



ayming

Projet debrief

Draft au 6 septembre 2019

Feature engineering

- 55 colonnes dont plusieurs ne portent aucune info (valeur unique)
 - browserSize, browserVersion, cityId, criteriaParameters, flashVersion, language, latitude, longitude, mobileDeviceBranding, mobileDeviceInfo, mobileDeviceMarketingName, mobileDeviceModel, mobileInputSelector, networkLocation, operatingSystemVersion, screenColors, screenResolution, socialEngagementType, visits
 - Reste 37 colonnes
- Feature engineering
 - Passage en numérique : fullVisitorId, pageviews, hits, visitNumber, transactionRevenue, visitId
 - Dates : champs date et visitStartTime transformé en date et heure (fonction (as.POSIXct)
 - Extraction de nouvelles colonnes : mois, jour de semaine, ..
 - Passage des dates elles memes en numérique
- Nouvelles features
 - Nb de session par joueur
 - Délai entre session
 - sessionId : Split en 2 colonnes numériques (de part et d'autre du « - »)
 - gcId : Extraction des 3 premiers et 3 derniers caracteres (présente des régularités. On pourrait pousser plus loin le parsing)
 - referralPath : détection des paths avec google (on pourrait pousser plus loin le parsing)
 - Keywords : features indicatrice de « youtube, google, matching, Content targeting, vertical targeting,, provided, store, shop
- Traitement des NA
 - « ex na » pour les strings
 - Médian pour les numériques
 - 0 pour la target transaction revenue
- Factorisation des strings, réduction du nombre de modalités à 1000 (information mutuelle)
- Passage de la target en $\log(x+1)$

55 colonnes



37 colonnes

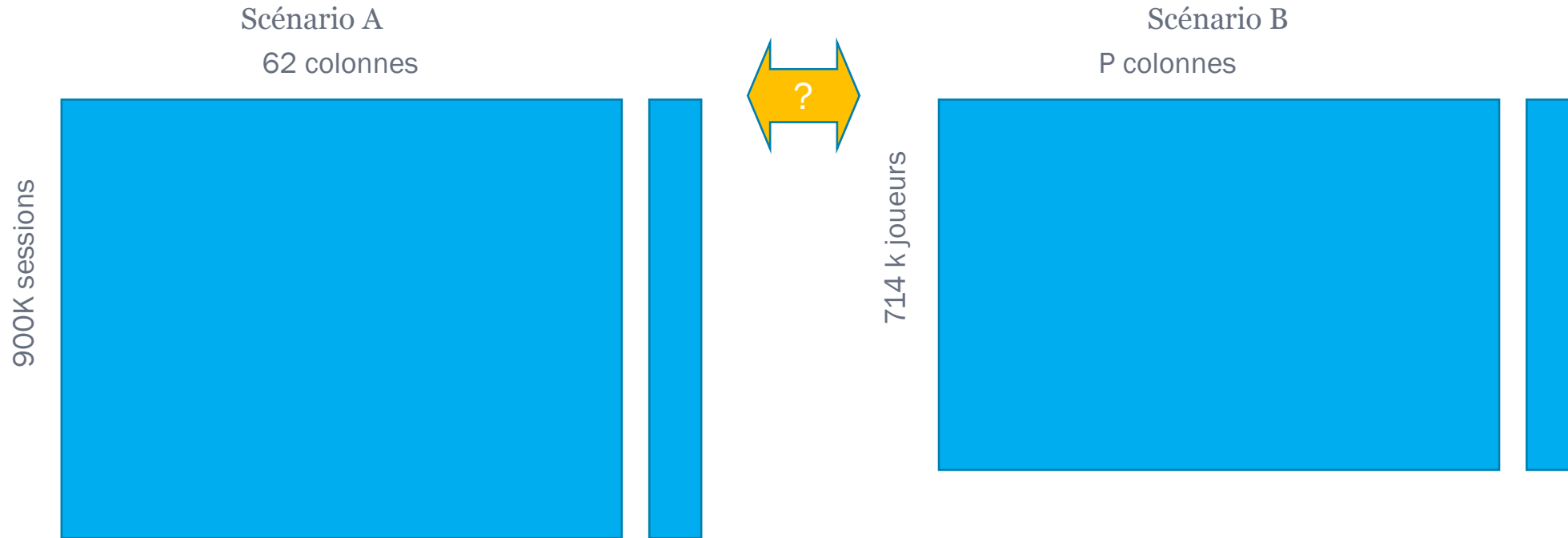


62 colonnes

2 approches possibles : prédiction de montant par session puis agrégation, ou construire un dataset par joueur puis prédiction directe du montant par joueur

