

OPENCLASSROOM

Parcours data scientist en alternance (@Synapse-Medicine)
Implémentez un modèle de scoring

16 décembre 2022

Maxime Dupouy

Livrable 7 Implémentez un modèle de scoring



Introduction

Objectifs du livrable

Présentation BDD et features
selection/engineering

2 - Modélisation

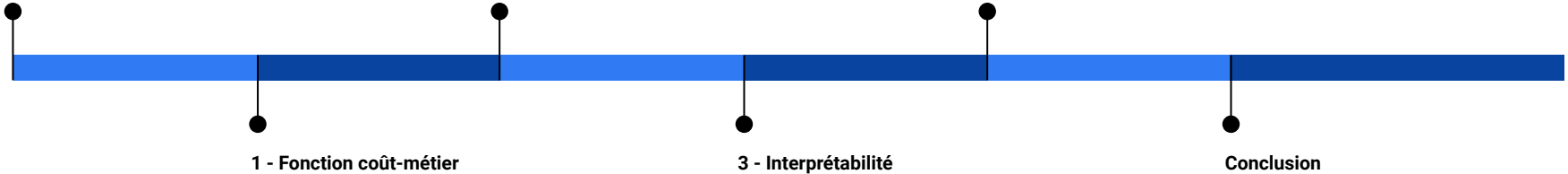
Méthodologie

Résultats

4 - Dashboard

Structuration

Démonstration



Introduction

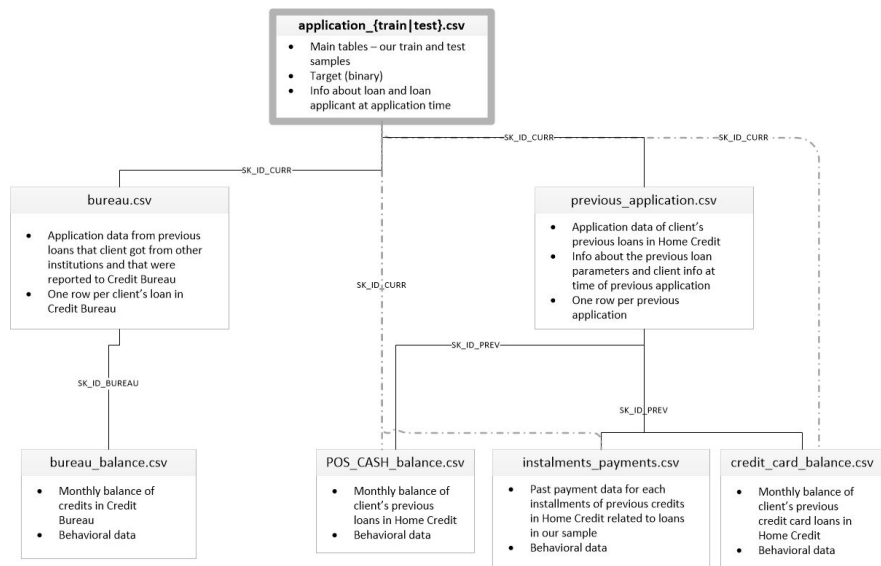
Objectifs

Implémentez un modèle de scoring

Missions	Objectifs
<ul style="list-style-type: none">- Construire un modèle de scoring qui donnera une prédiction sur la probabilité de faillite d'un client de façon automatique.- Construire un dashboard interactif permettant d'interpréter les prédictions faites par le modèle, et d'améliorer la connaissance client des chargés de relation client.	<ul style="list-style-type: none">- Présenter son travail de modélisation à l'oral- Réaliser un dashboard- Rédiger une note méthodologique- Utiliser un logiciel de version de code- Déployer un modèle via une API dans le Web

Introduction

BDD



Donnée provient de Home Credit Bank et est hébergée sur [Kaggle](#).

Le but de Home Credit est d'élargir l'inclusion financière de la population non bancarisée.

Ils utilisent une variété de données alternatives, comme des informations sur les télécommunications et les transactions, pour prédire les capacités de remboursement de ses clients.

