



# **Tarification des contrats d'assurance automobile**

Guy Tsang, Axel Gardahaut & Léo Dutertre-Ladurée

# INTRODUCTION

## CONTEXTUALISATION ET PRÉSENTATION DES DONNÉES

1

### • NATURE DES DONNÉES

Les variables sont anonymisées :

- Variables concernant l'assuré
- Variables concernant la région de l'assuré
- Variables concernant la voiture de l'assuré
- Variables calculées

### • BASE D'APPRENTISSAGE

Nombre de variables : 57

Nombre d'individus : 416 648

### • VARIABLE CIBLE

Réclamation dans un délai d'un an :

- Proportion de 1 dans la base d'apprentissage : 3.67%
- Proportion de 1 dans la base de test : 3.60%

### • MÉTRIQUE CHOISIE

Coefficient normalisé de Gini :

$$\text{Gini} = 2 \times \text{AUC} - 1$$

# PLAN



## I/. DÉMARCHE

- A) Benchmark
- B) Pré-traitement des données
- C) Sélection des variables
- D) Traitement des données

## II/. MODÈLES

- A) Régression pénalisée
- B) XGboost
- C) LightGBM
- D) Stratégie de validation

## III/. RÉSULTATS

- A) Comparaison des modèles
- B) Sélection du meilleur modèle & Performance

# I/. DÉMARCHE

## A) BENCHMARK

- Seuil à absolument dépasser avec :
  - Le traitement de la base de données
  - La sélection de variables
  - L'utilisation de modèles alternatifs
  - L'hyperparamétrisation des modèles
- Benchmark par LightGBM
  - Sans hyperparamétrisation (définis selon des valeurs usuelles en pratique)
  - Première étude de l'importance des variables du jeu de données
  - Coefficient Normalisé de Gini en validation croisée 5 blocs : 0.2719