DRAFT

3



Point d'attention : Séparation train test

Attention à splitter au hasard par joueur et non
au hasard par session dans le jeu d'évaluation
les nouveaux joueurs n'ont aucun historique
connu

Scénario A1

- directement une régression ?

Scénario A2

- d'abord une classif pour mettre de coté les 0 puis une régression ?

Point d'attention : feature engineering

Possibilités infinies de création de feautres
agrégées par joueur

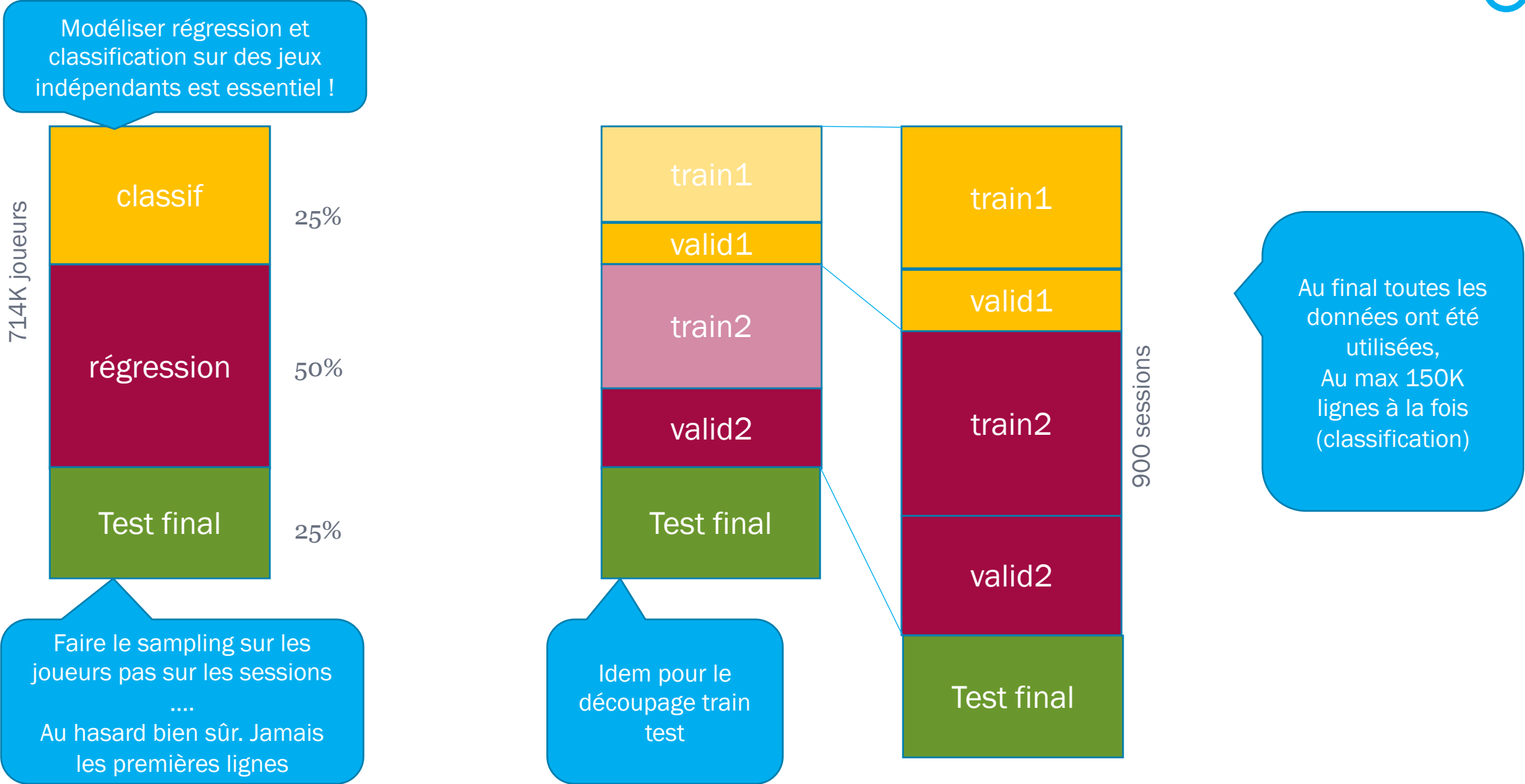
Quel point d'attention si stratégie
de modèles successifs ?