Introduction

Features selection/engineering

WILL KOEHRSEN · 4Y AGO · 80 715 VIEWS

398 Copy & Edit 686

Automated Feature Engineering Basics

Python · Home Credit Default Risk Feature Tools · Home Credit Default Risk

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 307511 entries, 0 to 307510
Data columns (total 17 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   SK_ID_CURR                               307511 non-null  int64
1   EXT_SOURCE_1                             307511 non-null  float64
2   EXT_SOURCE_2                             307511 non-null  float64
3   EXT_SOURCE_3                             307511 non-null  float64
4   DAYS_BIRTH                               307511 non-null  int64
5   AMT_CREDIT                               307511 non-null  float64
6   AMT_ANNUITY                              307511 non-null  float64
7   DAYS_EMPLOYED                             307511 non-null  int64
8   AMT_GOODS_PRICE                          307511 non-null  float64
9   DAYS_ID_PUBLISH                          307511 non-null  int64
10  OWN_CAR_AGE                              307511 non-null  float64
11  BUREAU_MAX_DAYS_CREDIT                   307511 non-null  float64
12  BUREAU_MAX_DAYS_CREDIT_ENDDATE           307511 non-null  float64
13  BUREAU_MAX_DAYS_ENDDATE_FACT             307511 non-null  float64
14  PREV_SUM_MIN_AMT_PAYMENT                 307511 non-null  float64
15  PREV_MEAN_MIN_AMT_PAYMENT                307511 non-null  float64
16  TARGET                                   307511 non-null  float64

