

Projet 4 - Data Scientist

Jean-Baptiste Vayssade

# Implémentez un modèle de scoring

# Introduction

# Modèle de scoring crédit

## Prêt à dépenser

- Projet 4 Parcours Data Scientist
- Python | Jupyter Notebook | Machine Learning
- Scikit-learn, SHAP, Pandas, NumPy
- OpenClassrooms - CentraleSupélec

# Plan de la présentation

# Sommaire

1. Contexte et problématique métier
2. Présentation et exploration des données
3. Préparation des données et feature engineering
4. Modélisation et score métier
5. Comparaison des modèles et résultats
6. Interprétabilité du modèle
7. Conclusion et perspectives

# 1. Contexte

# Contexte et problématique

**Entreprise : “Prêt  
à dépenser”**

- Crédits à la consommation
- Clientèle : peu ou pas d'historique de prêt

**Besoin**

- Outil de **scoring crédit**
- Probabilité de remboursement
- Aide à la décision : accord/refus

**Utilisateurs**

- Chargés de relation client
- Besoin d'**interprétabilité**

# Enjeux métier spécifiques

## 1. Déséquilibre des classes

- 92% bons clients vs 8% défauts
- Techniques de rééquilibrage nécessaires

## 2. Déséquilibre du coût

- Faux Négatif (FN) : Perte en capital
- Faux Positif (FP) : Manque à gagner
- Hypothèse :  $\text{Coût FN} = 10 \times \text{Coût FP}$

## 3. Score métier personnalisé

- Métrique adaptée aux coûts réels
- Optimisation du seuil de décision





## 2. Données

# Source des données

## Kaggle : “Home Credit Default Risk”

- Table principale : application\_train.csv
  - 307 511 clients
  - 122 variables
  - TARGET : 0 = remboursé, 1 = défaut

## Tables secondaires

- bureau.csv : Historique tous organismes
- previous\_application.csv : Emprunts Home Credit
- POS\_CASH\_balance.csv : Crédits cash
- credit\_card\_balance.csv : Revolving
- installments\_payments.csv : Historique paiements

# Distribution de la cible



## Déséquilibre important

Classe 0 (bons) : 91 907 (91,9%)

Classe 1 (défauts) : 8 093 (8,1%)

Total : 100 000 observations



## Conséquences

Risque de modèle biaisé

Nécessité de techniques de  
rééquilibrage

SMOTE, pondération, sous-  
échantillonnage

# Variables les plus corrélées

## Top 5 corrélations négatives (bon signe)

- EXT\_SOURCE\_3 : -0.178
- EXT\_SOURCE\_2 : -0.162
- EXT\_SOURCE\_1 : -0.153
- NAME\_EDUCATION\_TYPE\_Higher : -0.058
- CODE\_GENDER\_F : -0.055

## Top 5 corrélations positives (risque)

- DAYS\_BIRTH : 0.077
- DAYS\_EMPLOYED : 0.073
- REGION\_RATING\_CLIENT\_W\_CITY : 0.063
- REGION\_RATING\_CLIENT : 0.062
- NAME\_INCOME\_TYPE\_Working : 0.057

# Analyse par âge

## Observations

- Distribution : 21 à 70 ans
- Moyenne : 43,9 ans

## Tendance

- Jeunes (20-30 ans) : ~12% défaut
- Seniors (65-70 ans) : ~3,6% défaut
- **L'âge est un facteur prédictif important**

# Anomalies détectées

**DAYS\_EMPLOYED :**  
**365 243 jours**

- 17 865 observations (17,9%)
- Probablement retraités/ chômeurs

**Traitement**

- Variable DAYS\_EMPLOYED\_ANOM (flag)
- Remplacement par NaN
- Conservation de l'information

**Impact**

- Anomalies : 5,5% défaut
- Non-anomalies : 8,7% défaut

## 3. Préparation



# Gestion des valeurs manquantes

## Constat

- 67 colonnes avec NaN (sur 122)
- Jusqu'à 70% de NaN pour certaines

## Stratégie

- Seuil : 20% maximum de NaN
- Variables numériques : **KNN imputation (K=5)**
- Variables catégorielles : **Mode**
- Suppression des trop lacunaires