Une question importante avant de démarrer... → laquelle ?

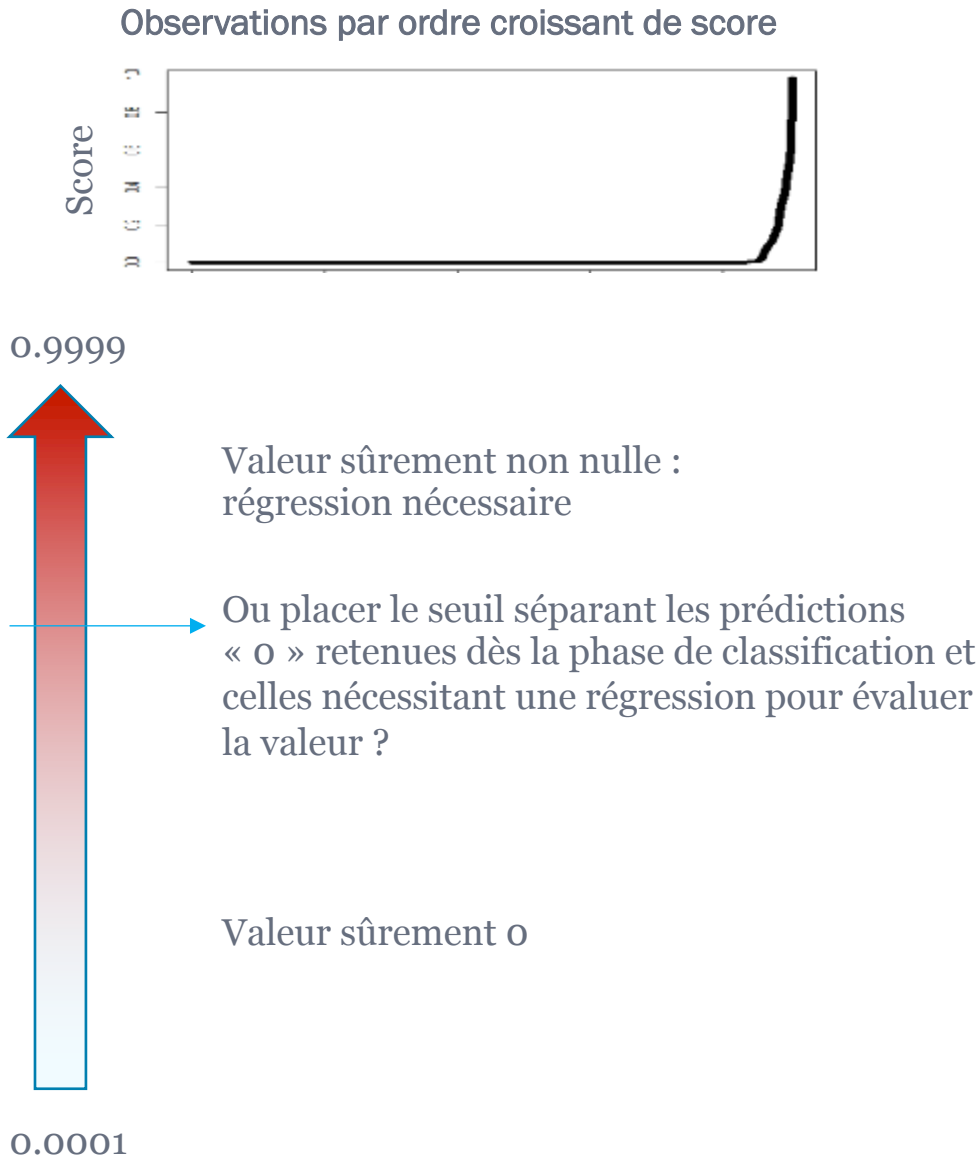
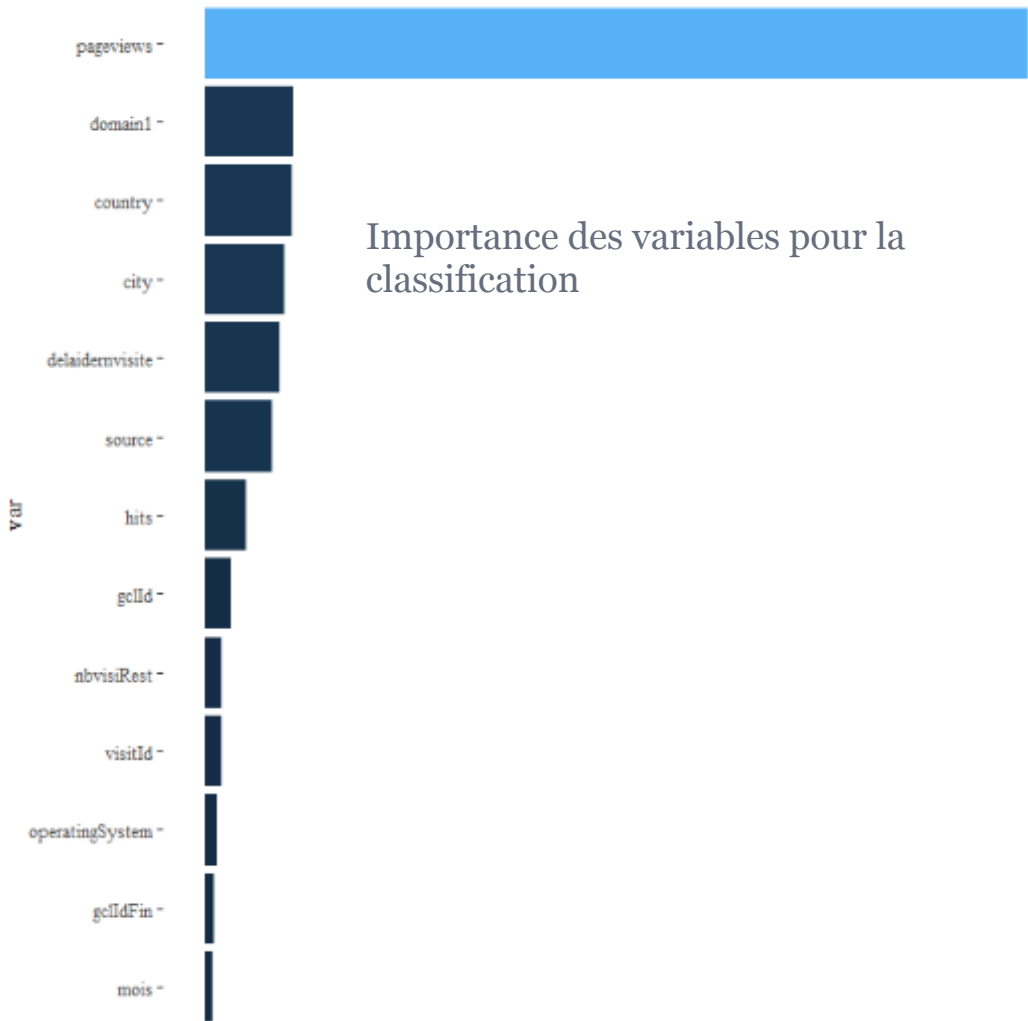
DRAFT

