

# Construction d'un modèle de scoring

Prédire la solvabilité d'un client

Projet 4 - Sylwia Bankowska - Avril 2022



# Sommaire

1. Compréhension de la problématique métier
2. Description du jeu de données
3. Transformation du jeu de donnée
4. Comparaison et synthèse des résultats pour les modèles utilisés
5. Interprétabilité du modèle
6. Conclusion

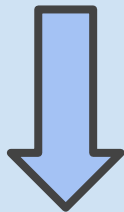
# Problématique métier

---

# Objectifs du projet

**Développer**

**un algorithme de classification**



pour identifier les clients solvables ou non

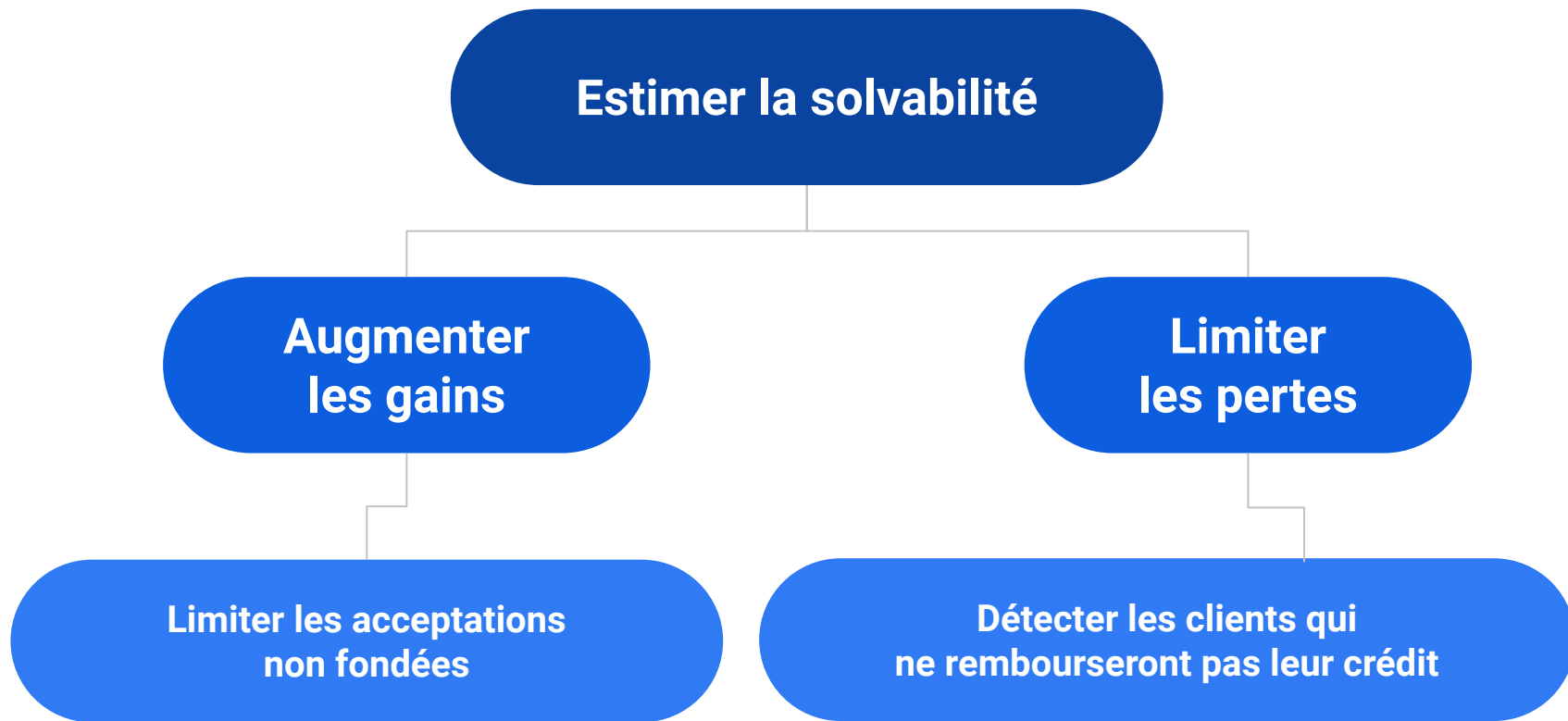
**Fournir**

**une mesure de l'importance  
des variables**



qui ont poussé le modèle  
à donner une probabilité à un client

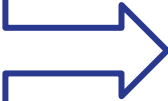
# Enjeux



# Matrice de confusion

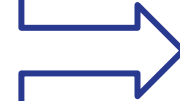
	CLIENT A RISQUE	CLIENT SOLVABLE
REFUS DE CRÉDIT	True Positive	False Positive
ACCEPTATION DE CRÉDIT	False Negative	True Negative