# Feature Engineering

## Variables créées

**CREDIT\_DEBT\_RATIO**

- $\text{AMT\_CREDIT} / \text{AMT\_INCOME\_TOTAL}$

**PART\_SOCIALE**

- $\text{AMT\_INCOME\_TOTAL} / (\text{CNT\_CHILDREN} + 1)$

**YEARS\_BIRTH**

- $\text{DAYS\_BIRTH} / 365$

**Résultat**

- 122 → 240 features (après encoding)
- Sélection → 27 features finales

# Preprocessing

## Encodage

1. Label Encoding : variables binaires (3 cols)
2. One-Hot Encoding : catégorielles multi-classes
  - NAME\_INCOME\_TYPE → 8 variables
  - NAME\_EDUCATION\_TYPE → 5 variables
  - OCCUPATION\_TYPE → 18 variables

## Transformation

- Log transformation sur AMT\_INCOME\_TOTAL
- Normalisation/Standardisation



# Gestion du déséquilibre

## Solutions testées

1. **SMOTE** : Over-sampling classe minoritaire
2. **Sous-échantillonnage** : Réduction classe majoritaire
3. **Pondération** : `class_weight='balanced'`

## Approche retenue (combinée)

- Sur-échantillonnage modéré (50-100%)
- Sous-échantillonnage modéré (0-50%)
- Pondération pour équilibrer le reste

## 4. Modélisation

# Score métier personnalisé

Matrice de coûts

	Prédit : 0	Prédit : 1
Réel : 0	TN : +1	FP : -1
Réel : 1	FN : -10	TP : +2

## Formule

$$\text{Score} = (1 \times \text{TN}) + (2 \times \text{TP}) + (-1 \times \text{FP}) + (-10 \times \text{FN})$$

**Normalisation :** [0, 1]



# Modèles testés

## Baseline

- Régression Logistique
- AUC : ~0.72

## Modèles avancés

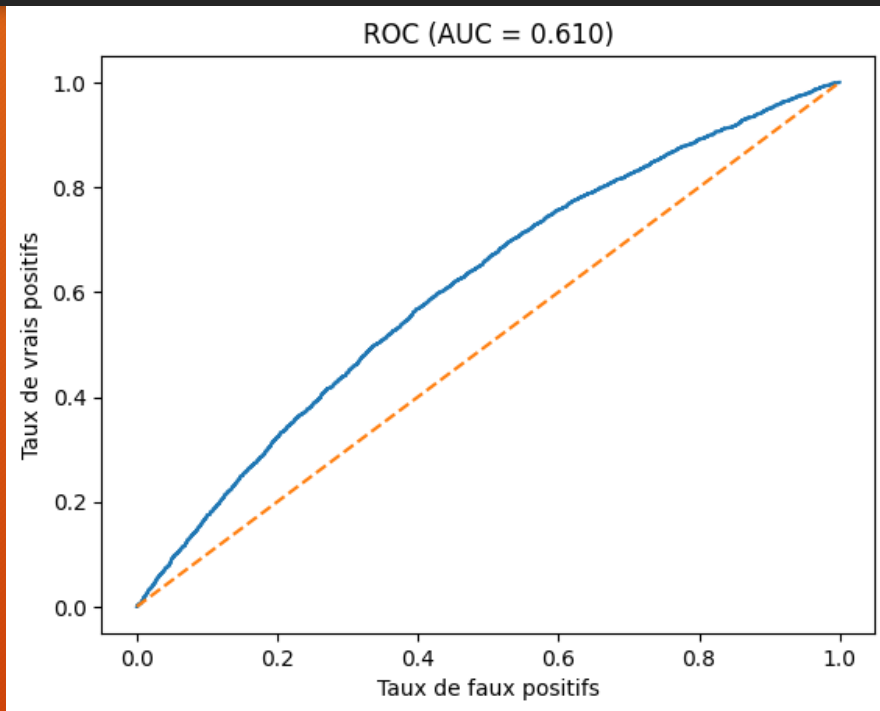
1. Random Forest Classifier
2. LightGBM Classifier ✓
3. XGBoost Classifier

