

## **P7 : Implémentez un modèle de scoring**

11/08/2022

DUBART Maxime

# Implémentez un modèle de scoring

## **Création d'un outil de « scoring crédit »**

Probabilité de défaut de paiement et classification des demandes

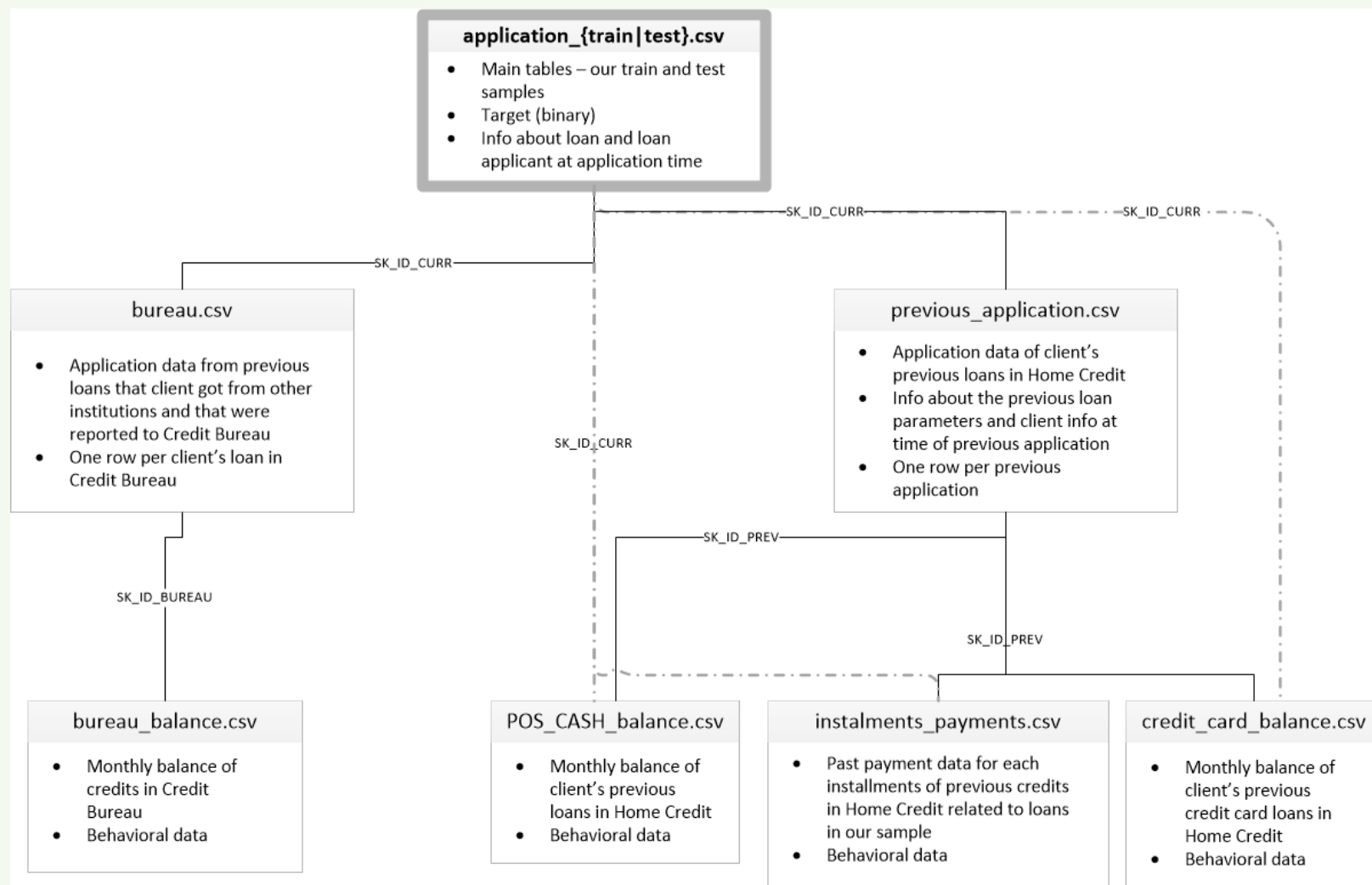
## **Développement d'un dashboard interactif**

Permettre plus de transparence vis-à-vis des décisions d'octroi de crédit  
pour les clients

