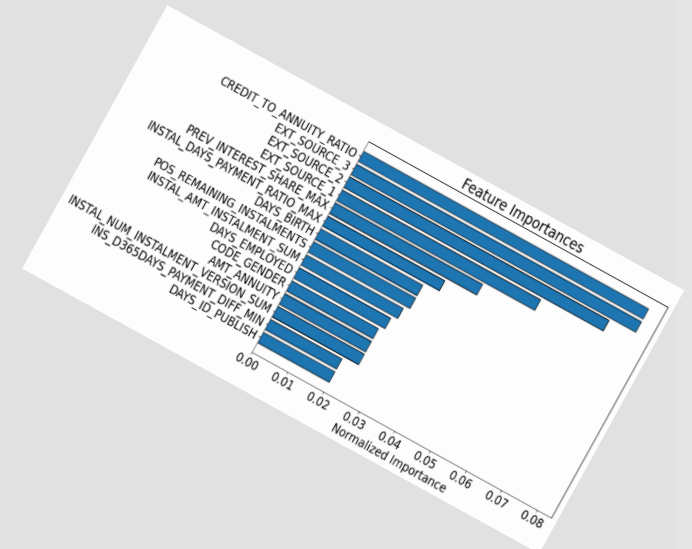
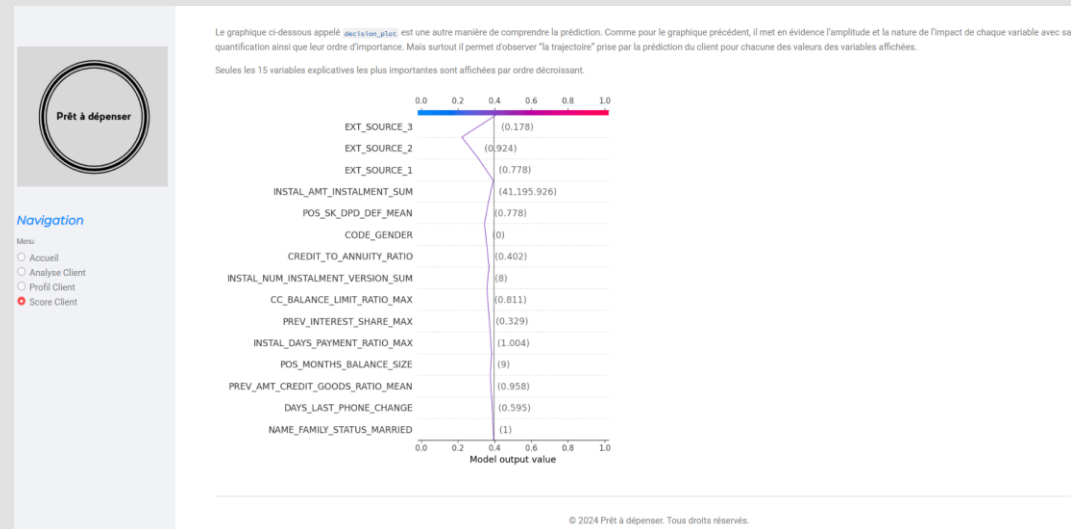
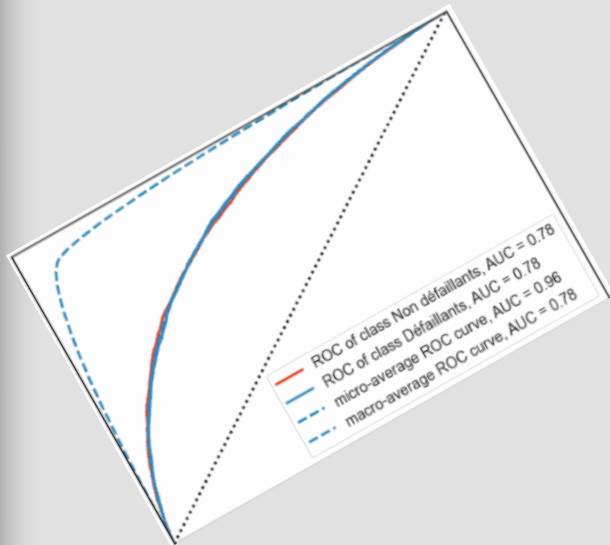


Implémenter un modèle de Scoring

Prêt à dépenser



- 1 Contexte
- 2 Traitement des données
- 3 Modélisation
- 4 Dashboard
- 5 Conclusions

- 1 Contexte
- 2 Traitement des données
- 3 Modélisation
- 4 Dashboard
- 5 Conclusions

Prêt à dépenser:



L'entreprise souhaite créer un outil de "scoring crédit" basé sur un algorithme de classification pour prédire la probabilité de remboursement d'un client et prendre des décisions d'accord ou de refus de crédit.

Mission:

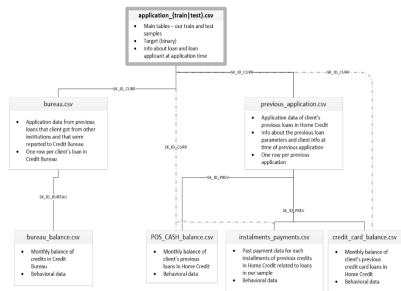
- Construire un **modèle de scoring automatisé** pour prédire la **probabilité de faillite** d'un client.
- Développer un **Dashboard interactif** pour les gestionnaires de la relation client, permettant d'interpréter les prédictions du modèle et d'améliorer leur connaissance des clients.

Objectifs:

Les chargés de relation client ont souligné l'importance de la **transparence** dans les décisions d'octroi de crédit. Pour répondre à cette demande, l'entreprise prévoit de développer un tableau de bord **interactif** permettant d'expliquer les **décisions de crédit** de manière transparente et de donner aux clients un **accès facile** à leurs informations personnelles.