Limiter les  
pertes



TP / FN

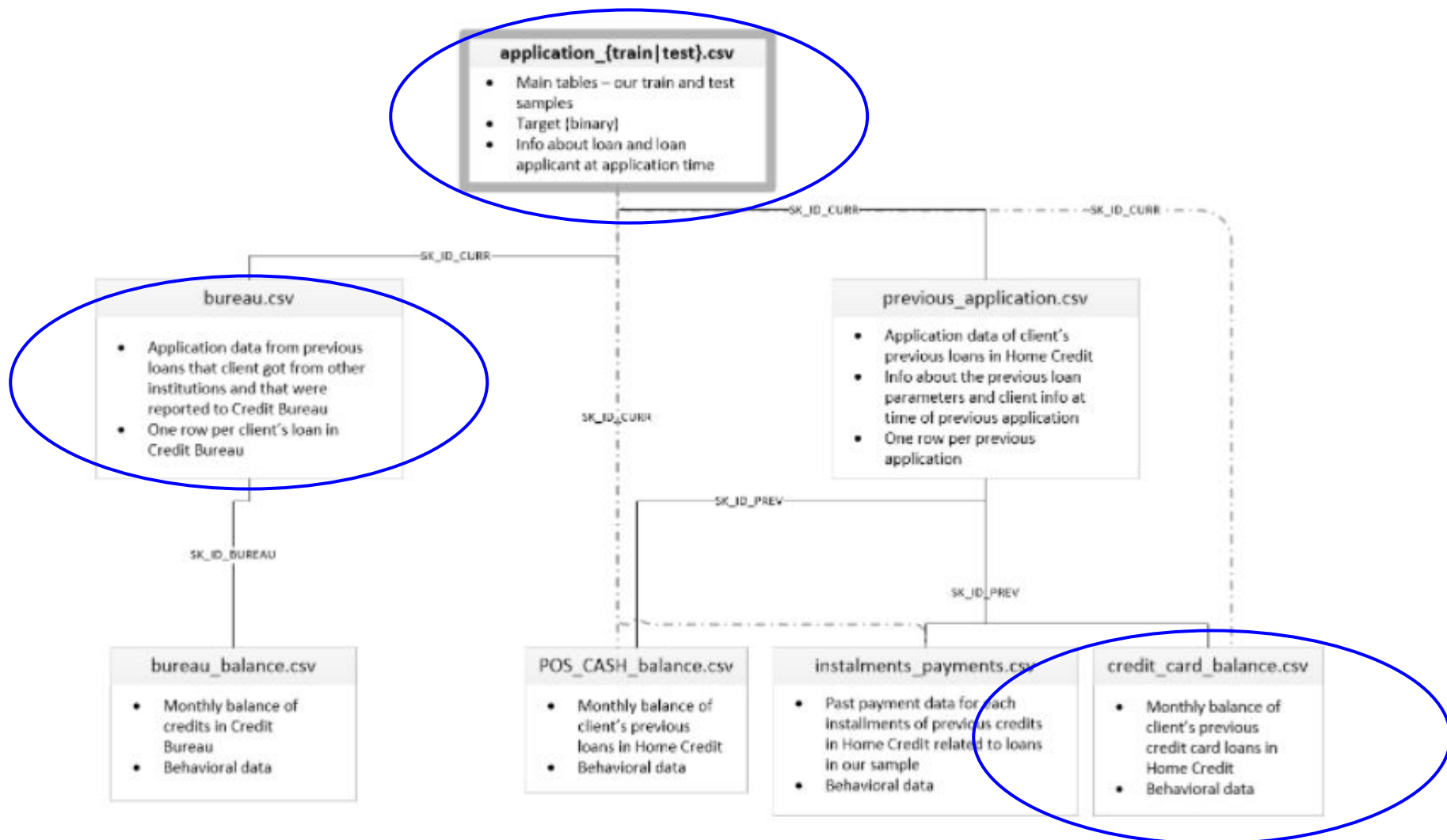
Augmenter les  
gains



FP / TN

# Description du jeu de données

---





# Agrégation de données

## **application\_train**

- 307511 lignes
- 122 variables



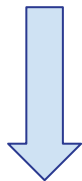
## **bureau**

- 1 716 428 lignes
- 17 variables



## **credit\_card\_balance**

- 3 840 312 lignes
- 23 variables

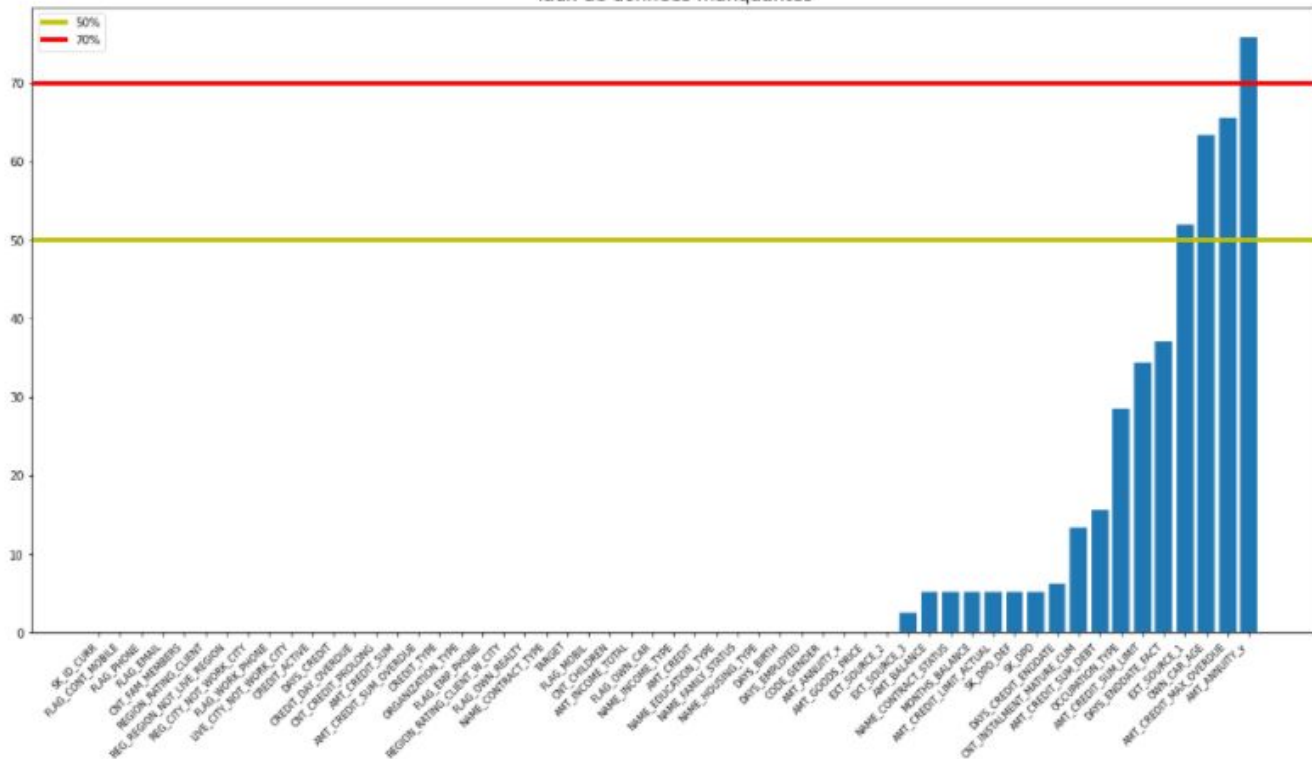
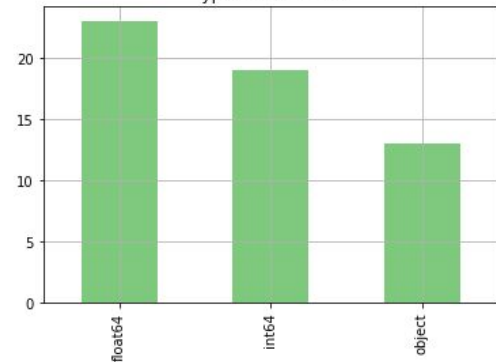


