
IMPLÉMENTATION DU MODÈLE DE SCORING

MISE EN PRODUCTION DE L'API ET DU DASHBOARD INTERACTIF



SOMMAIRE

Prêt à dépenser

- Présentation
 - Projet
 - Jeu de données
- Modélisation
 - Métriques
 - Démarche
 - Tracking des modèles
 - Synthèse des résultats
- Pipeline de déploiement
 - Architecture
 - Tests unitaires
- DataDrift
- Dashboard

- Notes :
 - Données : <https://www.kaggle.com/c/home-credit-default-risk/data>
 - Kernel de départ adapté ensuite : <https://www.kaggle.com/code/jsaguiar/lightgbm-with-simple-features/script>
 - GitHub : https://github.com/elisepoupi/P7_scoring
 - Dashboard : <https://epoupi-pret-a-depenser.herokuapp.com/>
noms à utiliser pour tests : personnages De Game of Thrones (ex: Robert, Ellaria...)

PRÉSENTATION PROJET

Prêt à dépenser

- Mettre en œuvre un outil de scoring crédit pour calculer la probabilité qu'un client rembourse son crédit et déterminer l'octroi ou non du crédit
 - **Maximisation des bénéfices de l'entreprise**
- Communiquer les résultats sur un dashboard à destination des chargés de relation client et du client.
 - **Rapidité de l'API**
 - **Clarté de l'analyse**
 - **Marge de négociation pour le chargé de relation client**
- S'appuyer sur des sources variées : données d'autres institutions financières, données comportementales...
 - **Utilisation d'un maximum de variables**

Prêt à dépenser

Formulaire de connexion

Veuillez entrer votre nom :