dtypes: float64(13), int64(4)
memory usage: 42.2 MB
```

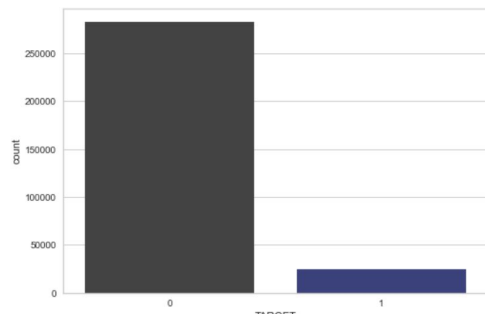
Partie “sous traitée” sur des notebooks déjà réalisés par des data scientist - [lien](#)

Obtention d’une BDD composée de :

- 15 colonnes les plus “impactantes” (sélection avec “Feature Tool” et “Deep Feature Synthesis”)
- 307 511 observations
- TARGET :
 - 0 = client sans défaut de paiement
 - 1 = client en défaut de paiement

1 - Fonction coût-métier

Problématique et résolution



Répartition des classes 0 et 1

	Prédiction clients non défaut (0)	Prédiction clients en défaut (1)
Clients non défaut (0)	TN	FP
Clients en défaut (1)	FN	TP

Matrice de confusion

Objectifs:

Éviter de se tromper lorsque le client est à défaut de paiement, tout en acceptant le plus de crédits possible.

Métriques sélectionnées :

- recall
- précision

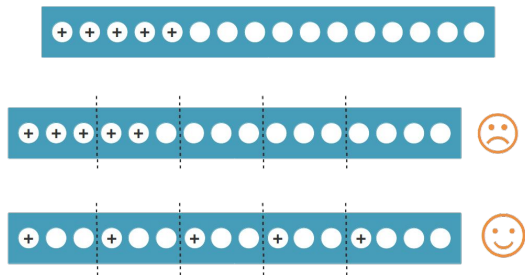
Une métrique existe pour avoir un ratio entre les deux : le $F_{\beta 2}$ score.

$$F_{\beta} = (1 + \beta^2) \cdot \frac{precision \cdot recall}{(\beta^2 \cdot precision) + recall}$$

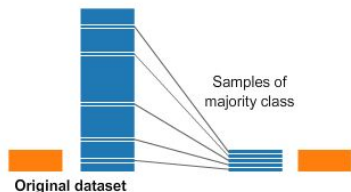
2 - Modélisation

Méthodologie

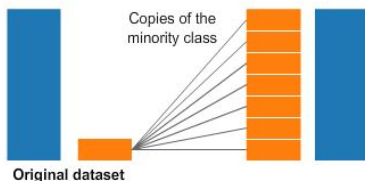
Stratification



Undersampling



Oversampling



Méthode pour évaluer de manière robuste :

- Stratified K fold

Contrebalancer le déséquilibre des classes :

- Oversampling (SMOTE)
- Undersampling (NearMiss)
- Variation threshold

Modèles entraînés:

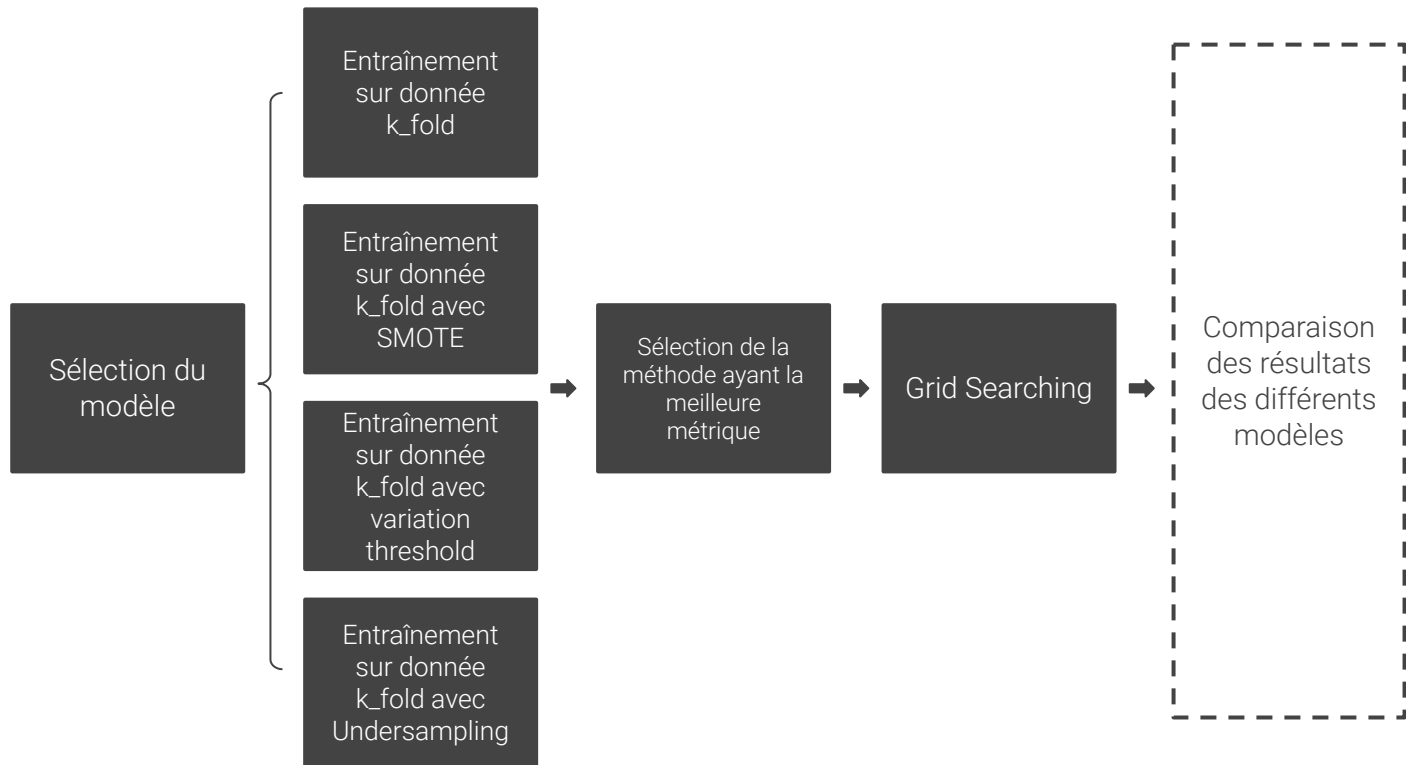
- Logit
- RandomForestClassifier base
- RandomForestClassifier notebook
- XGBoostClassifier
- LGBMClassifier

Evaluation sur :

- $F \beta 2$ score
- recall, précision, accuracy

2 - Modélisation

Méthodologie



Résumé

2 - Modélisation

Résultats

Baseline :

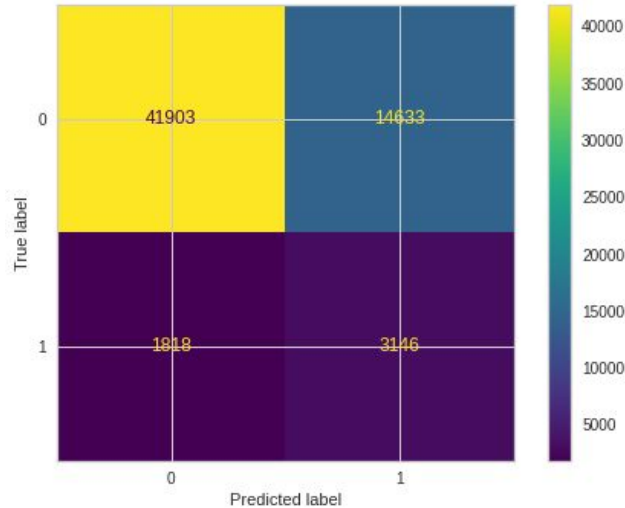
- $F \beta_2 = 0.086$
- accuracy = 0.85

	F β_2 score	Accuracy
Logit (SMOTE)	0.39	0.68
RandomForestClassifier (finetuné)	0.26	0.33
RandomForestClassifier finetuné notebook	0.399	0.68
