

ESAIP

Rapport Data Engineering & Analysis

Majeur IA



Mathis Herbreteau - Mathias Le Pottier - Samuel Pasquier
12/12/2025

Table des matières

Contexte et objectifs	2
Préparation des données et méthodologie de construction du dataset.....	2
Modélisation et métrique métier	3
Résultats obtenus	3
Conclusion	4
Annexe	4

Contexte et objectifs

Ce projet s'inscrit dans le cadre d'une mission pour une société financière spécialisée dans les crédits à la consommation pour des clients ayant peu ou pas d'historique de prêt. L'objectif principal est de développer un outil complet de Credit Scoring capable de prédire automatiquement la probabilité de défaut d'un client et de classer chaque demande en crédit accordé ou refusé.

Préparation des données et méthodologie de construction du dataset

Pour constituer notre jeu de données, nous avons exploité huit tables relationnelles (données applicatives, historiques des bureaux de crédit, soldes, etc.), fusionnées selon la stratégie suivante :

1. **Prétraitement et Nettoyage** : Nous avons traité les anomalies, telles que des durées d'emploi aberrantes, en les remplaçant par des valeurs nulles (*NaN*).
2. **Ingénierie des fonctionnalités (Feature Engineering)** : Afin de synthétiser l'information financière, quatre indicateurs clés ont été créés :
 - CREDIT_INCOME_PERCENT : ratio crédit/revenu.
 - ANNUITY_INCOME_PERCENT : ratio annuité/revenu.
 - CREDIT_TERM : ratio annuité/crédit.
 - DAYS_EMPLOYED_PERCENT : pourcentage de jours travaillés.
3. **Encodage des variables catégorielle** : Les variables catégorielles ont ensuite été encodées (Label Encoding pour les binaires, One-Hot Encoding pour les autres).
4. **Stratégie d'agrégation** : Compte tenu des relations "un-à-plusieurs" (ex: un client possède plusieurs crédits passés), nous avons appliqué une agrégation en cascade. Les tables de détails (comme *installments_payments*) ont d'abord été groupées par identifiant de prêt (SK_ID_PREV) pour extraire des statistiques (moyenne, somme, min, max). Ces résultats ont ensuite été regroupés par client (SK_ID_CURR) pour obtenir des indicateurs globaux.

Vérification et Gestion des données

Après la fusion, l'intégrité des ensembles d'entraînement (Train) et de test (Test) a été validée pour assurer la cohérence des colonnes. Deux points critiques ont été adressés par la suite :

- **Déséquilibre des classes** : La cible est fortement déséquilibrée (92 % de clients sans défaut de paiement contre 8 % en défaut). Ce biais a été traité par une pondération des classes (`class_weight='balanced'`) lors de l'entraînement.
- **Valeurs manquantes** : Aucune suppression massive n'a été effectuée afin de préserver l'information. Les modèles basés sur des arbres de décision (LightGBM, XGBoost) gérant nativement ces absences, nous avons choisi de conserver ces données en l'état.

Modélisation et métrique métier

Définition de la métrique métier : Afin de disposer d'une évaluation pertinente, nous avons élaboré une fonction de coût personnalisée alignée sur les enjeux business. Celle-ci pénalise dix fois plus un "Faux Négatif" (octroi d'un crédit à défaut, entraînant une perte de capital) qu'un "Faux Positif" (refus d'un client solvable, représentant un manque à gagner).

Cette stratégie vise à minimiser le risque financier global, acceptant une sélectivité accrue au détriment du volume de crédits accordés.

Choix des modèles et Stratégie d'entraînement Deux approches ont été retenues :

1. **Modèle de référence (Baseline) :** Une Régression Logistique a été établie pour fixer un seuil de performance minimal. Ce modèle ne gérant pas les valeurs manquantes nativement, il a été intégré dans un pipeline incluant une étape d'imputation par la médiane.
2. **Modèles avancés (Gradient Boosting) :** Nous avons sélectionné LightGBM et XGBoost. Si XGBoost est réputé pour sa robustesse et sa précision, LightGBM se distingue par sa rapidité d'exécution et son efficacité mémoire sur de grands volumes de données.

Optimisation des Hyperparamètres Pour chaque modèle avancé, nous avons comparé une version avec les paramètres par défaut et une version optimisée. L'optimisation de LightGBM a été réalisée via une recherche bayésienne avec le framework Optuna. Après 50 itérations (*trials*), la configuration optimale retenue est la suivante :

Résultats obtenus

Performance du meilleur Modèle (LightGBM) Le modèle final sur le jeu de validation affiche les performances suivantes :

Un ROC AUC de 0.786, ce qui nous indique que le modèle a une bonne capacité à discriminer les instances positives et négatives, notamment en prenant en compte les faux positifs et les faux négatifs ainsi qu'un seuil de 0.51.

Nous avons obtenu ce seuil après avoir analysé le score métier obtenu pour tous les seuils entre 0 et 1 avec un pas de 0.01. On retrouve ~0.5 car nos classes sont équilibrées.