# Implémentez un modèle de scoring

## Fichiers de données

7 fichiers



# Implémentez un modèle de scoring

## Fichiers de données

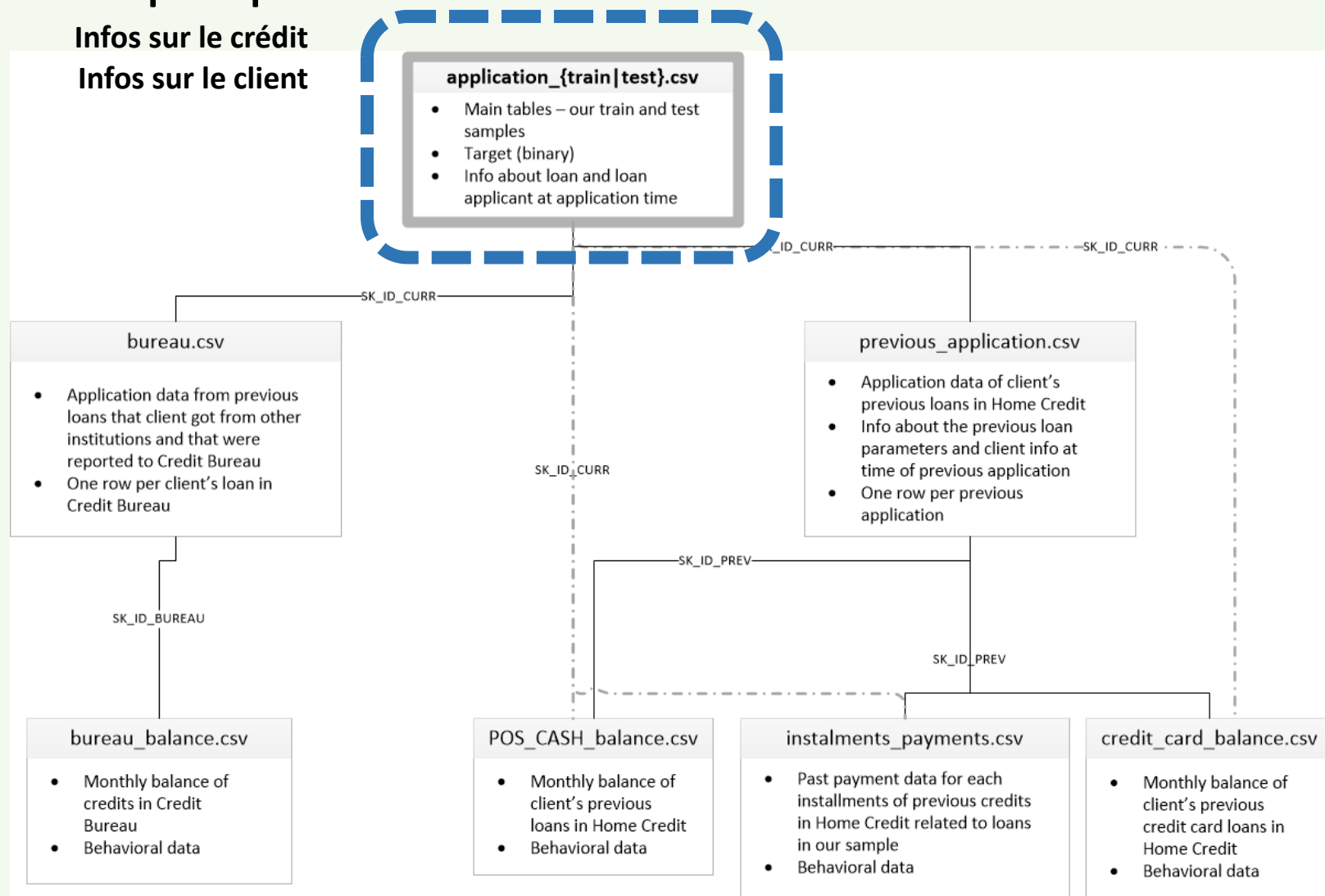
7 fichiers

8% de positifs (i.e. défaut)

### Fichier principal

Infos sur le crédit

Infos sur le client



# Implémentez un modèle de scoring

## Fichiers de données

7 fichiers

8% de positifs (i.e. défaut)

