

Segmentez des clients d'un site e-commerce

Problématique et objectifs

La société Olist souhaite réaliser une segmentation des clients

Objectifs:

- Proposer un clustering de clients adapté
- Proposer une fréquence de maintenance pour le modèle de segmentation proposée

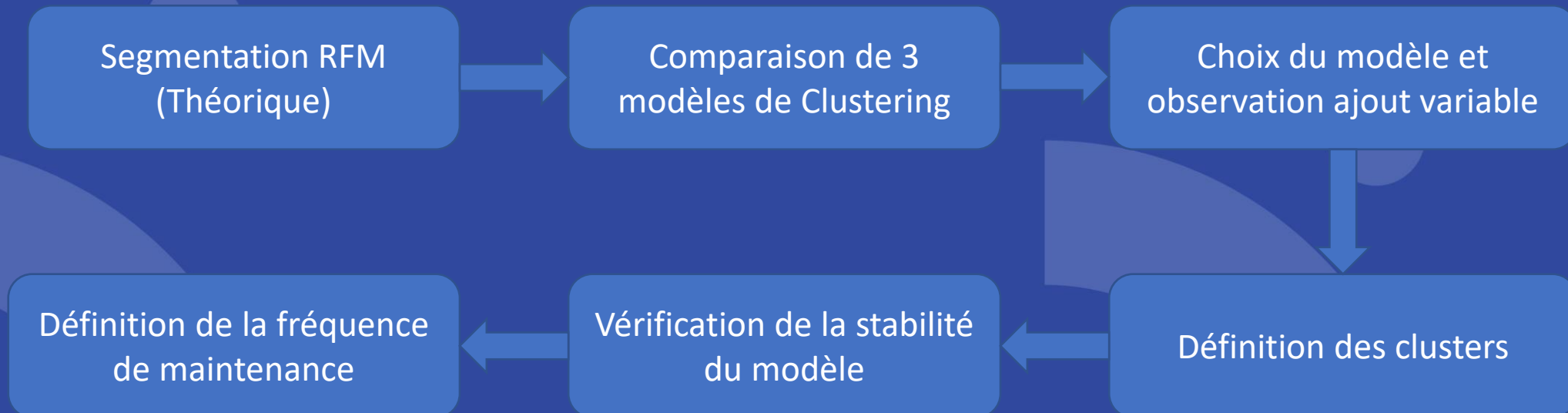
Objectif secondaire:

- Observer l'effet des variables sur le modèle

Stratégie choisie

La stratégie choisie pour traiter le dataset est la segmentation RFM qui se base sur:

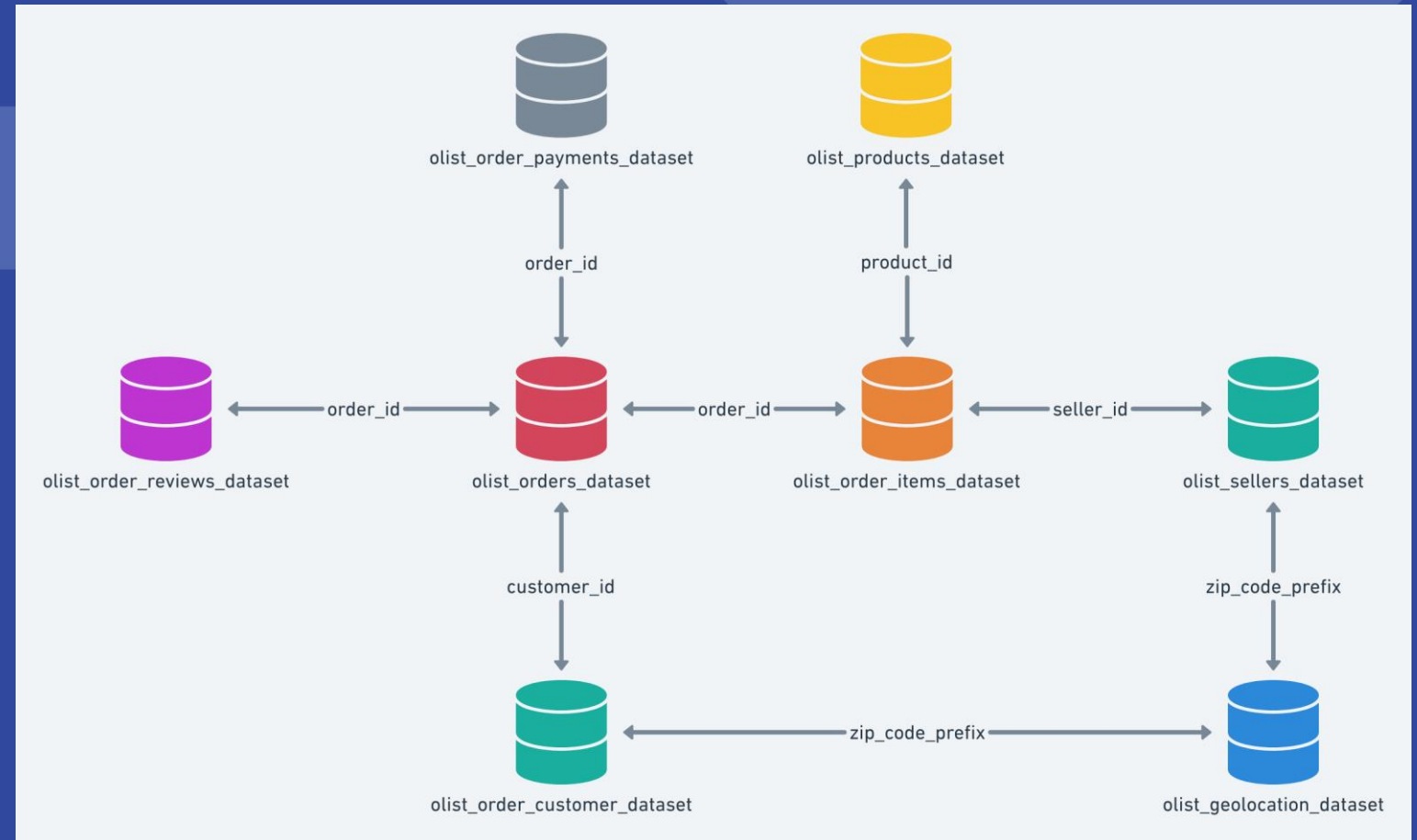
- La Récence du dernier achat
- La Fréquence d'achat
- Le Montant total des achats



Jeu de données

9 fichiers:

- Liste des clients
- Liste des commandes
- Liste des notes et commentaires
- Liste des vendeurs
- Liste des produits
- Liste des paiements
- Liste des géolocalisation
- Liste des produits par commandes
- Liste des catégories dans les 2 langues



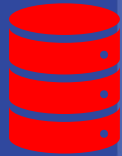
Nettoyage et traitements des données

Sélections des fichiers et variables d'intérêt:



Customers:

- Customer id
- Customer unique id
- City
- State
- Zip Code



Orders:

- Order id
- Order purchase time
- Order status



Payments:

- Payments id
- Payment sequential
- Payment installments
- Payment value



Reviews:

- Order id
- Reviews Score



Items:

- Order id
- Product id
- Order items



Products:

- Product id
- Category

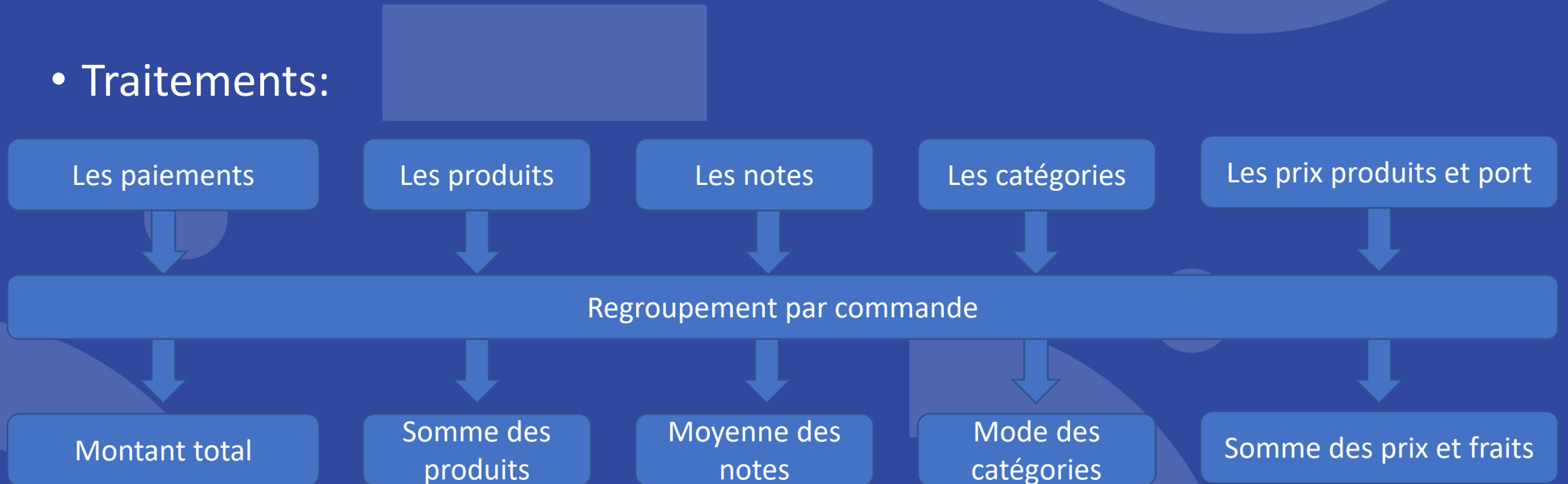
- Imputation

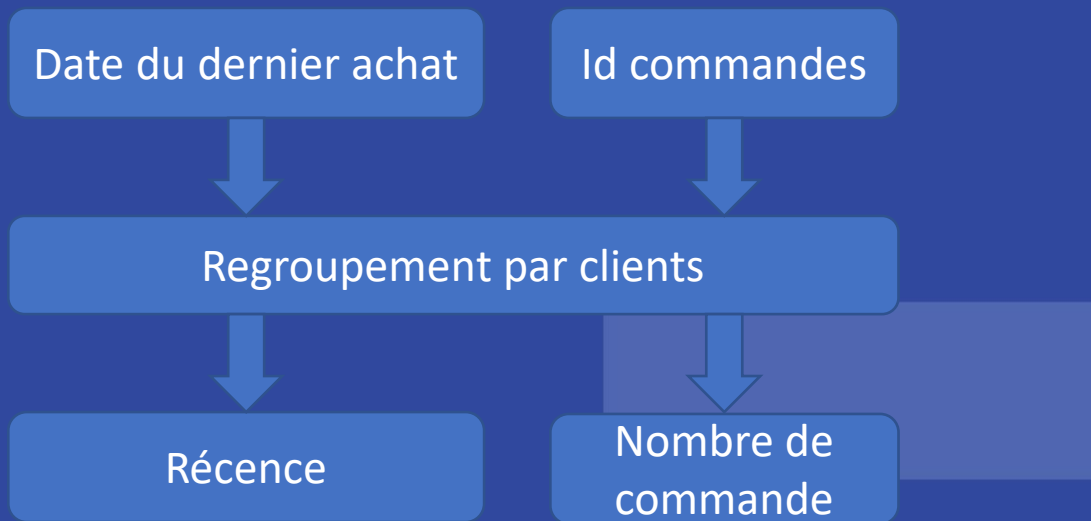
Variable note imputée par 3

Les variables prix produit, et frais de port imputées par le mode

La variable catégorie par 'unknown'

- Traitements:

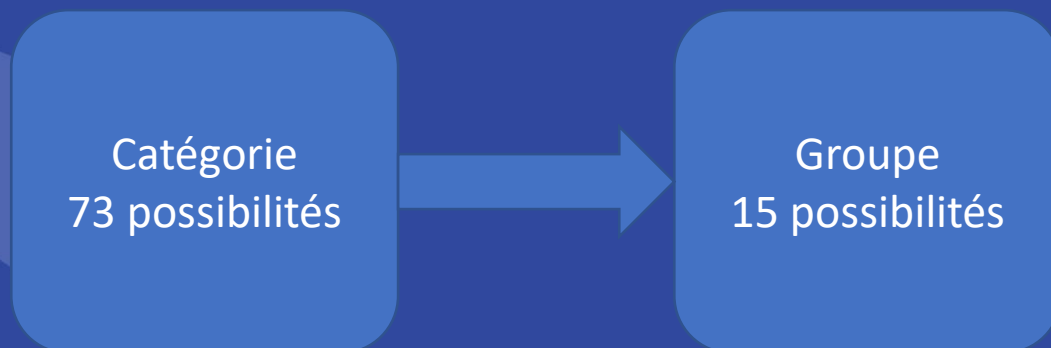




On se suppose que la date du test est le 01/10/2018 étant donné que la date la plus proche étant fin septembre

Nous avons nos:

- R -> Récence
- F -> Nombre de commande
- M -> Montant total



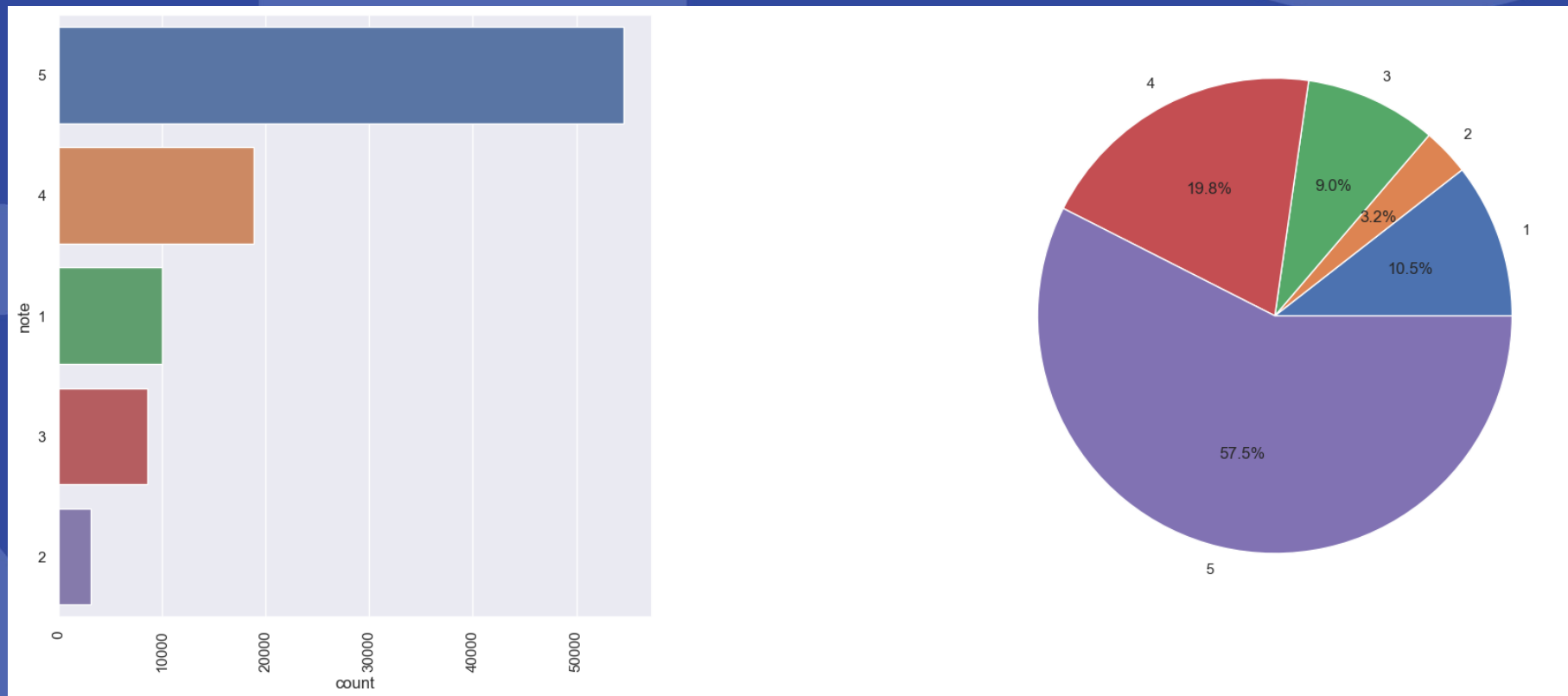
Répartition des catégories dans 15 groupes

Il y avait des catégories qui semblaient représenter la même gamme de produits

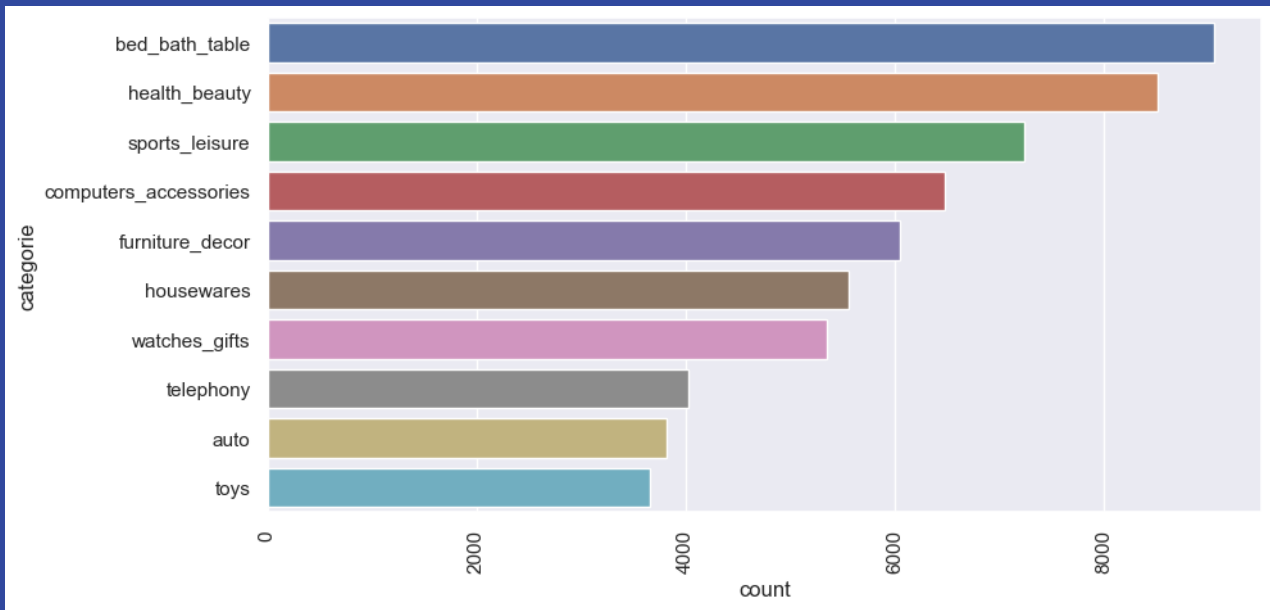
Préparation pour One-hot-Encoder

Analyse exploratoire

Répartition des notes:

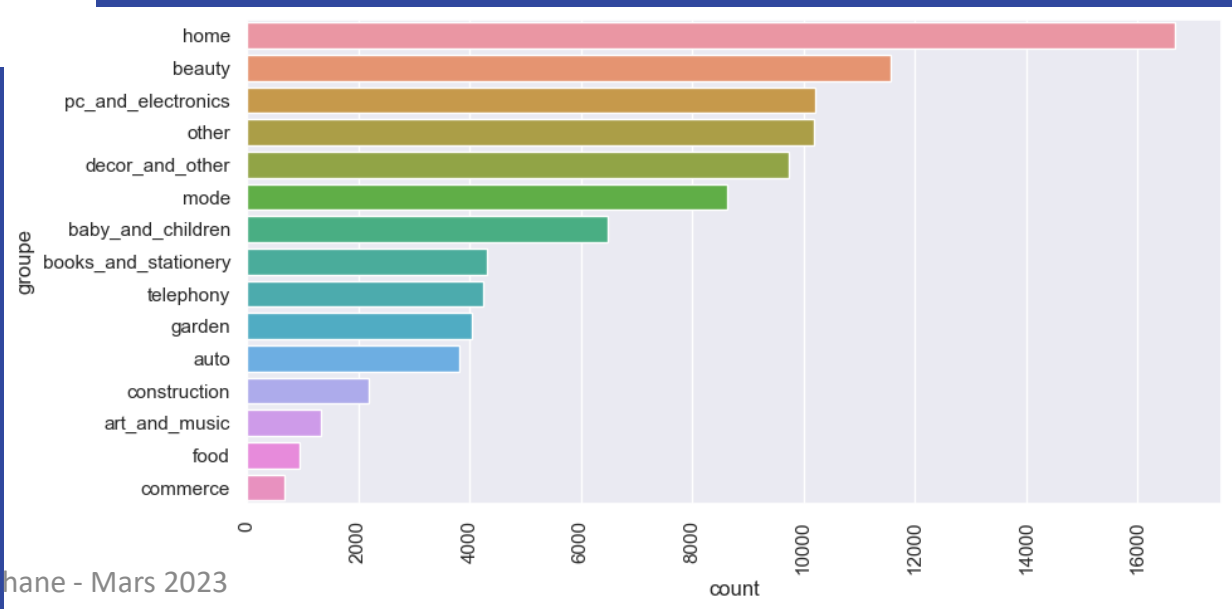


Répartition des catégories et des groupes:

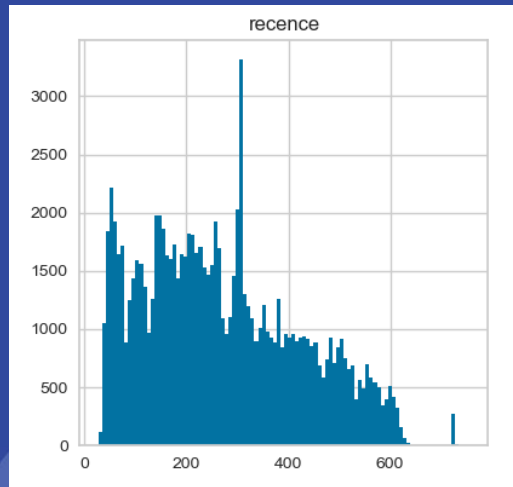


Les commandes concernent en majorité les accessoires de maison, la beauté et l'électronique

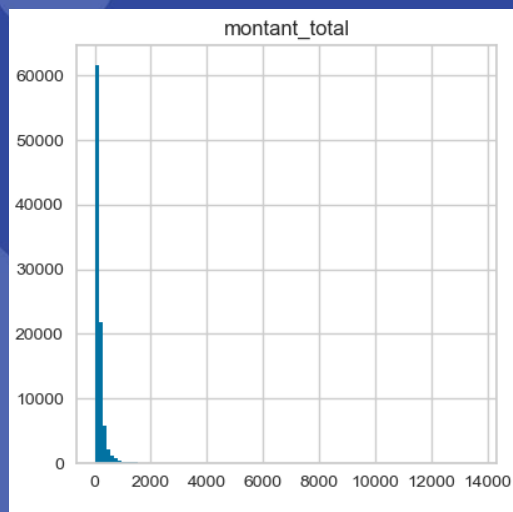
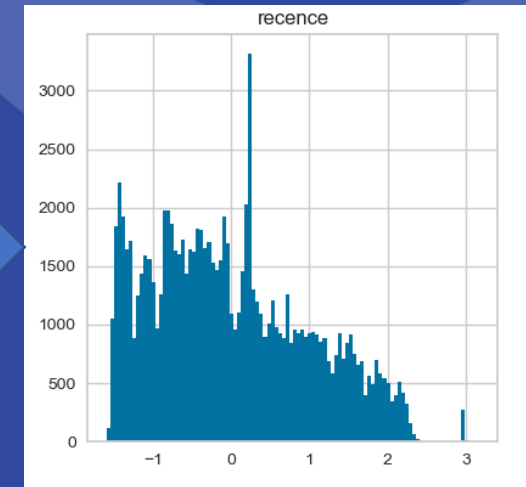
Les groupes représentent bien les 10 catégories les plus représentées



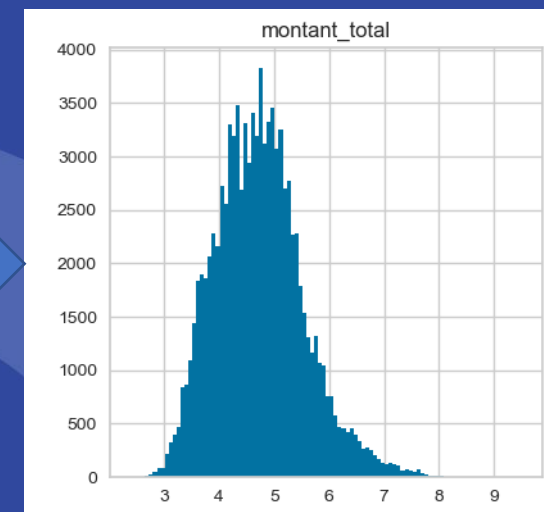
Transformations et encodage



Standard Scaler



Transformation log



groupe
Home
Beauty
Mode
....
Auto

Une colonne

One-hot-Encoder

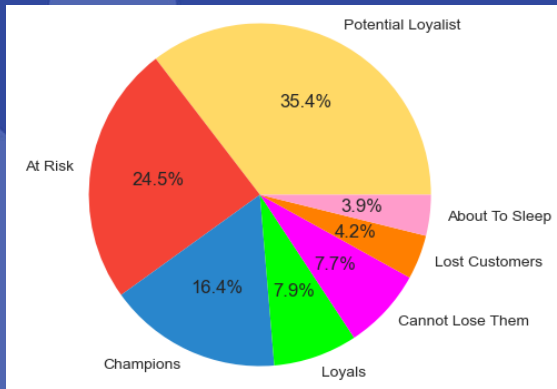
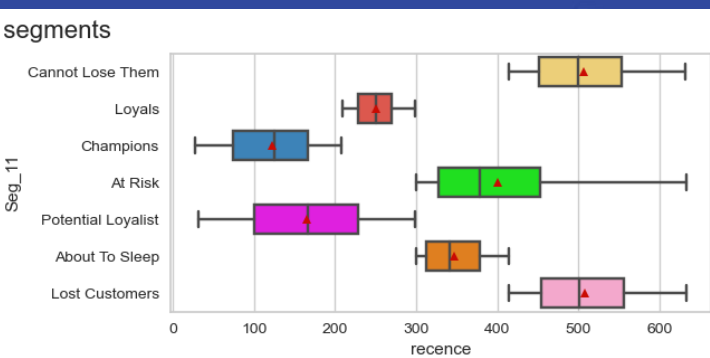
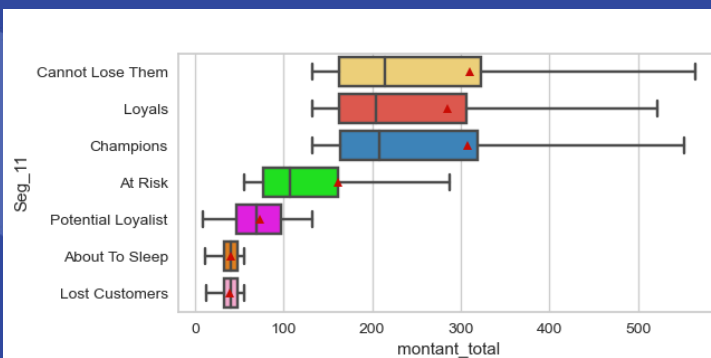
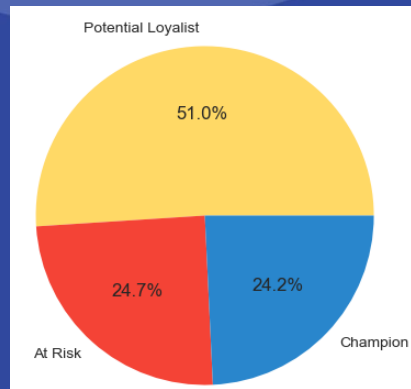
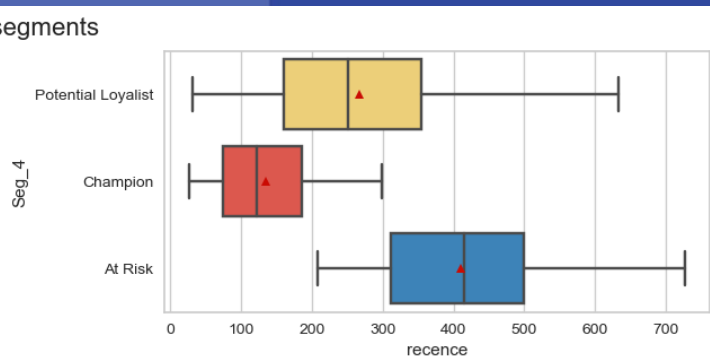
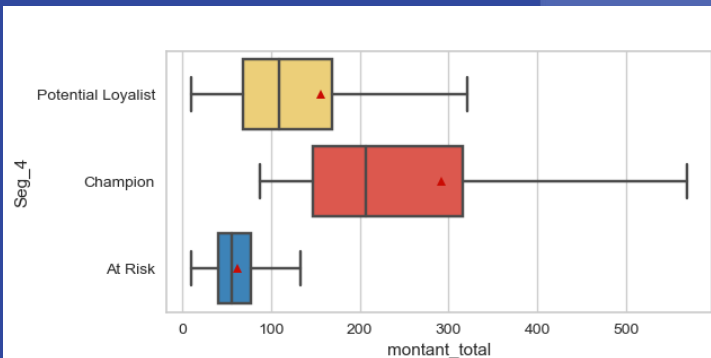
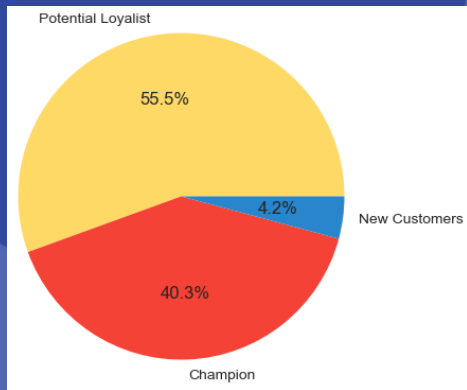
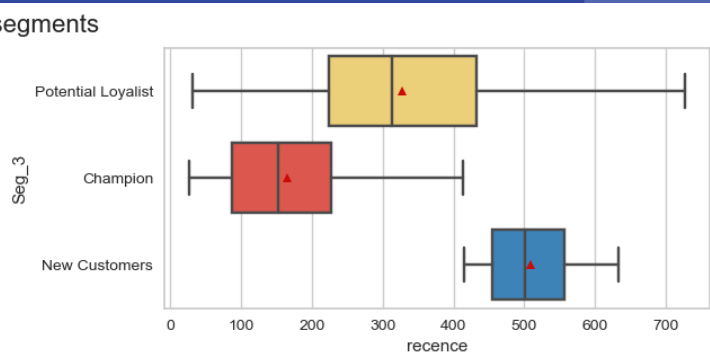
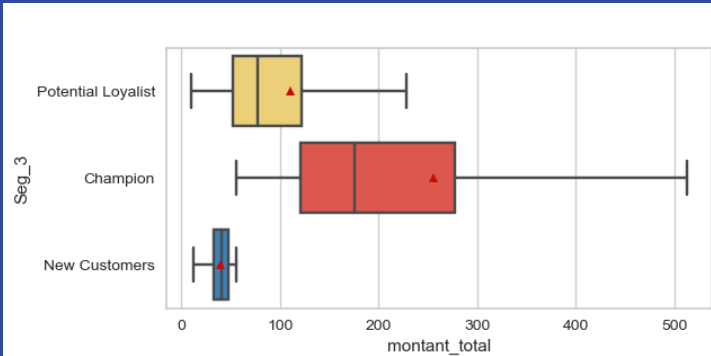
Home	Beauty	Mode	Auto
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	0
0	0	0	1