## **DATASET FINAL**

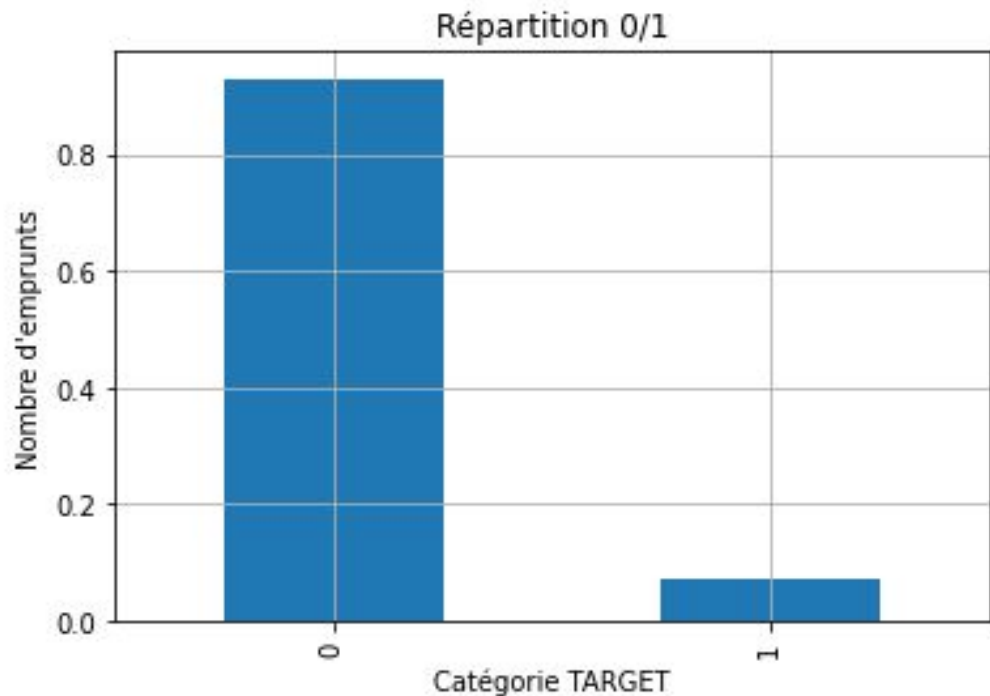
- 3 912 016 lignes
- 55 variables

- EDA : 10% => 391 202
- MOD : 40% => 152 285

# Analyse univariée



# Variable cible

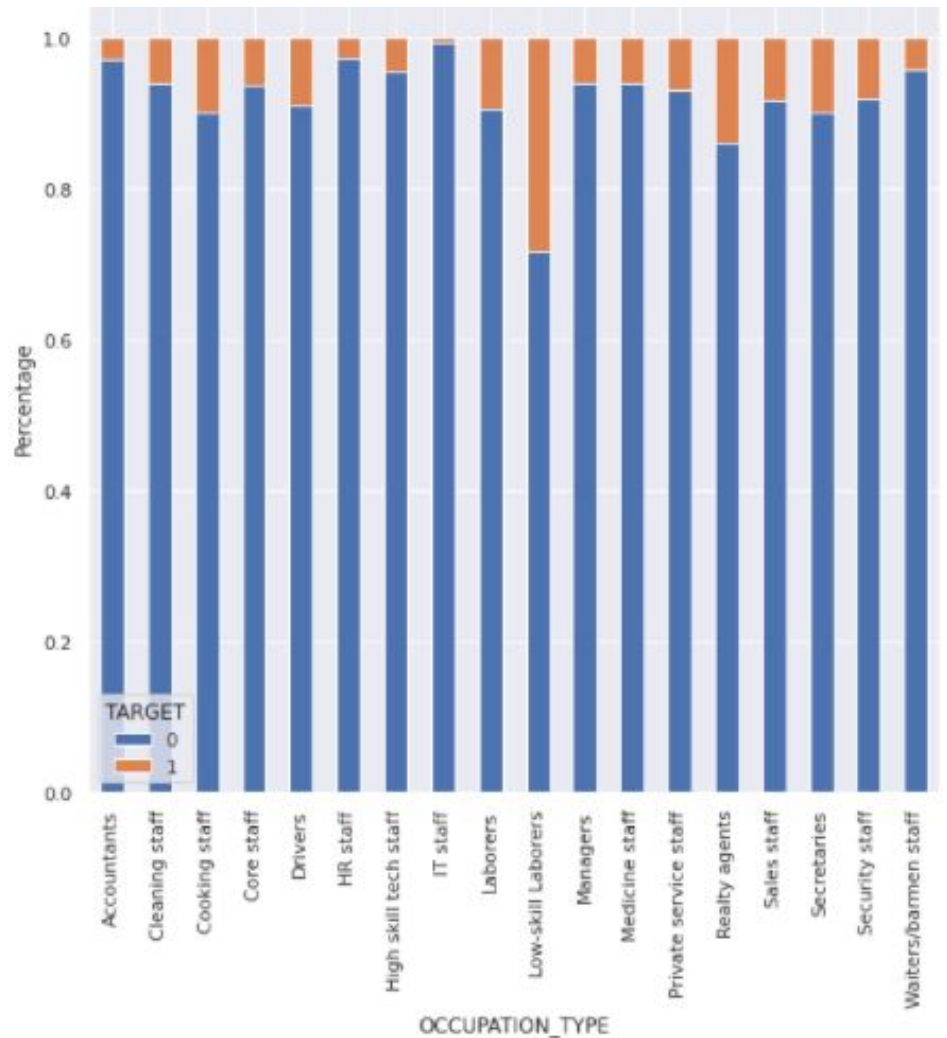
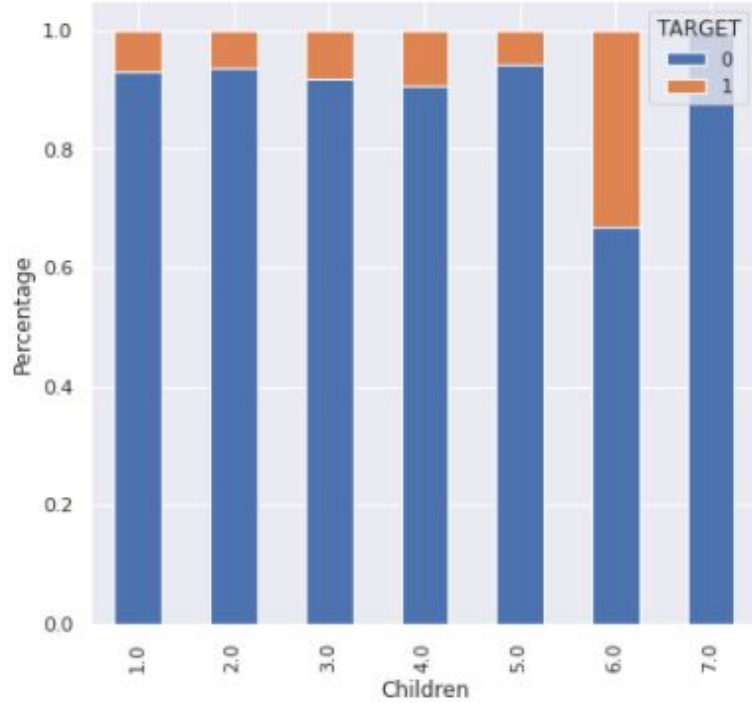


- Classe majoritaire
  - 93% de clients 0
- Classe minoritaire
  - 7% de clients 1

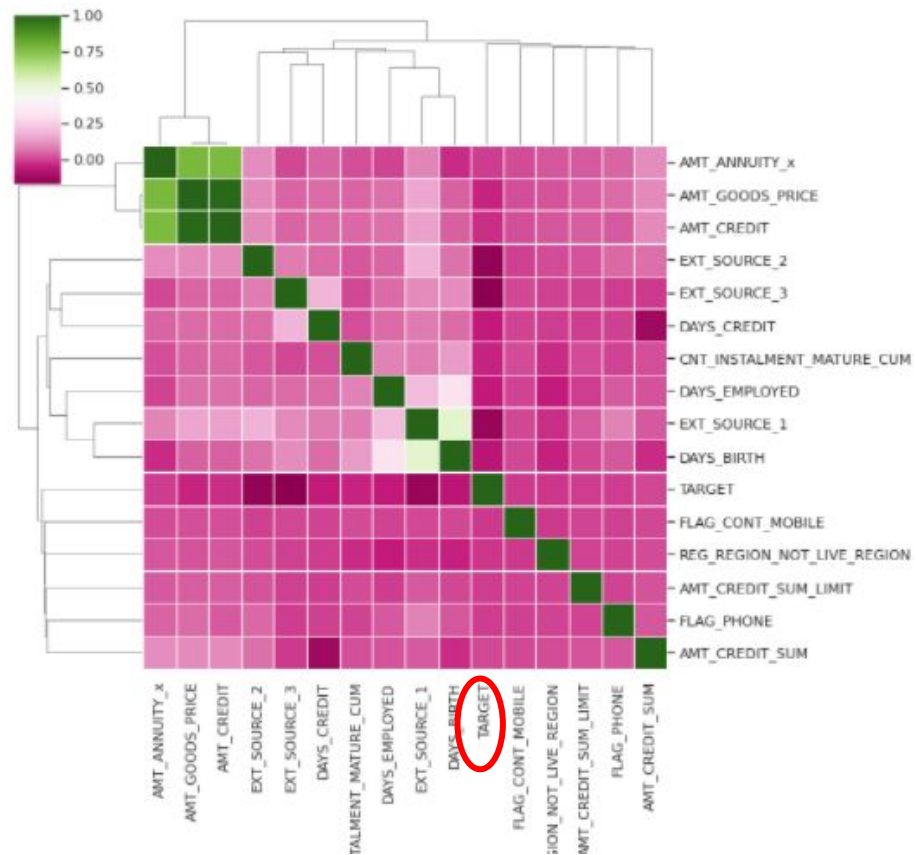
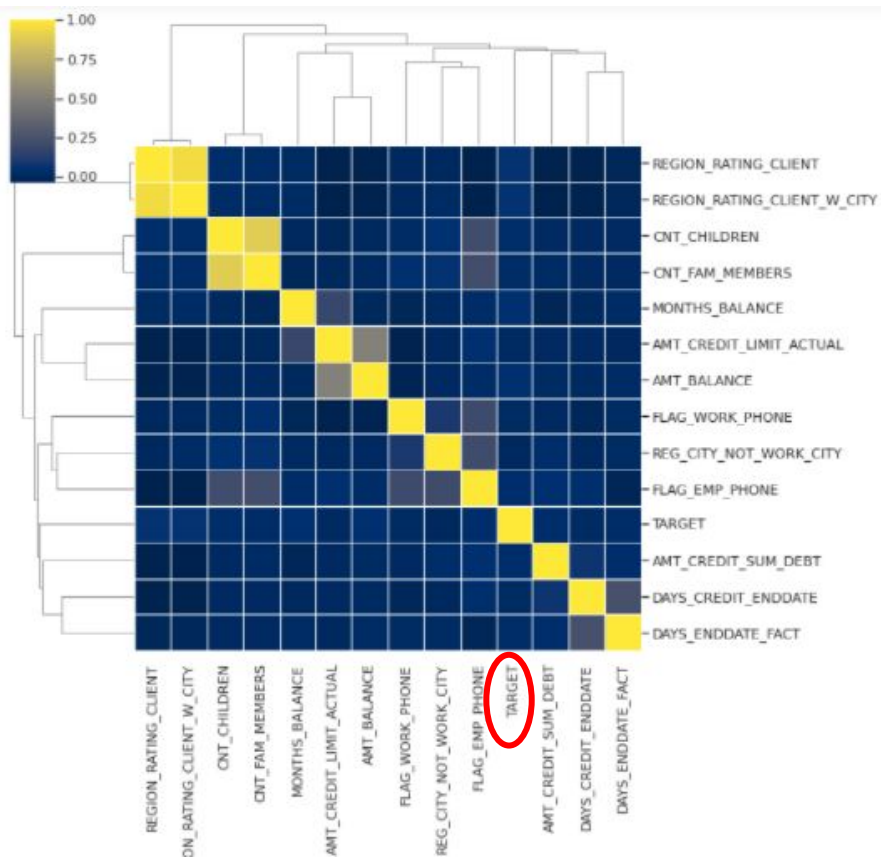
# Anomalies