### Fichier principal

Infos sur le crédit  
Infos sur le client

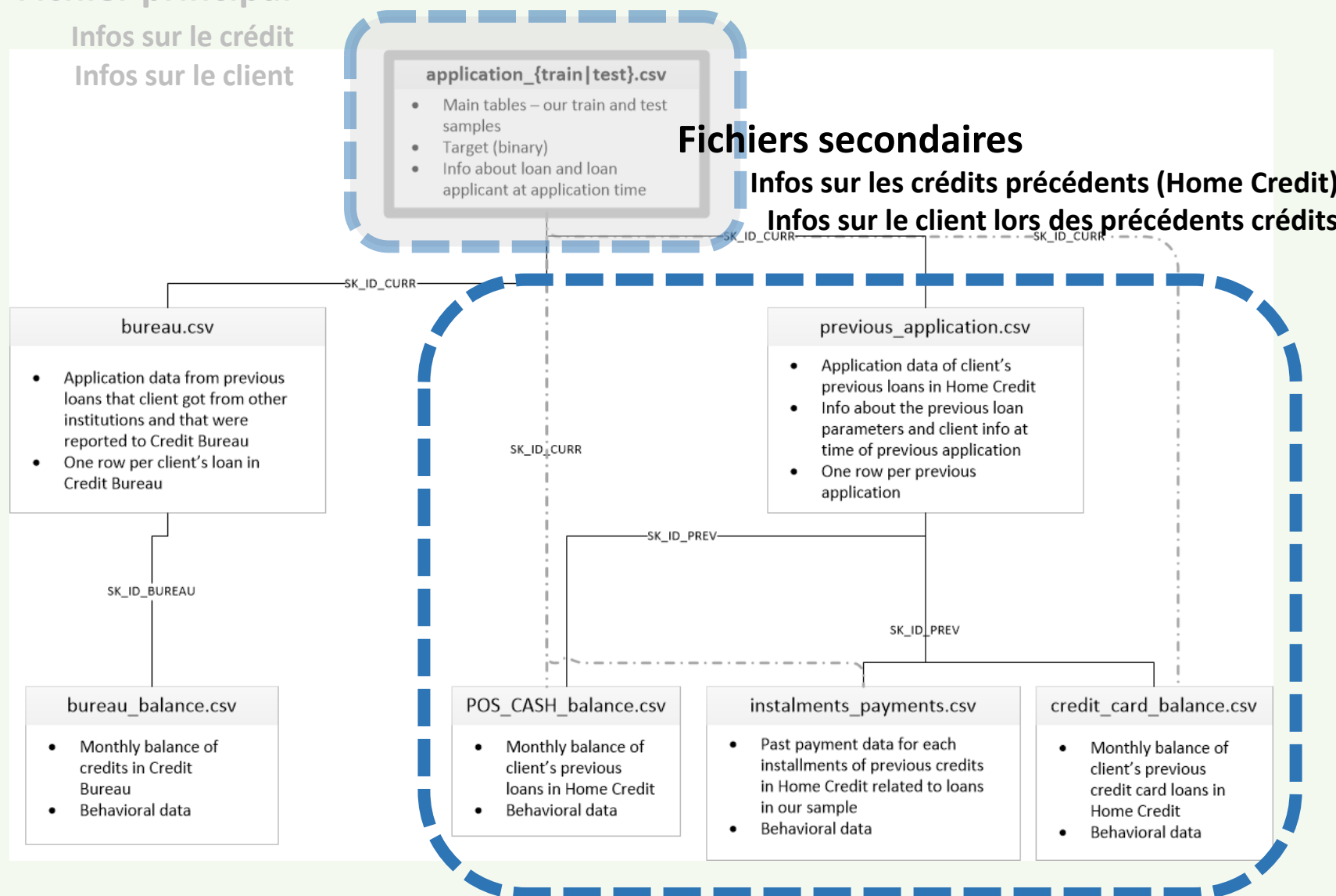
#### application\_{train|test}.csv

- Main tables – our train and test samples
- Target (binary)
- Info about loan and loan applicant at application time

### Fichiers secondaires

Infos sur les crédits précédents (Home Credit)

Infos sur le client lors des précédents crédits



# Implémentez un modèle de scoring

## Fichiers de données

7 fichiers

8% de positifs (i.e. défaut)

## Fichiers secondaires

Infos sur les crédits précédents (Autres institutions)  
Infos sur le client lors des précédents crédits

## Fichier principal

Infos sur le crédit  
Infos sur le client

### application\_{train|test}.csv

- Main tables – our train and test samples
- Target (binary)
- Info about loan and loan applicant at application time

## Fichiers secondaires

Infos sur les crédits précédents (Home Credit)

Infos sur le client lors des précédents crédits

### bureau.csv

- Application data from previous loans that client got from other institutions and that were reported to Credit Bureau
- One row per client's loan in Credit Bureau

SK\_ID\_BUREAU

### bureau\_balance.csv

- Monthly balance of credits in Credit Bureau
- Behavioral data

SK\_ID\_CURR

### POS\_CASH\_balance.csv

- Monthly balance of client's previous loans in Home Credit
- Behavioral data

SK\_ID\_PREV

### instalments\_payments.csv

- Past payment data for each installments of previous credits in Home Credit related to loans in our sample
- Behavioral data

SK\_ID\_PREV

### credit\_card\_balance.csv

- Monthly balance of client's previous credit card loans in Home Credit
- Behavioral data

### previous\_application.csv

- Application data of client's previous loans in Home Credit
- Info about the previous loan parameters and client info at time of previous application
- One row per previous application

# Implémentez un modèle de scoring

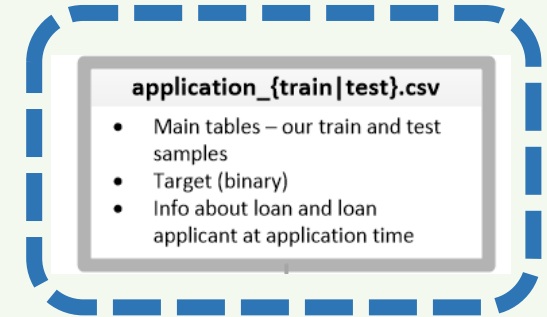
## Fichiers de données

7 fichiers

8% de positifs (i.e. défaut)

### Fichier principal

Infos sur le crédit  
Infos sur le client



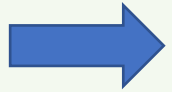
### 121 colonnes – différents types d'informations

- Sur le logement du client (surface, localisation, etc.)
- Sur les demandes à propos du client au Bureau des Crédits
- Informations administratives (e.g., quels documents ont été fournis)
- Informations sur le travail du client (localisation, type, salaire, etc.)
- Informations sur les hobbies du client
- Informations sur la situation du client (e.g., âge, statuts marital, # enfants etc.)
- Informations sur la demande de crédit (montant, mensualités, type de crédit)
- Informations provenant de sources externes