EDA

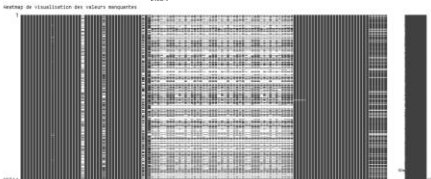
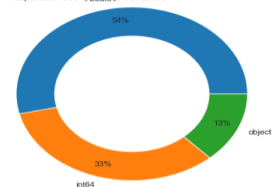
Traitement des 8 fichiers de prêt.



Stat, typage, données manquantes.

Nombre par type de variable	% des types de variable
float64	65
int64	4,2
object	16

Répartition des typage variables



PRE-PROCESSING

Nettoyage des données

Typage, valeurs manquantes, imputation

Feature engineering

- Création nouveaux types de features (min, max, sum, ...)
- Encodage (cat encoding, ...)
- Elimination des fortes colinéarités
- Fusion des datasets nettoyés

Données test et entraînement

FEATURES SELECTION

LightGBM

- Recherche des features importances avec 2 itérations afin d'éliminer les 0 importances.

Permutation importance

- Permutation-based Importance avec sklearn

- Permutation importance avec eli5

Conservation des variables les plus fréquentes

MODELISATION

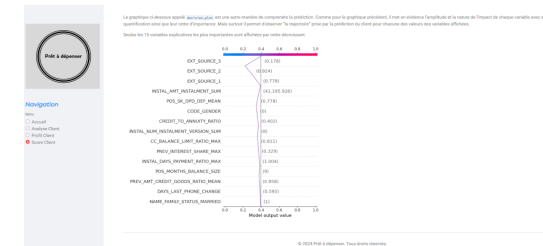


- Séparation du dataset
- Gestion du déséquilibre (SMOTE, ...)
- Sélection des métriques (métrier, ...)
- Optimisation bayésienne des différents modèles (selon les métriques, ...):
- Détermination de la probabilité optimale.

Modèle Final

DASHBOARDING

Développement



Déploiement

- 1 Contexte
- 2 **Traitement des données**
- 3 Modélisation
- 4 Dashboard
- 5 Conclusions

2 Traitement des données – jeu datasets

application_{train|test}.csv

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

application_train.csv
(307511,122)

application_test.csv
(48744,121)

previous_application.csv
(1670214,37)

- Demandes précédentes de prêts immobiliers

POS_CASH_balance.csv
(10001358,8)

- Instantanés mensuels du solde des points de vente précédents et des prêts en espèces

instalments_payments.csv
(13605401,8)

- Historique de remboursement des crédits précédents.

credit_card_balance.csv
(3840312,23)

- Reçu mensuel du solde des cartes de crédits précédents

bureau.csv
(3840312,23)

- Antécédents des crédits des clients

bureau_balance.csv
(3840312,23)

- Soldes mensuels des crédits précédents

8 datasets:
307511 clients et 530
variables

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

bureau_balance.csv

- Monthly balance of credits in Credit Bureau
- 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

POS_CASH_balance.csv

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

instalments_payments.csv

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

credit_card_balance.csv

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

Nettoyage

Typage des données

Valeurs aberrantes, valeurs différentes sur les données train/test

Valeurs manquantes

Imputation

- Numérique ordinal en object
- Optimisation mémoire du dataframe (int64 en int8, ...)

Transformation des valeurs aberrantes de l'EDA

Transformation ou suppression valeurs uniques du train set ≠ test set

Suppression des nan avec conservation des variables importantes de l'EDA

1ère étape
Imputation par la médiane*2ème étape*
Imputation par valeurs nulles 0
Imputation Qualitatives constante XNA*3ème étape*
Imputation NaImputer

Feature engineering

Variables métiers, statistiques

Encodage

Manuel:
Ajout variables métiers*Quantitatives:*
StandardScaler()**Automatique:*
Ajouts variables Statistiques (min, max, sum, ...)*Qualitatives:*
LabelEncoder (binaire)
pd.get_dummies (> 2 valeurs .uniques)

Assemblage

Merge

Fusion des 6 datasets avec les données train/test

Colinéarité

Colinéarité

Suppression des variables fortement colinéaires

Automatique

- Encodage features catégoriques
- Création de features statistiques: quantitatives et qualitatives (min, max, sum, mean, count,...)
- Fonction 'Feature Engineering' intégrant ces caractéristiques sur chaque dataset

