

# Projet 5 : Segmentation du comportement clients >

## SOMMAIRE

### > Choix des caractéristiques

- Objectifs et enjeux
- Présentation des variables et nettoyage des données
- Protection des données
- Piste de modélisation et choix des caractéristiques
- Présentation du modèle collaboratif et choix des hyper-paramètres adaptés

### > Segmentation / prédiction

- Apprentissage non supervisé : méthode
- Segmentation et interprétation (modèle collaboratif)
- Segmentation > 3 scénarios, choix du nombre de clusters / interprétation des clusters / visualisation (modèle manuel )
- Apprentissage supervisé : méthode
- Présentation des algorithmes de classifications
- Performances de différents scénarios
- Prédiction finale

# Problématique client > Segmenter les utilisateurs ? Pourquoi ?

**Postulat : Objectif d'une entreprise > Faire du profit**

Méthodes directes ou indirectes : améliorer la visibilité ou l'image de l'entreprise  
investir pour créer les nouveaux besoins, améliorer l'environnement de travail pour  
augmenter la productivité etc...)

- **Coûts structurels (entretien fonctionnement, équipements etc...)**
- **Investissement (développement de nouvelles gammes de produits, recherche et développement...)**
- **Marketing(publicité, démarchage, réductions etc...)**

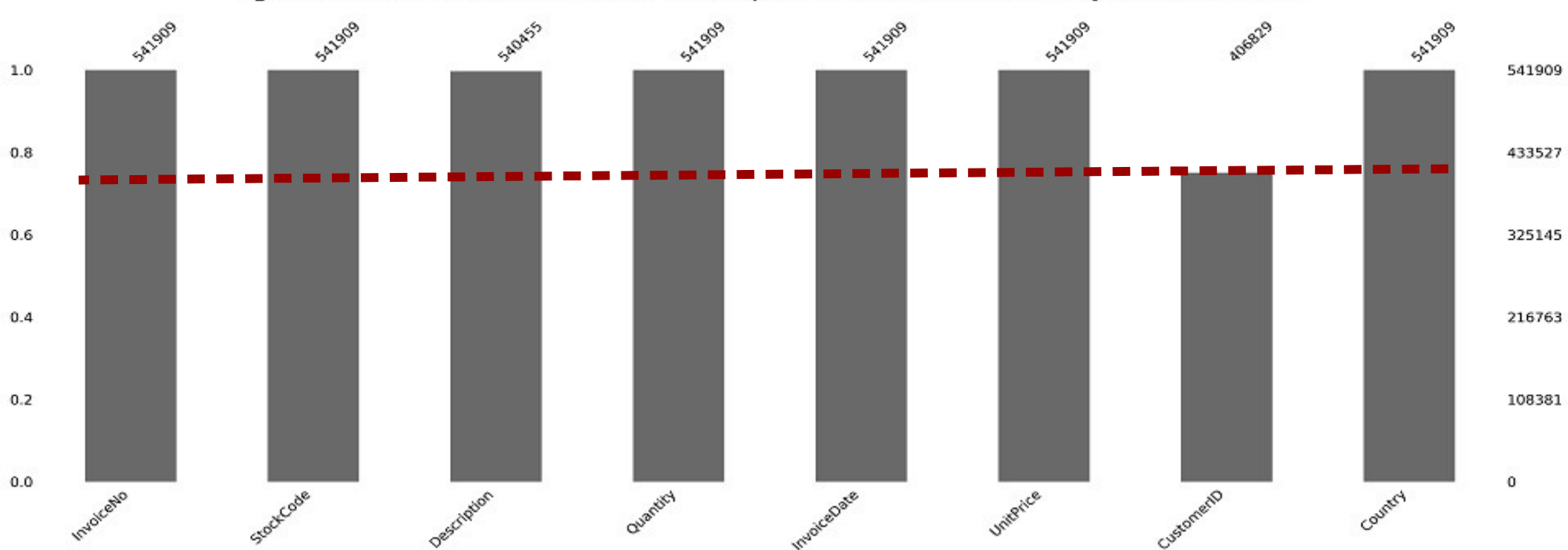
- Ces coûts sont très rarement indépendants des items proposés et de la nature des consommateurs.

Connaître les utilisateurs >  
Prédire leurs comportements  
futurs > établir des groupes  
d'utilisateurs caractéristiques >  
adapter la gouvernance, les  
stratégies de l'entreprise pour  
satisfaire les objectifs.

# Présentation des variables

- Variables peu nombreuses et de deux genres : qualitatives et transactionnelles
- Variables qualitatives difficilement utilisables
- Variables transactionnelles : prix par unité, nombre d'unités, origine et date.
- Plus de 4250 items pour près de 4000 utilisateurs
- La variable « utilisateur » n'a de sens que lorsque l'utilisateur est identifié

fig-1: Occurence en valeurs non manquantes des variables du jeux de donnée

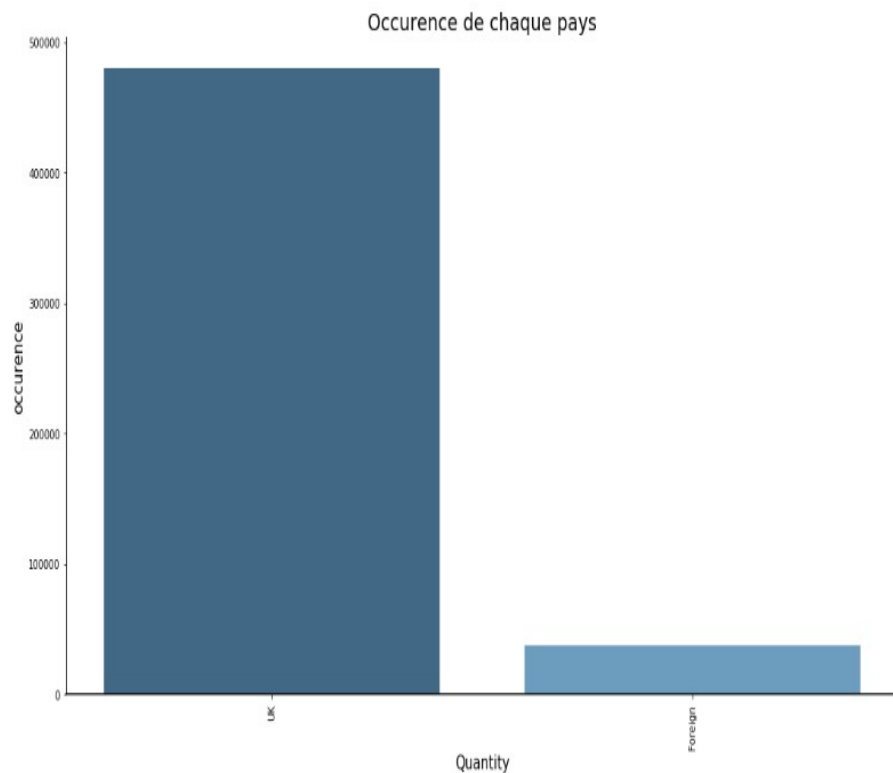


# Nettoyage des données

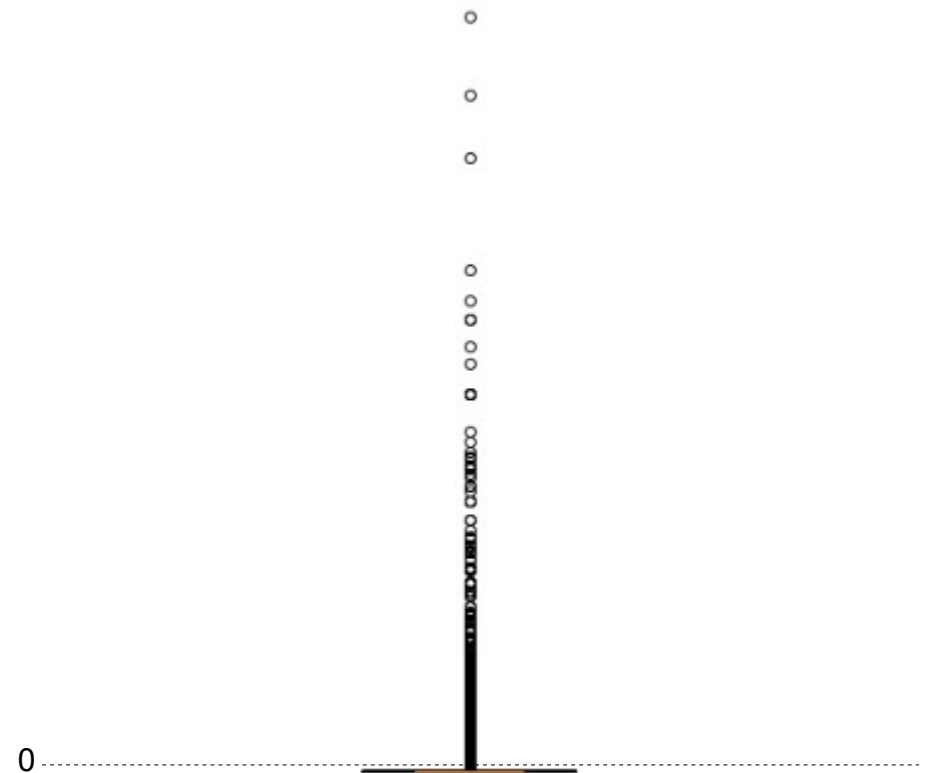
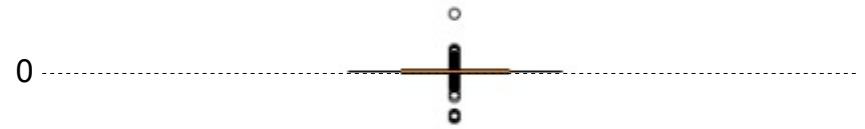
Présence de quantités négatives

Présence de données aberrantes

Traitement de la variable origine (prise en compte des clients anglais et des clients « étrangers »)