# Implémentez un modèle de scoring

## (i) Réduire la dimensionalité, suppression de :

- features avec > 40% de valeurs manquantes – excepté EXT\_SOURCE\_1
- features 'flag' *a priori* inutile pour la classification (e.g. enregistrement de tel document)
- features relatives à l'enregistrement dans la BDD (e.g. heure de la demande)



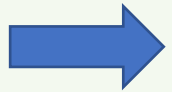
**Aboutit à la suppression de env. 70% des *features* (36 au lieu de 121)**



# Implémentez un modèle de scoring

## (i) Réduire la dimensionalité, suppression de :

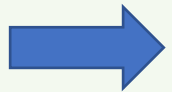
- features avec > 40% de valeurs manquantes – excepté EXT\_SOURCE\_1
- features 'flag' *a priori* inutile pour la classification (e.g. enregistrement de tel document)
- features relatives à l'enregistrement dans la BDD (e.g. heure de la demande)



**Aboutit à la suppression de env. 70% des *features* (36 au lieu de 121)**

## (ii) Adapté le jeu de données aux analyses :

- suppression des valeurs aberrantes (remplacées par NA)
- encodage des *features* catégorielles (one-hot encoding)
- création de variables composées (e.g., % de jours employés, valeurs du crédit relative au salaire, revenu moyen par membre du foyer, etc.)



**Au final, 143 *features* pour env. 300k clients**

# Implémentez un modèle de scoring

## Modèles de classification



### Régression logistique (Lasso/Ridge)

**Avantages :** Coefficients facilement interprétables, considération du déséquilibre de classe via poids, temps fit/prédiction, sélection variables possible, peu paramétrique

**Inconvénients :** Contrainte linéarité/additivité, imputation des valeurs manquantes (ici, par la moyenne)

# Implémentez un modèle de scoring

## Modèles de classification



### Régression logistique (Lasso/Ridge)

**Avantages** : Coefficients facilement interprétables, considération du déséquilibre de classe via poids, temps fit/prédiction, sélection variables possible, peu paramétrique

**Inconvénients** : Contrainte linéarité/additivité, imputation des valeurs manquantes (ici, par la moyenne)



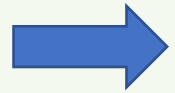
### XGBoost

**Avantages** : Pas de contrainte de linéarité, considération des valeurs manquantes

**Inconvénients** : Interprétation plus délicate (i.e. nécessite analyses postérieures), temps fit/prediction plus important, hautement paramétrique

# Implémentez un modèle de scoring

## Classes non balancées



**Peut-être problématique pour les algo. de ML**  
Pas vraiment pour la reg. log.



# Implémentez un modèle de scoring

## Classes non balancées



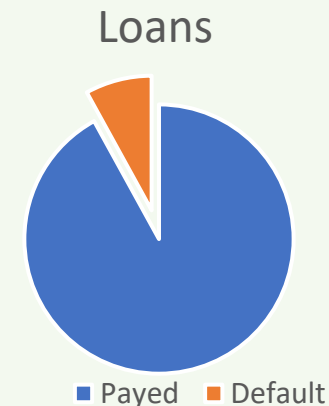
**Peut-être problématique pour les algo. de ML**

Pas vraiment pour la reg. log.



**Stratégies possibles**

- Sur-échantillonner la classe minoritaire
- Sous-échantillonner la classe majoritaire (1:2)
- Création d'échantillons synthétiques (e.g., SMOTE)



# Implémentez un modèle de scoring

## Classes non balancées



**Peut-être problématique pour les algo. de ML**

Pas vraiment pour la reg. log.



**Stratégies possibles**

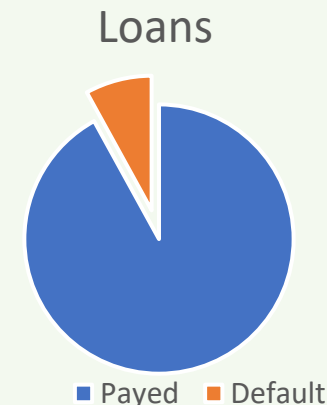
- Sur-échantillonner la classe minoritaire
- **Sous-échantillonner la classe majoritaire (1:2)**
- Création d'échantillons synthétiques (e.g., SMOTE)



*Forte augmentation de la taille du jeu de données*



*Via k-nn, problème des variables binaires*



# Implémentez un modèle de scoring

## Classes non balancées



➡ **Peut-être problématique pour les algo. de ML**  
Pas vraiment pour la reg. log.