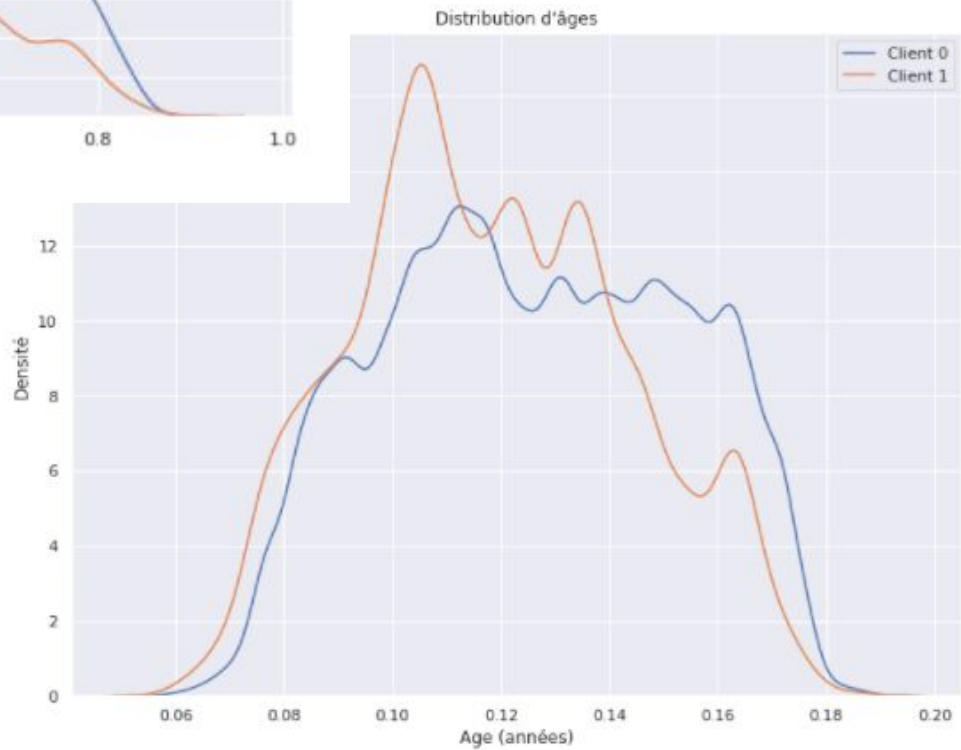
		Feature	Correction	Solution apportée
1	Chiffres négatifs	DAYS_BIRTH DAYS_EMPLOYED	✓	<ul style="list-style-type: none"><li>• Conversion en chiffres positifs :<ul style="list-style-type: none"><li>◦ diviser par 365</li><li>◦ multiplier par -1</li></ul></li></ul>
2	Présence des nan		✓	<ul style="list-style-type: none"><li>• Imputation par médiane / mode</li><li>• Suppression de features de plus de 60% de nan</li></ul>
3	Présence des inconnus	ORGANIZATION_ TYPE CODE_GENDER	✓	<ul style="list-style-type: none"><li>• "XNA"</li><li>• Imputation par mode</li></ul>
4	Présence des doublons		✓	<ul style="list-style-type: none"><li>• Suppression</li></ul>
5	Présence des outliers		✓	<ul style="list-style-type: none"><li>• Remplacement par médiane</li></ul>

# Analyse bivariable



# Corrélations



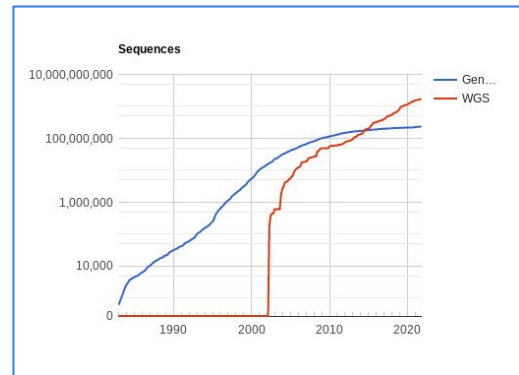
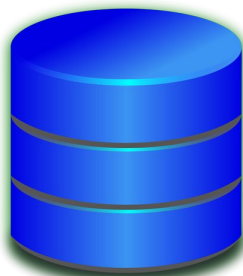


# Transformation du jeu de données

---



# Preprocessing



1

**Imputation des NaN**

Identification et élimination des nan

2

**Encodage**

One-Hot

3

**Feature engineering**

Création de nouvelles variables métier

4

**Feature selection**

Détermination de l'importance et choix de variables

5

**Normalisation**

Centrage d'une variable à zéro et de standardisation de la variance à 1

# Imputation des NaN



1.

- **Variables catégorielles**

- Remplissage par mode
- `fillna()`



2.

- **Variables numériques**

- Remplissage par médiane
- `SimpleImputer()`

# Encodage

id	color
1	red
2	blue
3	green
4	blue

One Hot Encoding

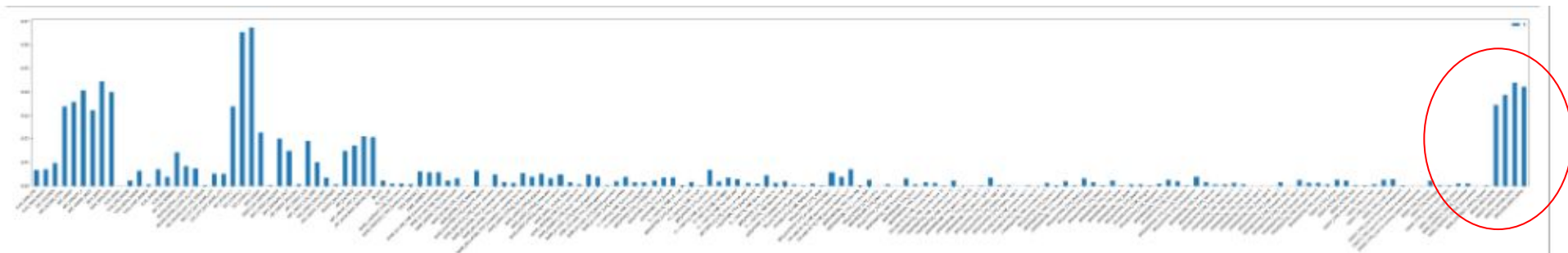
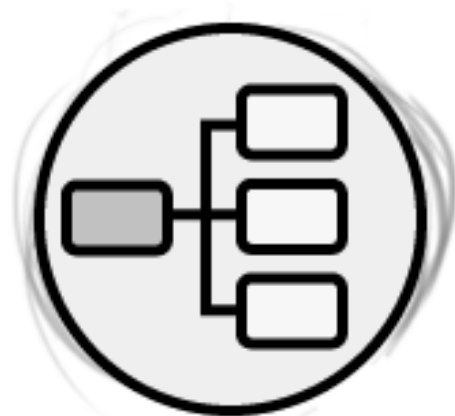
id	color_red	color_blue	color_green
1	1	0	0
2	0	1	0
3	0	0	1
4	0	1	0

```
NAME_CONTRACT_TYPE----- 2
CODE_GENDER----- 3
FLAG_OWN_CAR----- 2
FLAG_OWN_REALTY----- 2
NAME_INCOME_TYPE----- 6
NAME_EDUCATION_TYPE----- 5
NAME_FAMILY_STATUS----- 5
NAME_HOUSING_TYPE----- 6
OCCUPATION_TYPE----- 19
ORGANIZATION_TYPE----- 58
CREDIT_ACTIVE----- 4
CREDIT_TYPE----- 13
NAME_CONTRACT_STATUS----- 8
```

```
trainset.FLAG_OWN_CAR.unique(), trainset.FLAG_OWN_REALTY.unique()
(array(['Y', 'N'], dtype=object), array(['Y', 'N'], dtype=object))
```

# Feature engineering

- Ratio **montant emprunté** vs **prix du bien acheté** :  $\text{AMT\_CREDIT} / \text{AMT\_GOODS\_PRICE}$
- Ratio **annuités** vs **montant emprunté** :  $\text{AMT\_ANNUITY} / \text{AMT\_CREDIT}$
- Ratio **annuités** vs **revenus annuels** :  $\text{AMT\_ANNUITY} / \text{AMT\_INCOME\_TOTAL}$
- Ratio **ancienneté au travail** vs **âge** :  $\text{DAYS\_EMPLOYED} / \text{DAYS\_BIRTH}$