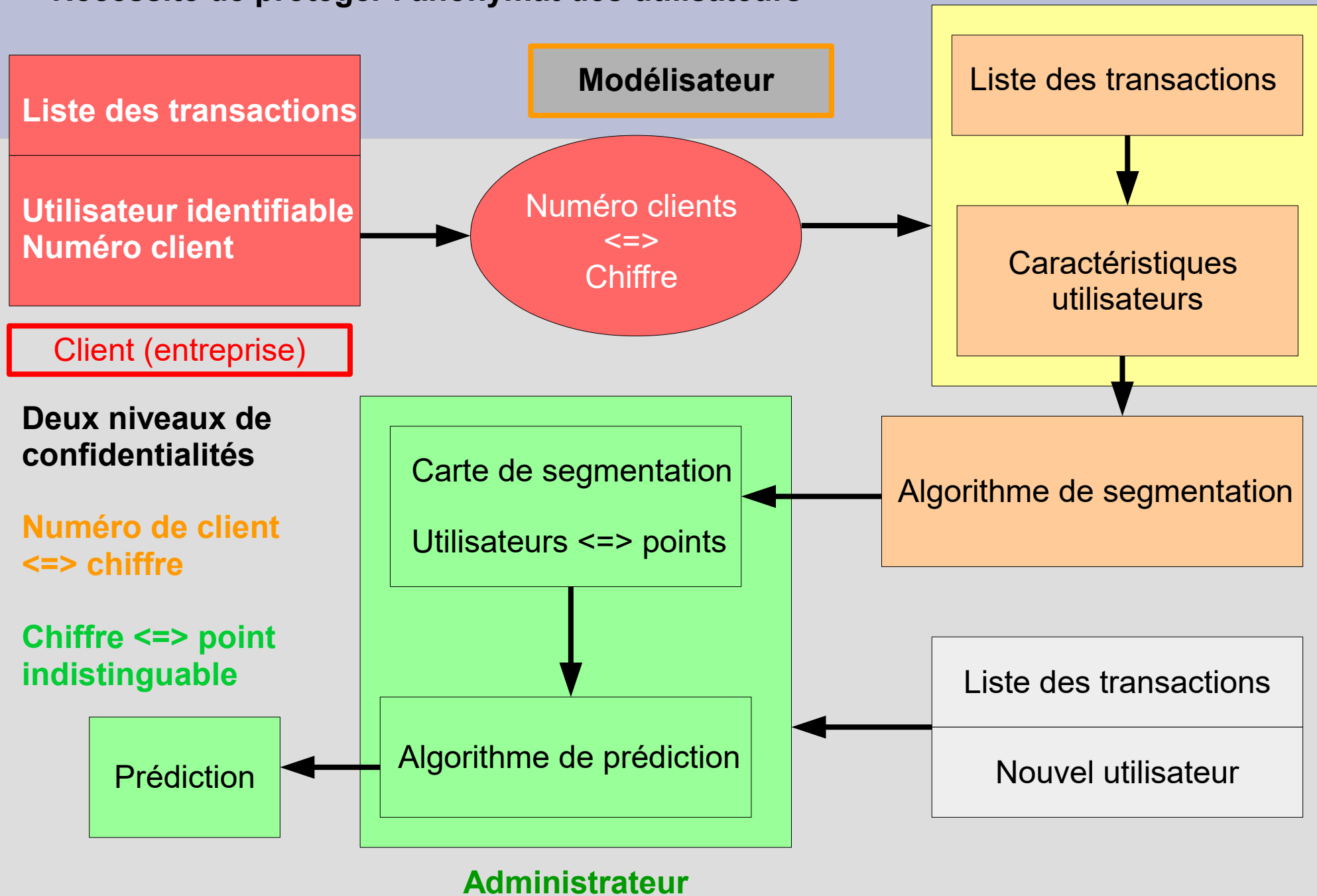


Distribution statistique de la quantité d'unités par transaction

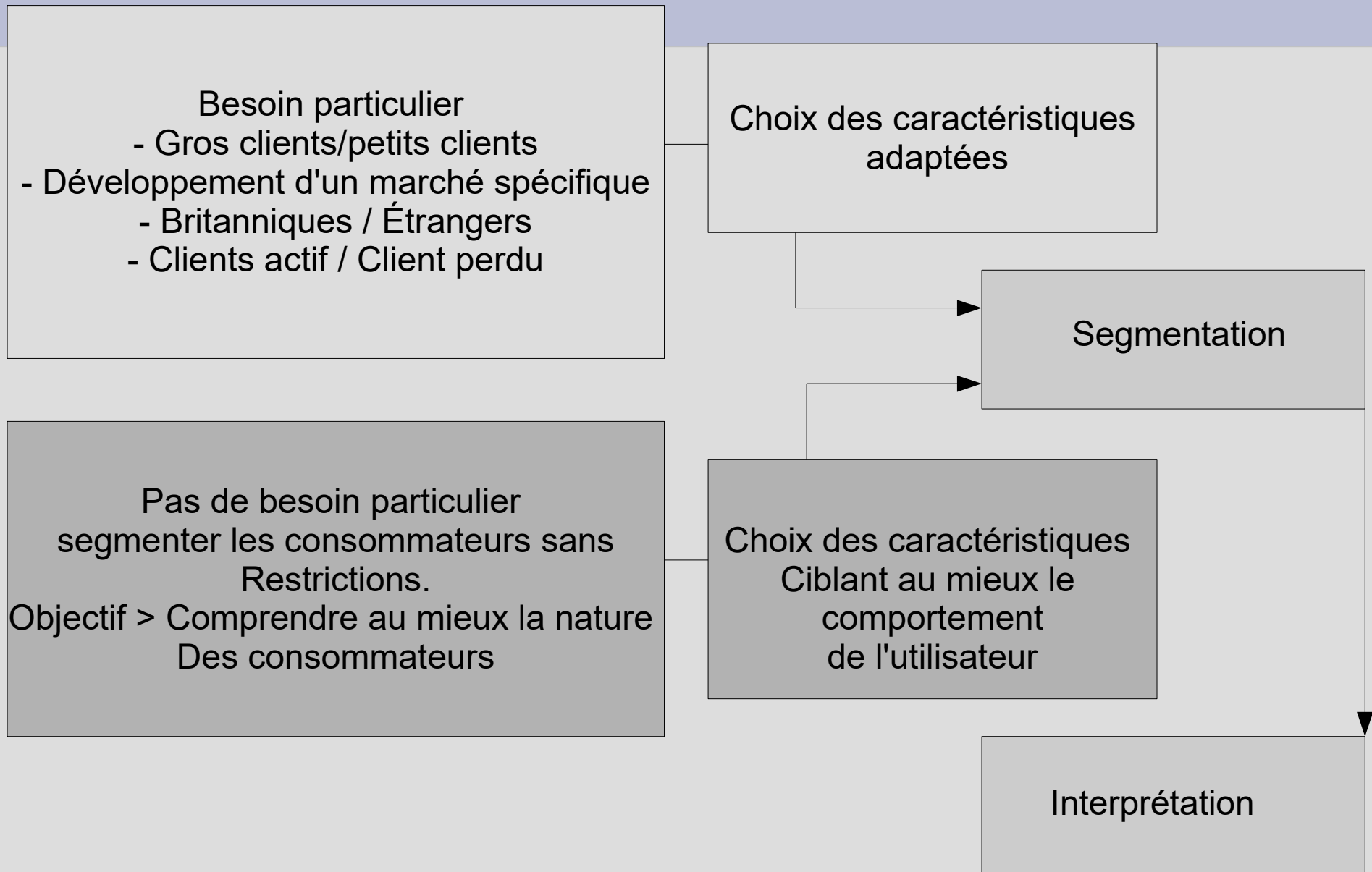


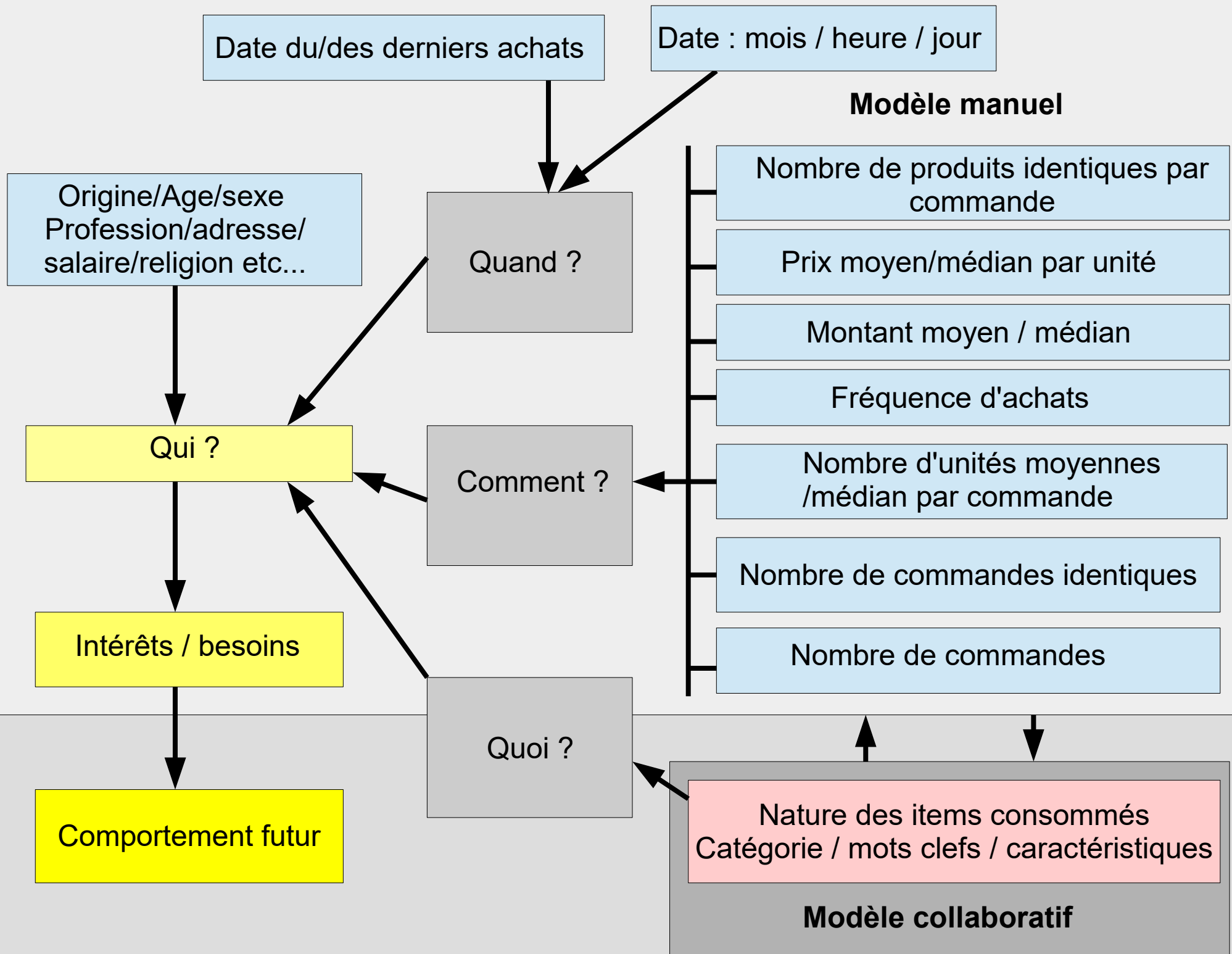
# Confidentialité

- Nécessité de protéger l'anonymat des utilisateurs

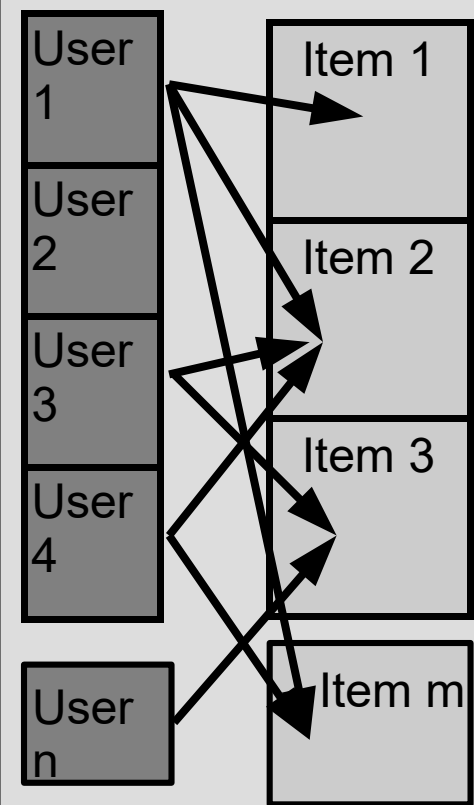


# Besoins de l'entreprise > Problématique spécifique / Comprendre les dynamiques d'achats





# • Caractéristiques Modèle collaboratif



Item 1	Item 2	Item 3

Item m

- 1
- 2
- 3
- 4
- 5
- 6

• Vecteurs  
Caractéristiques  
Items

• Vecteurs  
Caractéristiques  
utilisateurs

$\hat{R}_{i1u1}$ $R_{i1u1}$	$\hat{R}_{i2u1}$ $R_{i2u1}$	$\hat{R}_{i3u1}$ $R_{i3u1}$
$\hat{R}_{i1u2}$ $R_{i1u2}$	$\hat{R}_{i2u2}$ $R_{i2u2}$	$\hat{R}_{i3u2}$ $R_{i3u2}$
$\hat{R}_{i1u3}$ $R_{i1u3}$	$\hat{R}_{i2u3}$ $R_{i2u3}$	$\hat{R}_{i3u3}$ $R_{i3u3}$
$\hat{R}_{i1u4}$ $R_{i1u4}$	$\hat{R}_{i2u4}$ $R_{i2u4}$	$\hat{R}_{i3u4}$ $R_{i3u4}$

$\hat{R}_{imu1}$ $R_{imu1}$
$\hat{R}_{imu2}$ $R_{imu2}$
$\hat{R}_{imu3}$ $R_{imu3}$
$\hat{R}_{imu4}$ $R_{imu4}$

$\hat{R}_{i1un}$ $R_{i1un}$	$\hat{R}_{i2un}$ $R_{i2un}$	$\hat{R}_{i3un}$ $R_{i3un}$
--------------------------------	--------------------------------	--------------------------------

$\hat{R}_{imun}$ $R_{imun}$
--------------------------------

						User 1
						User 2
						User 3
						User 4

						User n
1	2	3	4	5	6	



Nombre de transactions utilisateur > item

$$p_{ui} = \begin{cases} 1 & r_{ui} > 0 \\ 0 & r_{ui} = 0 \end{cases}$$

## Présentation (mathématique) de l'algorithme et choix des hyper paramètres

$$\min_{x_*, y_*} \sum_{u, i} c_{ui} (p_{ui} - x_u^T y_i)^2 + \lambda \left( \sum_u \|x_u\|^2 + \sum_i \|y_i\|^2 \right)$$

$$c_{ui} = 1 + \alpha r_{ui}$$