➡ **Stratégies possibles**

- Sur-échantillonner la classe minoritaire
- **Sous-échantillonner la classe majoritaire (1:2)**
- Création d'échantillons synthétiques (e.g., SMOTE)

*Forte augmentation de la taille du jeu de données*

*Via k-nn, problème des variables binaires*

➡ **Entraînement sur jeux de données re-balancés ou non**

# Implémentez un modèle de scoring

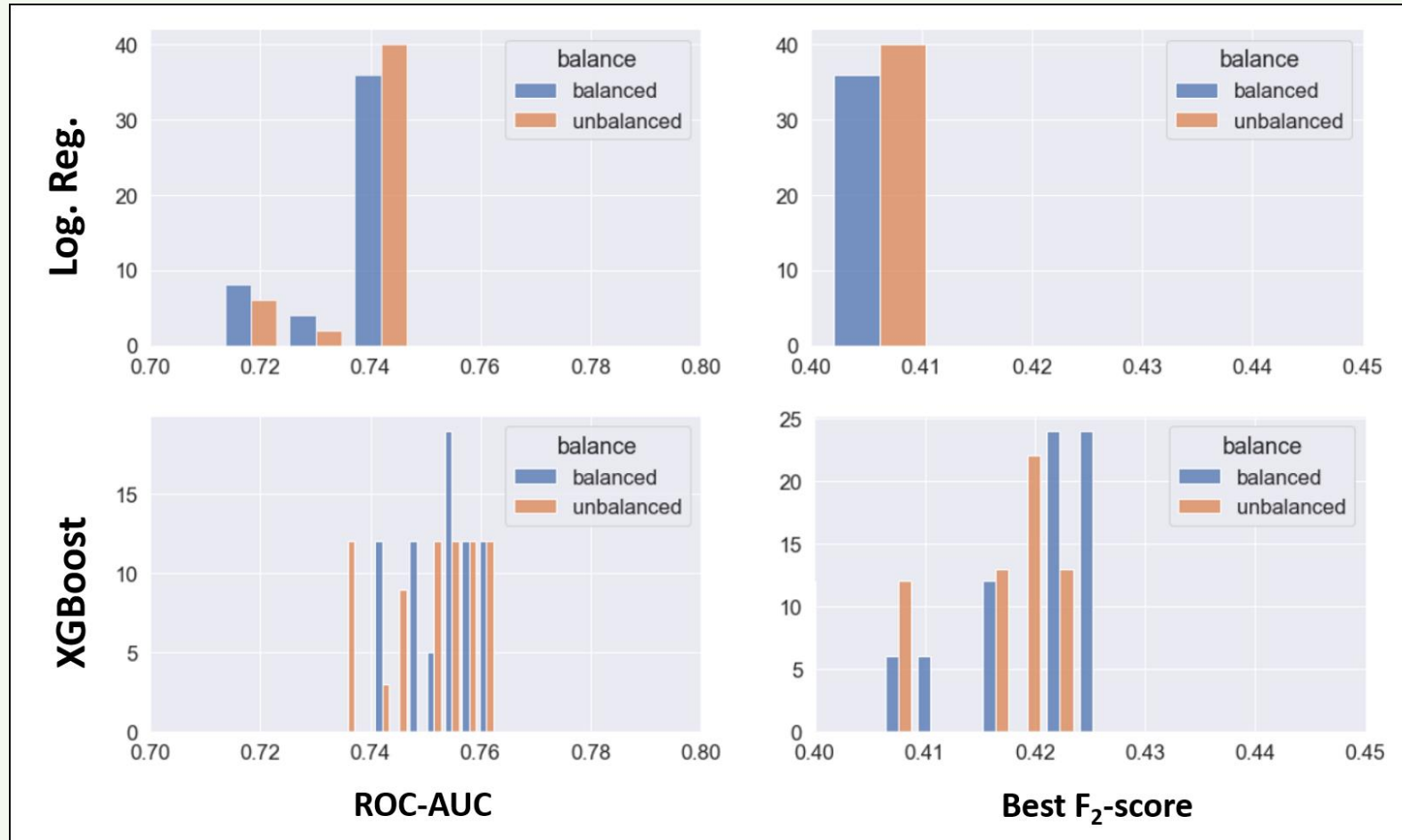
## Méthodologie

- ➡ **Datasets original ou rebalancé (downsampling)**
- ➡ **Choix des hyper-paramètres (recherche sur grille / cross-validation)**
  - Split train / validation : 70/30
  - 5-folds cross-validation
- ➡ **Métriques : ROC-AUC et  $F_{\beta}$ -Score**