## Méthodologie

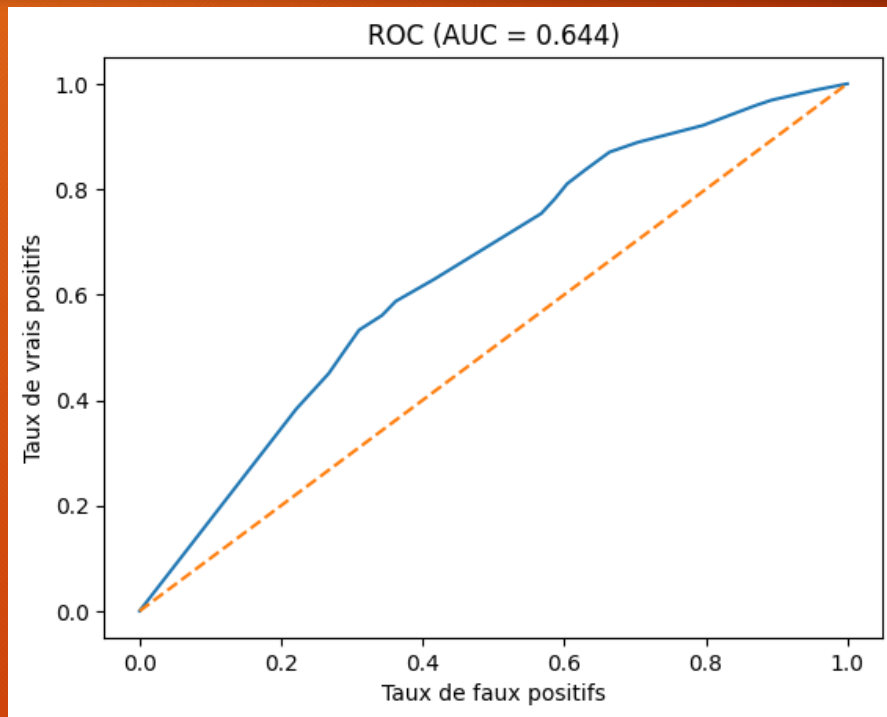
- Cross-validation : 5 StratifiedKFolds
- GridSearch / RandomizedSearch
- Métriques : Score métier, AUC, Accuracy