LightGBM non finetuné	0.42	0.69
XGBoost finetuné	0.418	0.73

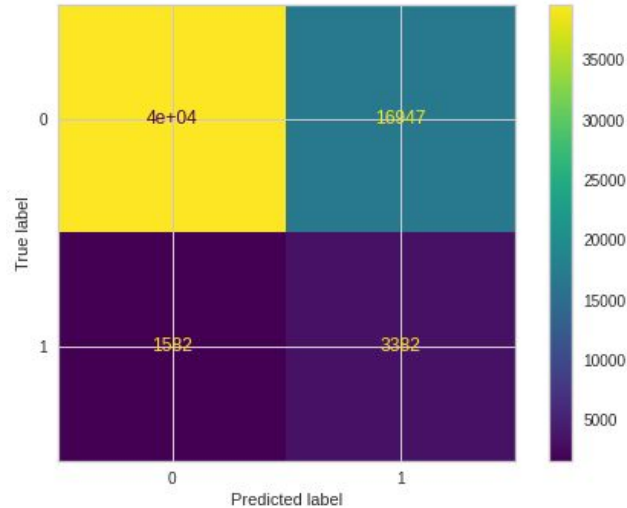
Benchmark

2 - Modélisation

Résultats



LightGBM

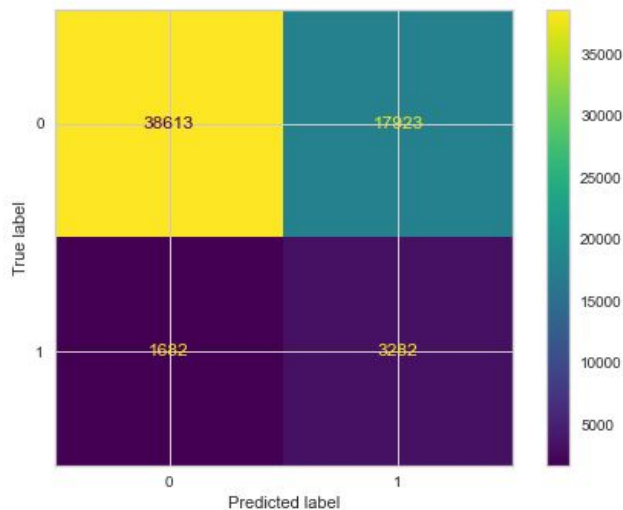


XGBoost

Matrice de
confusion

2 - Modélisation

Choix du modèle



random_forest_notebook	
f_beta_score_	0.399685
precision_	0.681218
recall_	0.661148
accuracy_	0.681218

Implémentation du **RandomForestClassifier** notebook, sans oversampling, undersampling, nouveau finetuning.

Choix sur F β 2 score **et** recall

3 - Interprétabilité

Globale

Réalisée avec scikit learn (feature importance)



Locale

Réalisée avec Shap (Force plot)



Descriptif des clients

Box plot afin de visualiser où se situe le client par rapport à la distribution globale



Implémentation

4 - Dashboard

 **FastAPI**



- gère la data
- gère modèle ML
- génère graphique

 **Streamlit**



- gère affichage
- gère input de données clients



Backend	→	Fastapi
Frontend	→	Streamlit
Déploiement	→	Heroku
Versionning	→	<u>Github</u>

4 - Dashboard



Conclusion



Modélisation :

- Résultats supérieurs à la baseline, mais reste assez faibles pour une utilisation brute
- Axes d'amélioration : réaliser la feature selection/engineering, tester modèle bagging, finetuner sur sample_weight

Dashboard :

- Répond au cahier des charges
- Axes amélioration : performance du site mitigée, graphiques interactifs

Global :

- S'appuyer sur des connaissances métiers et/ou de la donnée afin de réaliser certains choix

Points de difficultés :

- Ne pas avoir fait la feature selection : gain de temps mais mauvaise compréhension de la donnée
- Travail sur mac m1 (puce silicone) : problème de compatibilité (Heroku, xgboost, lgbm)



Je vous remercie pour votre attention

Dashboard



Versioning



*Désolé si le premier livrable ne correspond pas à ce que je présente, je ne comprends pas trop la consigne ...
Le site sera "destroy" à la fin de la soutenance pour des raisons évidentes de coût*