Régularisation

hyper-paramètres

Degrés de confiance

vecteurs utilisateurs

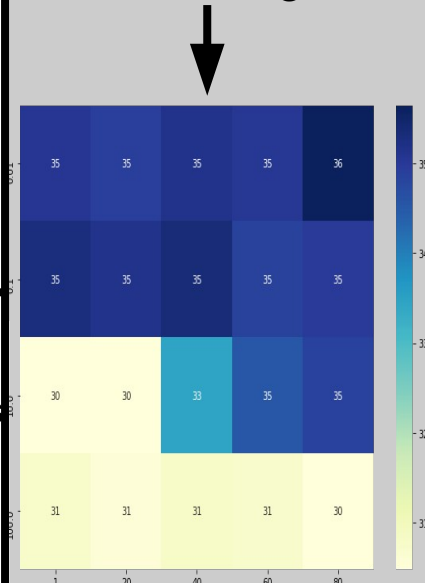
$$x_u = (Y^T C^u Y + \lambda I)^{-1} Y^T C^u p(u)$$

vecteurs items

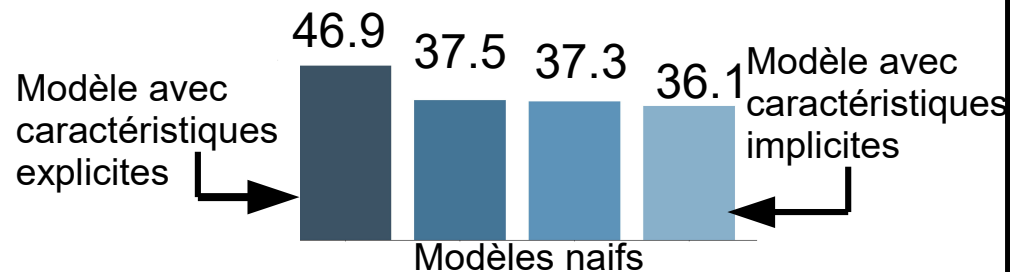
$$y_i = (X^T C^i X + \lambda I)^{-1} X^T C^i p(i)$$

- Meilleurs vecteurs utilisateurs > Meilleure prédiction
- Nécessité de définir un score de prédiction
- Définition des jeux de test et d'entraînement

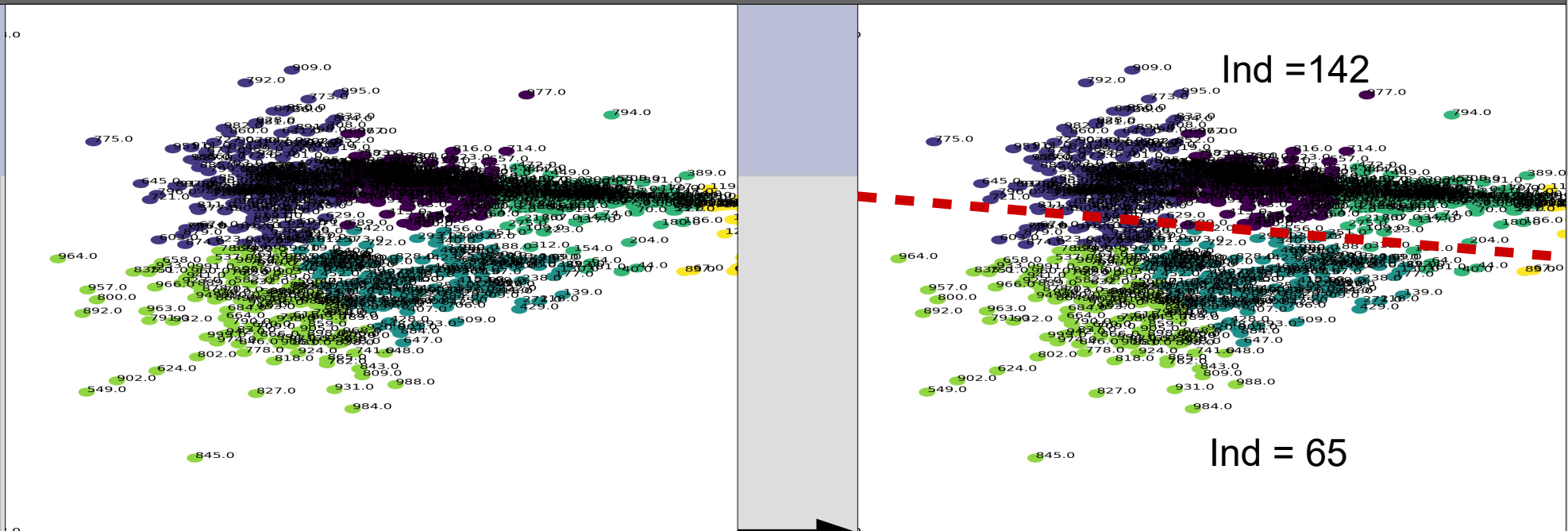
Détermination des hyper-paramètres par validation croisée sur grille



> Comparaison des modèles  
Algorithme de prédiction manuel  
distance item/barycentre des  
transactions dans l'espace items



# Clusters et interprétation : modèle collaboratif



Visualisation de l'historique des utilisateurs dans un cluster donné



Interprétation qualitative périlleuse  
Clients  $\Leftrightarrow$  entreprises concurrentes



Composante caractérisant le  
nombre de transactions effectuées

Définition d'un indice de  
popularité > « Ind »

Le modèle est peu performant  
pour prédire les comportements  
utilisateurs

Le modèle est peu performant  
pour prédire les comportements  
utilisateurs > mais fait émerger  
une caractéristique item basée sur  
la popularité

Prendre en compte la distribution statistique des transactions

Transactions :

- ~~Nombre de transactions~~
  - **Fréquence de transaction**
  - Montant moyen par transaction
  - Nombre d'unités moyennes par transaction
  - Prix moyen unitaire par transaction
- 
- Nombre de produit différent par commandes
  - Nombre de transaction identiques

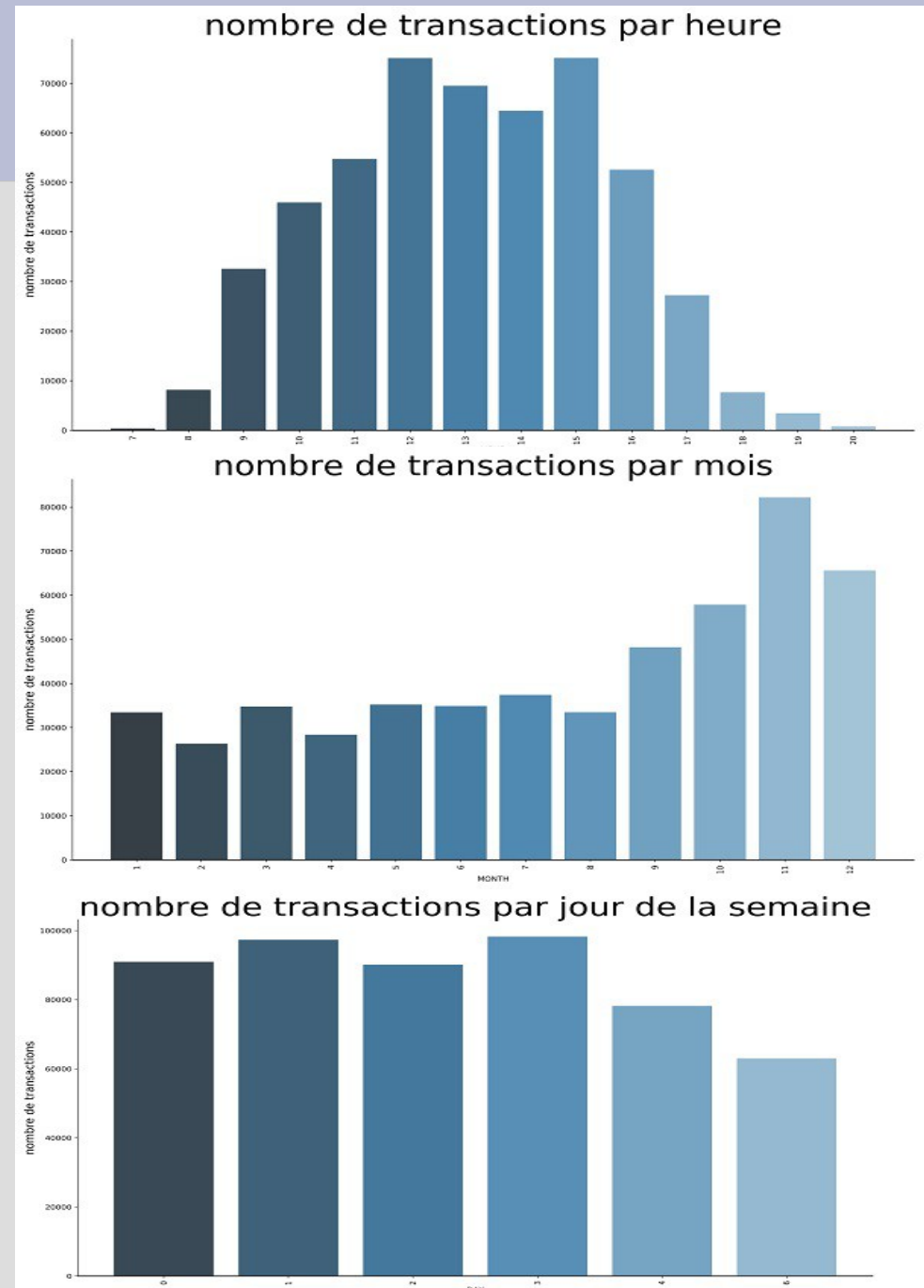
Prendre en compte les différences de consommation liées à l'origine