15 colonnes

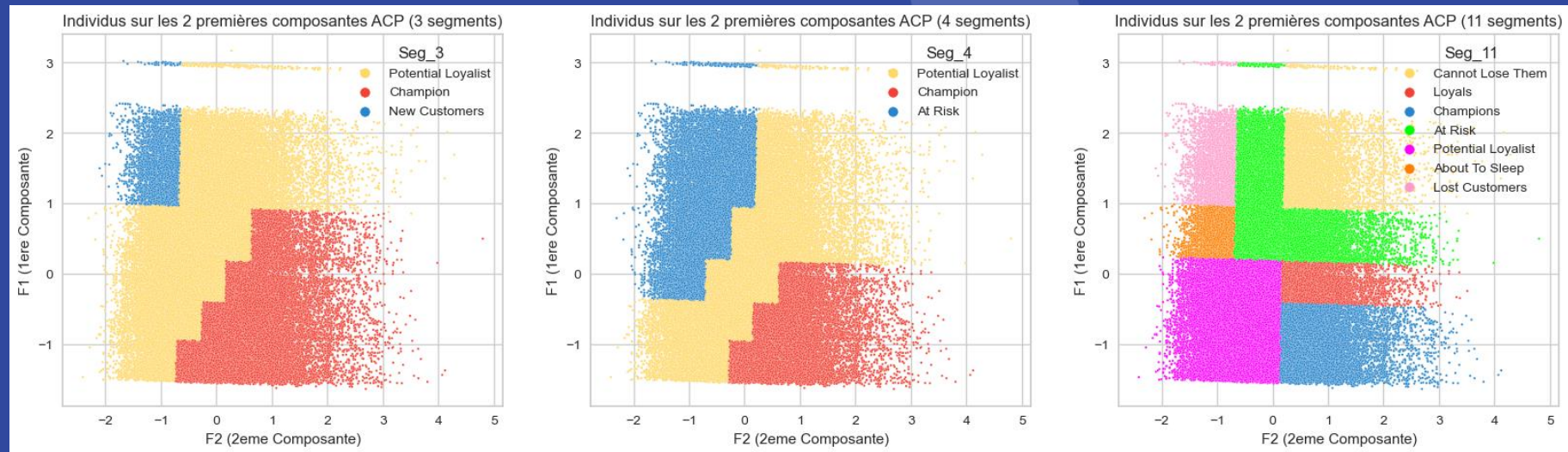
Les variables, nombre de commandes, moyenne de nombre de paiement, moyenne de nombre de produit, ne subissent aucune transformation pour préserver leur variance

Segmentation RFM

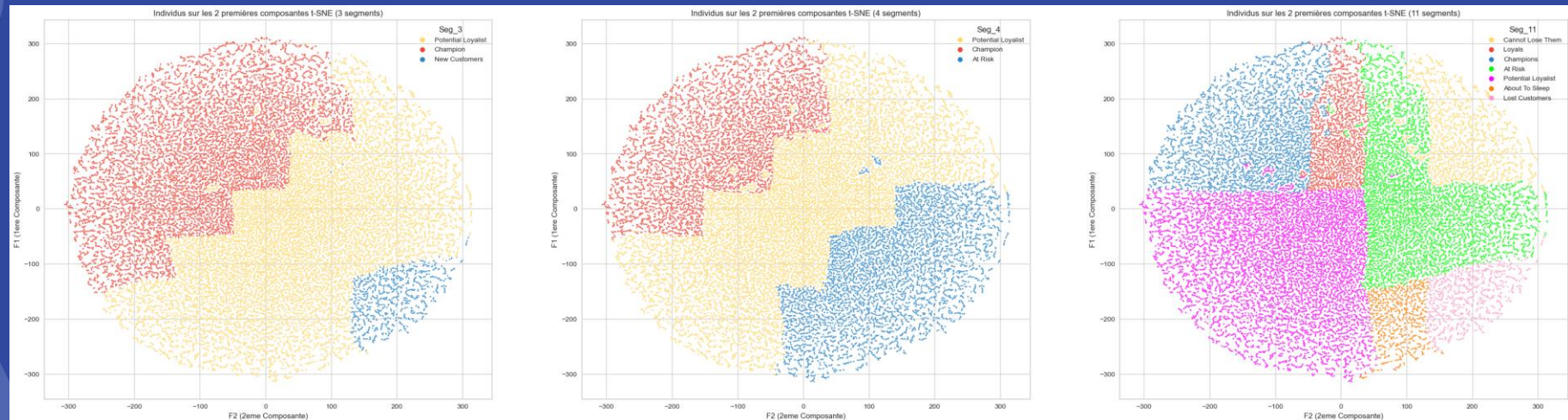
- La segmentation RFM basée sur un score utilisant la méthode des quantiles:
 - Répartition des clients selon la somme de leur score établi sur les variables Récence, Fréquence et Montant total
 - Score variant de 1 à 5 pour les variables soit un score total variant de 3 à 15
- La segmentation RFM basée sur la technique des 11 segments:
 - Score variant de 1 à 5 pour les variables, puis création de segments RFM (ex: 111, 545 ou 354)
 - Répartitions des segments dans 11 clusters



ACP



t-SNE



Comparaison des modèles

La comparaison est faite sur 50% du dataset choisie aléatoirement, à cause des limites matérielles.