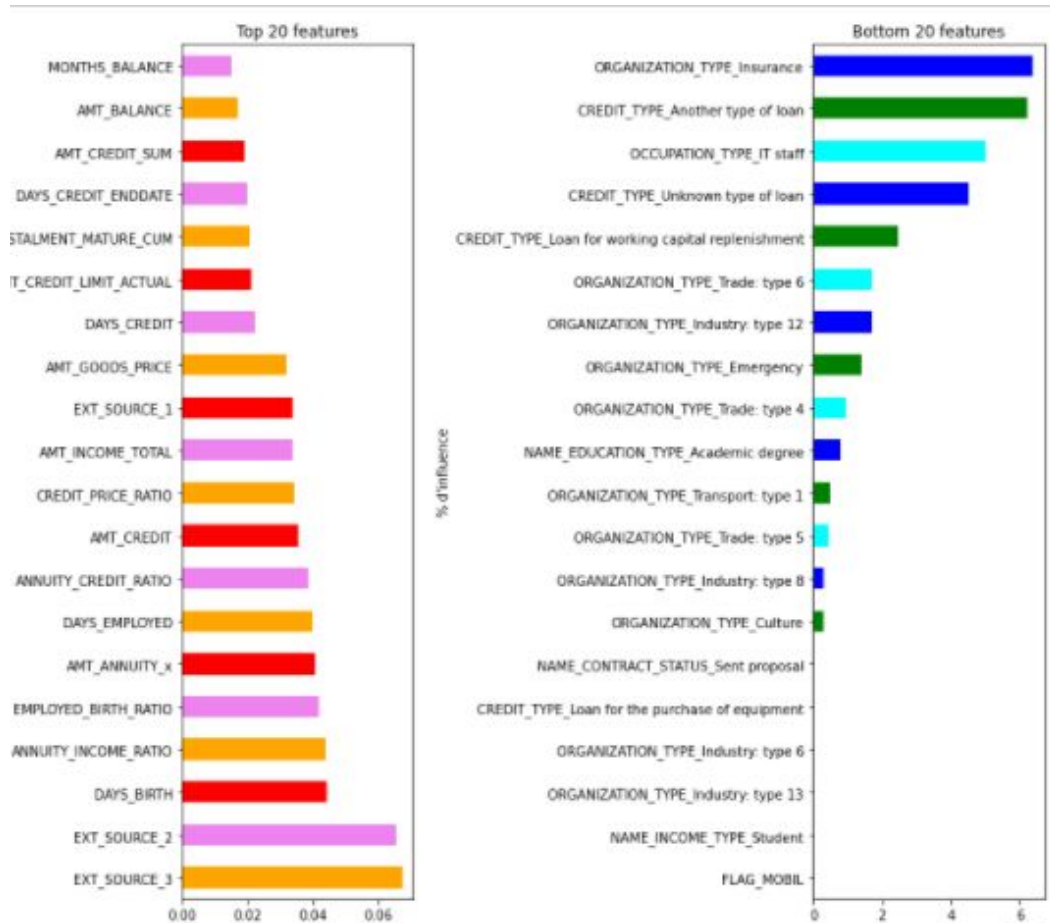
# Feature selection

## SelectFromModel

*Entraîne un estimateur puis sélectionne les variables les plus importantes pour cet estimateur.*

## SelectKBest

*Sélectionne les K variables X dont le score du test de dépendance (ici Chi2) avec y est le plus élevé*



