Credit repayment

Implémentation d'un modèle de scoring

Déploiement dans le

Cloud

DATA SCIENCE CREDIT REPAYMENT PREDICTION

Owner: Ph. Deflandre Création: 03/05/2025

Environnement de travail

Environnement Jupyter Lab sur Windows 11 (Edge Version 130.0.2849.56)

Environnement virtuel Imblearn version 0.13.0

Python version 3.11.7 XGBoost version 3.0.0

Numpy version 1.26.4 Mlflow version 2.21.0

Matplotlib version 3.8.0 Requests version 2.31.0

Pandas version 2.1.4 Json version 2.0.9

Seaborn version 0.12.2 LightGBM version 4.6.0

Scikit-learn version 1.5.2 Uvicorn version 0.34.0

Shap version 0.46.0 FastAPI version 0.115.12

Pydantic version 2.11.3 Streamlit version 1.45.0

,

Rappel de la mission

- Comparer des modèles de scoring pour déterminer la probabilité qu'un prêt soit remboursé.
- Comprendre les features qui vont être déterminantes dans la prédiction.
- Mettre en production le modèle de scoring le plus performant à l'aide d'une API. Réaliser une interface de test de l'API.
- Créer une interface utilisateur user friendly.
- Analyser un éventuel Data Drift.

Analyse de la base de données

000

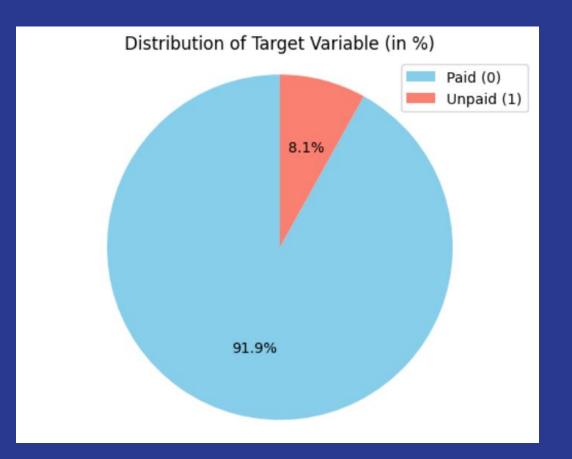
Source: https://www.kaggle.com/c/home-credit-default-risk/data

- application_train.csv
 - 121 features
 - 1 target
 - 307 511 lignes (crédits)

- > type de contrat
- informations personnelles
- > situation familiale
- > informations administratives
- > informations relatives au crédit

- application_test.csv
 - 121 features
 - 48 744 lignes (crédits)