Pays d'origine > 0 ou 1 (UK ou étranger)

Prendre en compte les nouveaux consommateurs / les clients perdus

- Date de la dernière commande

# Choix des variables et features engineering



# Choix des variables et features engineering

Prendre en compte l'impact du temps.

**Jour de la semaine et heure**

- > particulier / entreprise
- > grande ou petite entreprise

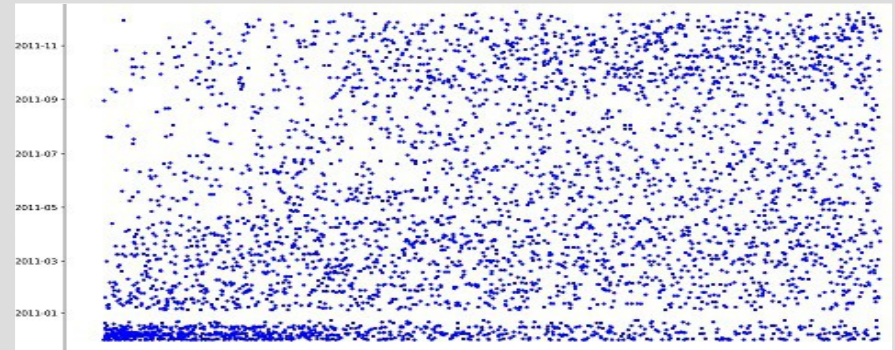
- Jour (heure) de maximum de consommation
- Probabilité de transaction le jour du maximum de consommation

Mois (périodes de vacances/nature des produits etc...)

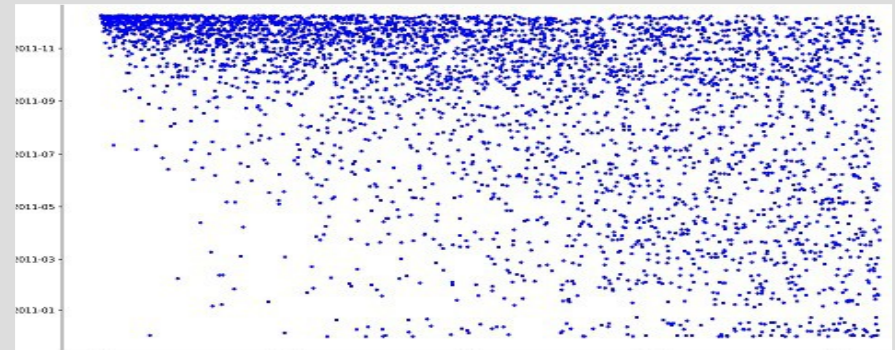
- > Entreprise (fonctionnement interne/catégorie etc...)
- > particulier (religion/préférences etc...)

- nombre de commandes par mois
- nombre de transaction au sein d'une commande par mois
- quantité moyen d'unités par transaction par mois
- prix moyen par transaction par mois

Date de la première commande par utilisateur



Date de la dernière commande par utilisateur



Imputation par la moyenne

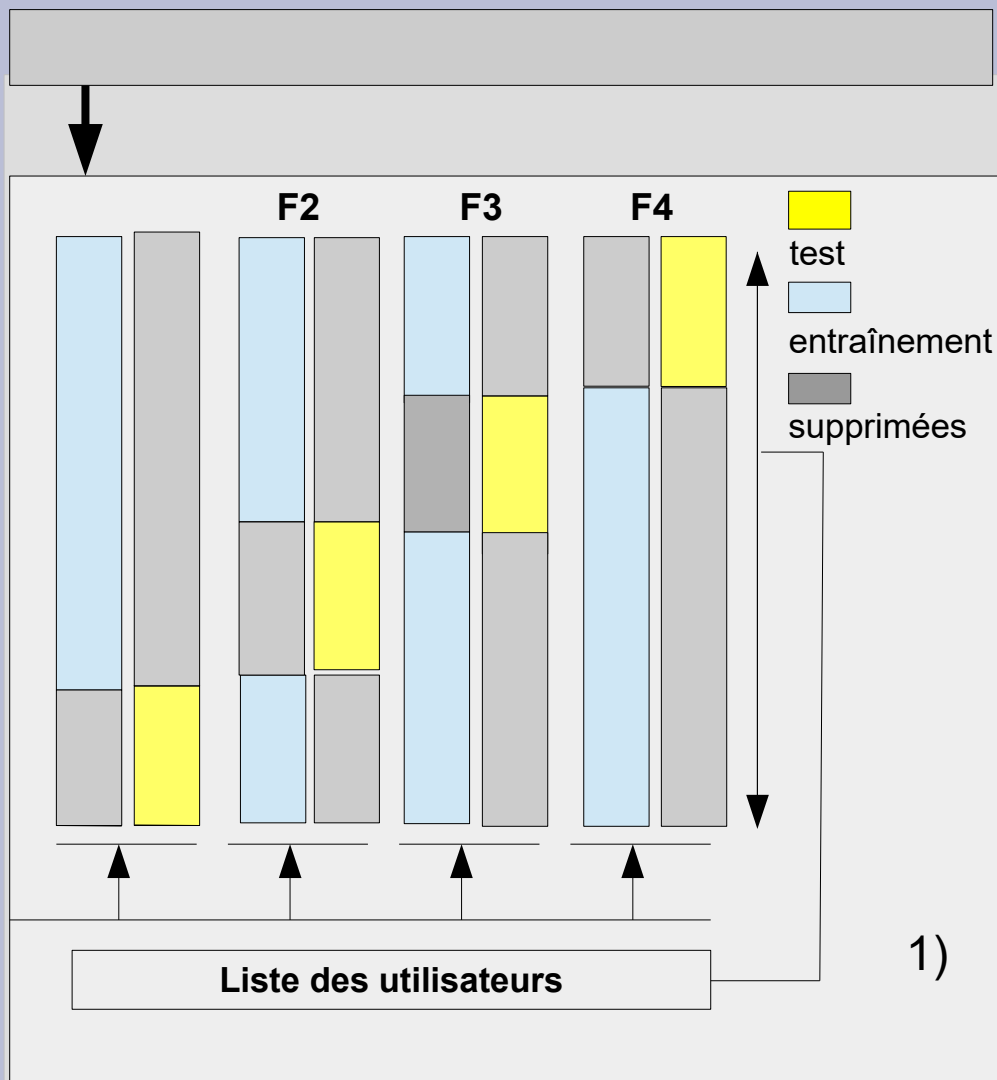
Imputation par régression

- **Technique** > effectuer une validation croisée sur grille pour trouver automatiquement la constante de régularisation > temps de calcul
- **Conceptuel** > comment connaître réellement qu'elle est la première commande ?

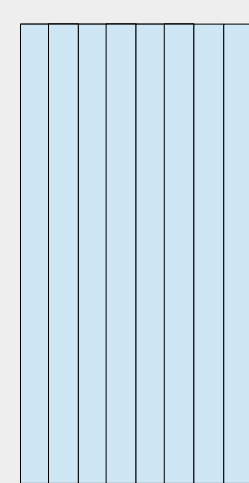
**- Établir plusieurs modèles de prédiction en fonction de la commande requête**

# Apprentissage non supervisé

Entrée > Jeu Original : liste des transactions



Génération des caractéristiques de clustering



caractéristiques

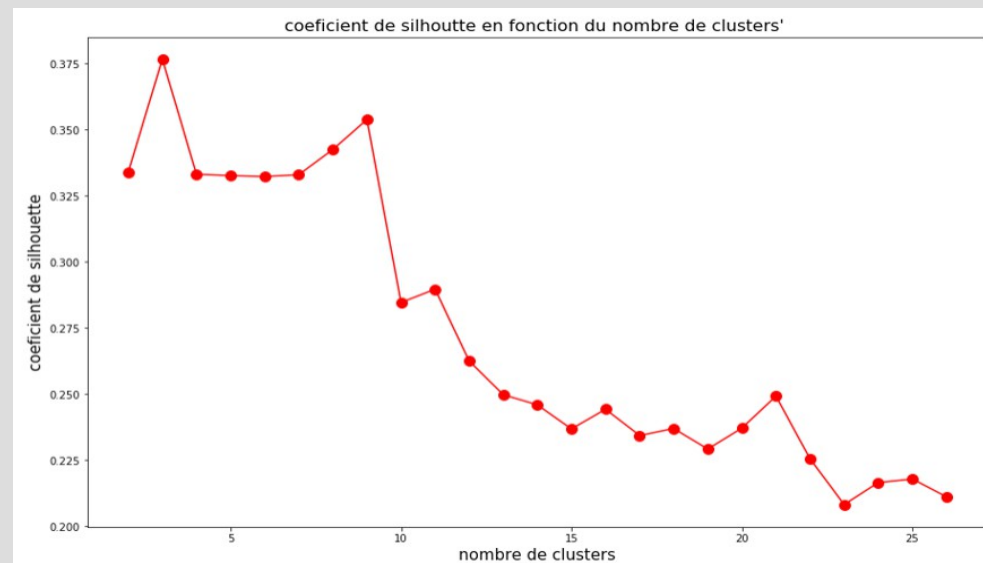
utilisateurs

2)

Segmentation  
K-means pour  
plusieurs K

Évaluation >  
structure du  
cluster :  
Coefficient de  
silhouette

3)



Sortie : Labels  
par client

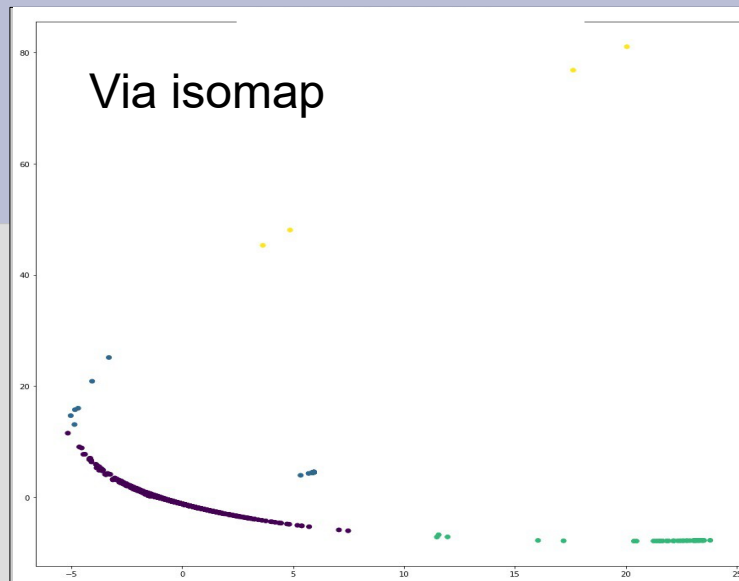
Interprétation  
qualitative

Étude des caractéristiques  
moyennes par cluster



# Segmentation > modèle manuel

Quel comportement souhaite t-on capturer ?