PRÉSENTATION JEU DE DONNÉES

Prêt à dépenser

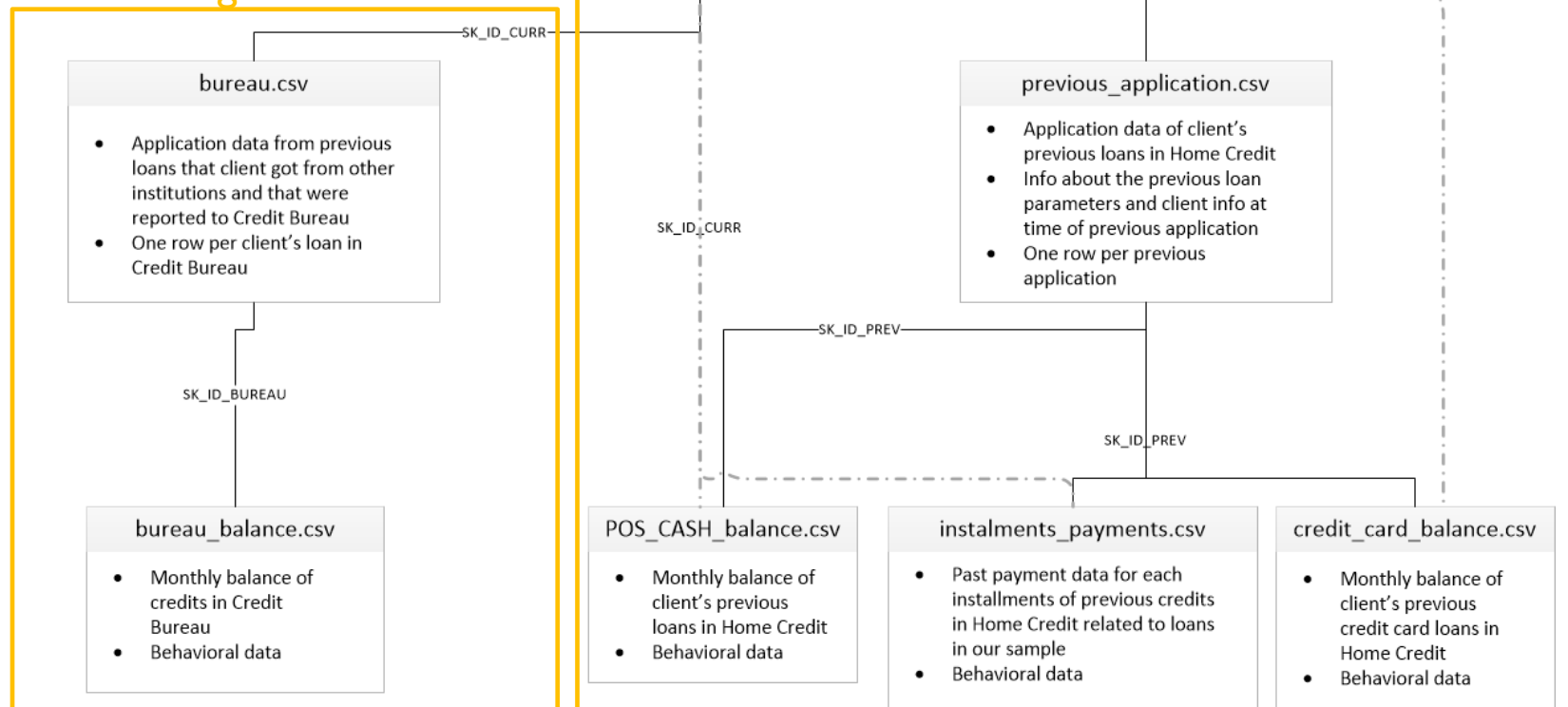
application_{train|test}.csv

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

120 variables
356 000 individus

Anciens prêts dans
d'autres organismes bancaires

Anciens prêts
chez prêt à dépenser



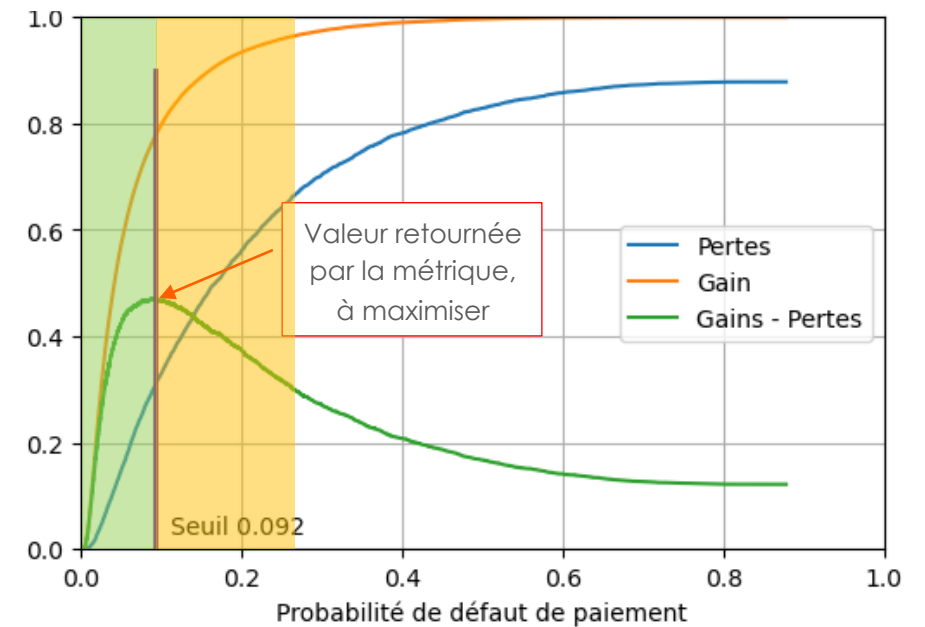
- Nombreuses valeurs manquantes
- 767 variables après regroupements

MODÉLISATION MÉTRIQUES

Prêt à dépenser

- Métriques standard suivies :
 - ROC AUC et Accuracy
- Métrique métier : maximiser les bénéfices de l'entreprise
 - Hypothèses :
 - Bénéfice d'un dossier réalisé sans défaut de paiement : 1
 - Perte d'un dossier réalisé avec défaut de paiement : 10
 - Formule : bénéfice potentiel maximum / bénéfice total des dossiers sans défaut de paiement

