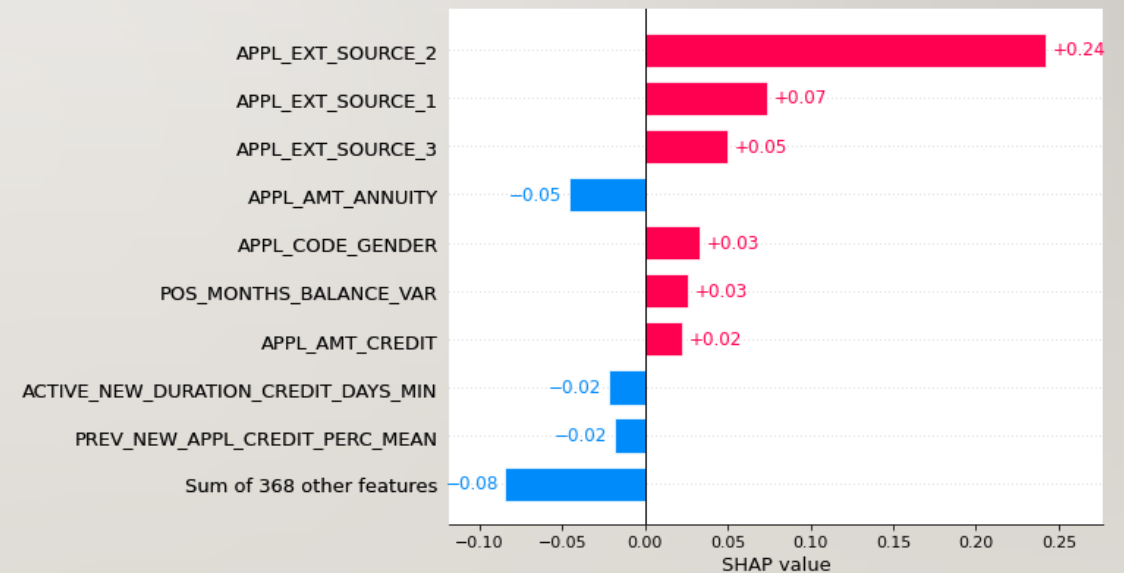
INTERPRÉTATION DU MODÈLE

Valeurs de Shapley (SHAP module)