Manuel

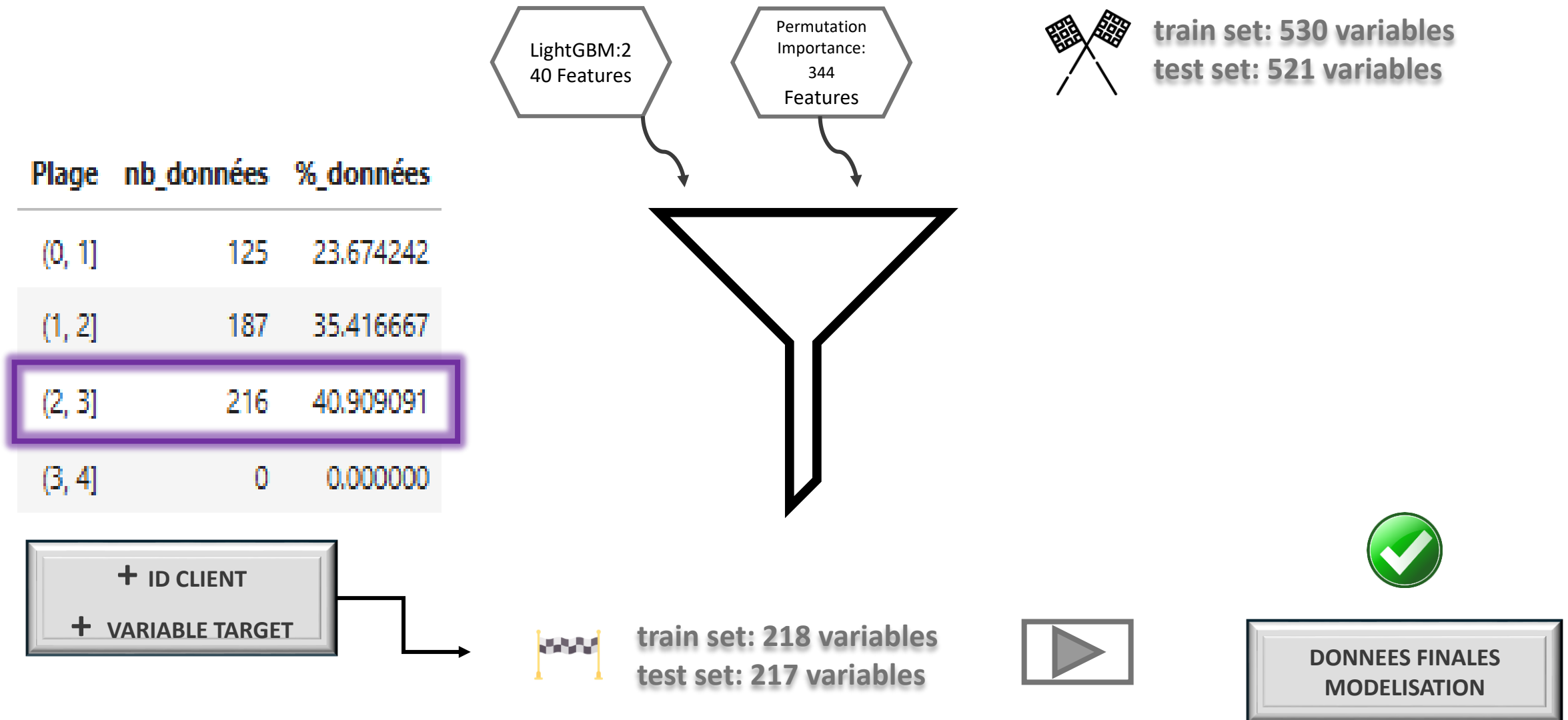
- Ratio/différence des variables métiers
- Conversion jours en années: ratio
- Nouveaux features créés par fusion d'autres features (addition, multiplication, ...)
- Données externes: ratio, moyenne, min, max

Dataset initial	Dimensions initiales	Dimensions Datasets après le process FE	Fusion des datasets, imputation nan et var. colinéaires
application_train/test	(307 511, 122) (48744,121)	(307507, 259) (48744, 258)	
cc_balance	(3840312, 23)	credit_balance_fe (103558, 389)	(307507, 294) (48744, 292)
installments_payments	(13605401,8)	installment_fe (339587, 56)	(307507, 301) (48744, 292)
POS_CASH_balance	(10001358, 8)	pos_cash_fe (337252, 38)	(307507, 322) (48744, 314)
previous_application	(1670214, 37)	previous_application_fe (338857, 258)	(307507, 499) (48744, 488)
bureau_balance	(27299925, 3)	bureau_bb_fe (817395, 12)	(307507, 495) (48744, 485)
bureau	(1716428, 17)	bureau_fe (305811, 60)	(307507, 530) (48744, 521)

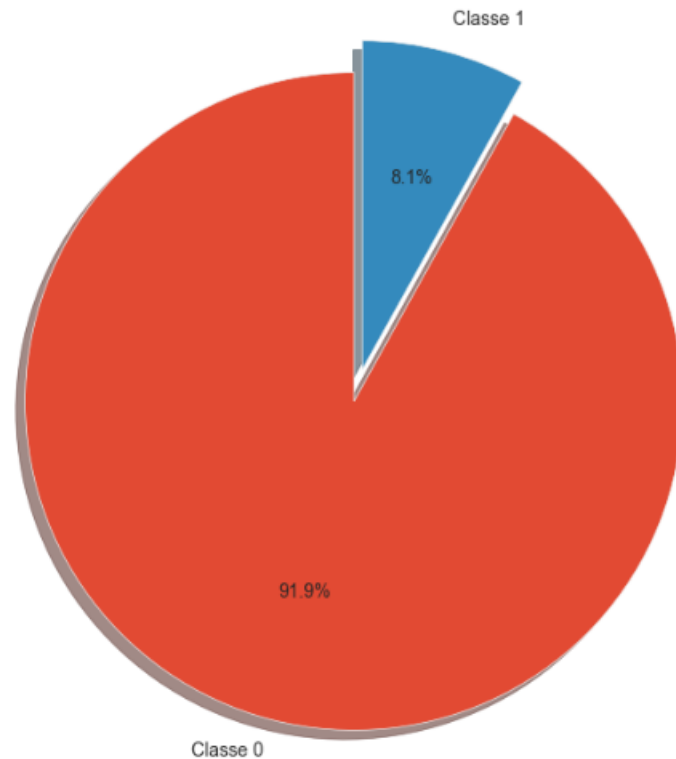
train set: 530 variables
test set: 521 variables



**FEATURE SELECTION
OBLIGATOIRE**



- 1 Contexte
- 2 Traitement des données
- 3 **Modélisation**
- 4 Dashboard
- 5 Conclusions

**Constat variable cible**

Fort déséquilibre entre les clients à risque (Target =1) ou sain (Target = 0)

**Modélisation**

Rééchantillonnage des classes déséquilibrées: SMOTE, ADASYN

Utilisation de plusieurs métriques standards (minimisation du déséquilibre des classes):