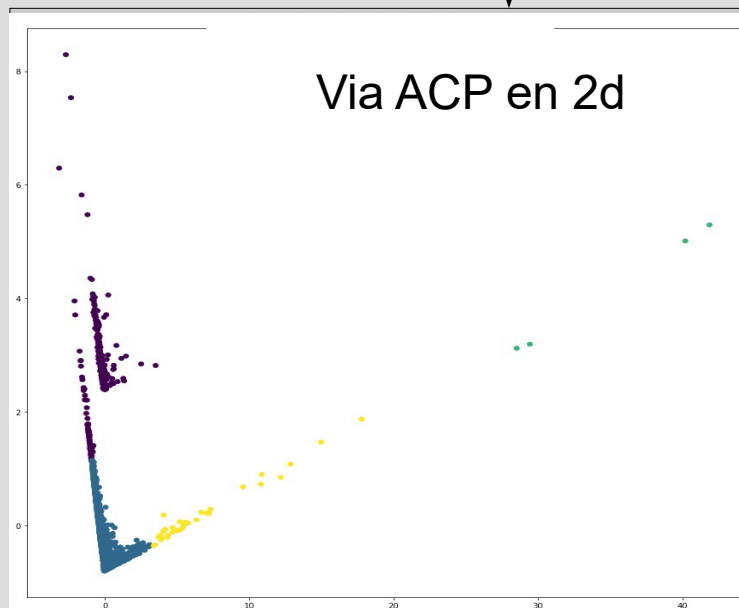


Taille des  
clusters /  
Nombre de  
clusters  
minimum /

maximum

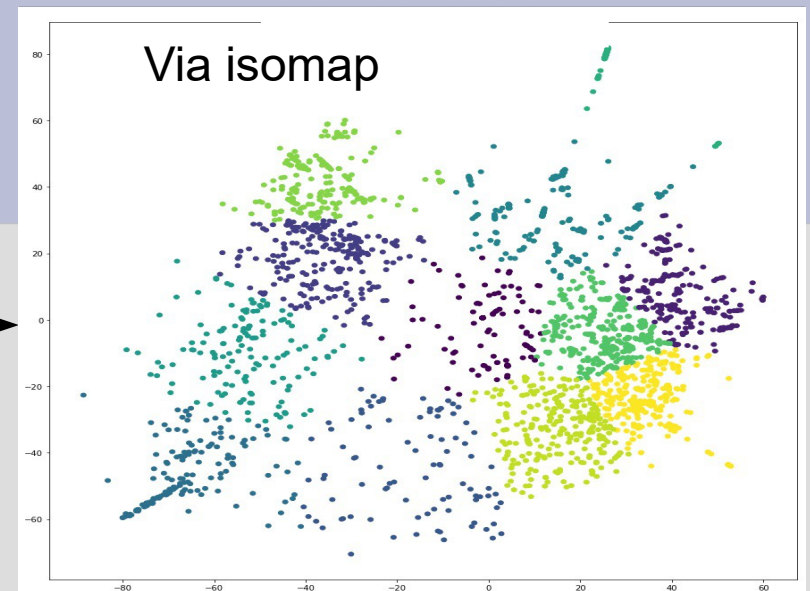
Dynamique  
linéaire / non  
linéaire

Standard Scaler

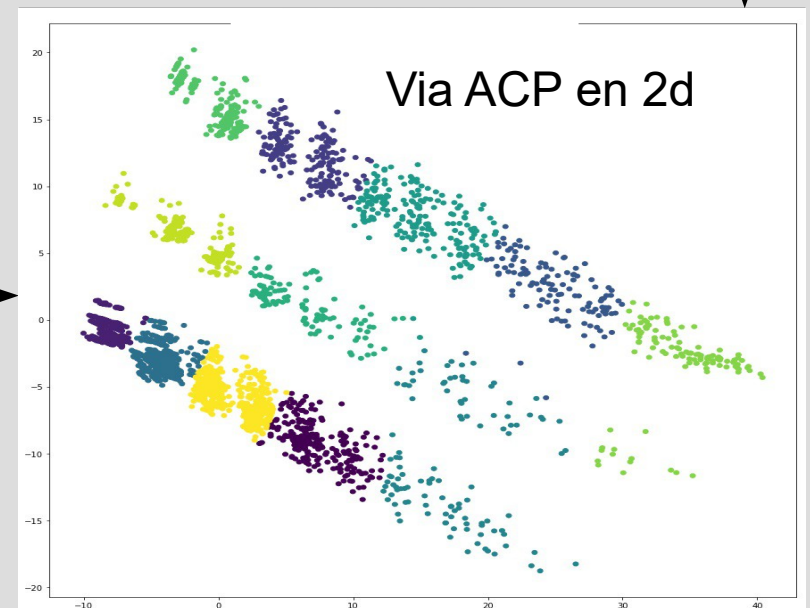


Données  
« écrasées »  
par les  
outliers

Utiliser un  
standardiser  
différent



Normal Quantil Transformer

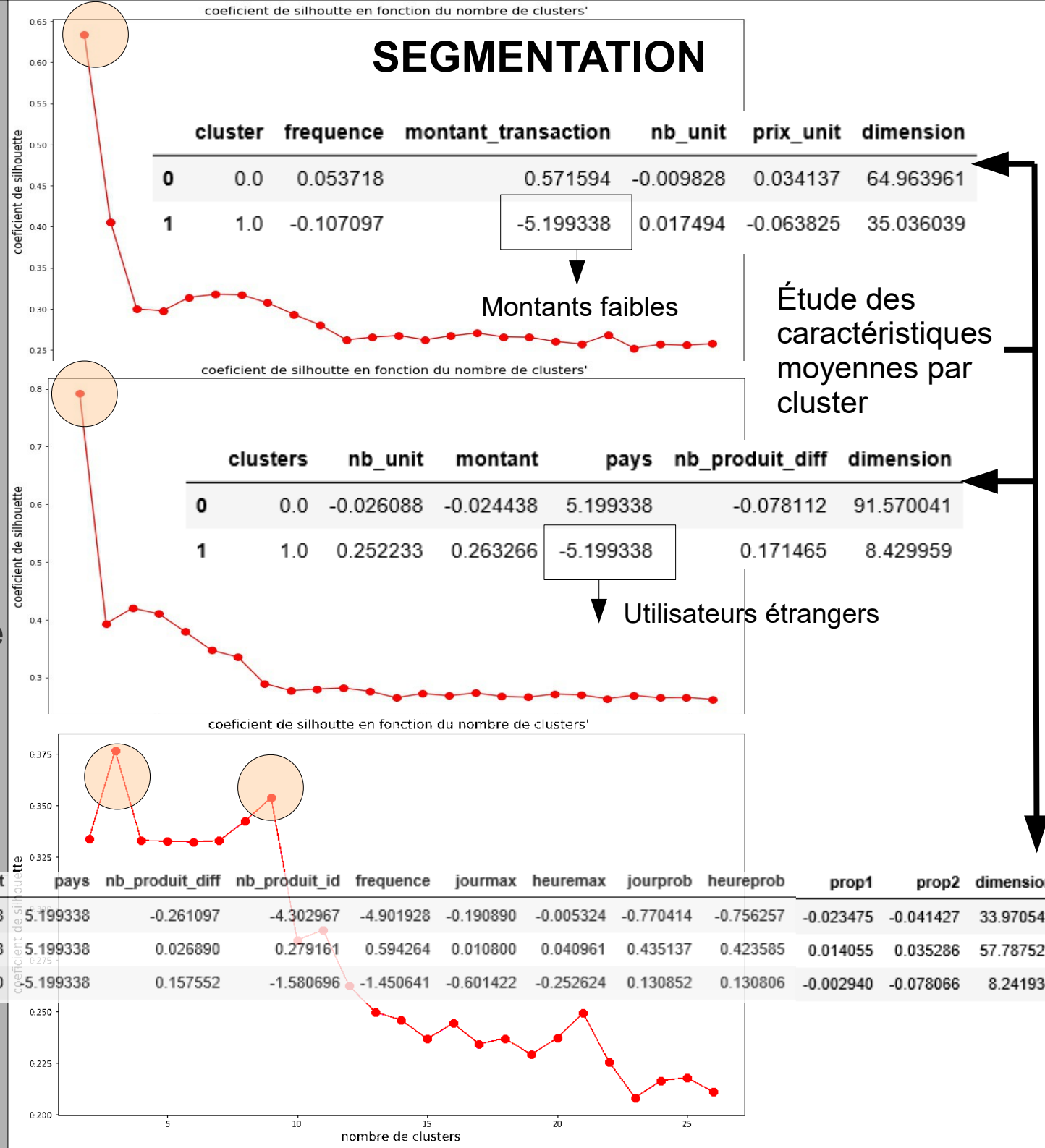


**Scénario 1 :**  
**Objectif >**  
Distinguer les petits  
clients des gros  
clients

**Scénario 2 :**  
**Objectif >** évaluer  
les différences de  
comportements des  
clients étrangers

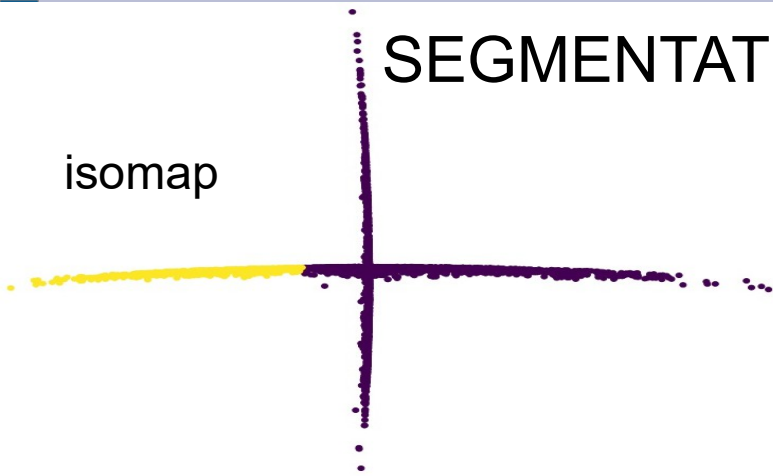
**Scénario 3 :** Ajout de  
caractéristiques  
temporelles,  
popularité, date de  
dernière transaction

clusters	last_trans	nb_unit	prix_unit	montant	pays	nb_produit_diff	nb_produit_id	frequence	jourmax	heuremax	jourprob	heureprob	prop1	prop2	dimension	
0	0.0	0.643874	-0.076650	0.004649	-0.043788	5.199338	-0.261097	-4.302967	-4.901928	-0.190890	-0.005324	-0.770414	-0.756257	-0.023475	-0.041427	33.970542
1	1.0	-0.375659	0.002486	0.006982	-0.014038	5.199338	0.026890	0.279161	0.594264	0.010800	0.040961	0.435137	0.423585	0.014055	0.035286	57.787527
2	2.0	-0.022331	0.249555	-0.066632	0.281910	5.199338	0.157552	-1.580696	-1.450641	-0.601422	-0.252624	0.130852	0.130806	-0.002940	-0.078066	8.241930