- Le clustering via k-Means:

Minimise la variance intra-cluster (l'inertie)

- Le clustering hiérarchique (Agglomerative Clustering):

Chaque individus est considéré comme un cluster, puis ils sont regroupés pour minimiser la variance intra-cluster en utilisant la méthode Ward

- Le clustering via DBScan:

Relie les points d'un même cluster en passant de voisin à voisin

Les métriques utilisées:

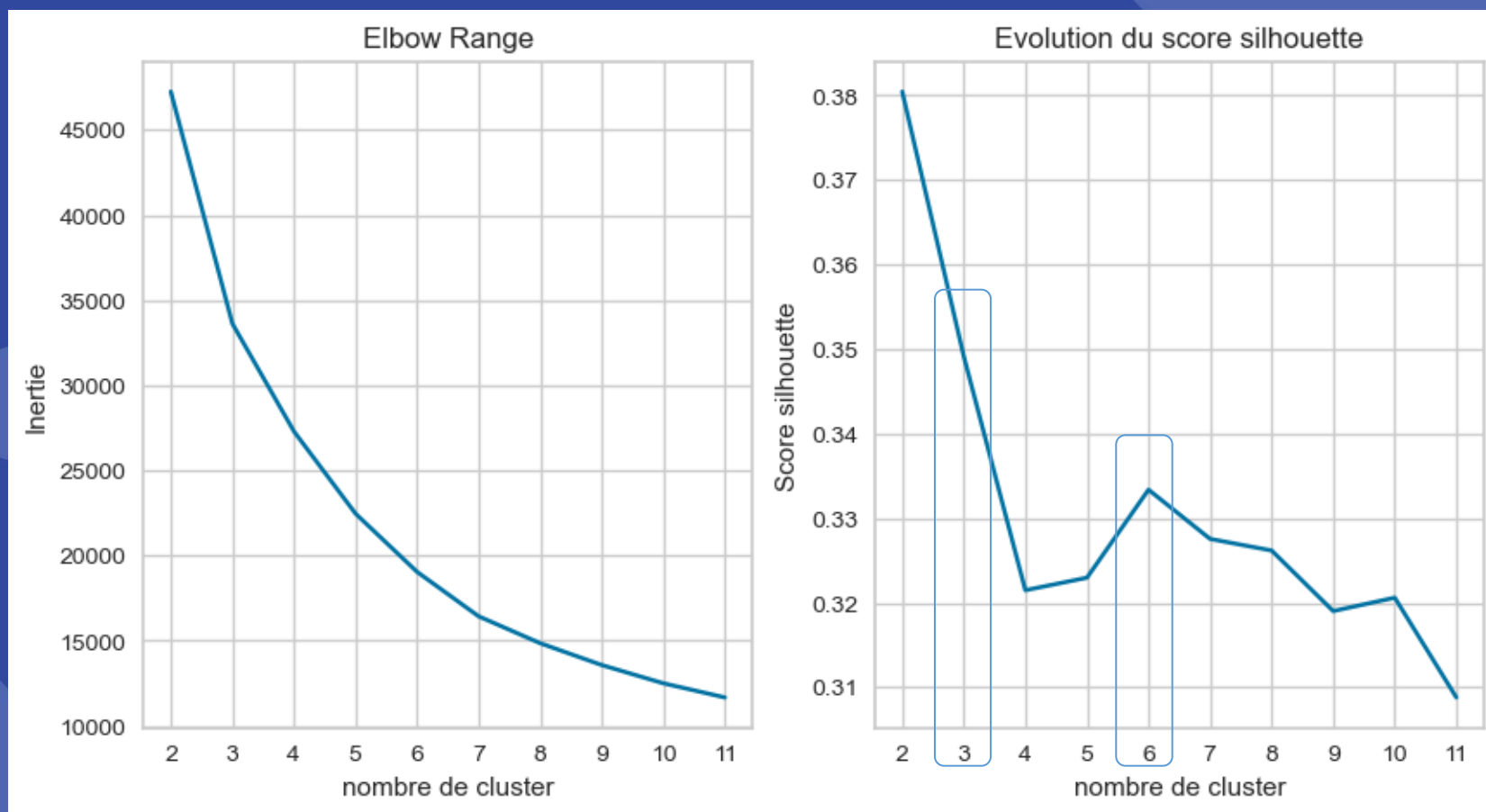
- L'inertie:
 - Détermine si les points d'un même cluster sont proches
 - Représente la variance intra-cluster
 - Si petit bon cluster
- Le score silhouette:
 - Détermine si chaque point est dans le bon cluster
 - Si grand bon cluster
- L'Indice de Rand ajusté:
 - Calcule une mesure de similarité entre deux regroupements en considérant toutes les paires d'échantillons et en comptant les paires qui sont attribuées dans le même groupe ou dans des groupes différents dans les regroupements prédits et vrais.
 - Il est ajusté pour avoir une valeur proche de 0 si clusterings aléatoires, proche de 1 si clusterings identiques et négatif si les clusterings ne s'accordent pas

Autres métriques:

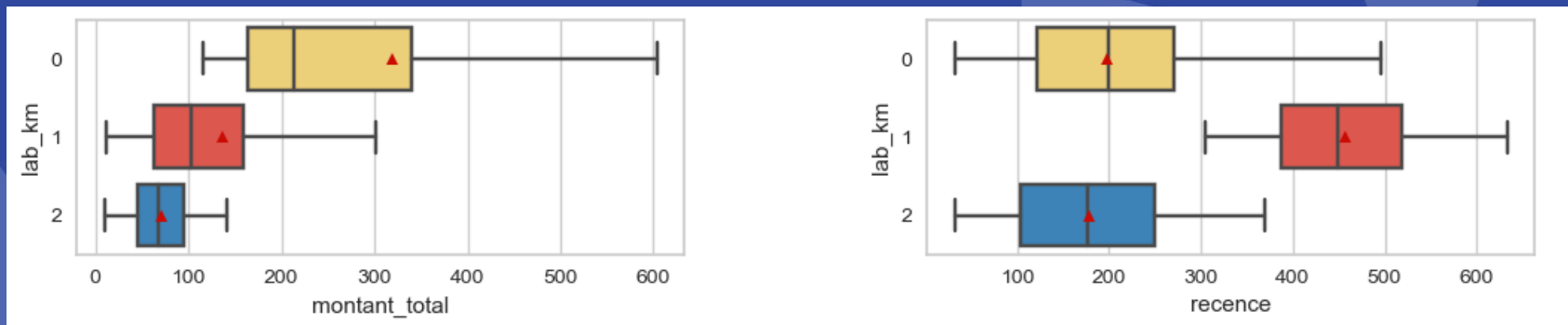
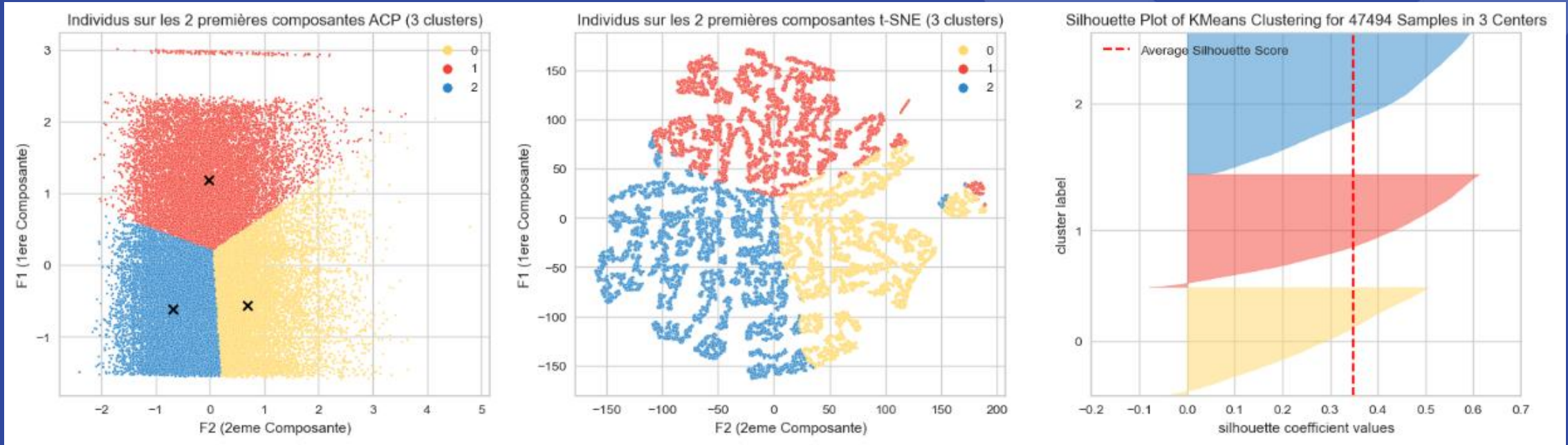
Détermine si les clusters sont bien éloignés les uns des autres et bien regroupés sur eux-mêmes

- Indice de Davies-Bouldin (Rapport entre la taille du cluster et la distance entre les clusters)
- Indice de Calinski-Harabasz (Variance inter-clusters / Variance intra-cluster)

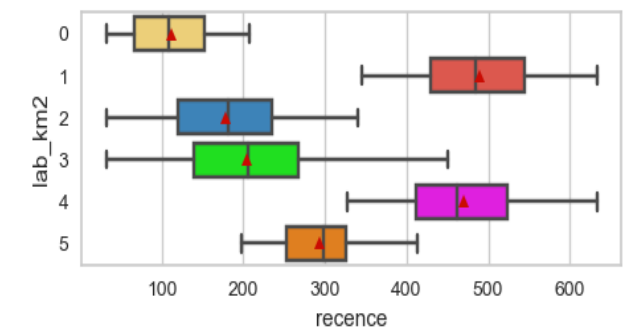
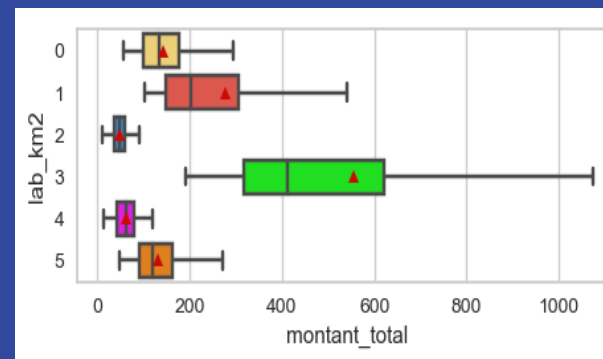
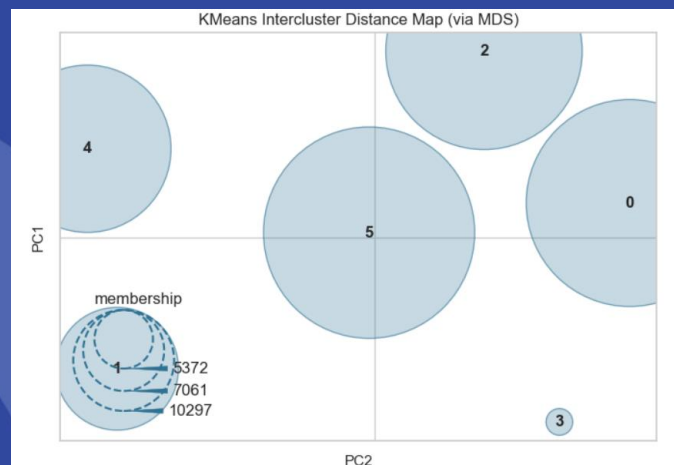
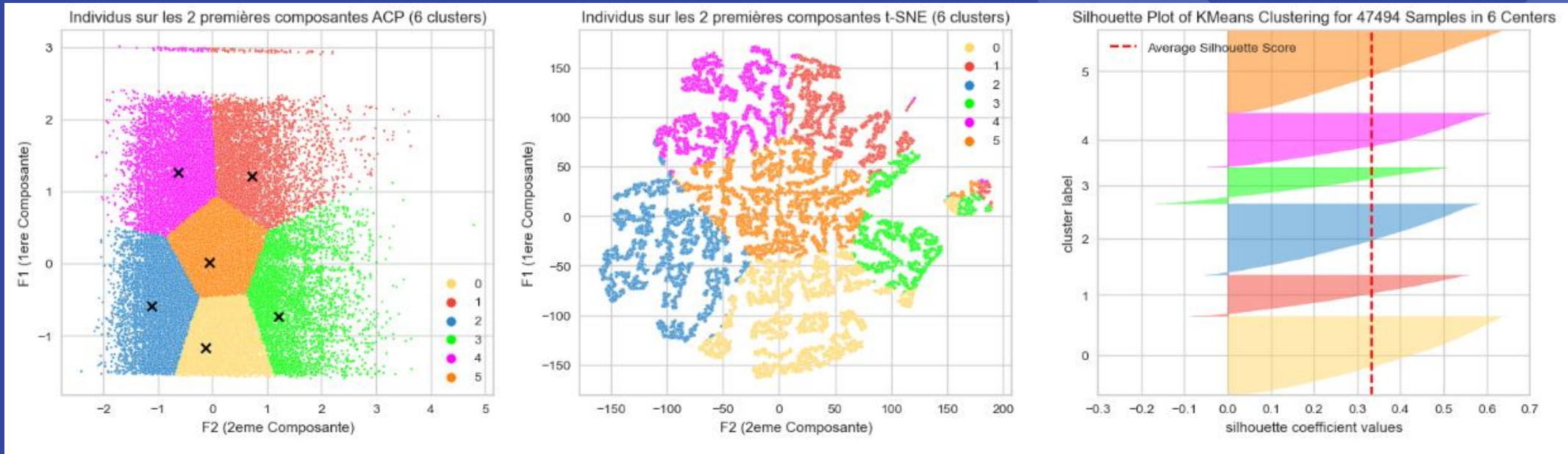
Clustering via k-Means