- Métrique métier choisie :



- Prêts accordés
- Marge de négociation laissée au chargé de relation client
Avec gains potentiels > 30% des gains possibles

MODELISATION

DÉMARCHE 1/2

Prêt à dépenser

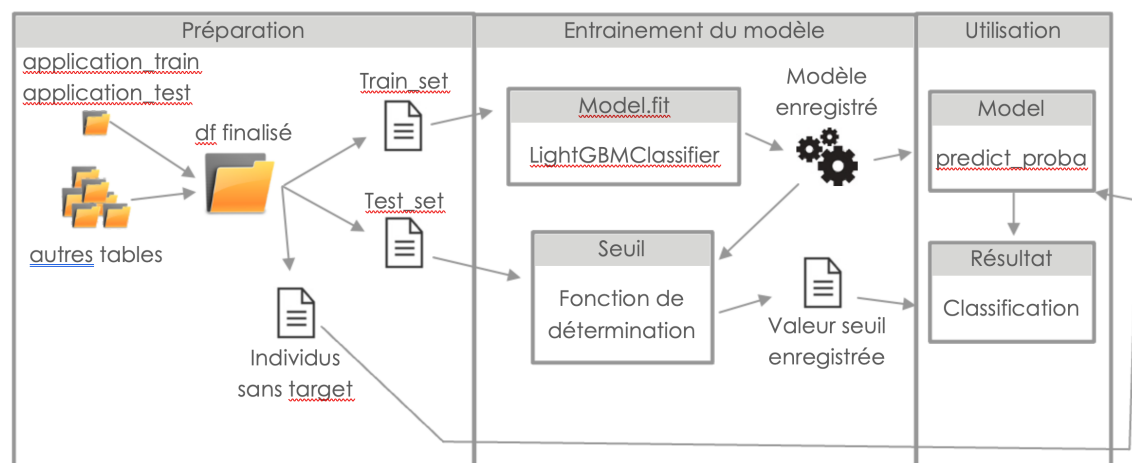
	Essai 1	Essai 2	Model 1	Essai 4	Model 2
Objectif	1 ^{er} résultat rapide	Pouvoir utiliser Lime ou une régression linéaire ou de l'over sampling	Améliorer le score métier	Réduire le nombre de variables pour explications claires	Négociation : marge plus importante et possibilité de modifier des données
Préprocessing		KNNImputer (trop de variables) ACP (NaN) ANOVA KNNImputer 200	Imputation manuelle des valeurs manquantes	Imputation manuelle des valeurs manquantes Sélection ANOVA	Imputation manuelle des valeurs manquantes Pipeline intégrant la transformation des données Under sampling
Valeurs NaN	Non traitées	Aucune	Aucune	Aucune	Aucune
Réduc. dimens.	Aucune	200 features ANOVA	Aucune	200 feat. ANOVA	Aucune
Modèle	LightBGMClassifier	LGBMClassifier	LGBMClassifier	LGBMClassifier	LGBMClassifier
Explication locale	Lime (NaN)	Lime imparfait (catégoriel : 1 variable par label existant)	Lime imparfait	Lime imparfait	Lime plus clair (variable catégorielle = 1 variable)
Score	0.470	0.439	0.470	0.451	0.404