4

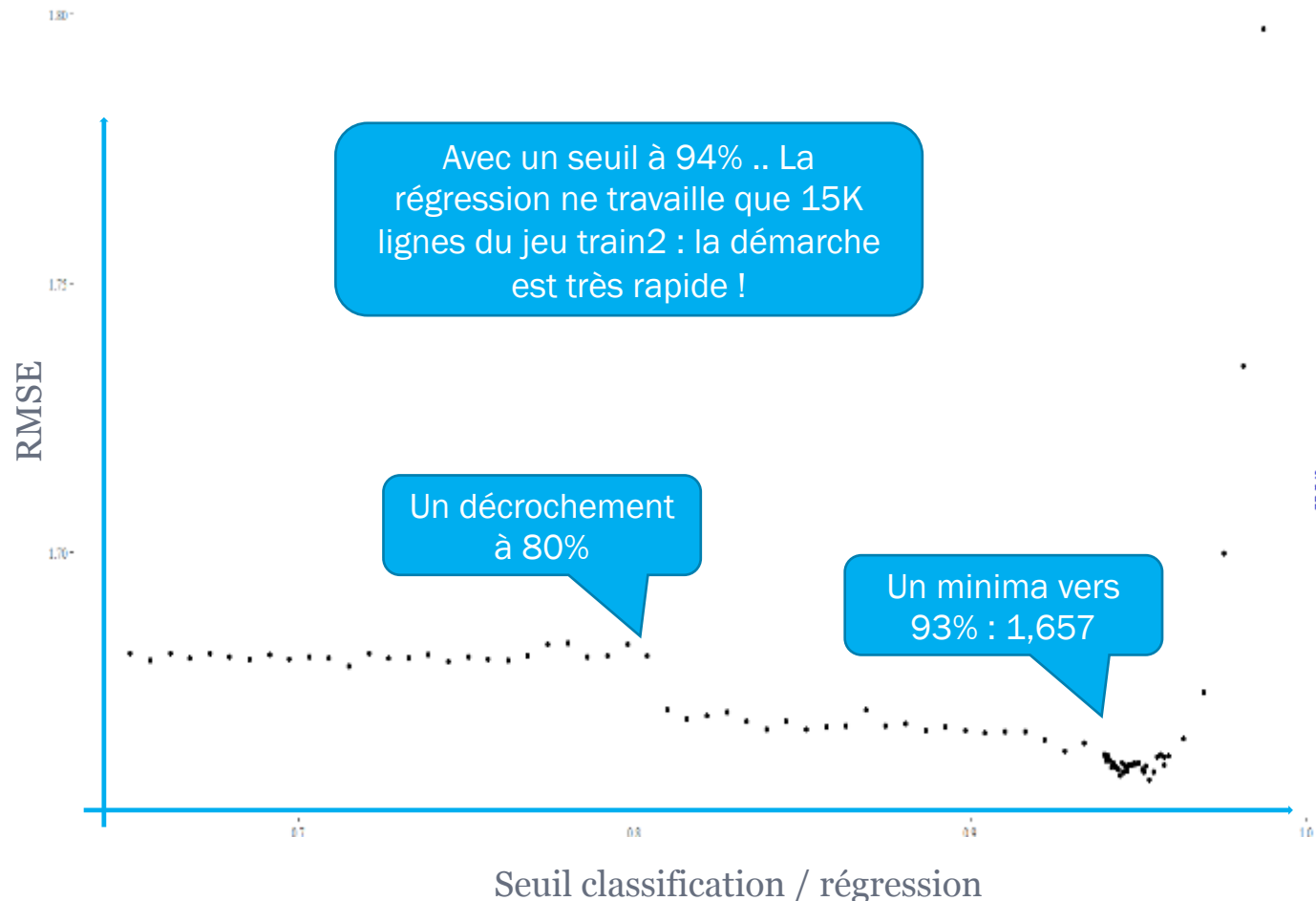
- Quelle serait la performance d'une prédiction naïve ?
- Exemple : prédiction aucun achat pour tout le monde
 - RMSE : 2,02
- Cette valeur est essentielle pour porter un jugement de valeur sur le résultat ...