# Séparation du jeu de données

## Trainset

**X\_train ,  
y\_train**

- 45 685 lignes
- 51 features

- 45 685 lignes
- 158 features

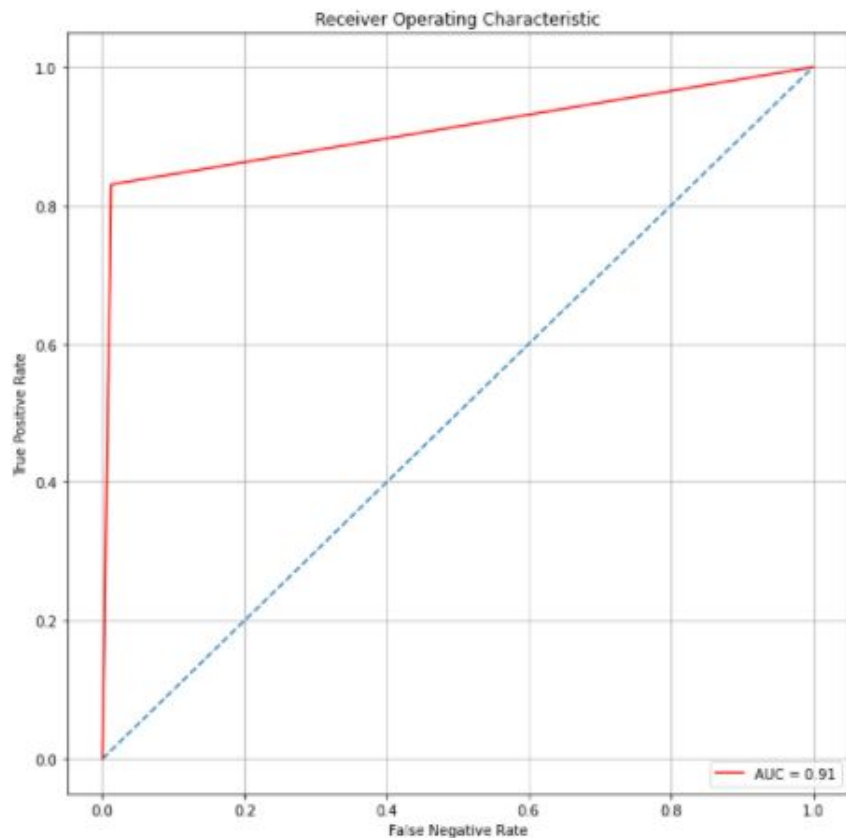
## Testset

**X\_test,  
y\_test**

- 11 422 lignes
- 51 features

- 11 422
- 156 features

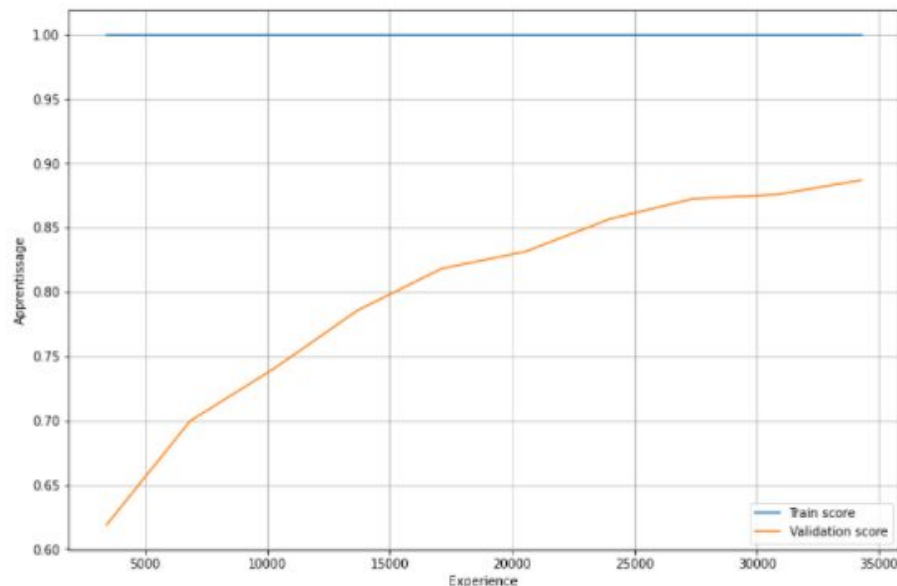
# Première modélisation



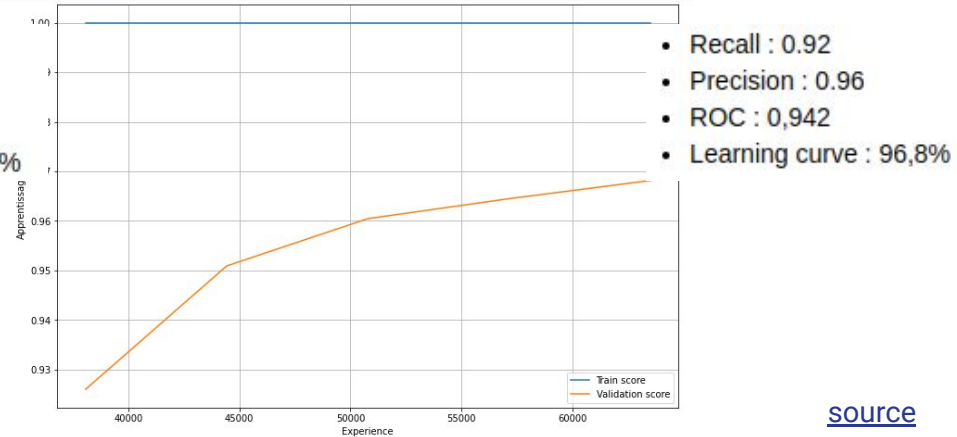
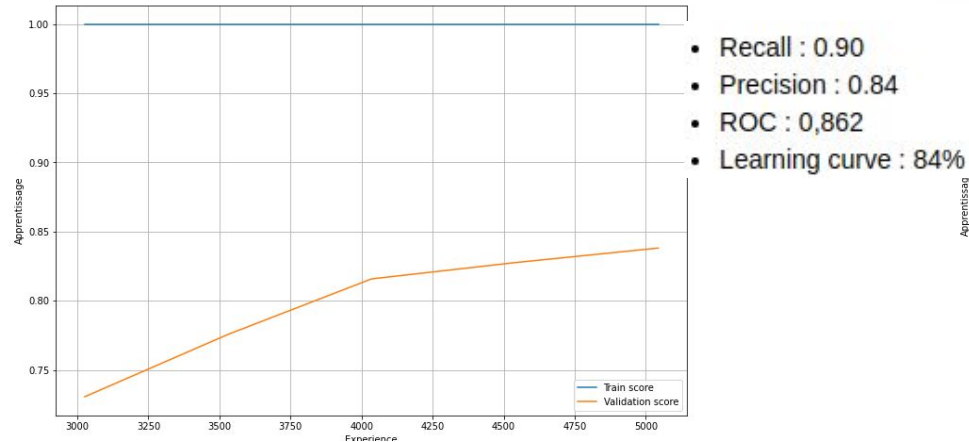
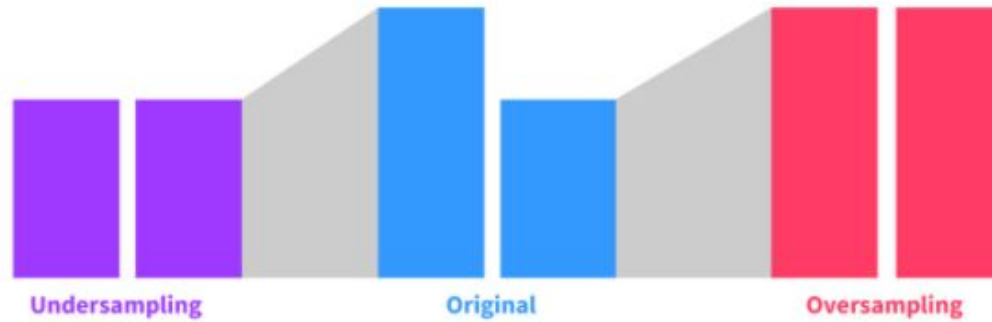
----- 1. Matrice de confusion -----

```
[[10513  131]
 [   132  646]]
```

	precision	recall	f1-score	support
0.0	0.99	0.99	0.99	10644
1.0	0.83	0.83	0.83	778
accuracy			0.98	11422
macro avg	0.91	0.91	0.91	11422
weighted avg	0.98	0.98	0.98	11422



# Équilibrage des classes



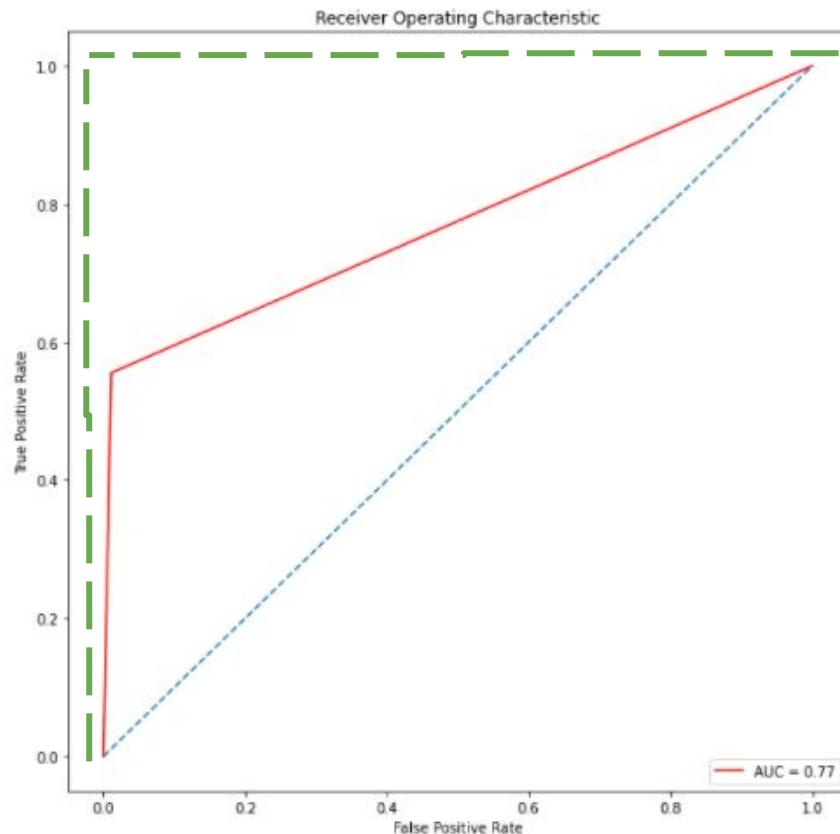


# Synthèse des résultats pour les modèles utilisés

---

# Métrique : ROC

- **modèle parfait AUC = 1**
  - 100% d'observations correctement classées
- **modèle naïf AUC = 0.5**
  - TP Rate = FP Rate
  - même proportion d'observations correctement & incorrectement classées
- **modèle KNeighbors Clf AUC**
  - supérieur à 0.5
  - proportion d'observations correctement classées en TP est supérieure aux observations incorrectement classées (FP)



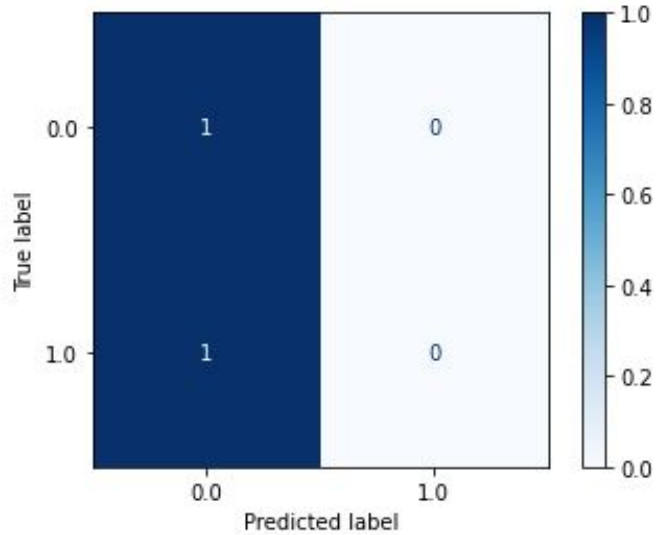
# KPI : Precision - Recall (Sensitivity)

		Predicted Class		
		Positive	Negative	
Actual Class	Positive	True Positive (TP)	False Negative (FN) Type II Error	<b>Sensitivity</b> $\frac{TP}{(TP + FN)}$
	Negative	False Positive (FP) Type I Error	True Negative (TN)	<b>Specificity</b> $\frac{TN}{(TN + FP)}$
		<b>Precision</b> $\frac{TP}{(TP + FP)}$	<b>Negative Predictive Value</b> $\frac{TN}{(TN + FN)}$	<b>Accuracy</b> $\frac{TP + TN}{(TP + TN + FP + FN)}$