MODELISATION

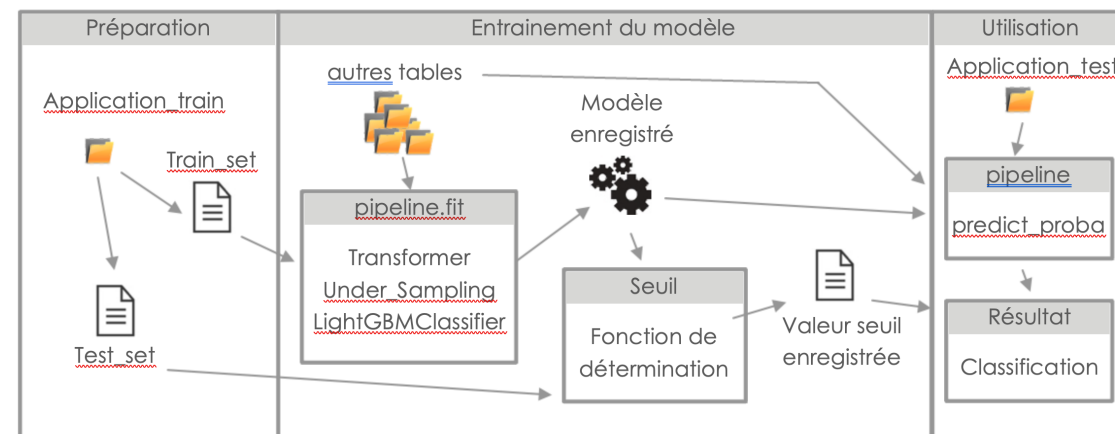
DÉMARCHE 2/2

Prêt à dépenser

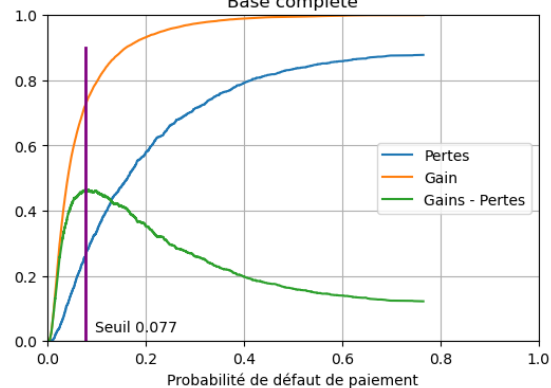
MODELE N° 1



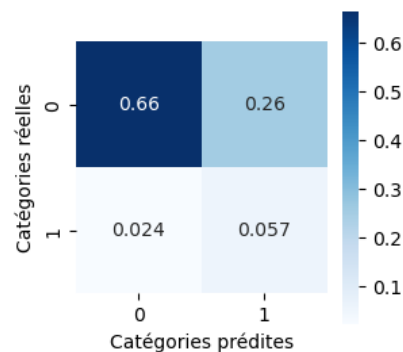
MODELE N° 2



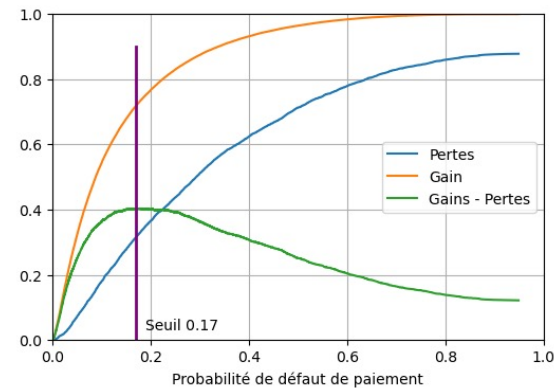
Détermination du seuil impliquant le refus du prêt
Base complète