Accuracy, ROC, F1 Score, F2 Score, FBeta, Custom metric

Expected	TP: Vrai Positifs	FN: Faux Négatifs
	FP: Faux Positifs	TN: Vrai Négatifs
Predicted		

Précision:

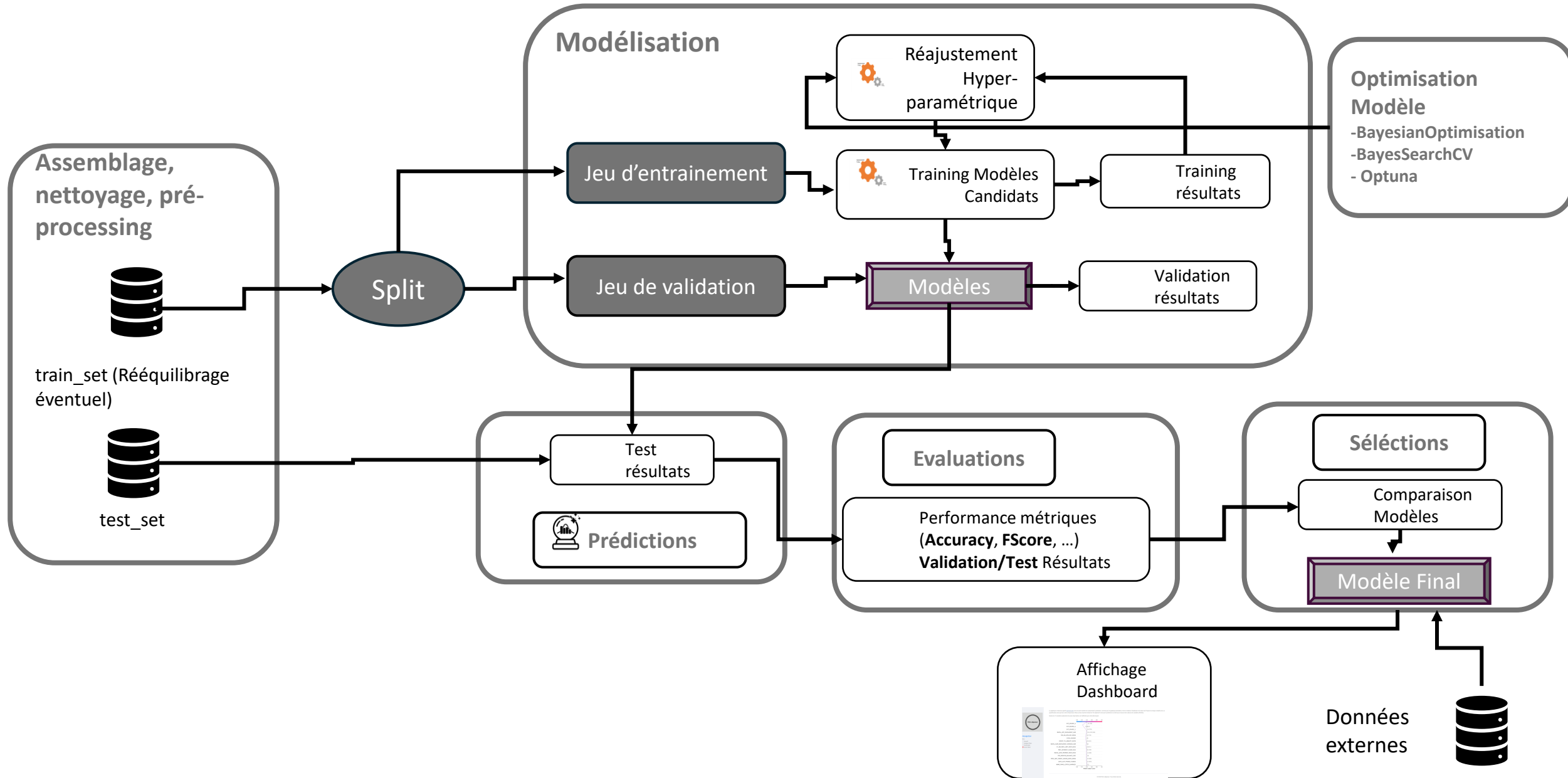
- défaillants = classe **positive**
- non-défaillants = classe **négative**

Minimisation des pertes de bénéfices (prédiction):

- un client non-défaillant s'il est défaillant ==> **minimiser le nombre de faux négatifs (erreur de type II)** (prédit non-défaillant mais client défaillant) ==> **maximiser le Recall**
- le client n'est pas défaillant ==> **minimiser les faux positifs (erreur de type I)** (classe 1 défaillant alors que non-défaillant) ==> **maximiser la Précision**

- Compromis FN/FP (si FN \rightarrow - P \rightarrow + et inversement)