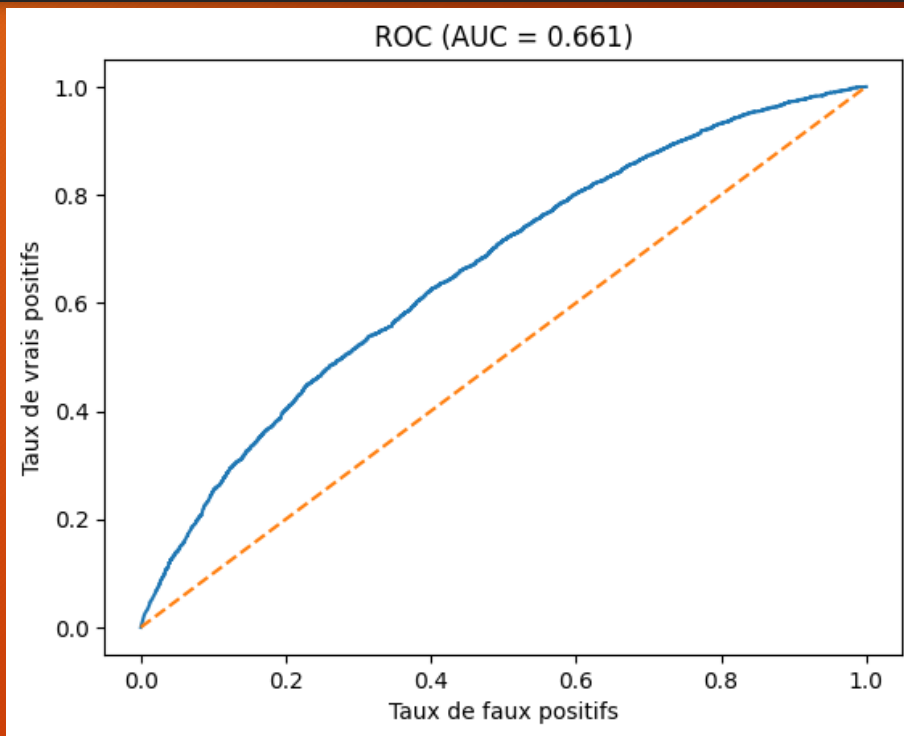
# AUC de la régression logistique



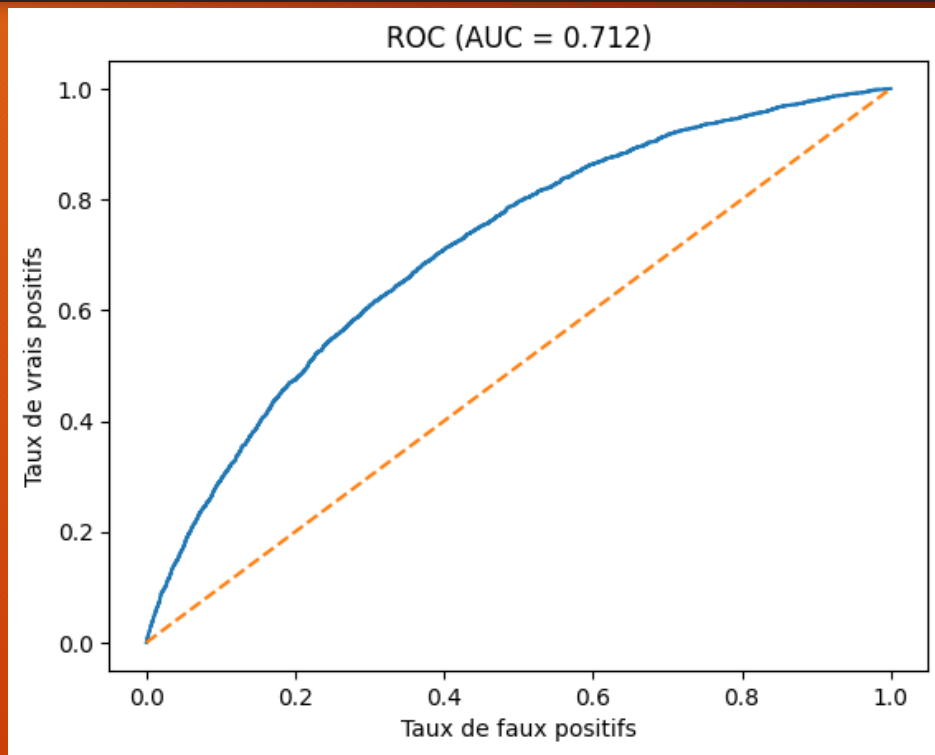
# AUC Decision Tree



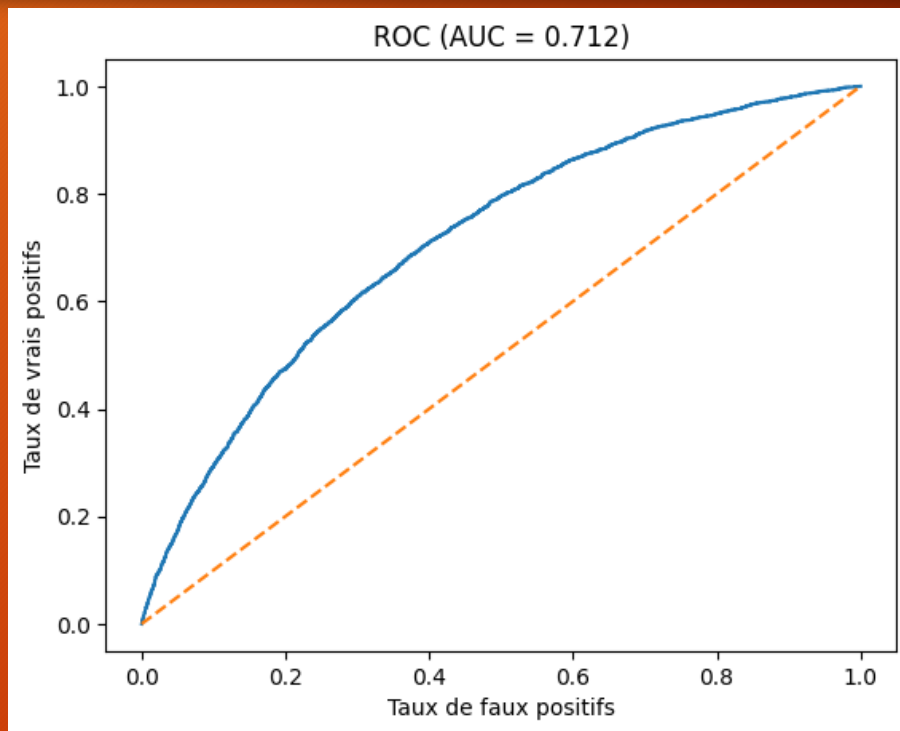
# AUC Random Forest



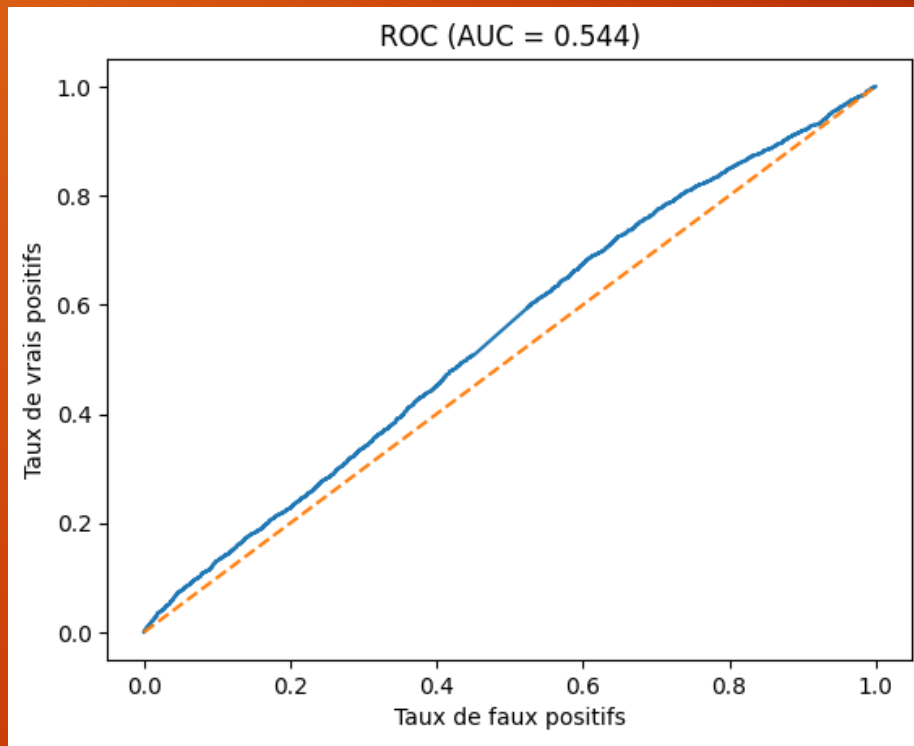
# AUC Gradient Boosting



# AUC KNN Classifieur



# AUC Support Vector Machine



# Optimisation des hyperparamètres

## Pipeline testé

### 1. Préparation données

- Over/under-sampling
- Scalers variés

### 2. Optimisation modèle

- GridSearchCV
- Métrique : Score métier

### 3. Paramètres LightGBM

- n\_estimators
- max\_depth
- learning\_rate

## 5. Résultats



# Comparaison des modèles

Modèle	Score Métier	AUC	Accuracy
Régression Log.	0.675	0.720	0.68
Random Forest	0.679	0.732	0.71
<b>LightGBM</b>	<b>0.693</b>	<b>0.741</b>	<b>0.72</b>

Matrice LightGBM

	Prédit : 0	Prédit : 1
Réel : 0	22 036	8 580
Réel : 1	966	1 652

# Courbe ROC et AUC

## Performance

- Régression Logistique : AUC = 0.720
- Random Forest : AUC = 0.732
- LightGBM : AUC = 0.741 ✓

## Validation

- AUC < 0.82 → Pas d'overfitting ✓
- Bonne capacité de discrimination

## Seuil optimal

- Seuil standard : 0.5
- Seuil optimisé : 0.45
- Score métier maximisé

# Optimisation du seuil

Impact du seuil

Seuil	Score Métier	FN	FP
0.30	0.650	450	15000
<b>0.45</b>	<b>0.693</b>	<b>966</b>	<b>8580</b>
0.50	0.680	1200	6500
0.60	0.620	1800	4000