# I/. DÉMARCHE

## B) PRÉ-TRAITEMENT DES DONNÉES

1

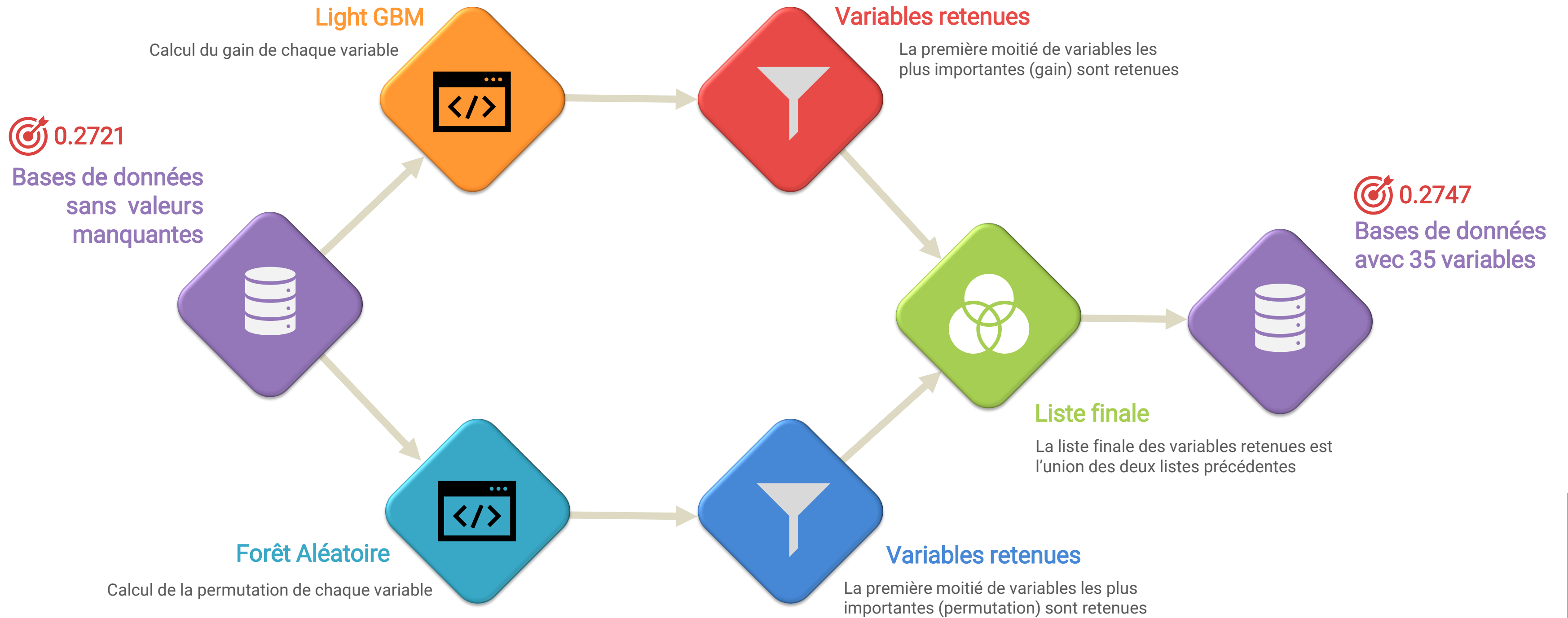




# I/. DÉMARCHE

## C) SÉLECTION DES VARIABLES

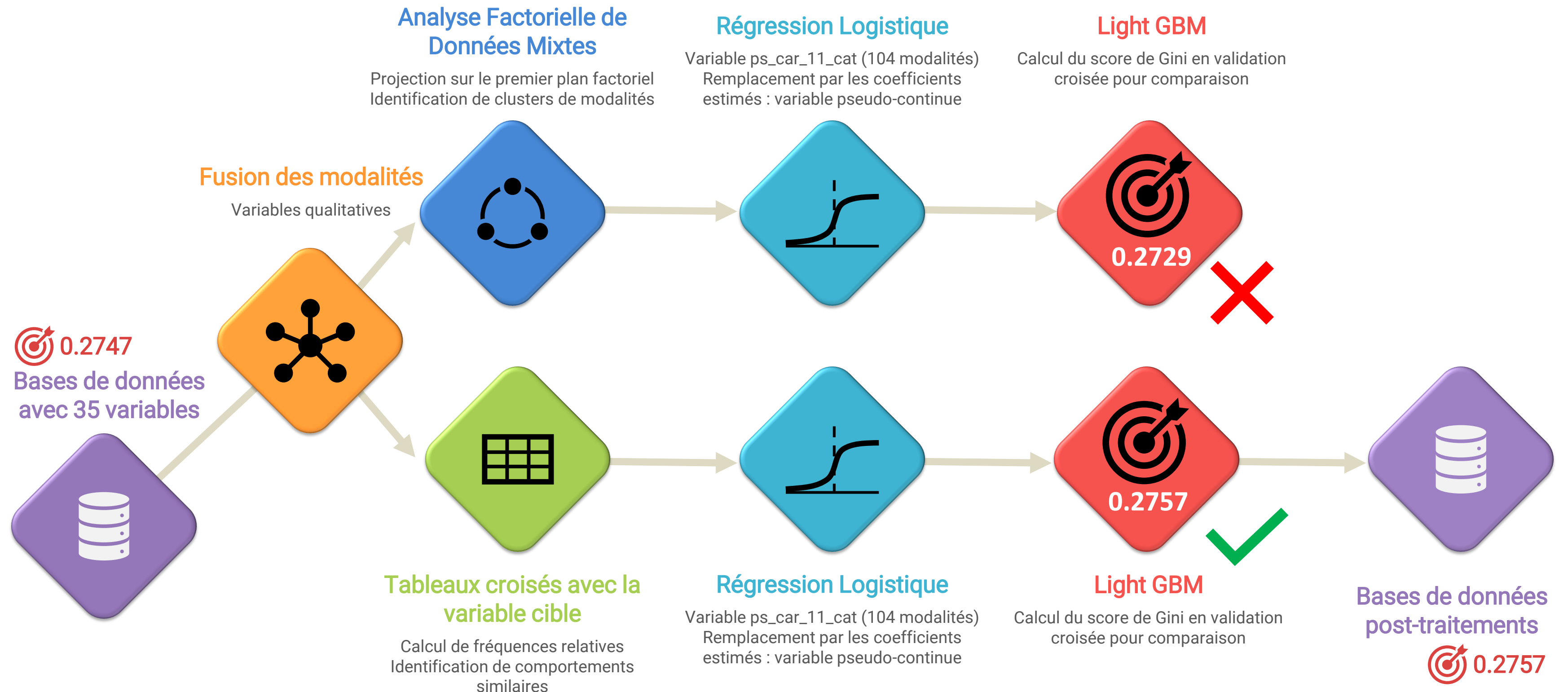
1



# I/. DÉMARCHE

## D) TRAITEMENT DES DONNÉES

1



PARTIE AXEL



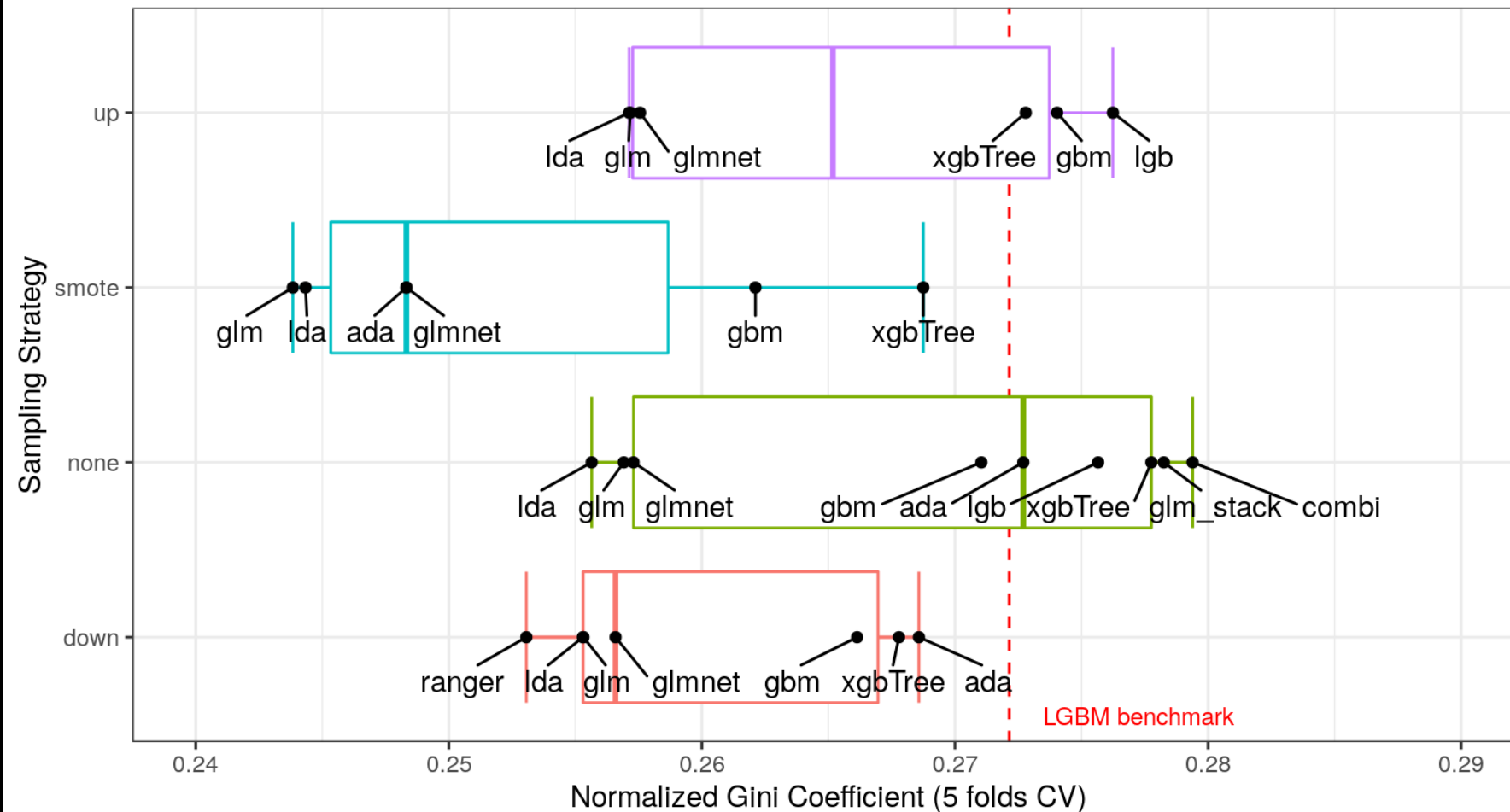