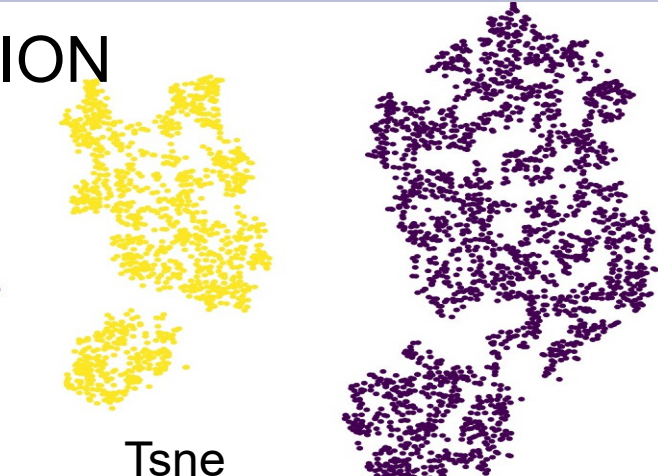


# SEGMENTATION

isomap

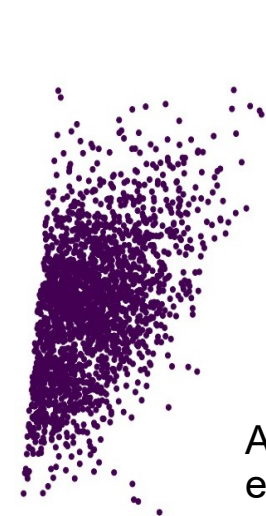


Tsne



**Scénario 1**

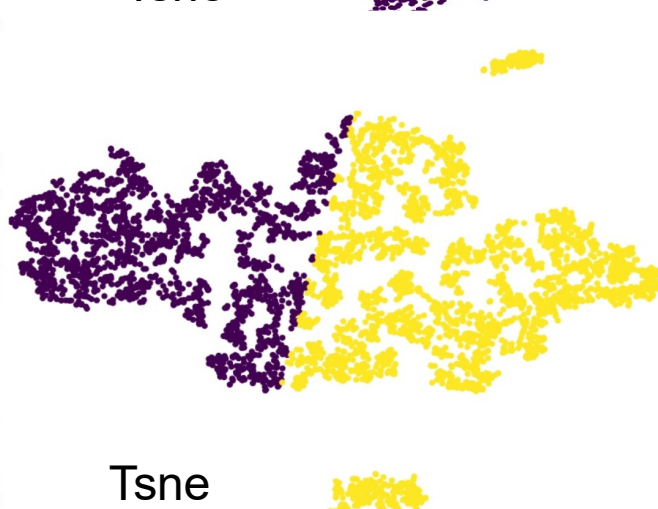
ACP variance  
expliquée : 89%



isomap



Tsne

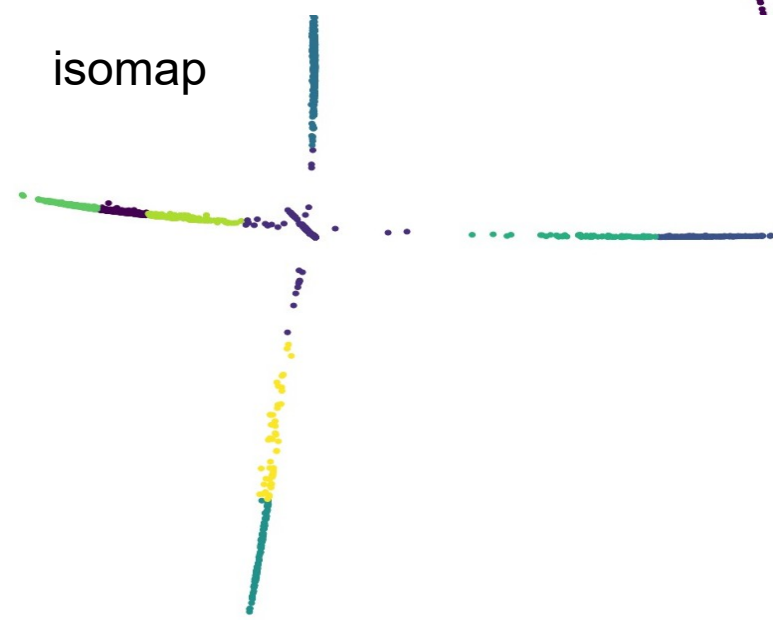


**Scénario 2**

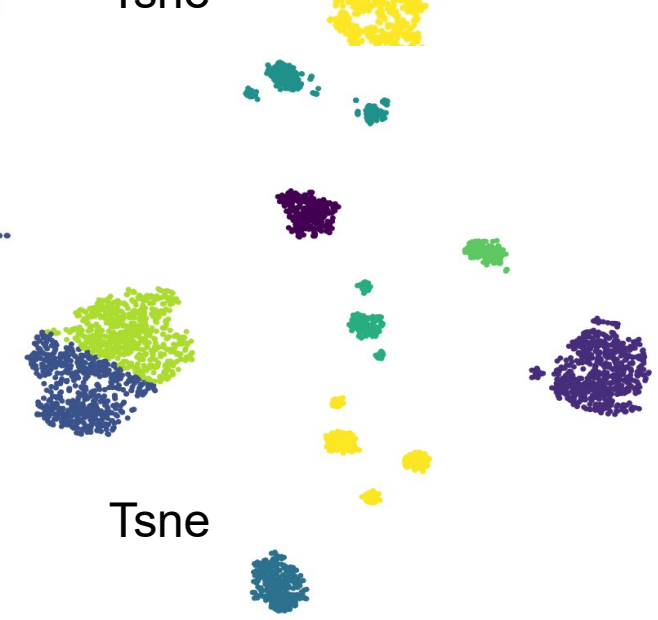
ACP variance  