Target

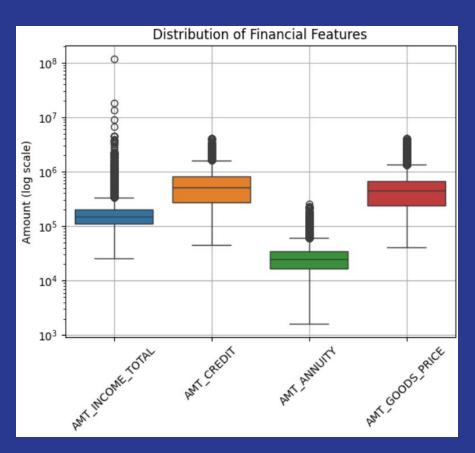


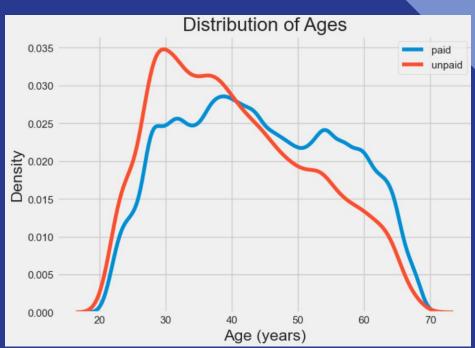
Classification binaire déséquilibrée

0 => le prêt est remboursé

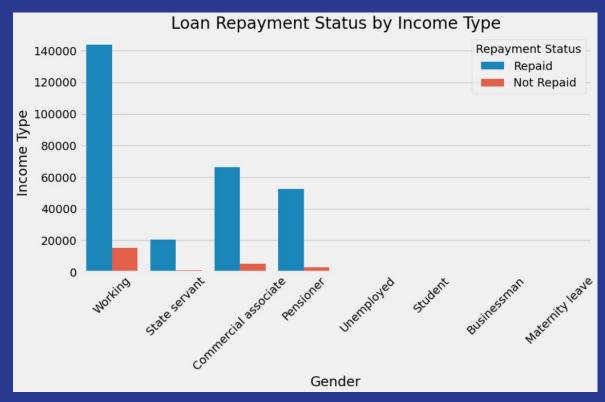
1 => le prêt n'est pas remboursé

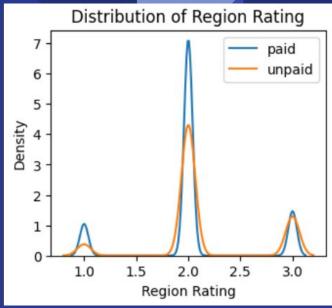
Features

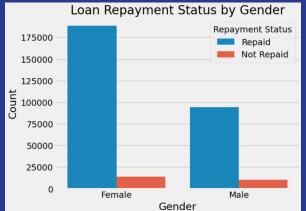




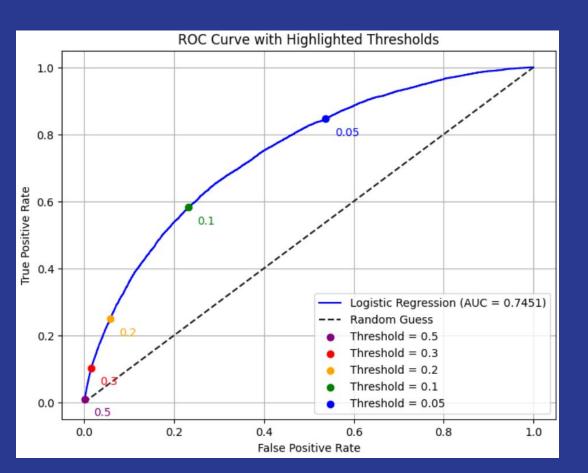
Features







Baseline - Logistic Regression classification



Premier modèle de référence

Comment estimer la qualité de prédiction du modèle ?

```
Train Accuracy: 0.919
Test Accuracy: 0.919
Test Balanced Accuracy: 0.504

Confusion Matrix:
[[70609 78]
[ 6132 59]]
```

Métrique business personnalisée

\$ E

Crédits non accordés qui auraient été remboursés

Train Accuracy: 0.919
Test Accuracy: 0.919

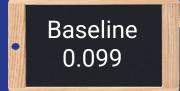
Test Balanced Accuracy: 0.504

Confusion Matrix:

[[70609 78] [6132 59]]

10 1

Crédits accordés non remboursés_



Feature engineering

Polynomial Features

$$(EXT_SOURCE_1)^a \times (EXT_SOURCE_2)^b \times (EXT_SOURCE_3)^c \times (DAYS_BIRTH)^d$$

 $(a, b, c, d) \in [0; 3]$

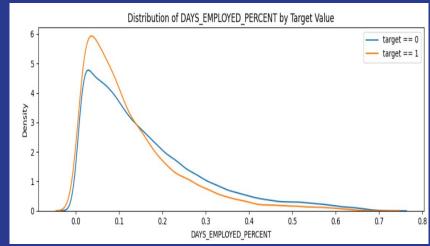
Features "métier"

$$taux \ d'endettement = \frac{montant \ du \ crédit}{revenus}$$

$$taux \ d'effort \ annuel = \frac{annuit\acute{e}}{revenus}$$

$$ratio \ de \ remboursement = \frac{annuit\acute{e}}{montant \ du \ crédit}$$

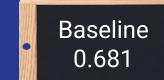
$$taux \ d'emploi \ relatif = \frac{jours \ travaill\acute{e}s}{\mathring{a}ge}$$