# Implémentez un modèle de scoring

## Scores de validation



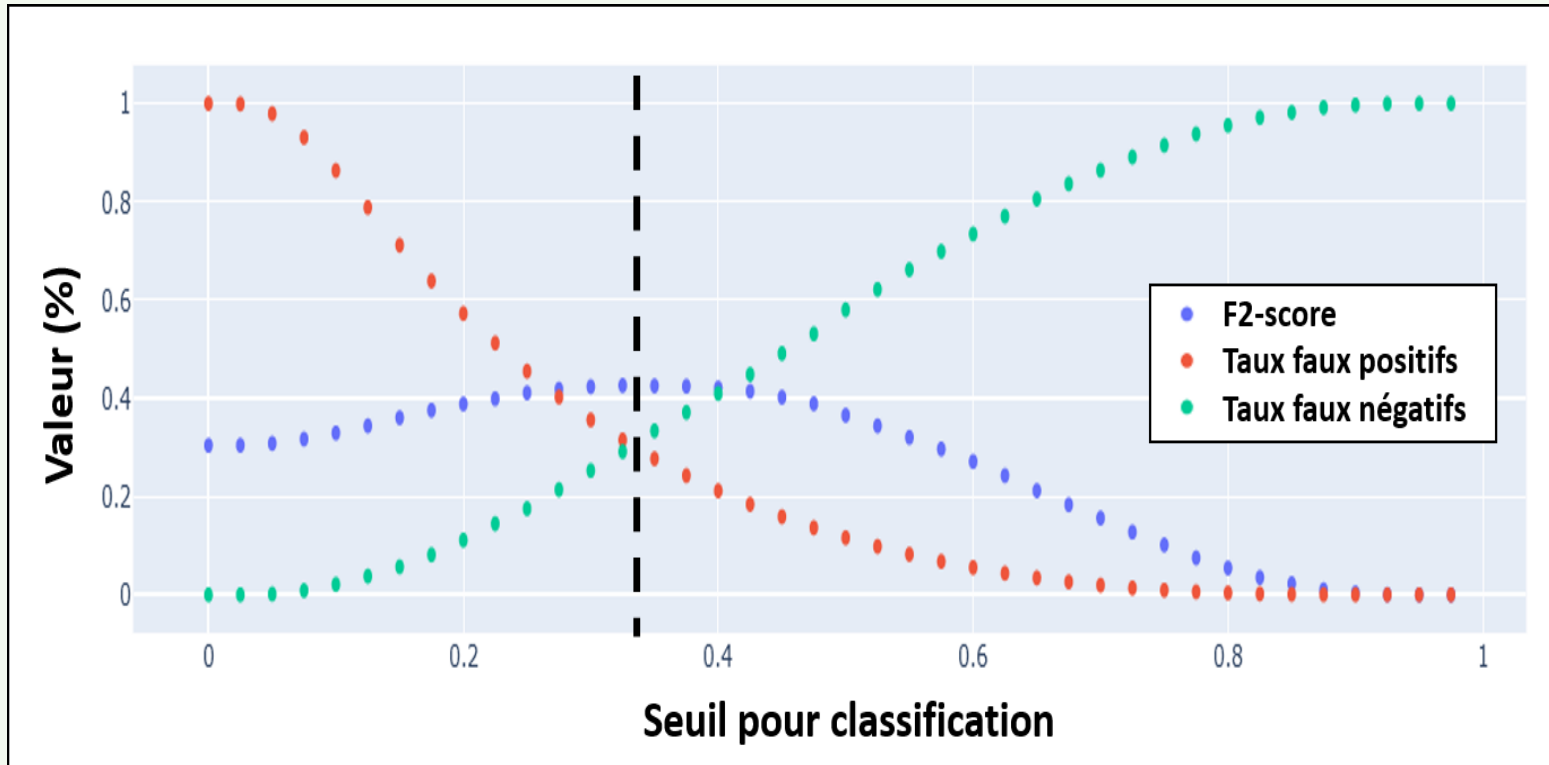
**XGBoost légèrement meilleur que la régression logistique**

**Peu d'effet du rééquilibrage des classes**

**Scores proches de ceux obtenus sur les meilleurs résultats Kaggle (avec l'ensemble des fichiers)**

# Implémentez un modèle de scoring

## Sélection du seuil



Seuil de 0.3/0.35 maximise le  $F_2$ -score

A ce seuil : env. 30% de FP et FN

# Implémentez un modèle de scoring

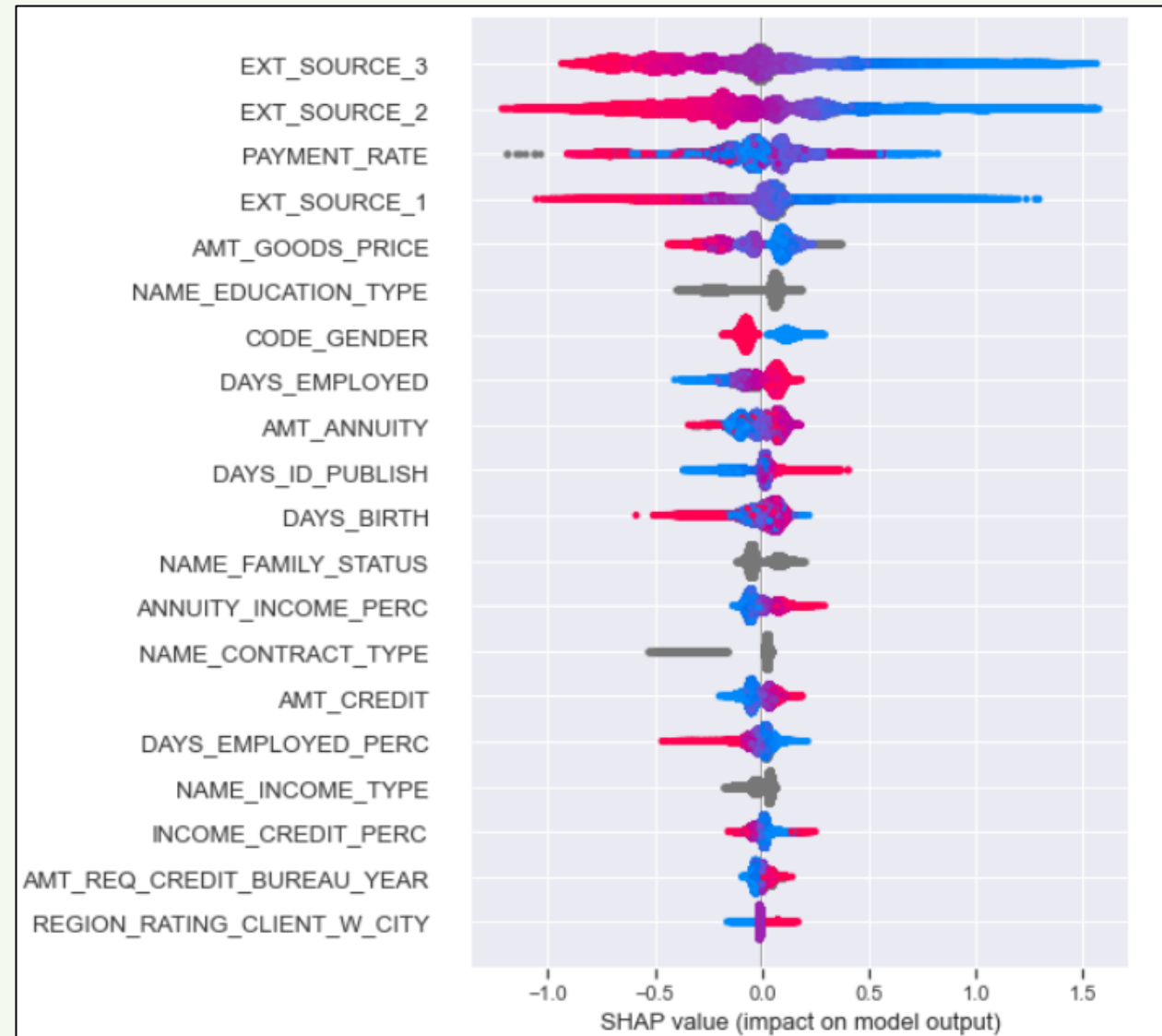
## Global features importances (SHAPs)