Pour la métrique `f1_score` nous obtenons 0.3, ce résultat montre que notre modèle a tendance ne pas accorder de crédit même s'il le pouvait car il ne veut pas prendre de risque du au fort coût d'un non remboursement.

On a également un recall de 0.68 donc notre modèle reconnaît 68% des bons clients.

La dernière métrique importante est la précision pour laquelle on obtient 0.19 ce qui est faible. La signification est que sur tous les mauvais payeurs identifiés seulement 19% sont réellement des mauvais payeurs.

Conclusion

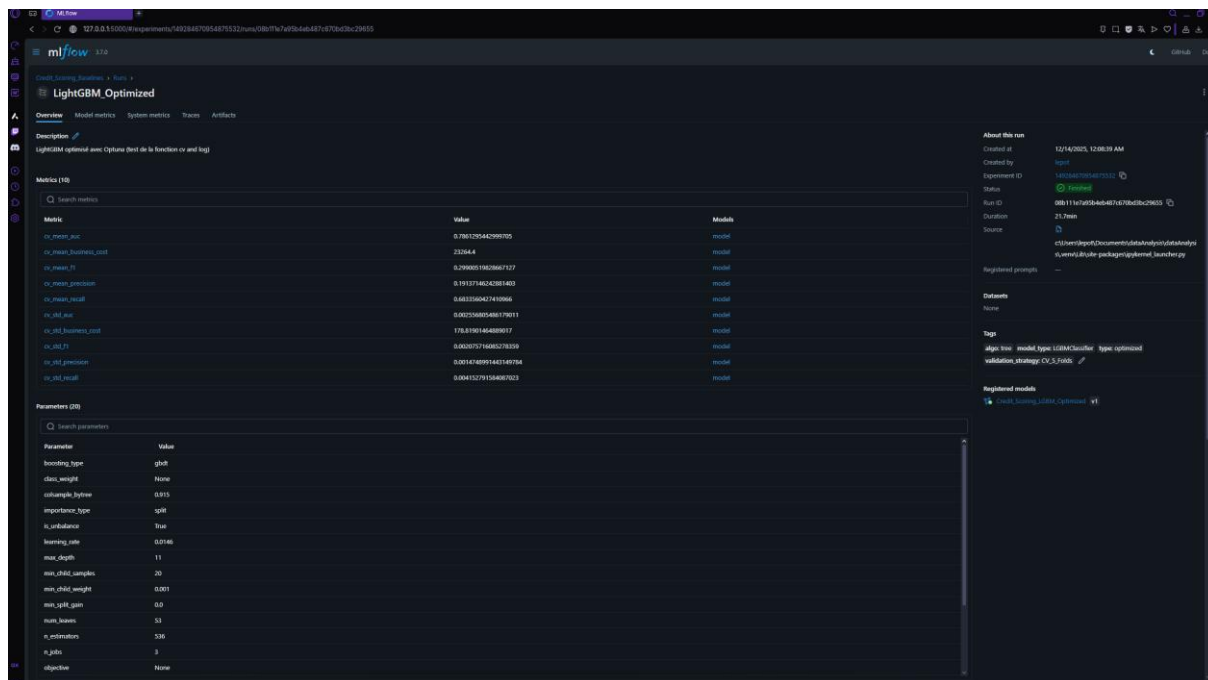
Ce projet a permis de développer un outil de Credit Scoring robuste, reposant sur l'algorithme LightGBM optimisé. Avec une ROC AUC de 0.786, le modèle démontre une capacité satisfaisante à discriminer les clients solvables des profils à risque.

L'alignement avec la stratégie financière de l'entreprise est respecté : en pénalisant fortement les défauts de paiement (coût x10), le modèle adopte un comportement prudent.

- Sécurité du capital : Le seuil de décision de 0.51 (cohérent avec l'équilibrage des classes) permet de sécuriser les actifs en évitant les profils douteux.
- Le compromis Risque/Volume : Cette prudence se traduit par un Recall de 0.68 (bonne détection des cibles) mais une Précision de 0.19. Cela indique que pour éviter les impayés, l'entreprise accepte de refuser un certain nombre de clients potentiellement solvables (manque à gagner) afin de garantir un taux de défaut minimal.

En somme, l'outil déployé répond efficacement à l'objectif prioritaire de minimisation du risque financier, fournissant une aide à la décision fiable pour l'octroi de crédits.

Annexe



Parameter	Value
boosting_type	gbdt
class_weight	None
colsample_bytree	0.915
importance_type	split
is_unbalance	True
learning_rate	0.0146
max_depth	11
min_child_samples	20
min_child_weight	0.001
min_split_gain	0.0
num_leaves	53
n_estimators	536
n_jobs	3
objective	None

Metric	Value	Model
cv_auc	0.7861295442999705	model
cv_auc_bounded_loss	23264.6	model
cv_auc_f1	0.29980519623667327	model
cv_auc_precision	0.19137146243261403	model
cv_auc_recall	0.682358427410966	model
cv_auc_auc	0.00235885488179011	model
cv_auc_bounded_loss	178.81961464889857	model
cv_auc_f1	0.00207514685278339	model
cv_auc_precision	0.0014748991542146784	model
cv_auc_recall	0.004152791584867623	model