Comparaison des modèles

Logistic Regression with features engineering

- class_weight hyperparamètre
- cross-validation
- > MLFlow
- utilité du Feature Engineering



| | Feature | Coefficient | Absolute_Coefficient |
|-----|---------------------------|-------------|----------------------|
| 7 | AMT_ANNUITY | 2.494649 | 2.494649 |
| 8 | AMT_GOODS_PRICE | -2.436208 | 2.436208 |
| 80 | DEF_60_CNT_SOCIAL_CIRCLE | 1.737929 | 1.737929 |
| 278 | ANNUITY_INCOME_PERCENT | 1.648104 | 1.648104 |
| 78 | DEF_30_CNT_SOCIAL_CIRCLE | 1.533950 | 1.533950 |
| 10 | DAYS_BIRTH | 1.145663 | 1.145663 |
| 246 | EXT_SOURCE_3_y | -0.954602 | 0.954602 |
| 252 | EXT_SOURCE_2 EXT_SOURCE_3 | -0.846523 | 0.846523 |
| 33 | EXT_SOURCE_3_x | -0.820314 | 0.820314 |
| 93 | FLAG_DOCUMENT_13 | -0.785311 | 0.785311 |
| | | | |





Comparaison des modèles

Random Forest

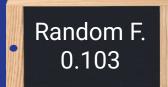
- class_weight hyperparamètre
- problème de surapprentissage

AdaBoost

- class_weight hyperparamètre
- optimisation des hyperparamètres avec GridSearchCV
- > entraînement lent du modèle

GradientBoosting

- optimisation des hyperparamètres avec GridSearchCV
- amélioration de l'AUC
- pas d'hyperparamètre d'équilibrage des classes







Comparaison des modèles

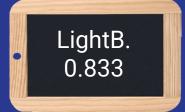
XGBoost

- scale_pos_weight hyperparamètre
- optimisation des hyperparamètres avec GridSearchCV

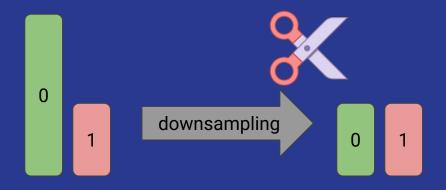


LightBoosting

- scale_pos_weight hyperparamètre
- > optimisation des hyperparamètres avec GridSearchCV
- > rapide