Sources extérieures d'info important beaucoup

Le taux de remboursement (i.e., annuité / total)

Le niveau d'éducation, le genre, la durée depuis laquelle le client travaille (inversée)

Low values High values



# Déployer le modèle & dashboard

## API & Dashboard

```
1 # -*- coding: utf-8 -*-
2
3 from flask import Flask, jsonify, request
4 import pandas as pd
5 import shap
6 import pickle
7 import sklearn
8 import xgboost
9 import requests
10 import numpy as np
11
12 app = Flask(__name__)
13
14 # Read Scaler
15 with open('standard_scaler.pkl', 'rb') as inp:
16     ss = pickle.load(inp)
17
18 # Read model
19 with open('xgb_model.pkl', 'rb') as inp:
20     xgb = pickle.load(inp)
21
22 # Read explainer
23 with open('shap_xgb_explainer.pkl', 'rb') as inp:
24     xgb_explainer = pickle.load(inp)
25
26 # Read data (X)
27 df = pd.read_csv('./us_data_subset.csv')
28 df_unscaled = df.copy()
29
30 # Get targets & ids as separate lists
31 target = df.pop('TARGET')
32 sk_id_curr = df.pop('SK_ID_CURR')
33
34 # Scaled df
35 df_ss = pd.DataFrame(ss.transform(df), columns = df.columns)
36 probs = xgb.predict_proba(df_ss)[:,1]
37
38 # Get fn, fp, tp, tn
39 def fp_fn_tp_tn(pp_, y_):
40     thresh = np.arange(0,1,0.025)
41
42     output = []
43
44     for t in thresh:
45         yp_ = (pp_ > t).astype('int')
46         cm = sklearn.metrics.confusion_matrix(y_, yp_)
47         tn_ = cm[0,0]
48         fp_ = cm[0,1]
49         fn_ = cm[1,0]
50         tp_ = cm[1,1]
51         prec_ = tp_ / (tp_ + fp_)
52         recall = tp_ / (tp_ + fn_)
```