- Test de minimisation fonction cout pour différentes métriques

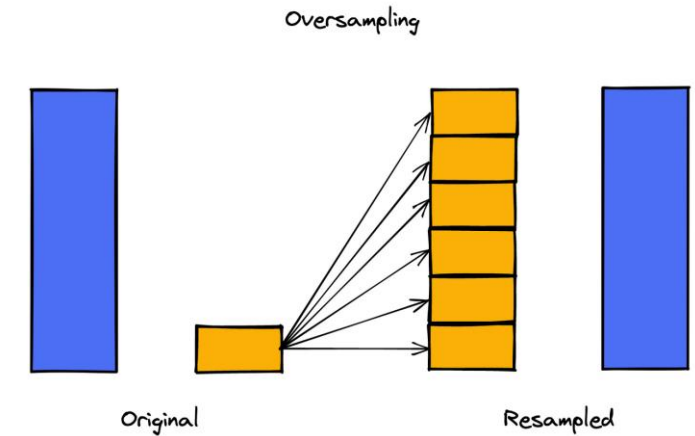
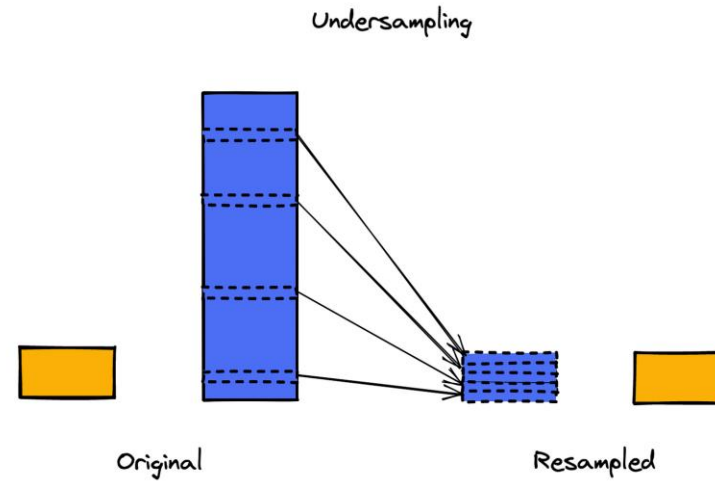


Méthodes utilisées

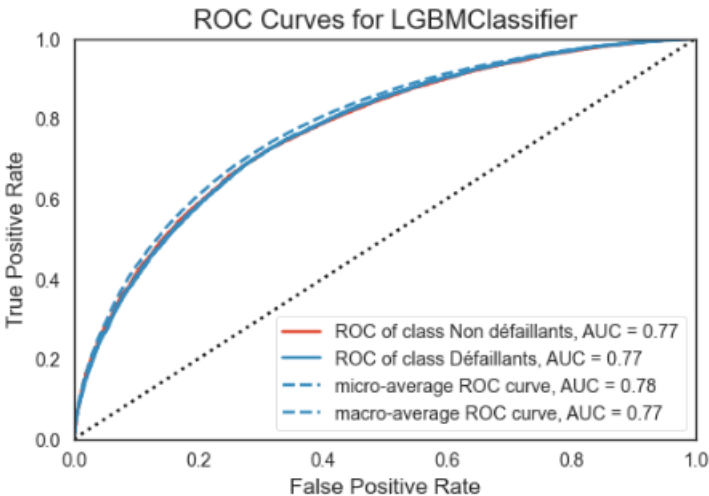
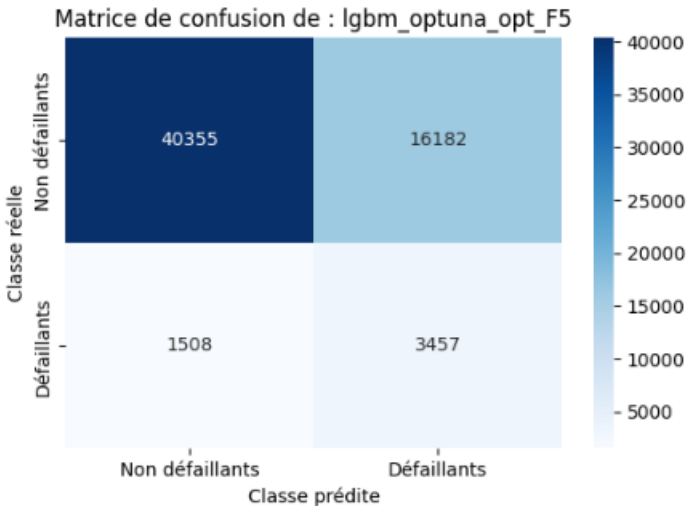
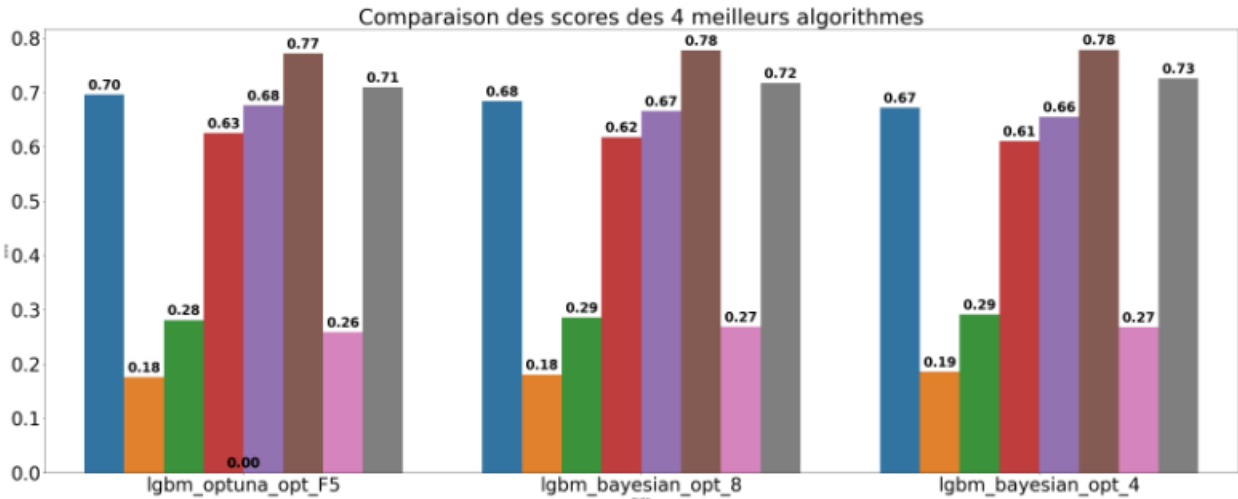
UNDERSAMPLING - *SMOTE*