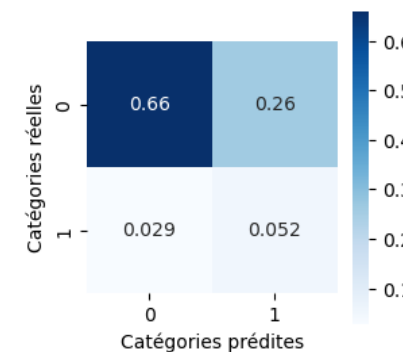
Matrice de confusion



Détermination du seuil impliquant le refus du prêt



Matrice de confusion



MODELISATION

TRACKING DES MODÈLES

Prêt à dépenser

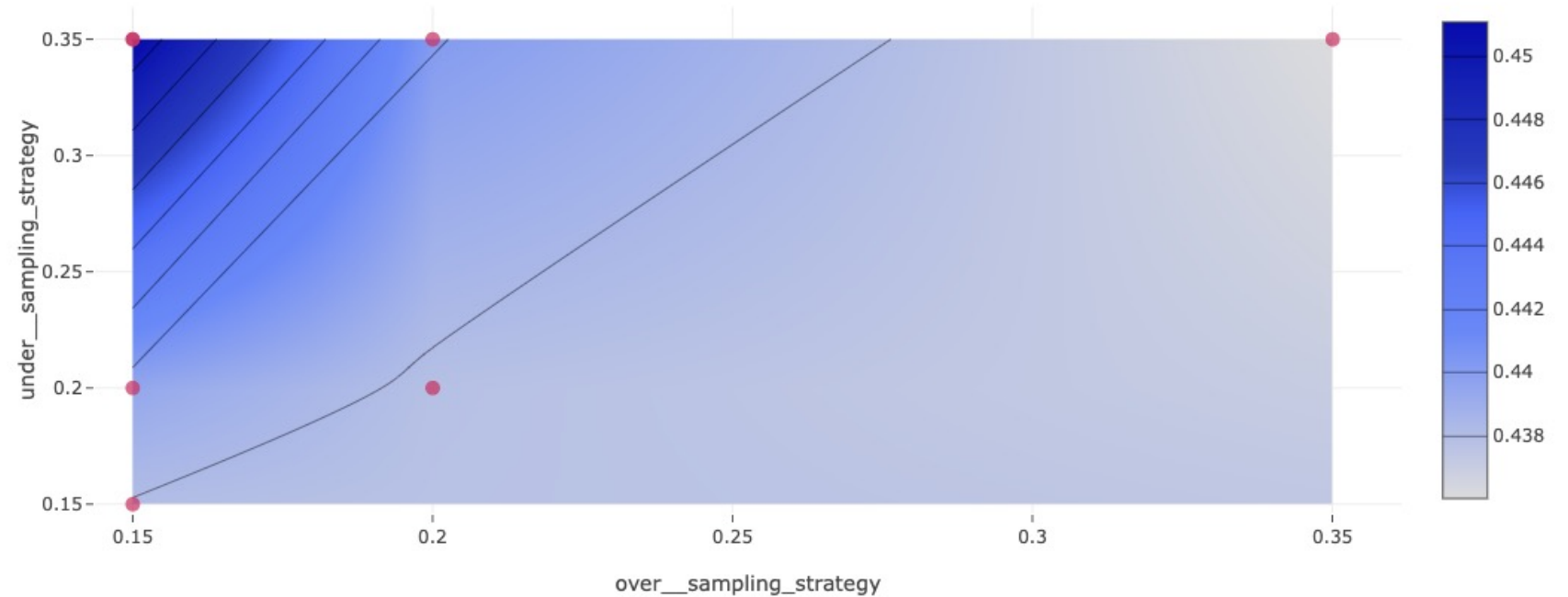
Exemple de visualisation MLFlow

X-axis:
over__sampling_strategy

Y-axis:
under__sampling_strategy

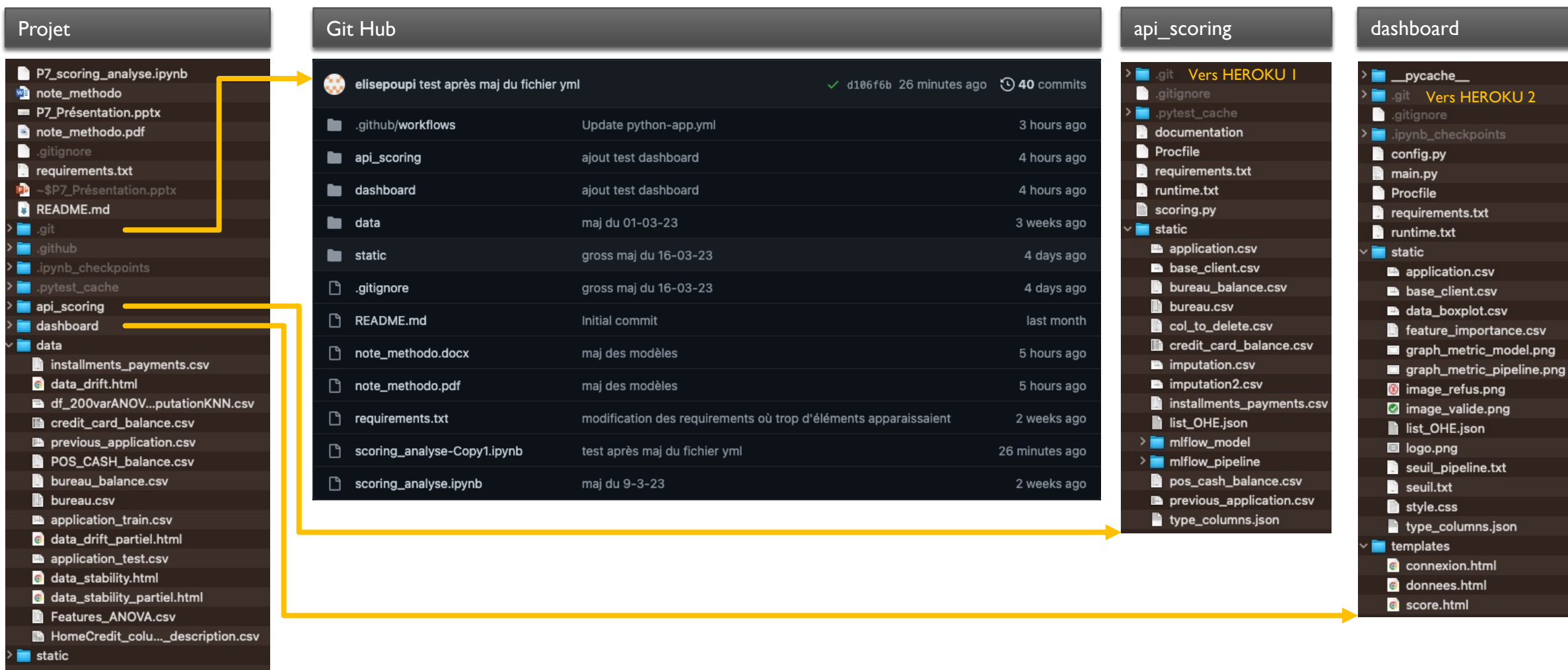
Z-axis:
custom_score

Reverse color:
☐



PIPELINE DE DÉPLOIEMENT ARCHITECTURE

Prêt à dépenser



PIPELINE DE DÉPLOIEMENT GIT

Prêt à dépenser

Visualisation des commits sur GitHub

The screenshot shows the GitHub interface for the repository 'P7_scoring' by user 'elisepoupi'. The repository is marked as 'Private'. At the top, there are buttons for 'Unwatch' (1), 'Fork' (0), and 'Star' (0). Below these are tabs for 'Code', 'Issues', 'Pull requests', 'Actions', 'Projects', and 'Wiki'. A dropdown menu shows the 'main' branch. The commit history is displayed, grouped by date. The most recent commits are from March 20, 2023:

- test après maj du fichier yml** (commit d106f6b) - elisepoupi committed yesterday ✓
- Update python-app.yml** (commit ab5ab83) - elisepoupi committed yesterday ✓ (marked as 'Verified')
- ajout test dashboard** (commit e1f1b4b) - elisepoupi committed yesterday ✓
- déclaration variable df transformer** (commit 3bb9e7b) - elisepoupi committed yesterday ✗
- maj des modèles** (commit 8ad80bd) - elisepoupi committed yesterday ✗