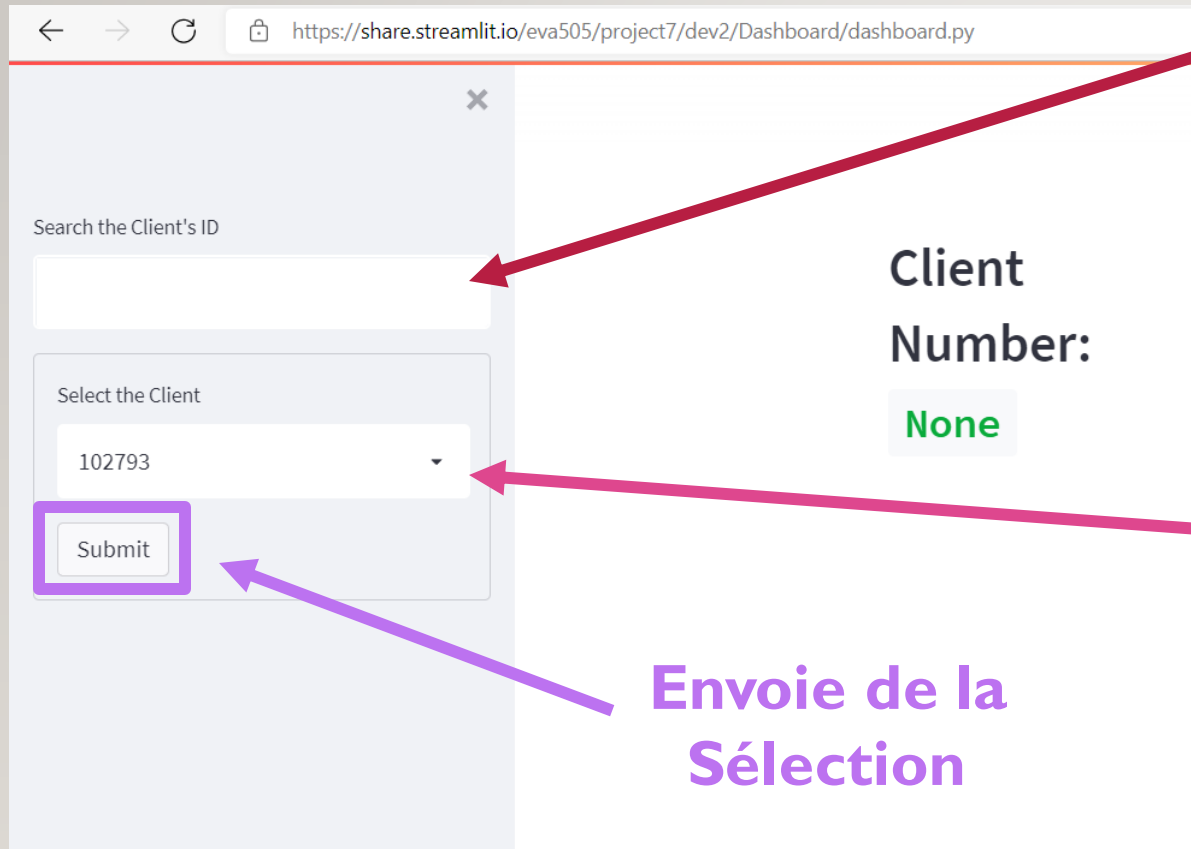
Un individu qui a
remboursé son crédit



Un individu en
défaut de crédit



LE DASHBOARD – SÉLECTION CLIENT



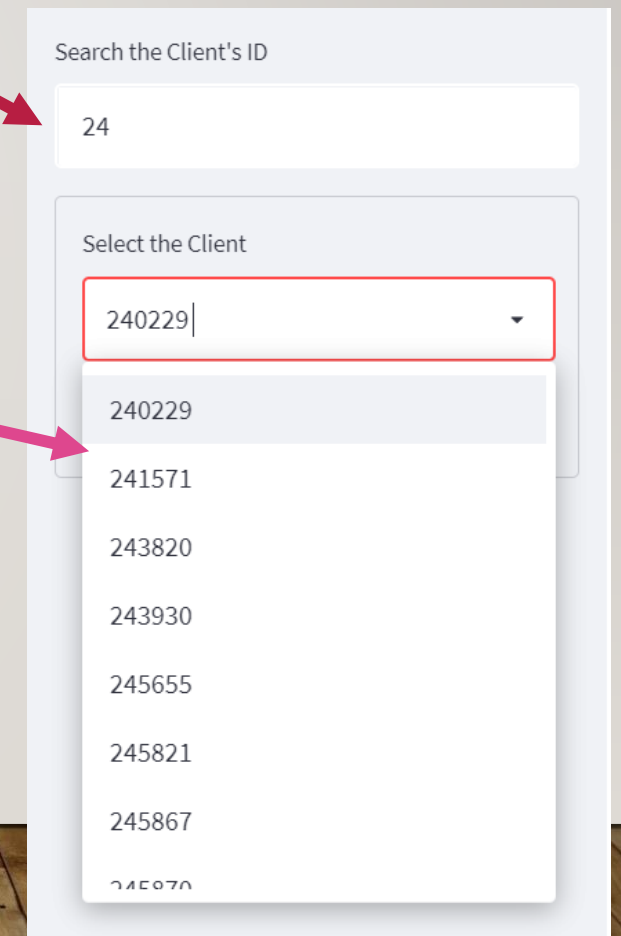
The screenshot shows a web browser window with the URL <https://share.streamlit.io/eva505/project7/dev2/Dashboard/dashboard.py>. The dashboard has a sidebar on the left with two sections: "Search the Client's ID" with a text input field, and "Select the Client" with a dropdown menu showing "102793" and a "Submit" button. The main area displays "Client Number:" followed by a green box containing the word "None".

Annotations with arrows point to the following elements:

- A red arrow points from the text "Recherche du identifiant du client (correspondant à son début)" to the "Search the Client's ID" input field.
- A red arrow points from the same text to the dropdown menu in the "Select the Client" section.
- A pink arrow points from the text "Sélection du identifiant du client" to the "Submit" button.
- A purple arrow points from the text "Envoie de la Sélection" to the "Submit" button.

**Recherche du identifiant du client
(correspondant à son début)**

**Sélection du
identifiant du client**



This image is a close-up of the "Select the Client" dropdown menu. The search bar at the top contains the text "24". The dropdown list is open, showing a list of client IDs. The first item, "240229", is highlighted with a red border. Below it, other visible items include "241571", "243820", "243930", "245655", "245821", and "245867".

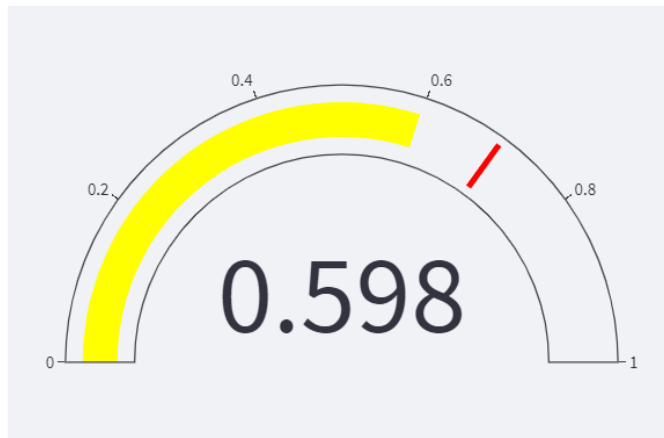
LE DASHBOARD – PRÉDICTION

Le Risque de Défaut de Crédit comme prédit par le modèle

Client
Number:

241571

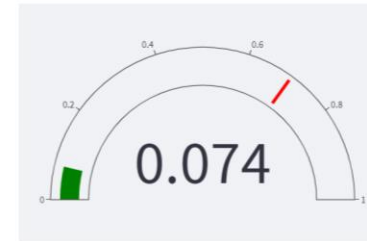
Default Risk



Client
Number:

100577

Default Risk



Très Faible
Risque

Client
Number:

429767

Default Risk



Risque
Elevée

LE DASHBOARD – LES FACTEURS DE DÉCISION



Les Facteurs le plus importantes dans la prédiction de modèle

Nombre des features à afficher

Défavorable

Favorable

LE DASHBOARD – COMPARAISON DES CLIENTS

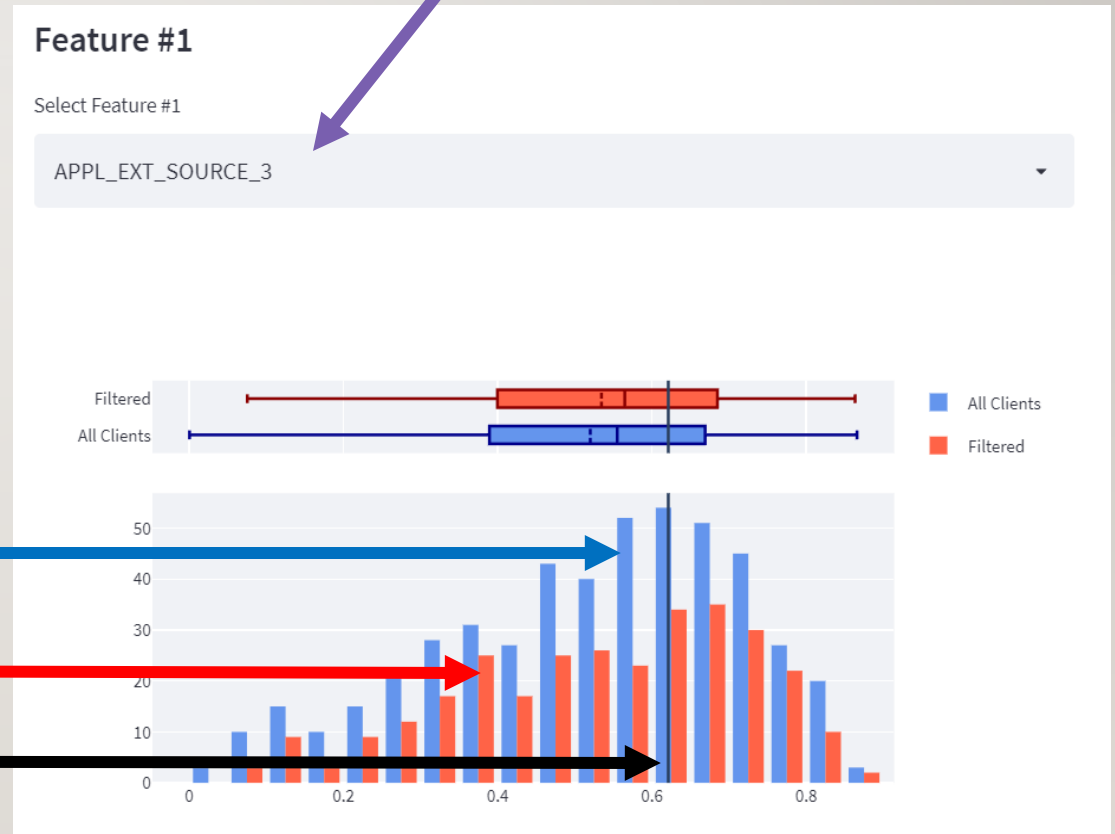
Choix des Filtres
pour faire la comparaison

Feature Comparisons among Clients

Filter Options

Gender	Defaulted	Income Range
<input checked="" type="checkbox"/> 1	<input type="checkbox"/> Yes	<input type="checkbox"/> +/- 10% of client's
<input type="checkbox"/> 0	<input checked="" type="checkbox"/> No	<input type="checkbox"/> +/- 20% of client's

Choix du Feature



Distribution de tous les clients
dans la base de données

Avec le filtre appliqué

Valeur du Client

POINTS D'AMÉLIORATION

Préparation et Sélection des Données

- ☐ nettoyage plus soigneux, i.e. détection des valeurs aberrantes, features catégorielles,...
- ☐ réduction de nombre des features utilisés
- ☐ features métiers ou polynomiales

Le Modèle - LGBM

- ☐ re-évaluation du modèle, i.e. mise à l'échelle, hyperparamètres,...

La fonction de coût

- ☐ analyse coûts-avantages, mauvais clients acceptés versus bons clients rejetés
- ☐ prendre en compte le montant du crédit



POINTS D'AMÉLIORATION

Evaluation du biais dans le modèle

- le genre est un feature le plus important dans le modèle, ce que problématique éthiquement

API et Dashboard

- retours des utilisateurs
- ajoute des fonctionnalités
- mise en place d'un propre base de données
- meilleur serveur