OVERSAMPLING - *ADASYN*

COMBINAISON OVERSAMPLING et
UNDERSAMPLING - *SMOTE*



Modèle	Jeu_donnees	FN	FP	TP	TN	Metrique	Optimisation	Class_weight	Rappel	Précision	F1	F5	F10	ROC_AUC
lgbm_optuna_opt_F1	train	1500	16369	3465	40168	F1	optuna	oui	0.697885	0.174700	0.279447	0.625803	0.677788	0.770713
lgbm_optuna_opt_F5	train	1508	16182	3457	40355	F5	optuna	non	0.696274	0.176027	0.281011	0.625205	0.676479	0.771812
lgbm_optuna_opt_F10	train	1514	15871	3451	40666	F10	optuna	non	0.695065	0.178605	0.284185	0.625499	0.675720	0.776303
lgbm_optuna_opt_10_train	train	1544	15089	3421	41448	F10	optuna	oui	0.689023	0.184819	0.291459	0.623592	0.670902	0.780011



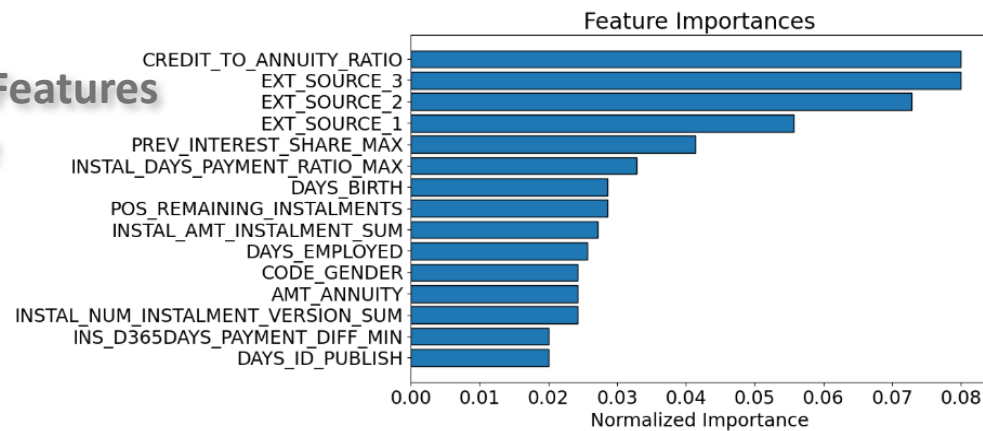
MODELE RETENU:
LGBM_OPTUNA_OPT_F5

MODELE LIGHTGBM:
Meilleurs Paramètres

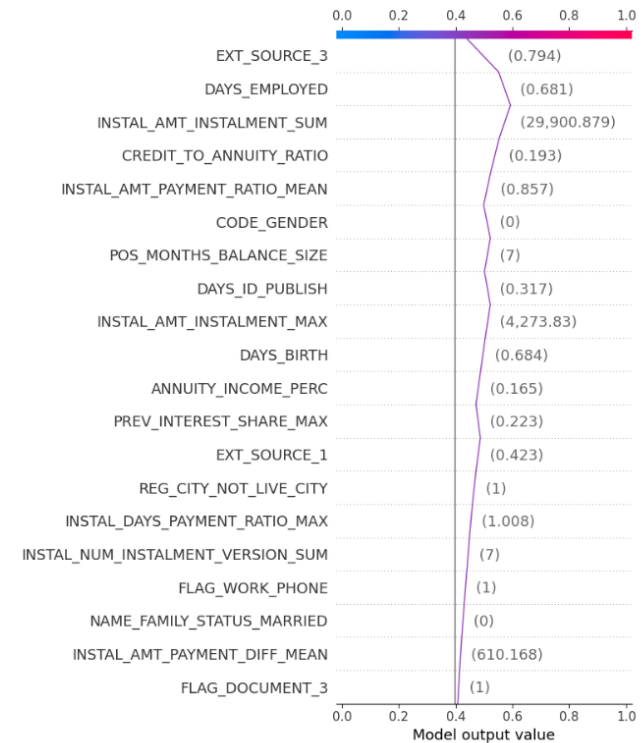
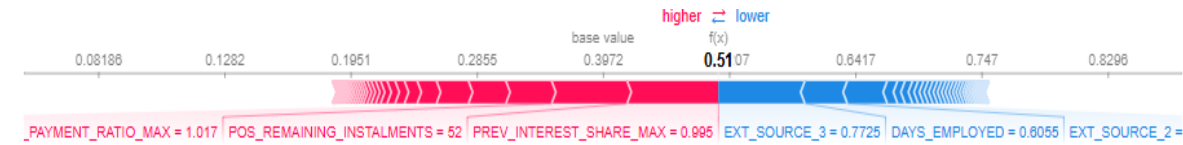
Hyperparamètre	Meilleurs paramètres
n_estimators	100 (par défaut)
learning_rate	0,1 (par défaut)
objective	binary
boosting_type	gbdt
class_weight	balanced
colsample_tree	0.883696173865355
max_depth	4
min_child_samples	37
min_child_weight	0.9053832802852111
num_leaves	8
reg_alpha	0.0013343227256418153
reg_lambda	1.1168060057563535e-06
subsample	0.876335534267455
subsample_freq	4

VUE GLOBALE

LightGBM: Features Importance



VUE CLIENT



- 1 Contexte
- 2 Traitement des données
- 3 Modélisation
- 4 Dashboard**
- 5 Conclusions