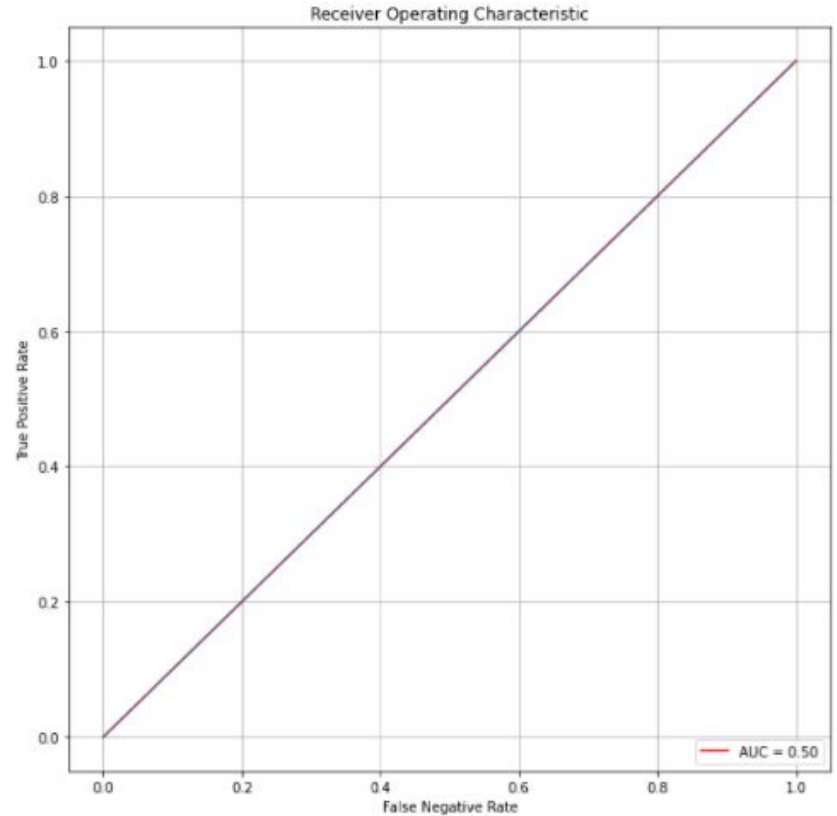
TP / **Actual** True Predictions

TP / **Total** True Predictions

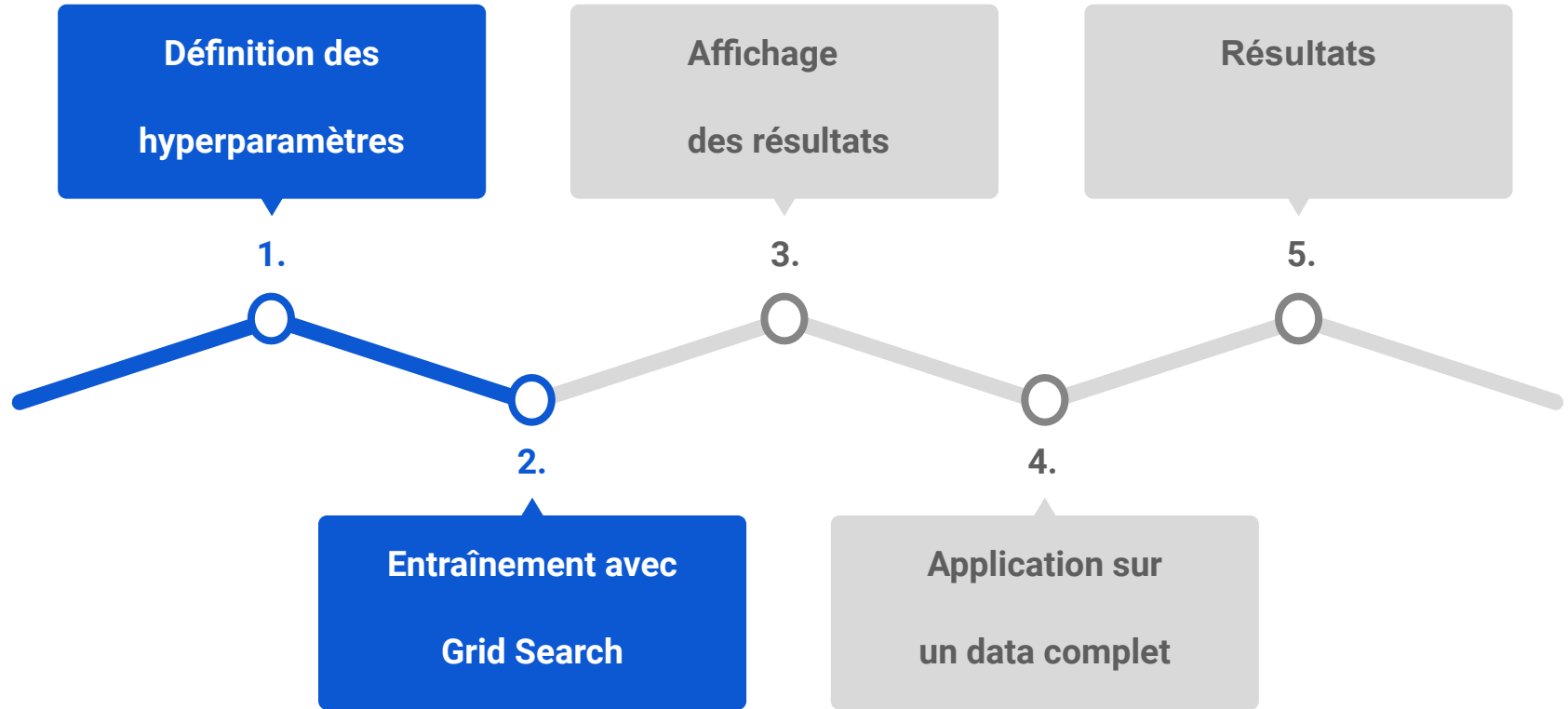
# DummyClassifier



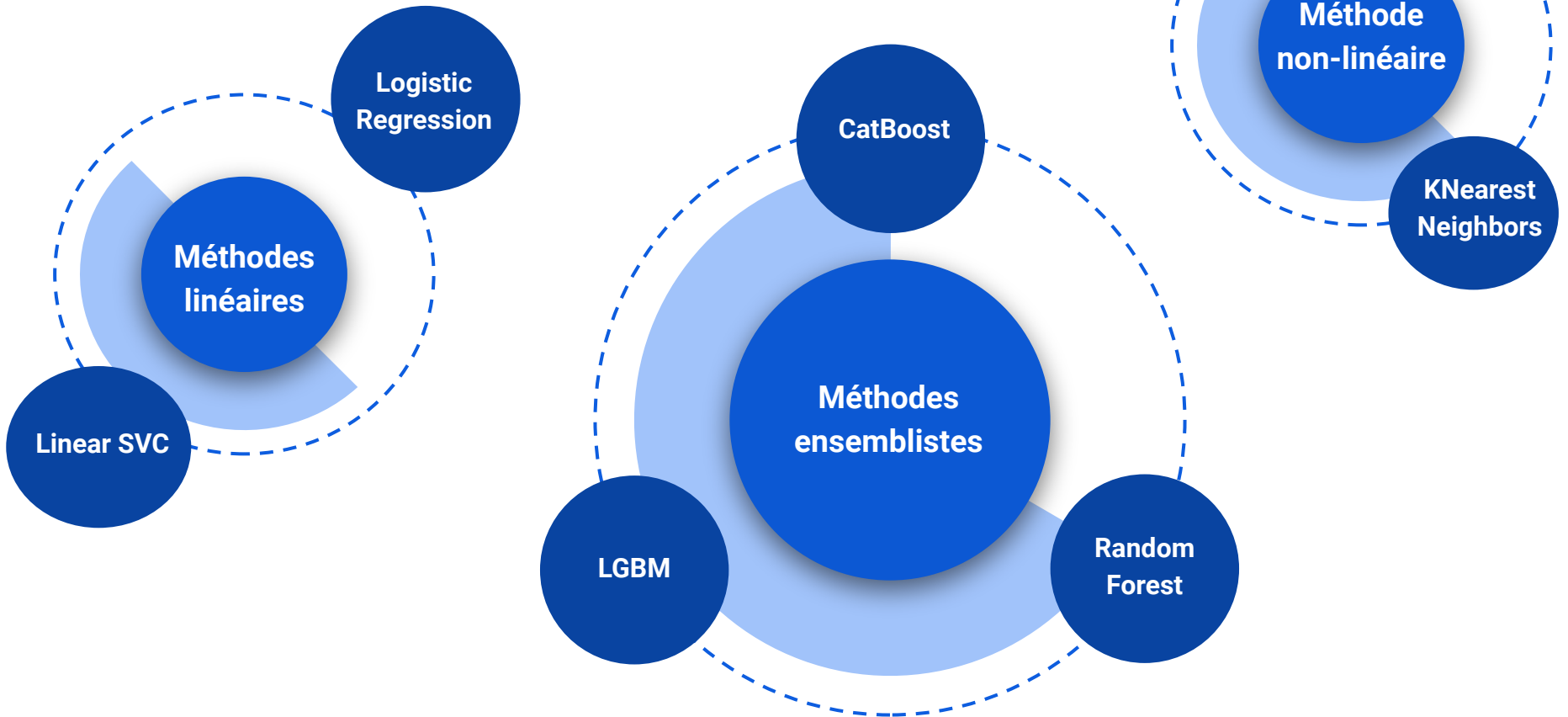
- $TN = 1$
- $FP = 1$



# Méthodologie



# Choix des algorithmes



# Résultats

	Recall TP	Precision FP	Accuracy	AUROC
Random Forest	0.77	0.86	0.82	0.821
Logistic Regression	0.68	0.70	0.70	0.695
SVM	0.69	0.70	0.70	0.698
KNN	0.91	0.78	0.82	0.823
CatBoost	0.96	0.96	0.96	0.960
LGBM	0.83	0.91	0.88	0.876
RandomForest Tuned	0.08	1.00	0.93	0.538
Logistic Regression Tuned	0.00	0.00	0.93	0.500
SVM Tuned	0.00	0.00	0.93	0.500
KNN Tuned	0.66	0.81	0.96	0.826
CatBoost Tuned	0.87	1.00	0.99	0.934
LGBM Tuned	0.33	0.97	0.95	0.663

# CatBoost Classifier

	precision	recall	f1-score	support
0.0	0.99	1.00	0.99	28271
1.0	1.00	0.87	0.93	2186
accuracy			0.99	30457
macro avg	0.99	0.93	0.96	30457
weighted avg	0.99	0.99	0.99	30457

[[282638]

[ 2871899]]

Receiver Operating Characteristic

True Positive Rate

False Positive Rate

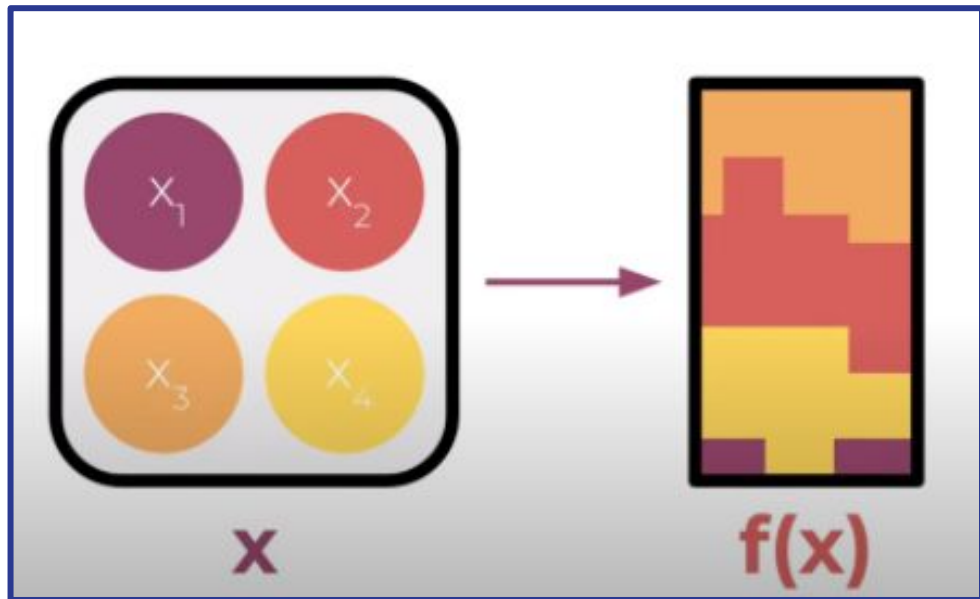
AUC = 0.93



# Interprétabilité du modèle

---

# SHAP

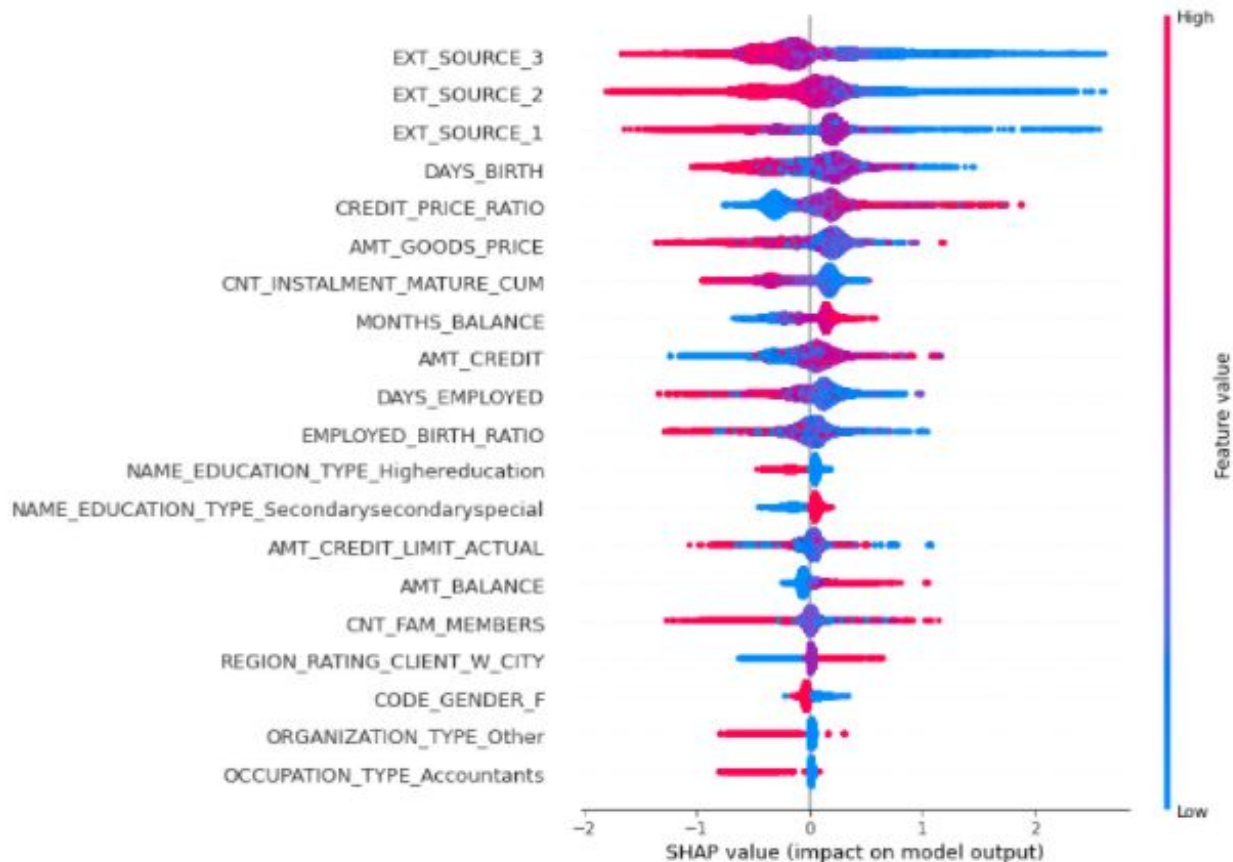