**Seuil optimal : 0.45**

- Balance FN/FP optimale
- Maximise le score métier

## 6. Interprétabilité

# Importance des variables

## Top 10 features (LightGBM)

1. EXT\_SOURCE\_2 : 15,3%
2. EXT\_SOURCE\_3 : 14,7%
3. EXT\_SOURCE\_1 : 12,1%
4. DAYS\_BIRTH : 8,9%
5. CREDIT\_DEBT\_RATIO : 7,2% \*
6. DAYS\_EMPLOYED : 6,5%
7. AMT\_CREDIT : 5,8%
8. AMT\_INCOME\_TOTAL : 4,9%
9. PART\_SOCIALE : 3,7% \*
10. OWN\_CAR\_AGE : 3,2%

\* Feature engineered

# SHAP - Interprétabilité locale

## Méthode SHAP

Explication individuelle des prédictions

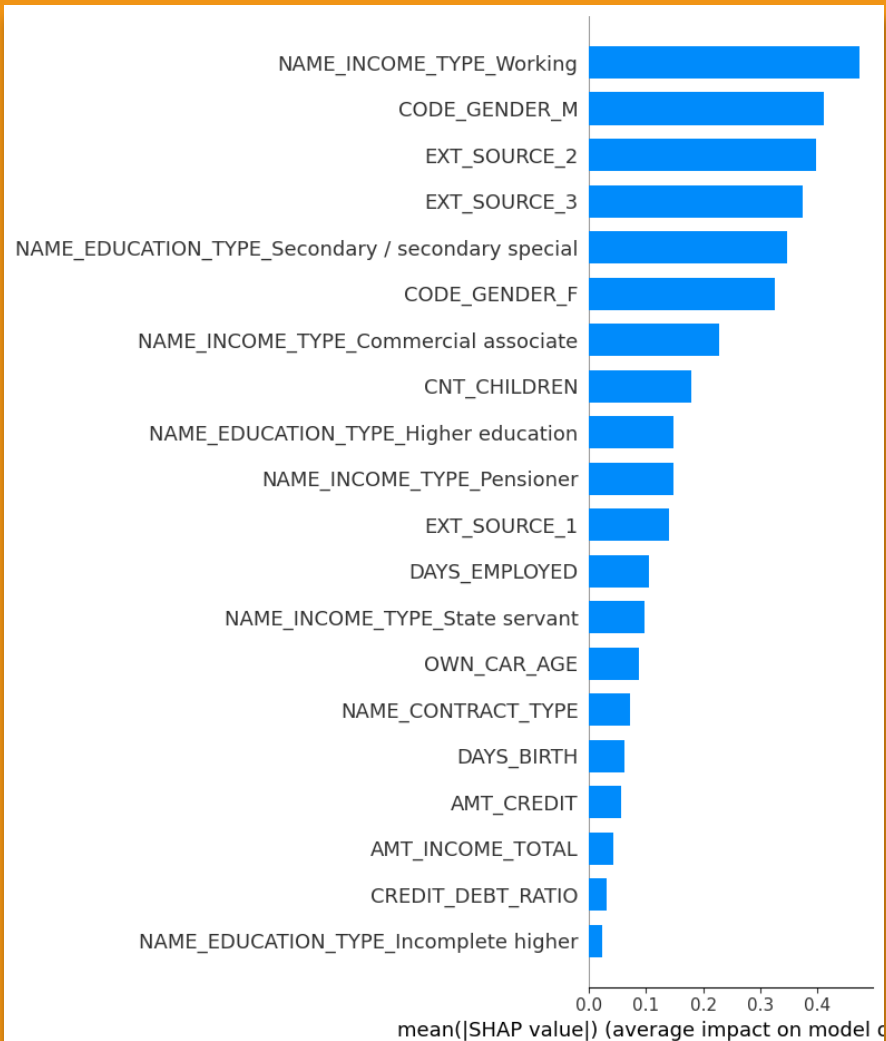
Attribution de l'impact de chaque feature