Downsampling - imblearn



XGBoost

- scale_pos_weight = 1
- optimisation des hyperparamètres avec GridSearchCV



Upsampling - SMOTE



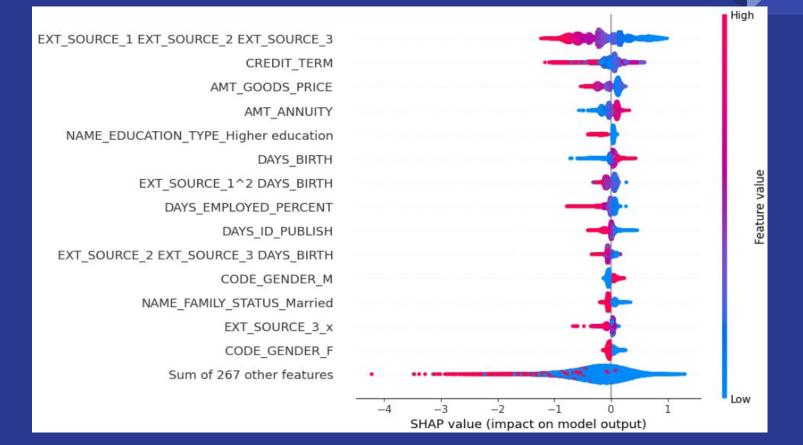
XGBoost

- scale_pos_weight = 1
- > hyperparamètres du downsampling



Interprétabilité globale - XGBoost





Interprétabilité locale - XGBoost

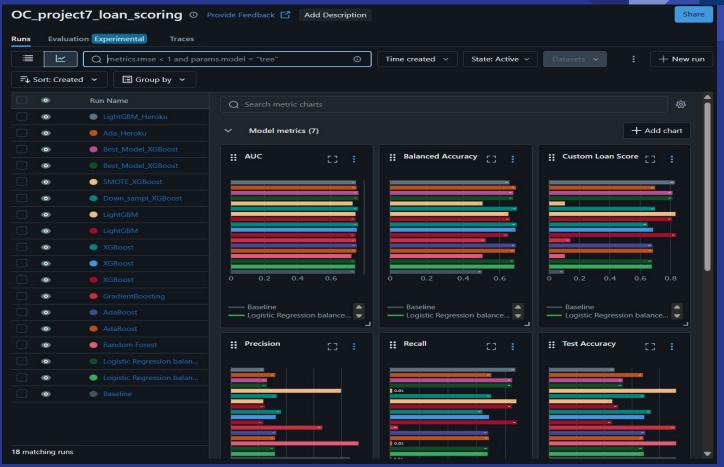


exemple prêt n°3



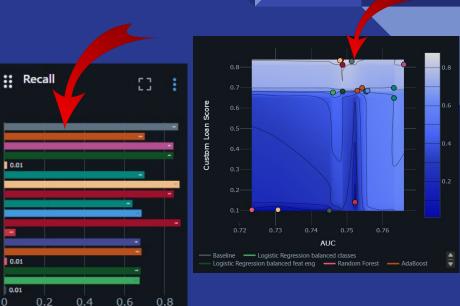
Tracking d'expérimentation - UI MLFlow





Choix du modèle à déployer en production

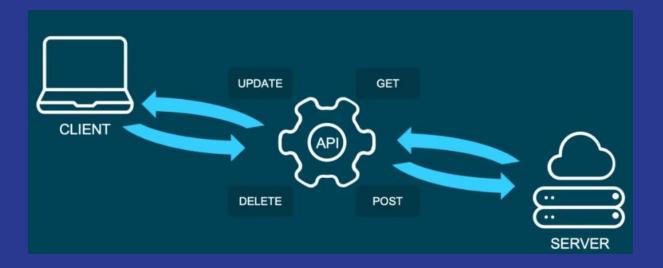




- combinaisons de plusieurs métriques
- taille des dependencies nécessaires sur Heroku!