expliquée : 92%



isomap

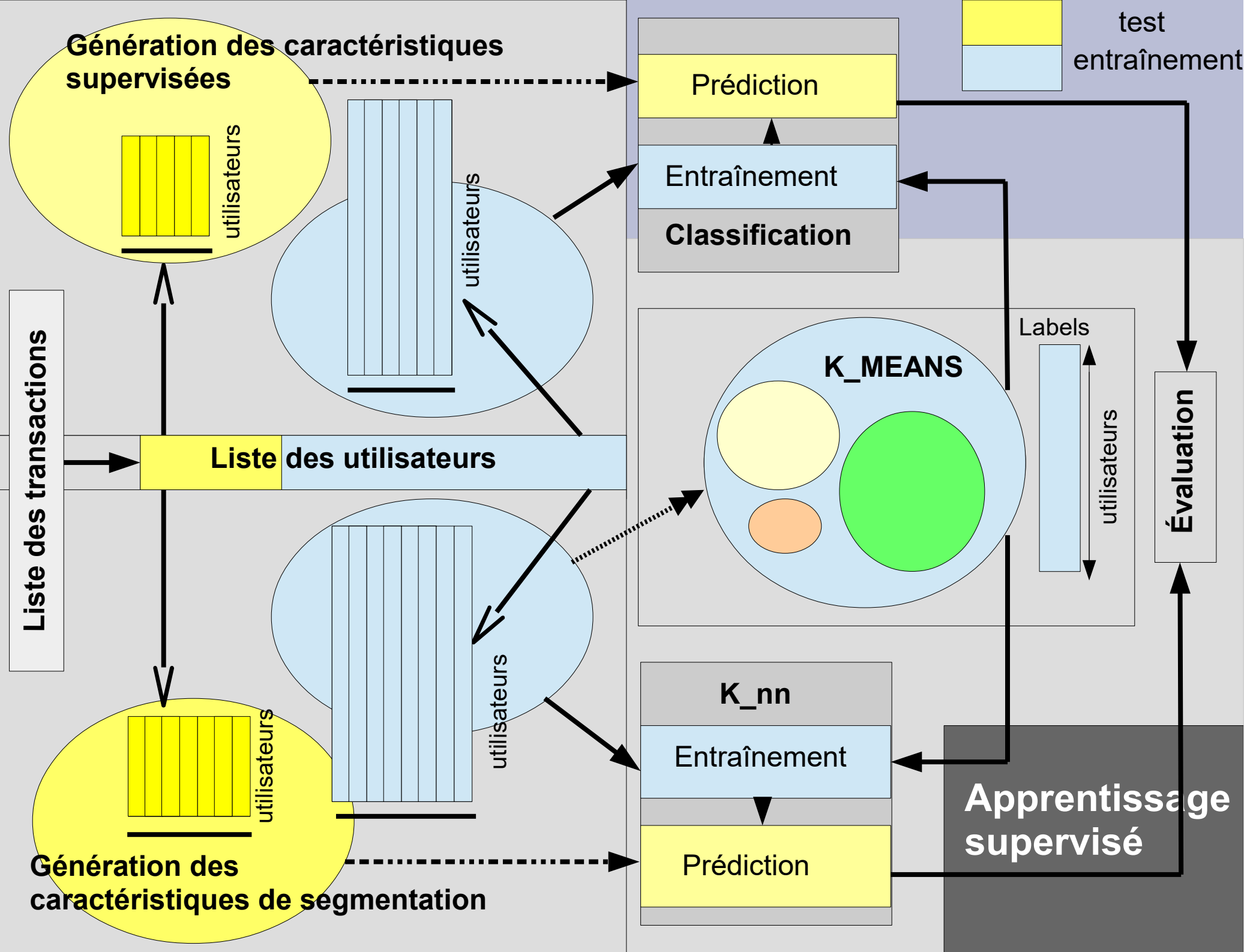


Tsne



**Scénario 3**

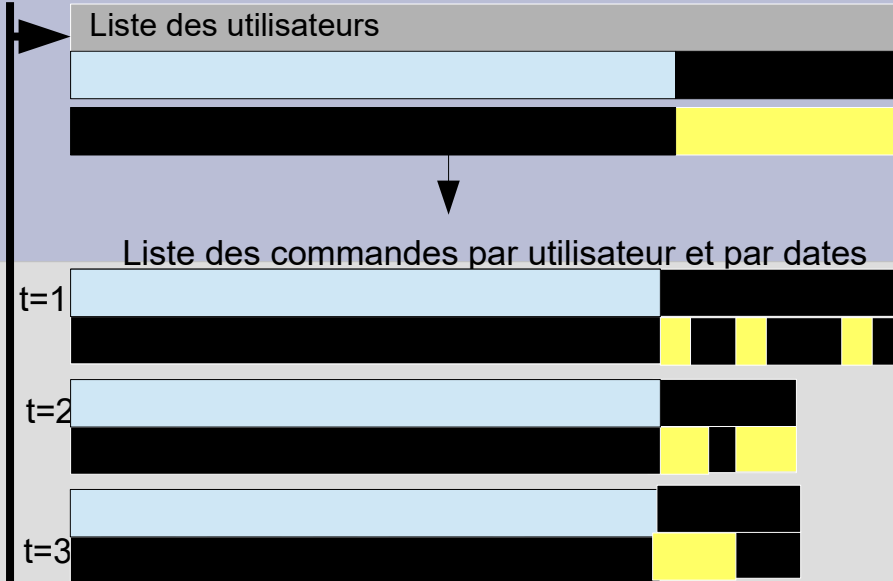




# Étude de la qualité de prédiction en fonction du scénario

entraînement  
test

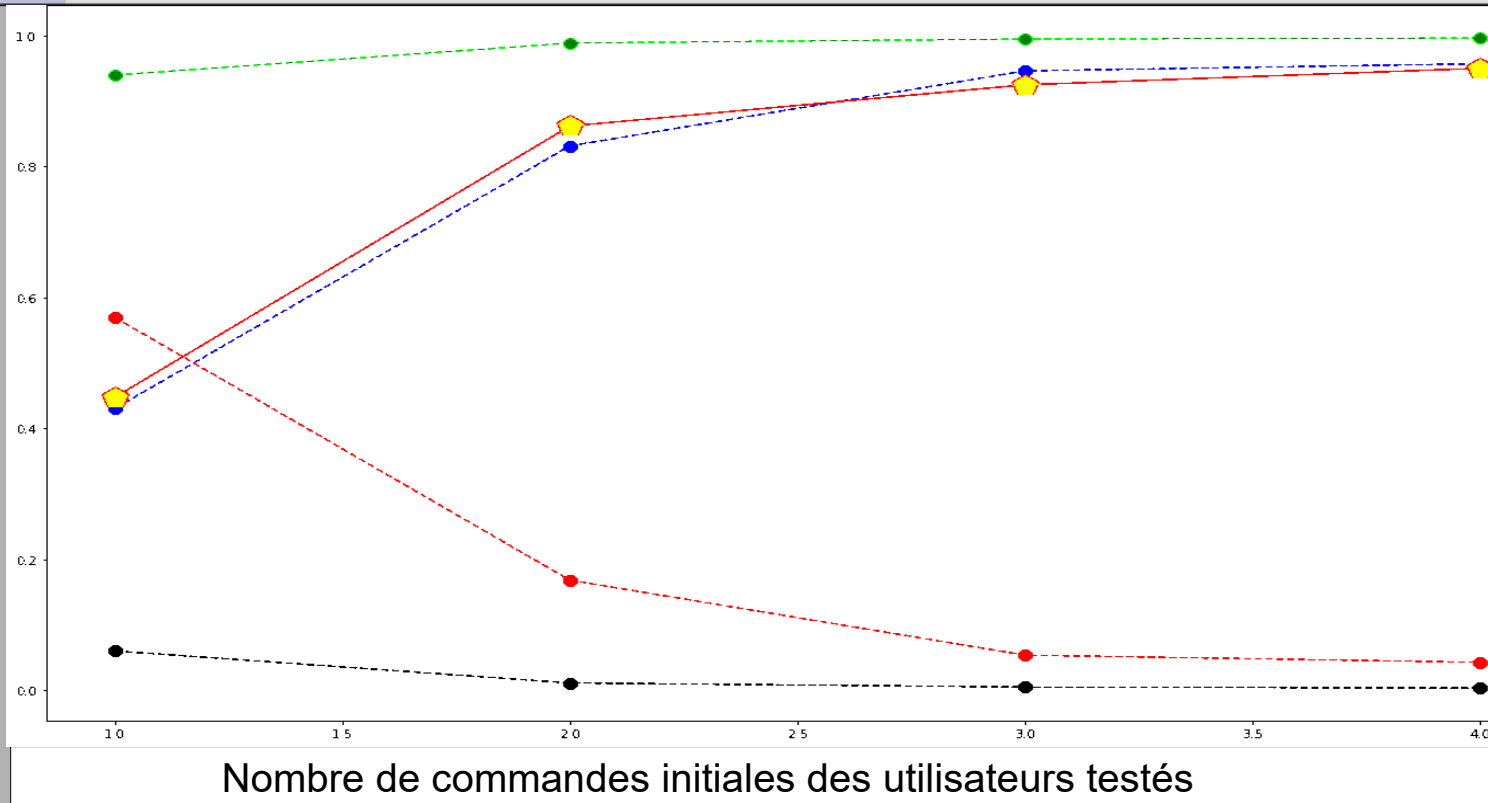
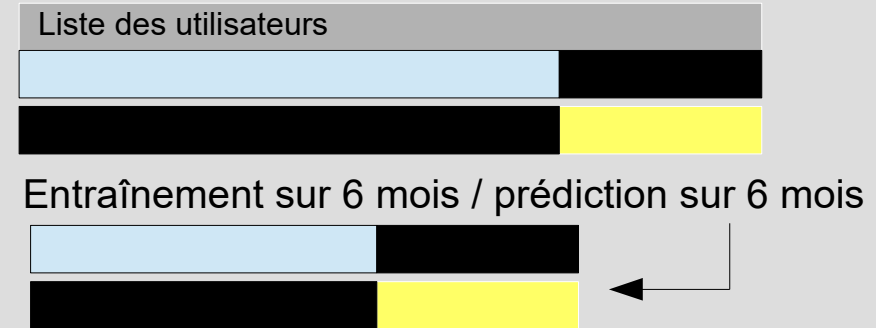
Liste des transactions