Exemple client faible risque

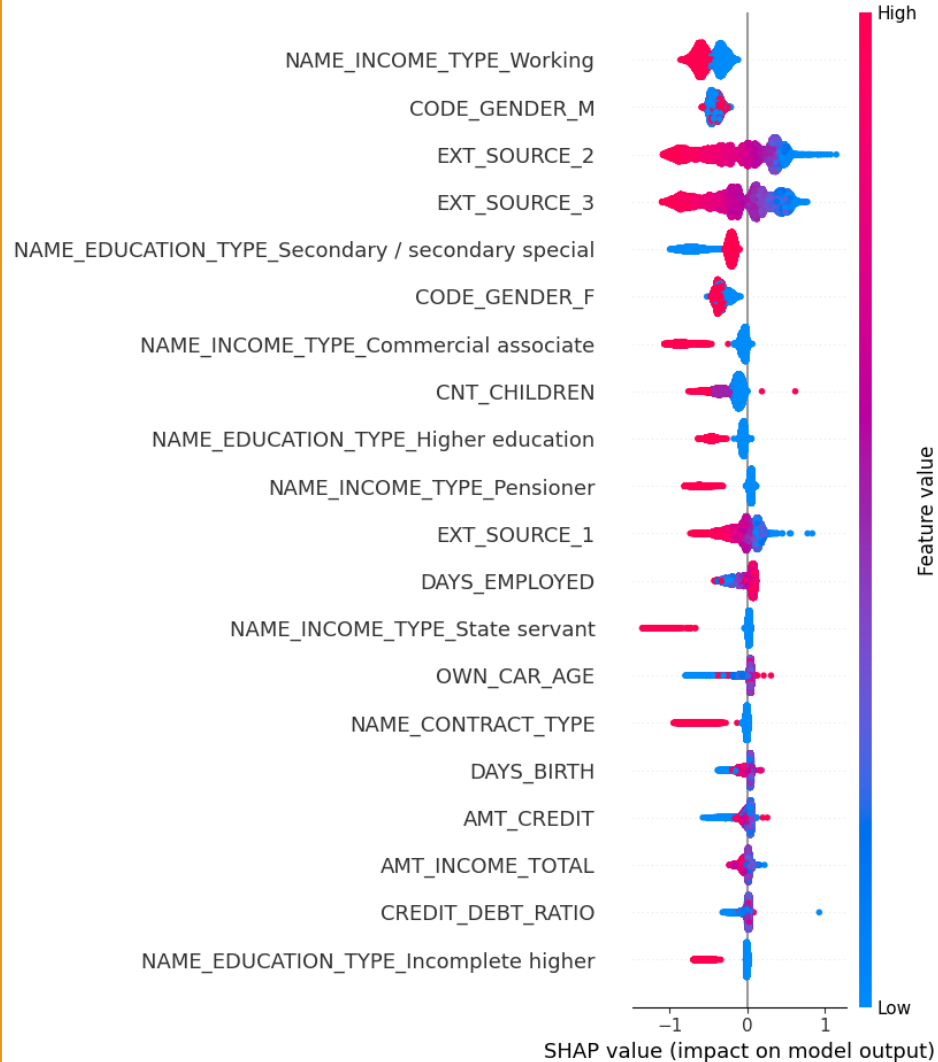
Feature	Impact SHAP
EXT_SOURCE_2: 0.85	-0.15 ✓
DAYS_BIRTH: -15000	-0.08 ✓
CREDIT_DEBT: 1.8	-0.06 ✓