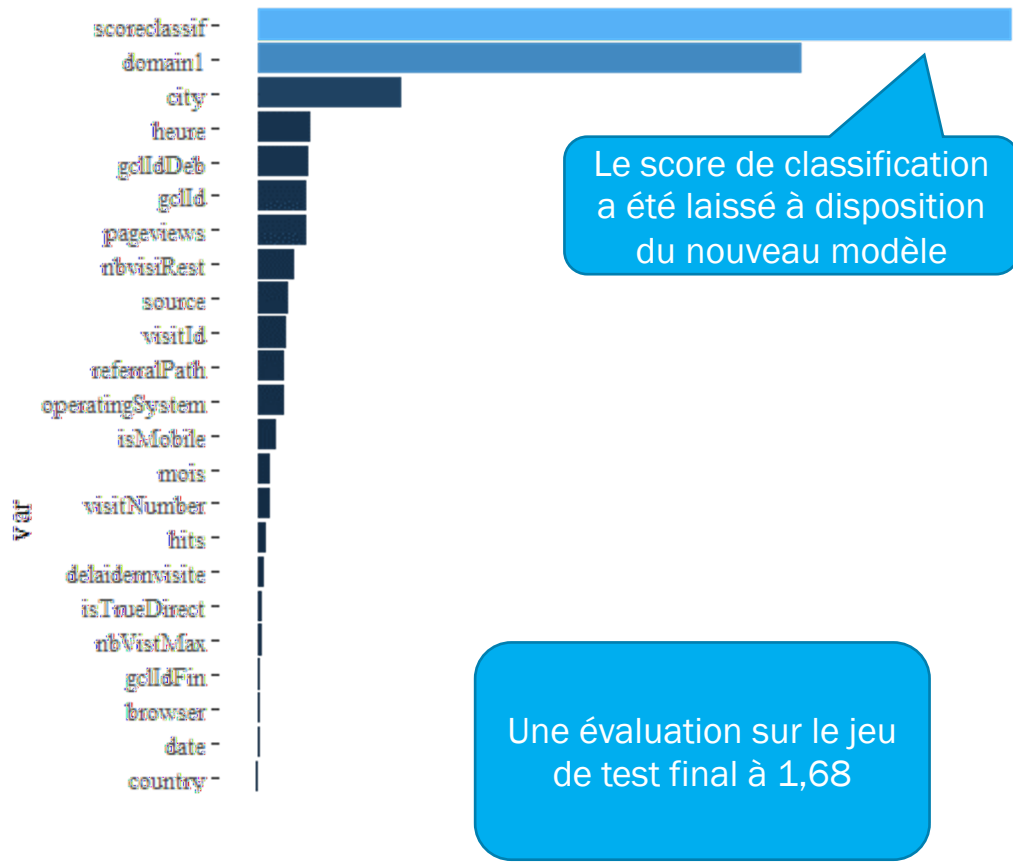
98% d'auc : Encouragement à poursuivre dans cette voie



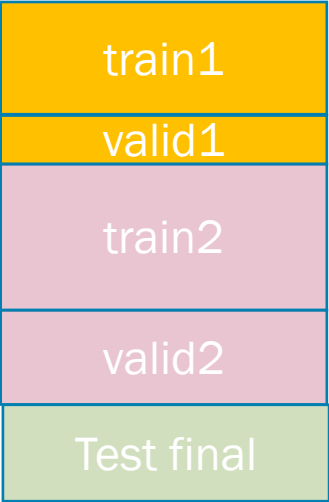
Pour déterminer le seuil .. La seule démarche est de tester !