Bienvenue sur l'application de Crédit 🏠

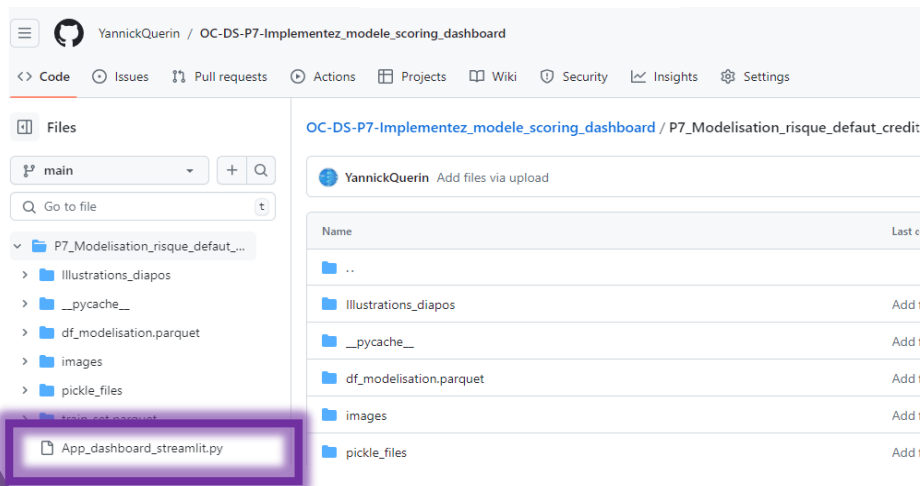
Cette application vous permet d'analyser le risque de défaut de paiement de vos clients, de visualiser le profil de vos clients et de comprendre comment les scores de risque sont calculés. Utilisez le menu à gauche pour naviguer entre les différentes sections.

© 2024 Prêt à dépenser. Tous droits réservés.

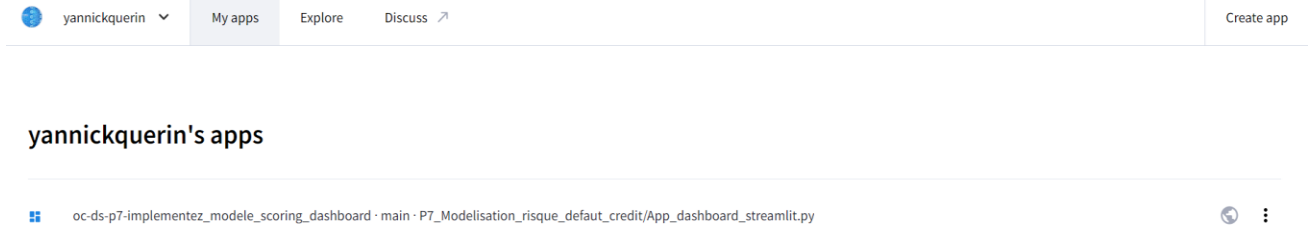
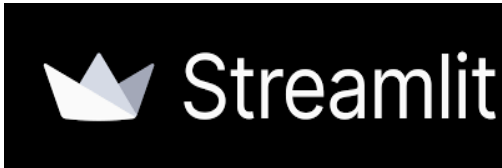
Navigation

Menu

- ☒ Accueil
- ☐ Analyse Client
- ☐ Profil Client
- ☐ Score Client

**Lien Github:**

https://github.com/YannickQuerin/OC-DS-P7-Implementez_modele_scoring_dashboard/blob/main/P7_Modelisation_risque_default_credit/App_dashboard_streamlit.py



Lien local: App_dashboard_streamlit.py

Lien Streamlit:

<https://yannickquerin-p07-dashboard.streamlit.app>



Navigation

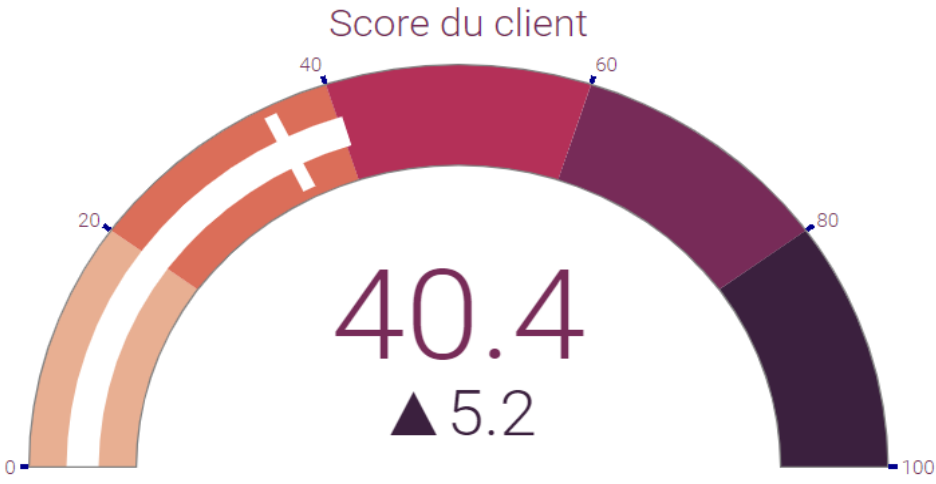
- Menu
- Accueil
 - Analyse Client
 - Profil Client
 - ☒ Score Client

Score Client 🏆

Veuillez sélectionner le numéro de votre client 🗨️

100001

Vous avez sélectionné l'identifiant n° : 100001



Le client dont l'identifiant est **100001** a obtenu le score de **40.4%**.

Il y a donc un risque de 40.4% que le client ait des difficultés de paiement.

Le client est donc considéré par '*Prêt à dépenser*' comme **solvable** et décide de lui **accorder** le crédit.