**Ajouter** les valeurs des contributions des variables sur toutes les prédictions, puis les **classer** par ordre d'importance.

**Expliquer** chaque prédiction individuellement.

**Estimer** la contribution de chaque variable dans la prédiction.

# Importance des variables



# Identification d'un TN



- Client solvable
- Probabilité de  $TARGET=1$  : 0.01
- Prédiction brute : 1,82
- Variables en bleu

# Identification d'un TP



- Client à risque
- Probabilité de TARGET==1 : 0.86
- Prédiction brute : -4.55
- Variables en rouge

# Conclusion

---

# Informations-clés

pour les conseillers de la banque

01	EXT_SOURCE_XXX	<ul style="list-style-type: none"><li>• Notes externes du demandeur d'un crédit<ul style="list-style-type: none"><li>○ favorable</li><li>○ défavorable</li></ul></li></ul>
02	DAYS_BIRTH	<ul style="list-style-type: none"><li>• L'âge du demandeur d'un crédit<ul style="list-style-type: none"><li>○ favorable</li><li>○ défavorable</li></ul></li></ul>
03	CREDIT_PRICE_RATIO	<ul style="list-style-type: none"><li>• Rapport entre le montant emprunté et le prix du bien acheté<ul style="list-style-type: none"><li>○ favorable</li><li>○ défavorable</li></ul></li></ul>