Liste des transactions

## Mise en place d'un scénario spécifique

> Restriction sur le choix des données utilisées pour l'apprentissage supervisé



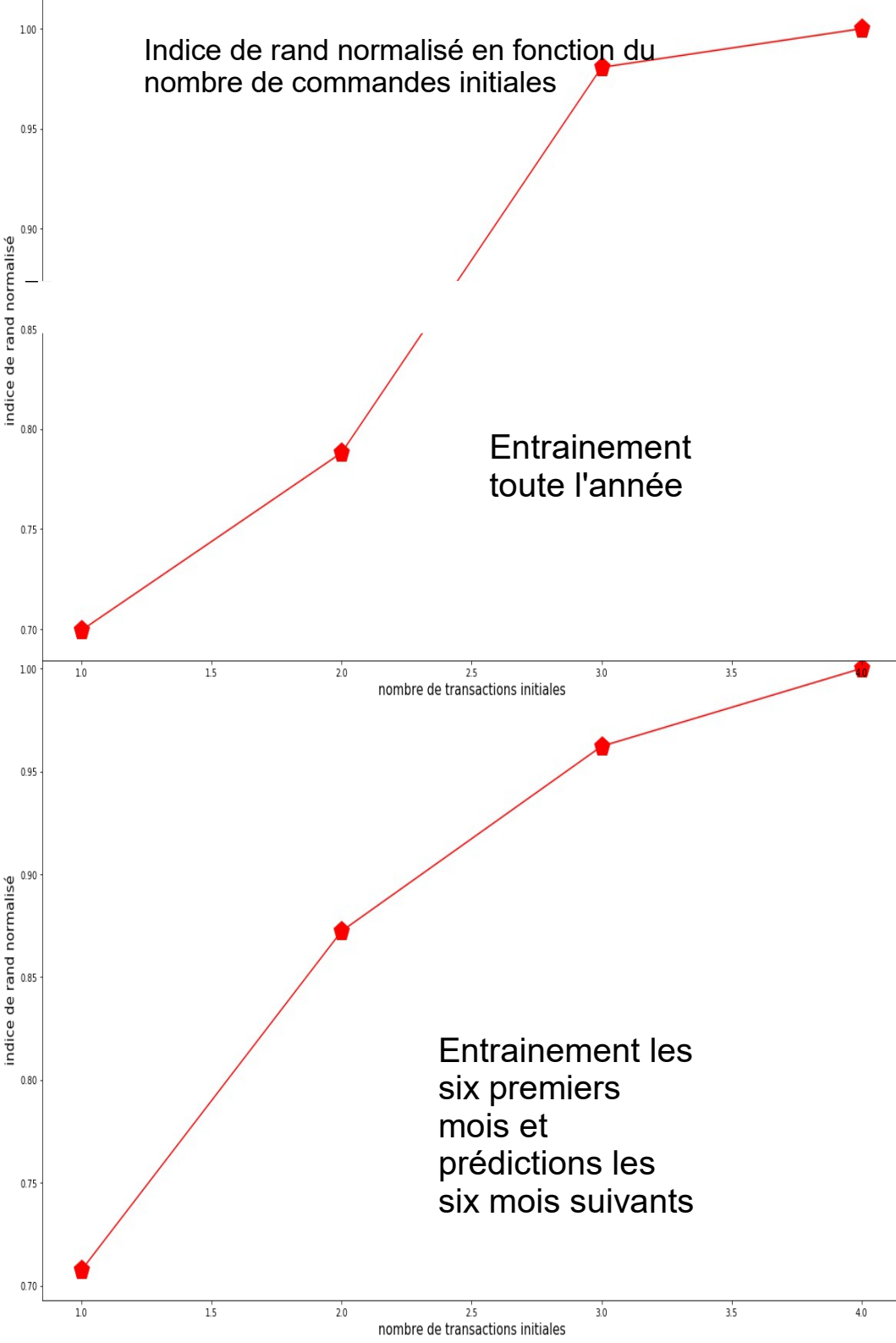
Vrais négatif

Indice de rand  
normalisé

Vrais positif

Faux négatif

Faux positif



# Apprentissage supervisé

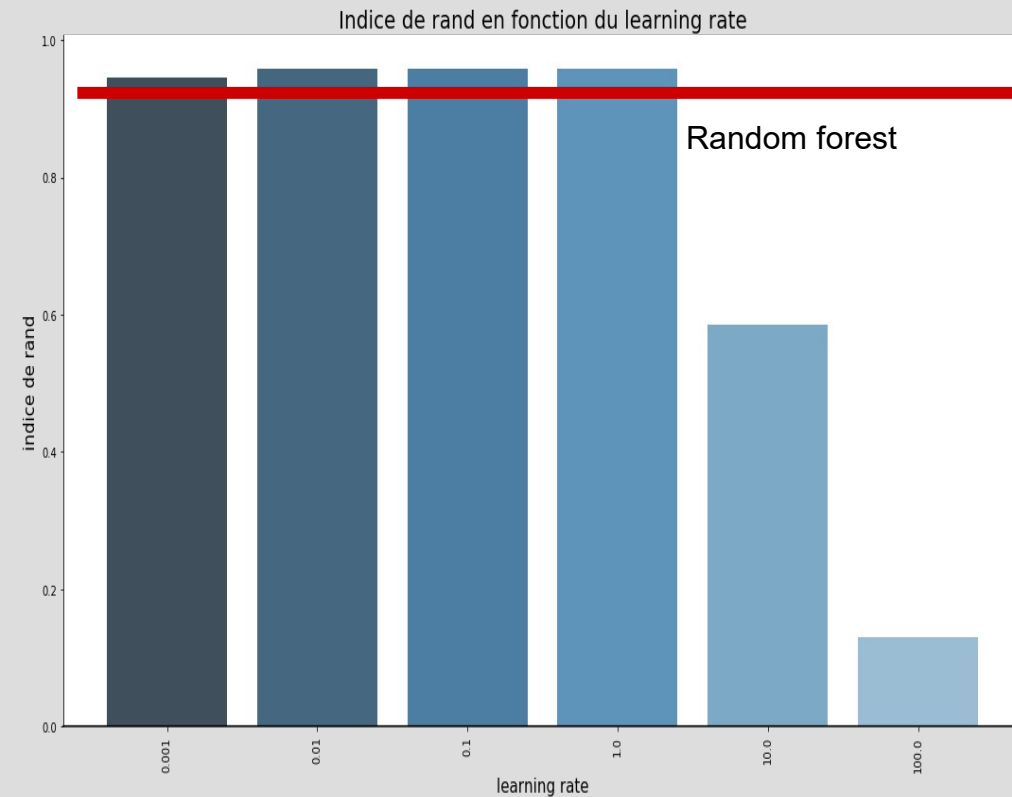
Comparaison des modèles

## Random forest

classifieur

Gradient boosting  
classifieur

- Détermination de l'hyper-paramètre « alpha » par grille



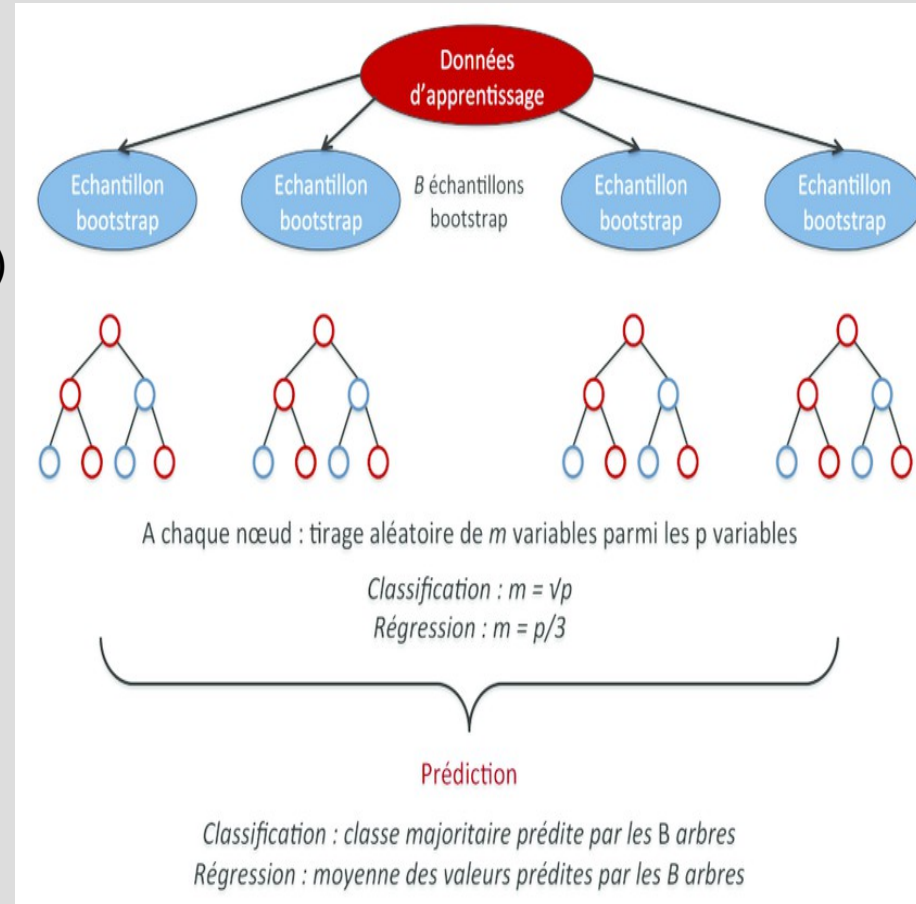
# Algorithme du bagging > présentation théorique

Méthode ensembliste > combinaison d'apprenants faibles pour obtenir des résultats plus performants

- Bootstrap > effectuer  $n$  échantillonnages avec remise  
effectuer la moyenne ou vote à la majorité des apprenants faibles appliqués aux  $n$  échantillons.  
Modèles individuels simples > forte variance / erreur répartie de manière normale  
conjugaison des modèles > réduction de la variance et donc prédiction plus fiable

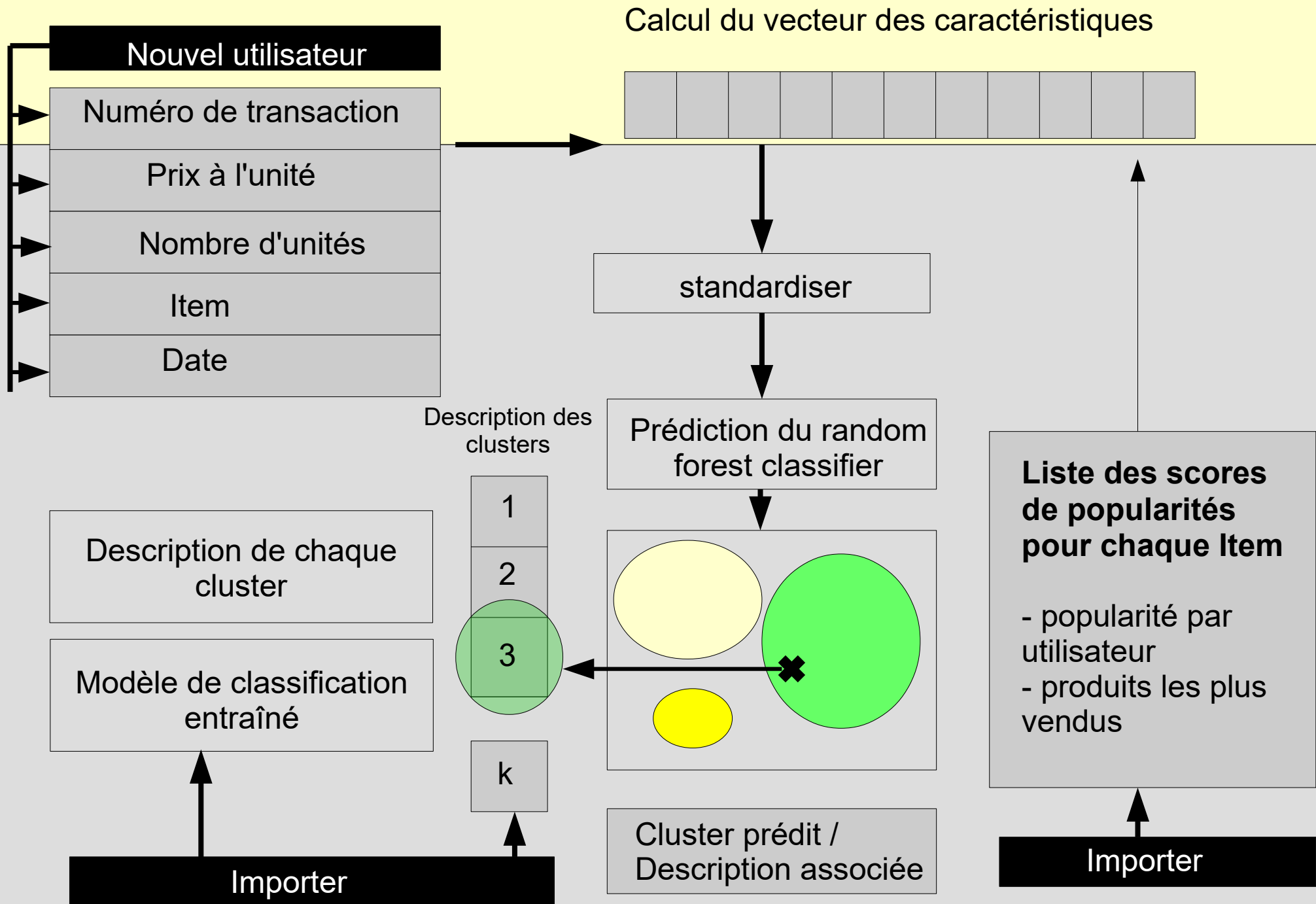
# Forêts aléatoires > présentation théorique

- Utilisation d'arbre de décisions :  
chaque branche > détermination de la caractéristique et de la valeur de séparation  
> utilisation d'heuristique (ex : index de Gini)  
  
nombre de branches optimal > compromis biais / variance
- Forêt > Agrégation d'arbres qui sur-apprennent individuellement
  - application du bootstrap pour n arbres différents
  - prise en compte d'un sous ensemble aléatoire de features à chaque nœud > réduction de la corrélation entre les arbres



- Avantages : mémoire > taille du dataset initial / pas d'overfitting / complexité raisonnable

# Prédiction



# Conclusion

Modèle collaboratif

Modèle manuel

Détermination des objectifs

Comprendre au mieux la nature des comportements client

Séparer les utilisateurs en jeux de test et d'entraînement

Étude des variables

Variables manuelles  
(transactionnelles/temporelles etc..)

Variables descriptives des items

Qualité de prédiction > Qualité des vecteurs utilisateurs générés

- détermination des jeux de test et d'entraînement
- détermination d'un score
- validation croisée sur grille

- définition des caractéristiques supposées pertinentes
- détermination d'un algorithme de prédiction pour la comparaison des modèles
- définition des jeux de test et d'entraînement

Segmentation via K-means / Choix du nombre de clusters via l'indice de silhouette

Interprétations qualitatives difficiles  
Deux caractéristiques interprétables émergent.

- Interprétation des scénarios via la moyenne des caractéristiques par cluster

Prédiction > modèles supervisés > différents scénarii (nombre de commandes initiales, échelle de temps...)

Évaluation > Indice de rand normalisé

Nouveau client > fonction python