# III/. RÉSULTATS

## A) COMPARAISON DES MODÈLES

1

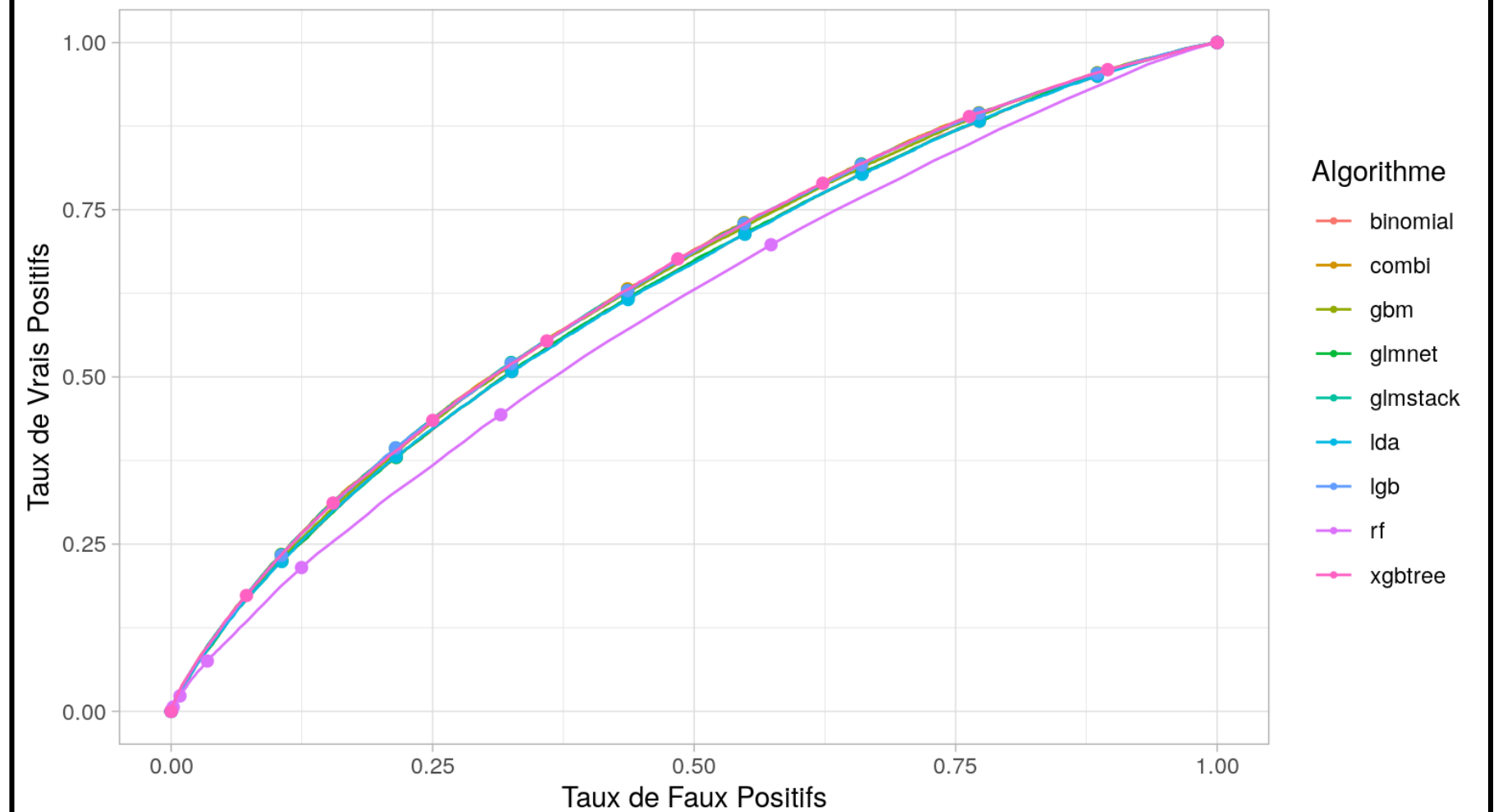
Comparaison des performances (validation croisée 5 blocs)

6 Modèles ont été retirés avant le construction des boxplots car trop mauvais  
 $\text{glm\_stack} = \text{stacking de XGB et LGB avec un GLM}$  ///  $\text{combi} = 0.1 * \text{pred\_LGB} + 0.9 * \text{pred\_XGB}$



Courbes ROC

Prédictions 5fCV des modèles entraînés sans resampling



	GLM	LDA	GLMNET	SVM	KPPV	ARBRE	FORÊT	GBM	ADA	LGB	XGB	STACK	COMBI
Up													
Down													
SMOTE													
Aucun													