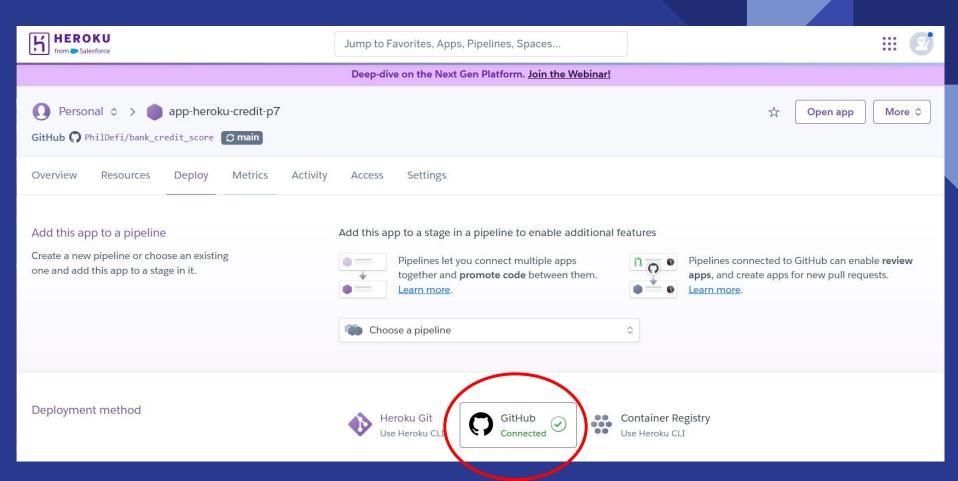
API

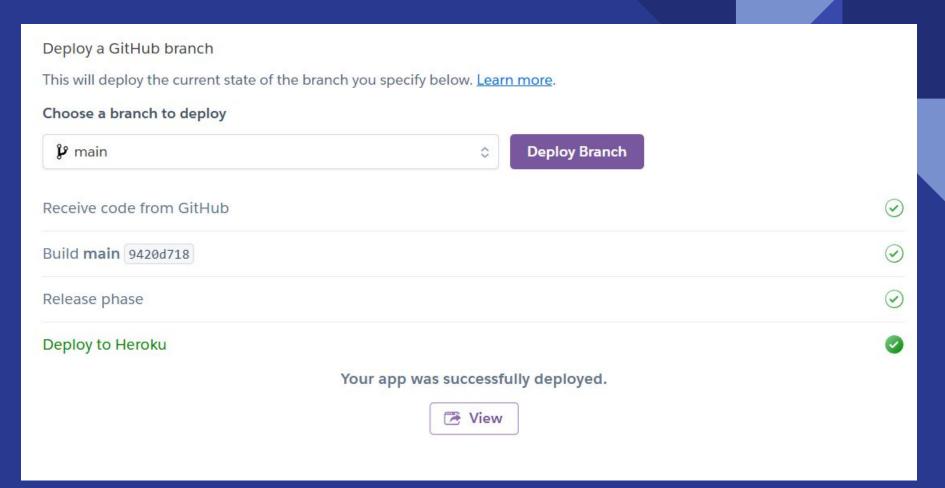


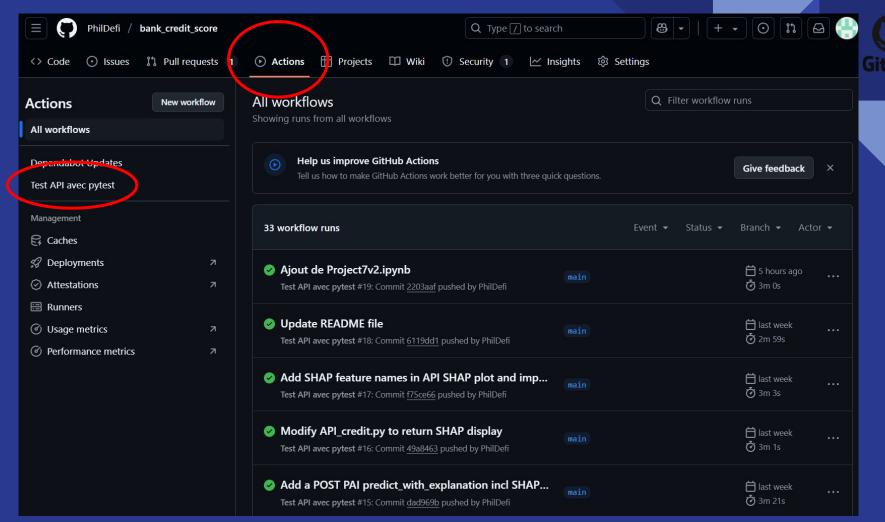


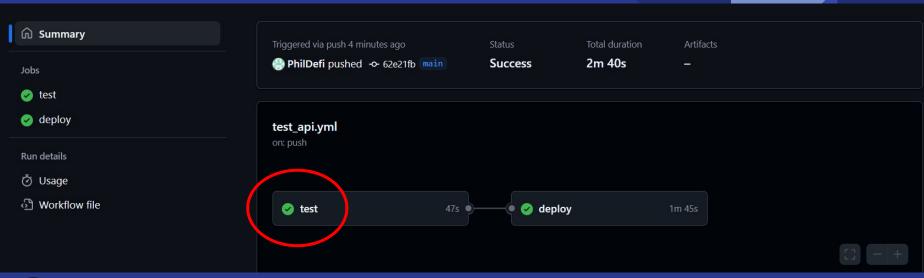
- FastAPI
- Tests de l'API pytest
- Déploiement sur Heroku
- Workflow : déploiement automatique sur Heroku depuis GitHub conditionné au succès des tests unitaires