Below these, commits from March 16, 2023 are shown:

- gross maj du 16-03-23** (commit 1cedb61) - elisepoupi committed 5 days ago ✓

Commandes git

Pull

Status

Add

Commit

Push

```
elisepoupinet@MacBook-Air-de-Elise Projet % git pull
remote: Enumerating objects: 9, done.
remote: Counting objects: 100% (9/9), done.
remote: Compressing objects: 100% (3/3), done.
remote: Total 5 (delta 2), reused 0 (delta 0), pack-reused 0
Unpacking objects: 100% (5/5), 827 bytes | 118.00 KiB/s, done.
From https://github.com/elisepoupi/P7_scoring
   e1f1b4b..ab5ab83  main       -> origin/main
Updating e1f1b4b..ab5ab83
Fast-forward
 .github/workflows/python-app.yml | 1 +
 1 file changed, 1 insertion(+)
elisepoupinet@MacBook-Air-de-Elise Projet % git status
On branch main
Your branch is up to date with 'origin/main'.

Changes not staged for commit:
  (use "git add <file>..." to update what will be committed)
  (use "git restore <file>..." to discard changes in working directory)
        modified:   scoring_analyse-Copy1.ipynb

no changes added to commit (use "git add" and/or "git commit -a")
elisepoupinet@MacBook-Air-de-Elise Projet % git add .
elisepoupinet@MacBook-Air-de-Elise Projet % git commit -m "test après maj du fichier yml"
[main d106f6b] test après maj du fichier yml
 1 file changed, 43 insertions(+), 5 deletions(-)
elisepoupinet@MacBook-Air-de-Elise Projet % git push -u origin main
Enumerating objects: 5, done.
Counting objects: 100% (5/5), done.
Delta compression using up to 8 threads
Compressing objects: 100% (3/3), done.
Writing objects: 100% (3/3), 39.98 KiB | 602.00 KiB/s, done.
Total 3 (delta 2), reused 0 (delta 0), pack-reused 0
remote: Resolving deltas: 100% (2/2), completed with 2 local objects.
To https://github.com/elisepoupi/P7_scoring.git
   ab5ab83..d106f6b  main -> main
branch 'main' set up to track 'origin/main'.
elisepoupinet@MacBook-Air-de-Elise Projet %
```

PIPELINE DE DÉPLOIEMENT PYTEST

Prêt à dépenser

En local

```
(LightGBM) elisepoupinet@MacBook-Air-de-Elise api_scoring % pytest scoring.py
===== test session starts =====
platform darwin -- Python 3.9.13, pytest-7.1.2, pluggy-1.0.0
rootdir: /Users/elisepoupinet/Documents/Data_Scientist/P7_scoring/Projet/api_scoring
plugins: anyio-3.5.0
collected 1 item

scoring.py . [100%]

===== 1 passed in 2.70s =====
(LightGBM) elisepoupinet@MacBook-Air-de-Elise api_scoring %
```

Automatisé sur push Git

pytest
python-app.yml

27 workflow runs

- ✓ test après maj du fichier yml
pytest #27: Commit d106f6b pushed by elisepoupi
- ✓ Update python-app.yml
pytest #26: Commit ab5ab83 pushed by elisepoupi
- ✓ ajout test dashboard
pytest #25: Commit e1f1b4b pushed by elisepoupi
- ✗ déclaration variable df transformer
pytest #24: Commit 3bb9e7b pushed by elisepoupi
- ✗ maj des modèles
pytest #23: Commit 8ad80bd pushed by elisepoupi
- ✓ gross maj du 16-03-23
pytest #22: Commit 1cedb61 pushed by elisepoupi
- ✓ maj suite erreur génération csv
pytest #21: Commit f5a9477 pushed by elisepoupi