Importance des variables pour la régression

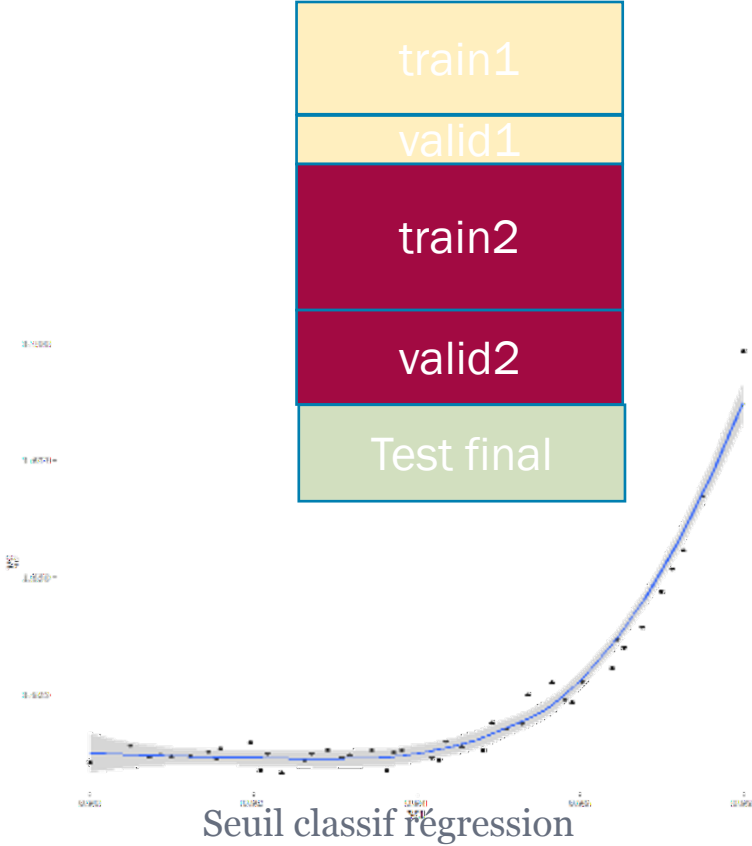


Classification



AUC 0.985

Régression au-delà d'un certain score
En dessous : 0
Au dessus : prédiction de la classification

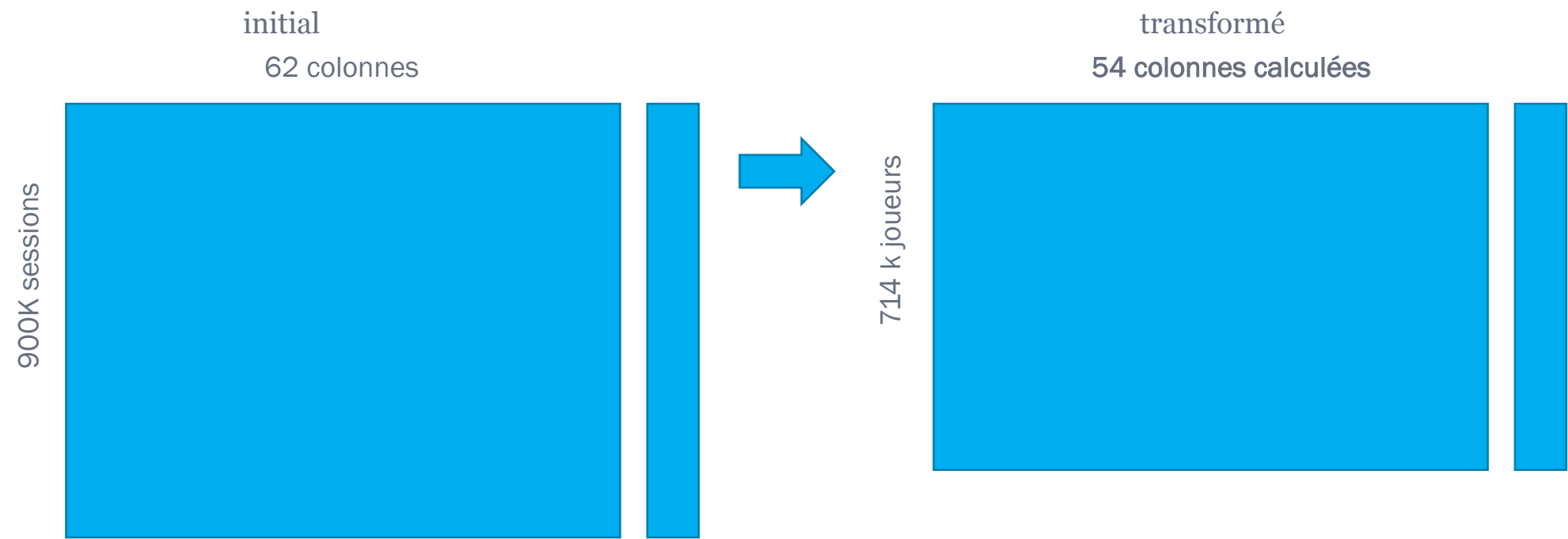


Régression 1.605

Jeu de test



Final 1.61



RMSE : 1,67 ... approche mise de côté

Pour aller plus loin

- Réapprendre avec toutes les données
- Intégrer le score par session dans la régression par visiteur (calculer des moyennes et des sommes de score ?)
- Combiner des modèles : comment ?