# Déployer le modèle & dashboard

## Outils utilisés

**API**

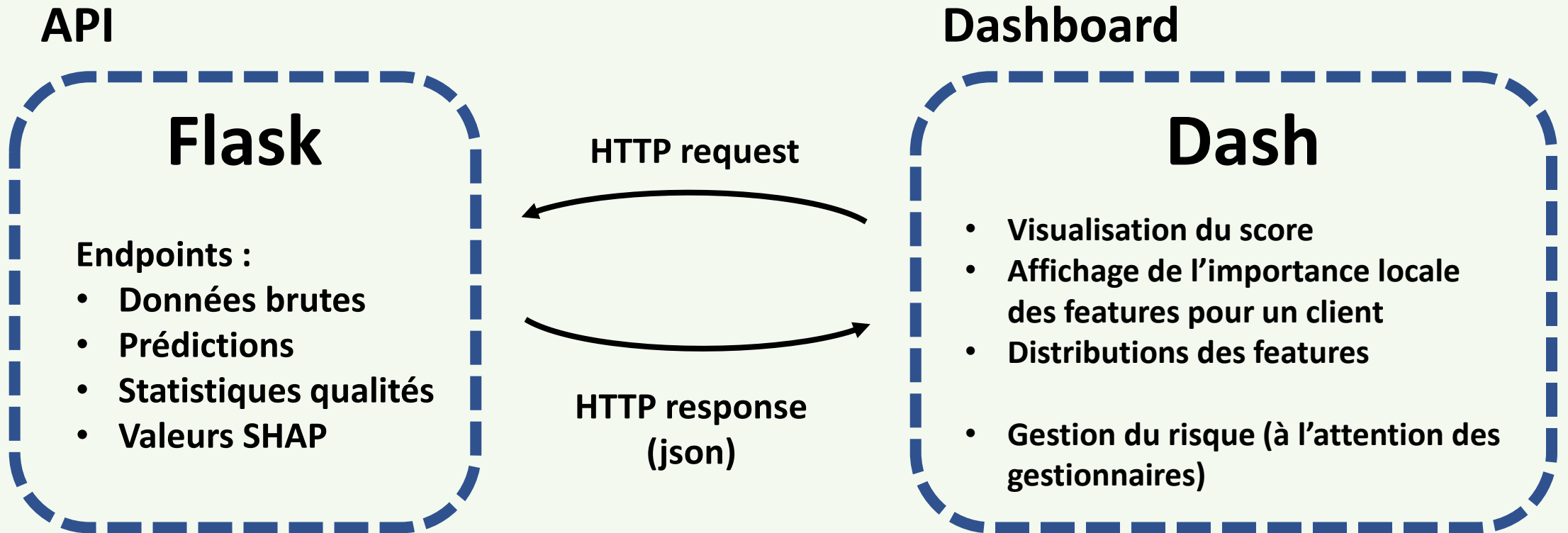
**Flask**

**Endpoints :**

- **Données brutes**
- **Prédictions**
- **Statistiques qualités**
- **Valeurs SHAP**

# Déployer le modèle & dashboard

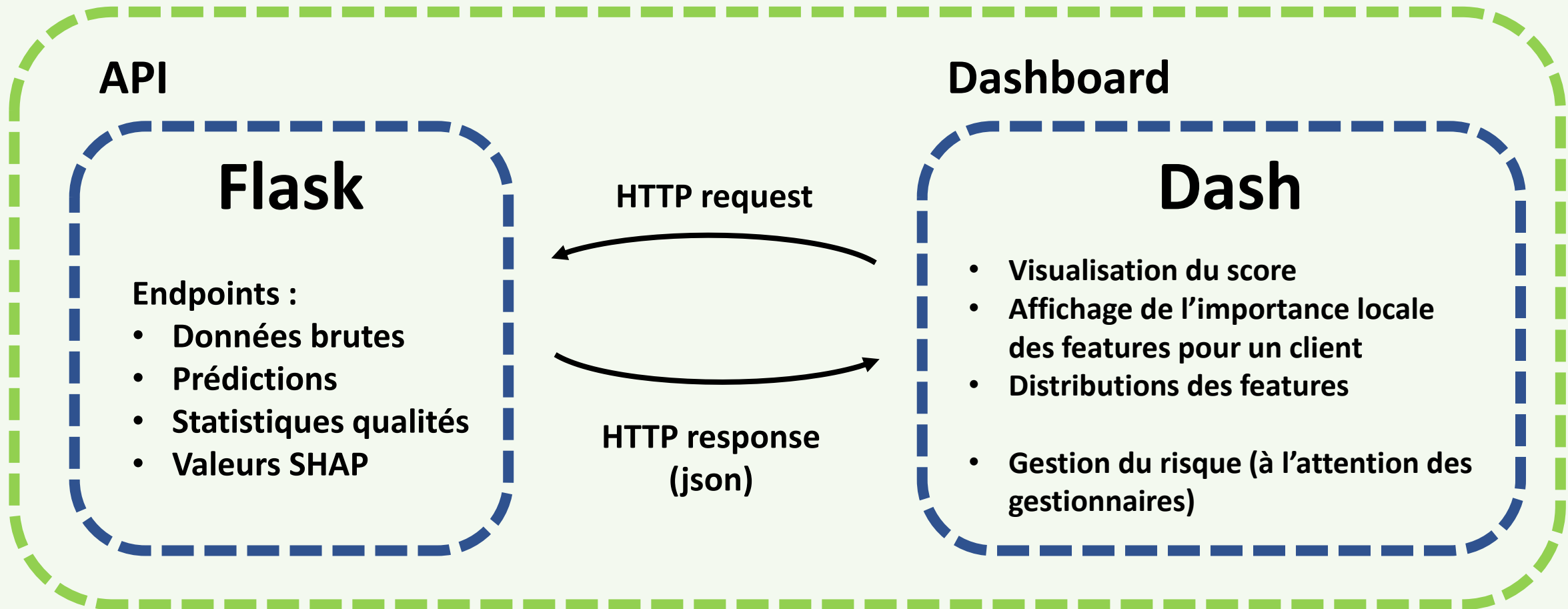
## Outils utilisés



# Déployer le modèle & dashboard

## Outils utilisés

Versionnées sur GitHub & déployées sur Heroku



# Déployer le modèle & dashboard

**<https://credits-ocr-dashboard.herokuapp.com/>**



**Merci**

# Implémentez un modèle de scoring

## Sélection du seuil