Résultat test automatisé sur GitHub

```
> ✓ Set up job
> ✓ Run actions/checkout@v3
> ✓ Set up Python 3.10
> ✓ Run actions/cache@v3
> ✓ Install dependencies
✓ Run pytest
  1 ▶ Run pytest ./api_scoring/scoring.py
  8 ===== test session starts =====
  9 platform linux -- Python 3.10.10, pytest-7.1.2, pluggy-1.0.0
 10 rootdir: /home/runner/work/P7_scoring/P7_scoring
 11 collected 1 item
 12
 13 api_scoring/scoring.py . [100%]
 14
 15 ===== 1 passed in 1.96s =====
 16 ===== test session starts =====
 17 platform linux -- Python 3.10.10, pytest-7.1.2, pluggy-1.0.0
 18 rootdir: /home/runner/work/P7_scoring/P7_scoring
 19 collected 1 item
 20
 21 dashboard/main.py . [100%]
 22
 23 ===== 1 passed in 0.63s =====
> ✓ Post Run actions/cache@v3
> ✓ Post Set up Python 3.10
> ✓ Post Run actions/checkout@v3
> ✓ Complete job
```

■ Tests réalisés :

- Api : test de la fonction de prédiction du score à partir d'un individu du df complet
- Dashboard : test de la fonction renvoyant la variable « SK_ID_CURR » à partir du nom du client

DATA DRIFT

Prêt à dépenser

- Data Drift recherché sur le df final :

- Toutes les tables apparaissent 😊
- Les variables catégorielles ont été encodées 😞
- Nombreuses variables _sum, _mean... diluent les résultats 😞

Dataset Drift

Dataset Drift is NOT detected. Dataset drift detection threshold is 0.5

767

Columns

113

Drifted Columns

0.147

Share of Drifted Columns

Data Drift Summary

Drift is detected for 14.733% of columns (113 out of 767).

Search

Column	Type	Reference Distribution	Current Distribution	Data Drift	Stat Test	Drift Score
> BURO_STATUS_C_MEAN_MEAN	num			Detected	Wasserstein distance (normed)	1.221259

- Data Drift détecté si Drift_score > 0.10 sur 50% des variables. Ici drift sur 0.147% des variables. Parmi elles :

- Data Drift sur les tables Bureau et Bureau_balance (provenance d'autres organismes bancaires) > 0.6
- Data Drift sur Payment_Rate (var. importante) : 0.58
- Data Drift sur EXT_SOURCE_1 : 0.16

- Dans l'idéal il faudrait refaire une analyse avec un df finalisé où les tables sont regroupées, sans encodage des variables catégorielles, et sans création de toutes les variables _mean, _sum, _var... pour vérifier la validité des résultats.

DASHBOARD

Prêt à dépenser

■ Connexion

Prêt à dépenser

Formulaire de connexion

Veuillez entrer votre nom :

■ Visualisation des données



Données de Ellaria – id [213818]

Ellaria

encore quelques clics et vous saurez si le prêt demandé pourra vous être accordé!
Veuillez vérifier les données puis cliquer sur 'calculer mon score' pour afficher le résultat:

Crédit

Type de contrat :
Montant du crédit :
Montant de l'annuité :
Montant des biens financés :
Accompagnateur :

...

Document 19 :
Document 20 :
Document 21 :

[Retour à l'accueil](#)

■ Visualisation des résultats



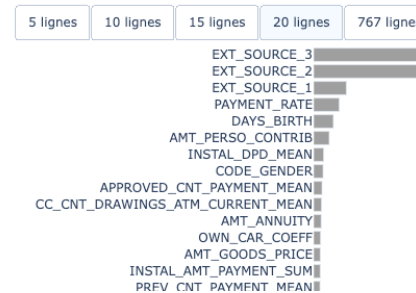
Score – Ellaria – id [213818]

Pour cette modélisation, les prêts sont accordés pour un score inférieur à 0.077.
Votre demande de prêt obtient le score de : 0.112

Le prêt est donc **refusé...**



Variables de décision selon leur importance



MERCI POUR VOTRE ATTENTION

PLACE AUX QUESTIONS