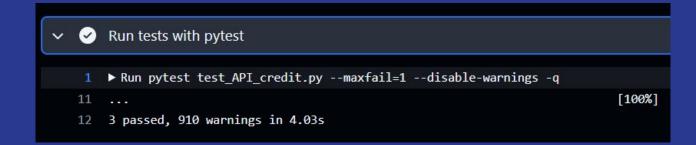


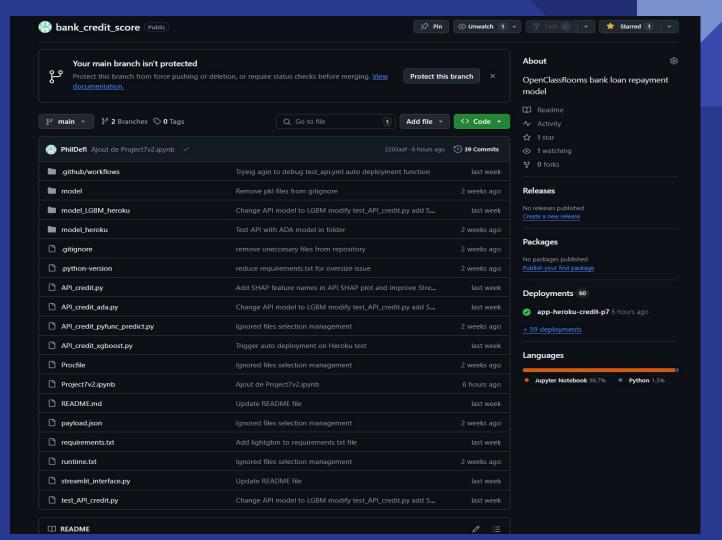










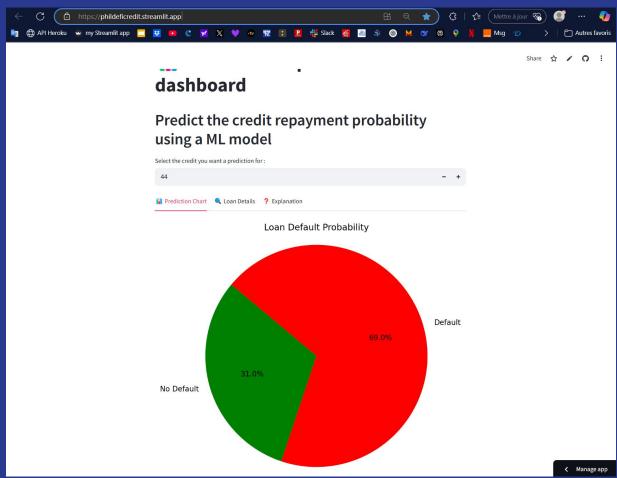


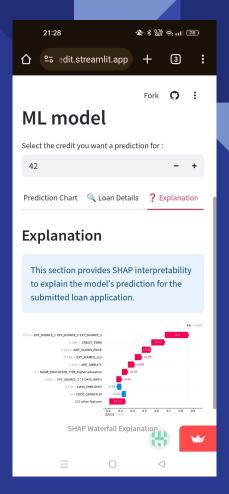


Architecture MLOps **Data scientist** pytest 💙 heroku **GitHub** API sync Streamlit https://phildeficredit.streamlit.app/ mlflow :: Precision

User interface







Data Drift

current reference count 307511.00 48744.00 0.00 (0.00%) 0.00 (0.00%) missing CODE_GENDER cat XNA 4.0 (0%) 0.0 (0%) M 105059.0 (34%) 16066.0 (33%) 202448.0 (66%) 32678.0 (67%)



 Count
 307511.00
 48744.00

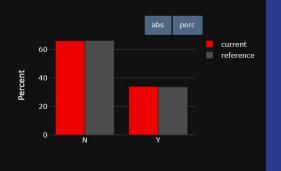
 FLAG_OWN_CAR missing
 0.00 (0.00%)
 0.00 (0.00%)

 Y
 104587.0 (34%)
 16433.0 (34%)

 N
 202924.0 (66%)
 32311.0 (66%)

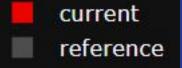
current

reference





application_test.csv



application_train.csv



Conclusion

Elaboration d'un modèle en tenant compte du déséquilibre des classes

Définition d'une métrique métier sur-mesure

Mise en place d'une API déployée dans le Cloud

Mise en place d'une interface utilisateur déployée dans le Cloud

Démonstrateur :

https://phildeficredit.streamlit.app/

Repository GitHub:

https://github.com/PhilDefi/bank_credit_score