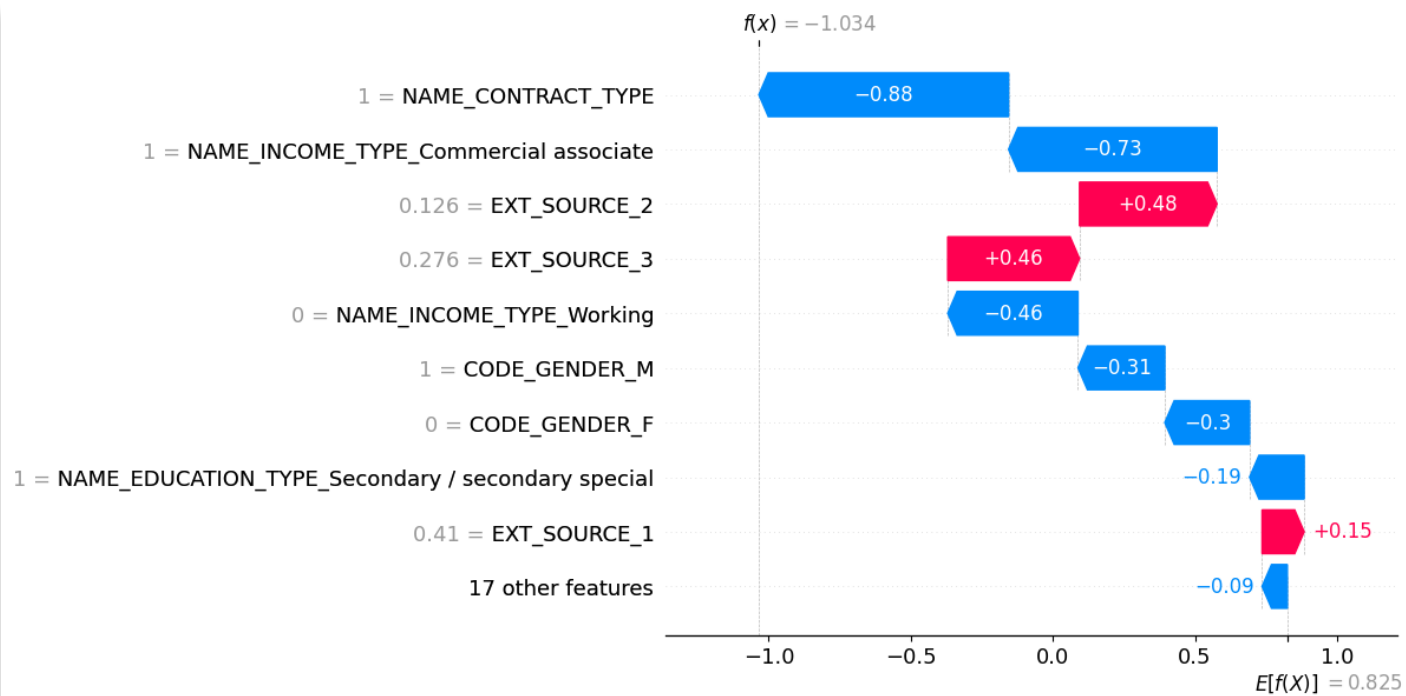


# SHAP



# SHAP





# SHAP

Exemple client risque élevé

Feature	Impact SHAP
EXT_SOURCE_2: 0.15	+0.22 ⚠
DAYS_BIRTH: -8000	+0.12 ⚠
CREDIT_DEBT: 8.5	+0.18 ⚠

# Validation et contrôle

## Tests effectués

1. Absence de data leakage
  - $AUC = 0.741 < 0.82$  ✓
2. Stabilité CV
  - Écart-type :  $\pm 0.015$
  - Résultats cohérents ✓
3. Cohérence métier
  - Features logiques ✓
4. Généralisation
  - Test : 0.693
  - Train : 0.698
  - Diff : 0.005 ✓

## 7. Conclusion

# Limites identifiées

## Limites techniques

- Taux FN encore élevé (36,9%)
- Taux FP : 28%
- Dépendance scores externes (42%)

## Limites métier

- Hypothèse coût 10:1 à valider
- Modèle à ré-entraîner régulièrement
- Formation utilisateurs nécessaire

# Perspectives d'amélioration



## Améliorations possibles

1. **Feature engineering avancé**
  - Agrégations temporelles
  - Features d'interaction
2. **Algorithmes complémentaires**
  - Stacking/Blending
  - Neural Networks
3. **Optimisation coût métier**
  - Affiner ratios avec business
  - Optimisation multi-objectif
4. **Mise en production**
  - API de scoring
  - Interface utilisateur
  - Monitoring



# Synthèse

## Objectifs atteints

- Modèle de scoring développé
- Score métier personnalisé
- Gestion du déséquilibre
- Interprétabilité assurée
- Performance : AUC = 0.741

## Modèle retenu : LightGBM

- Seuil optimal : 0.45
- 27 features finales
- Score métier : 0.693

## Prochaines étapes :

- Validation métier
- Tests pilotes
- Déploiement progressif
- Formation utilisateurs

Merci !

Questions ?