Prêt à dépenser

Navigation

Menu

- ☐ Accueil
- ☒ Analyse Client
- ☐ Profil Client
- ☐ Score Client

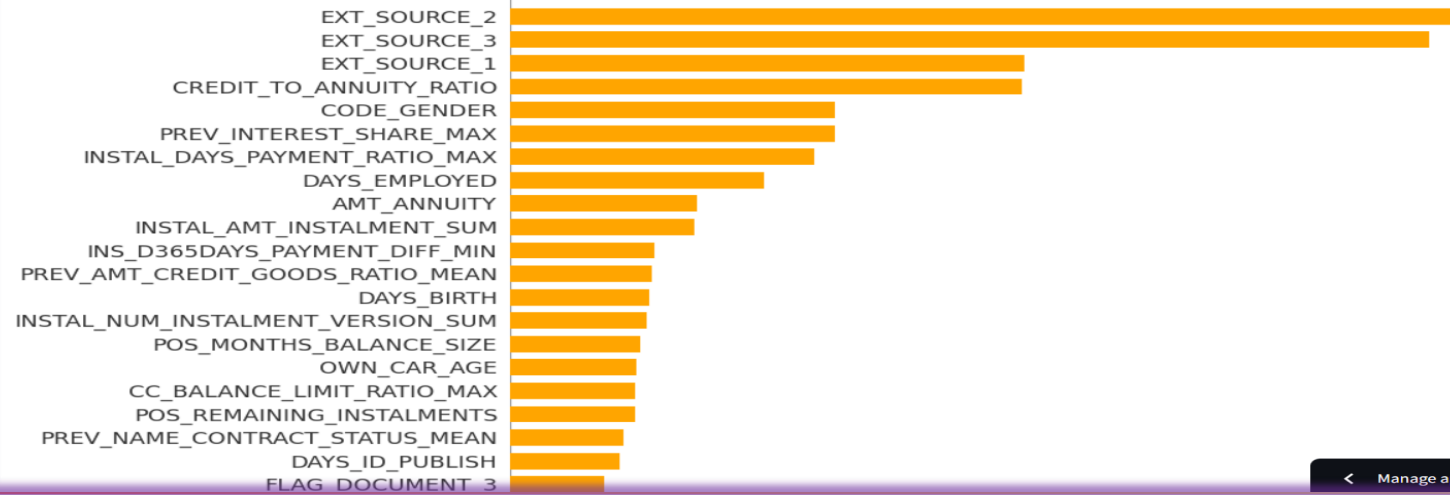
Prêt à dépenser

Navigation

Menu

- ☐ Accueil
- ☒ Analyse Client
- ☐ Profil Client
- ☐ Score Client

Interprétation Globale : Diagramme d'Importance des Variables

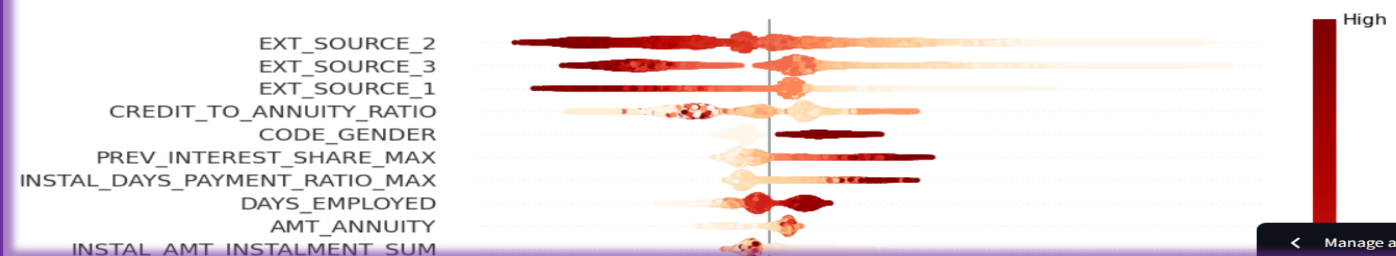


2. Quel est l'Impact de chaque caractéristique sur la prédiction ?

Le diagramme des valeurs SHAP ci-dessous indique également comment chaque caractéristique impacte la prédiction. Les valeurs de Shap sont représentées pour chaque variable dans leur ordre d'importance. Chaque point représente une valeur de Shap (pour un client).

Les points rouges représentent des valeurs élevées de la variable et les points beiges des valeurs basses de la variable.

Interprétation Globale : Impact de chaque caractéristique sur la prédiction



Prêt à dépenser

Profil Client 📄

1. Quel est le profil de votre client ?

Veuillez sélectionner le numéro de votre client à l'aide du menu déroulant 📌

100001

Vous avez sélectionné l'identifiant n° : 100001

Le client dont l'identifiant est **100001** a obtenu le score de **40.4%**.

Il y a donc un risque de 40.4% que le client ait des difficultés de paiement.

Profil socio-économique

Sexe de l'individu : **Féminin**

Situation familiale : **Married**

Nbre enfants : **0**

Niveau éducation : **Higher education**

Type revenu : **Working**

Ancienneté emploi : **6 ANS**

Revenus: **135000.0 \$**

Profil emprunteur

Type de prêt : **Cash loans**

Montant du crédit: **568800.0 \$**

Annuités: **20560.5\$**

Montant du bien: **450000.0 \$**

Type de logement : **House / apartment**

- 1 Contexte
- 2 Traitement des données
- 3 Modélisation
- 4 Dashboard
- 5 **Conclusions**

- Manque de connaissance précise du **secteur bancaire**: vérification au niveau **métier** des **variables retenus** lors du pré-process
- Définir de façon plus détaillée la **métrique d' évaluation** et de **fonction de cout** avec les équipes métiers
- Développer un dashboard avec une page **Données bancaires** coté décideur et une page **Données Clients** de façon à séparer, et sécuriser certaines données sensibles détenus uniquement du côté de la banque, sans pour autant certaines données au client.
- Entrevoir une section de **recommandation** qui permettrait au client de voir quelle variable aurait pu influencer sur son obtention ou pas du prêt en question.