Modèle construit



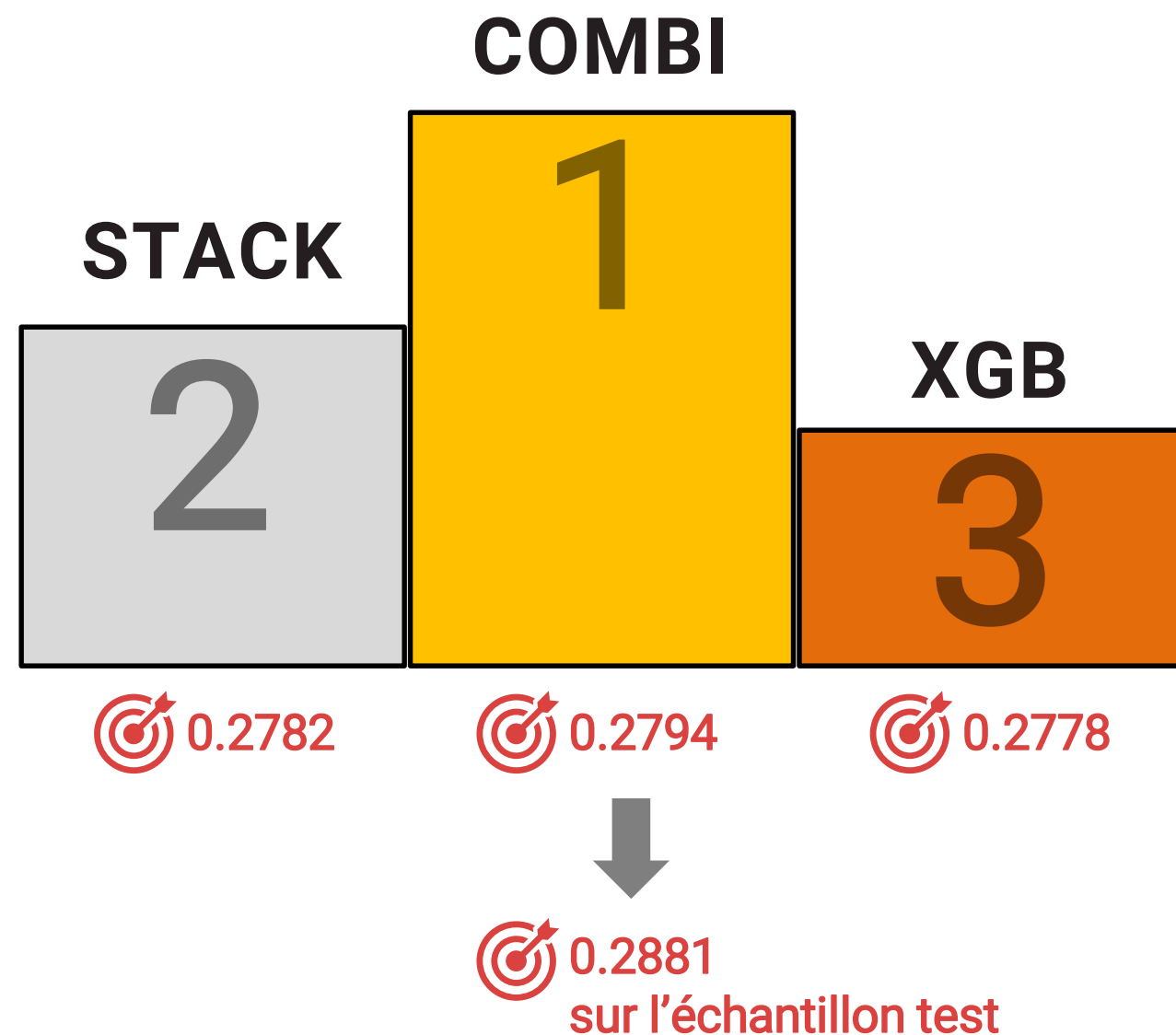
Modèle non construit



Modèle non pertinent

# III/. RÉSULTATS

## B) SÉLECTION DU MEILLEUR MODÈLE & PERFORMANCE



- Modèle combinaison de scores
  - Modèle issu d'une combinaison linéaire des prédictions CV du LightGBM (10%) et ceux du XGB (90%)
  - Pour information : AUC sur la base test égale à 0.6440
- Classification à partir du modèle retenu
  - Cut obtenu en optimisant le score F1 sur la base train
  - Matrice de confusion et performances de classification

		Observé	
		0	1
Prédiction	0	152 542	4 827
	1	19 596	1 599

	F1	Rappel
Validation	0.1176	0.2482
Test	0.1158	0.2488