- 3 clusters:

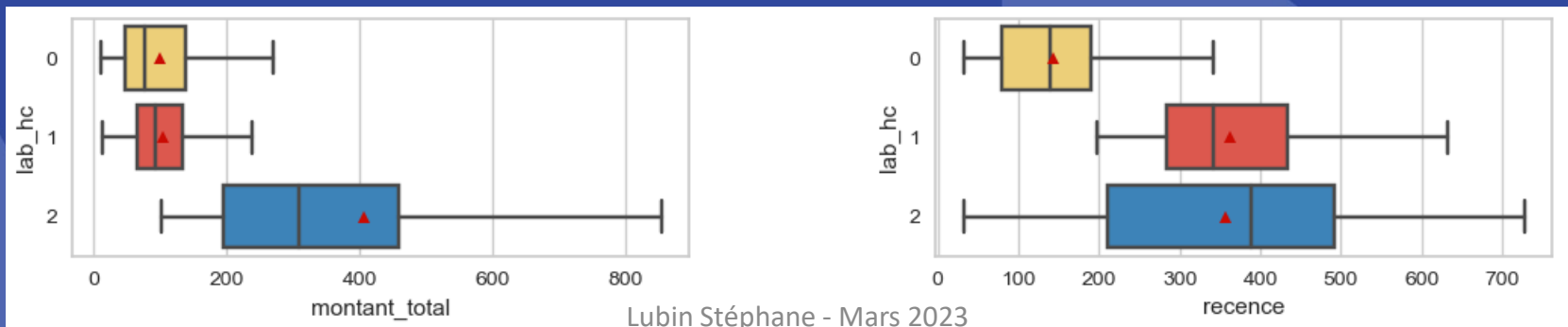
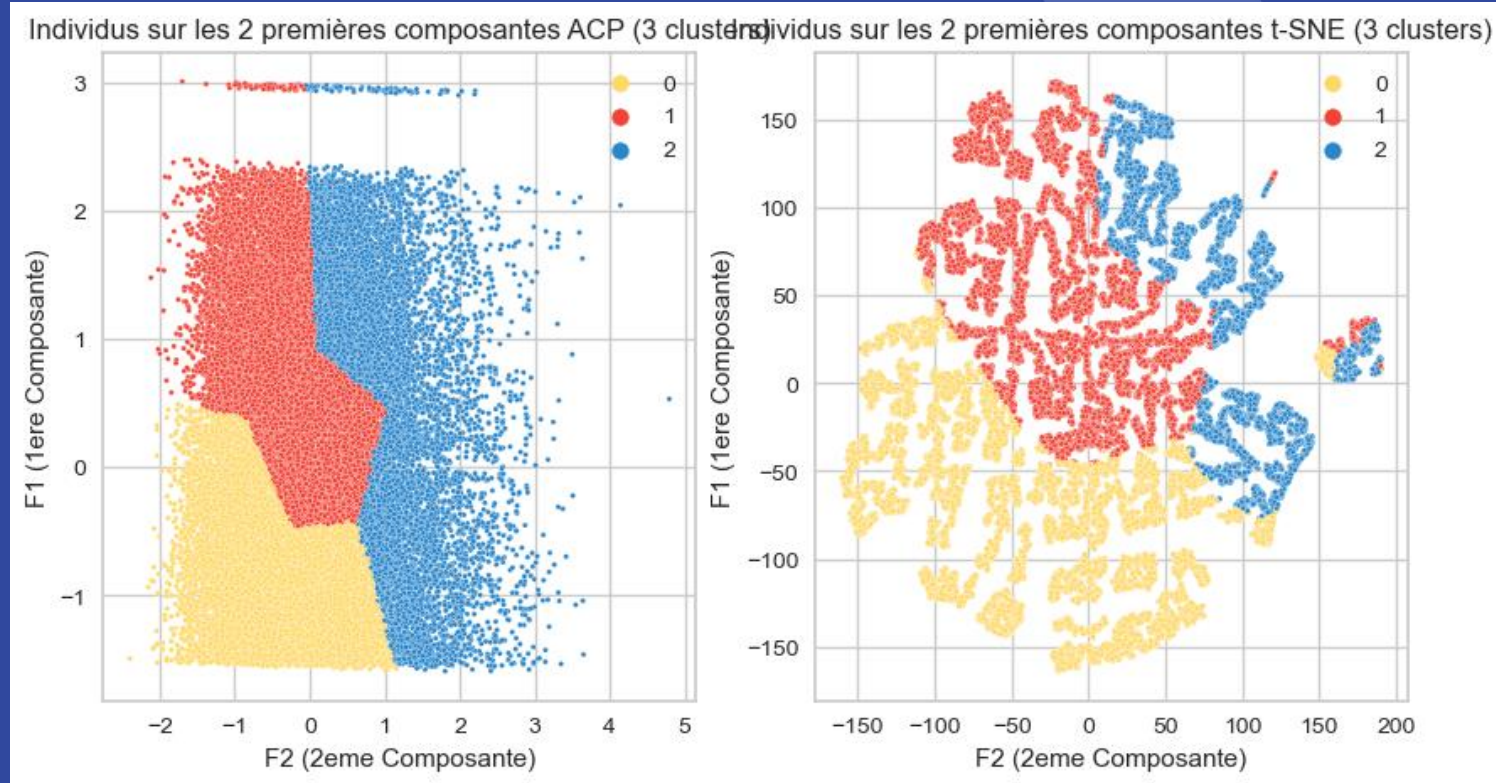


- 6 clusters:

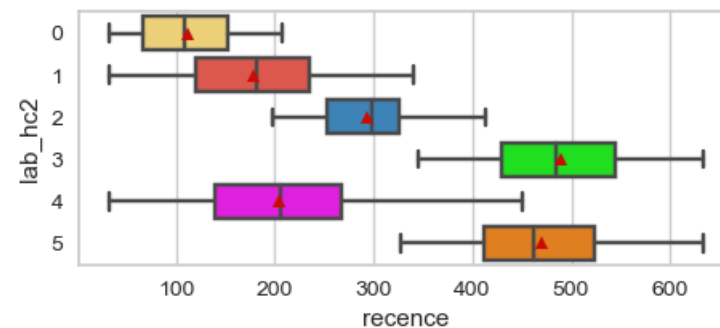
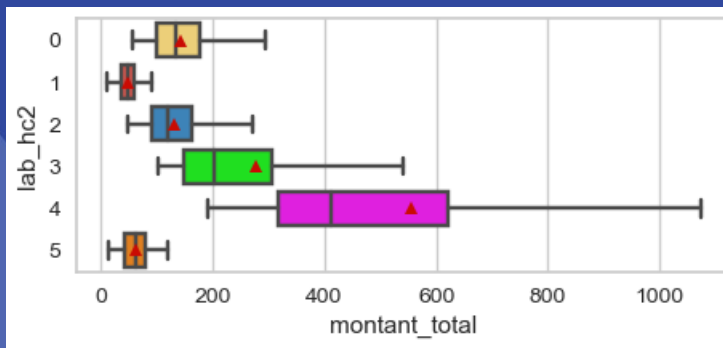
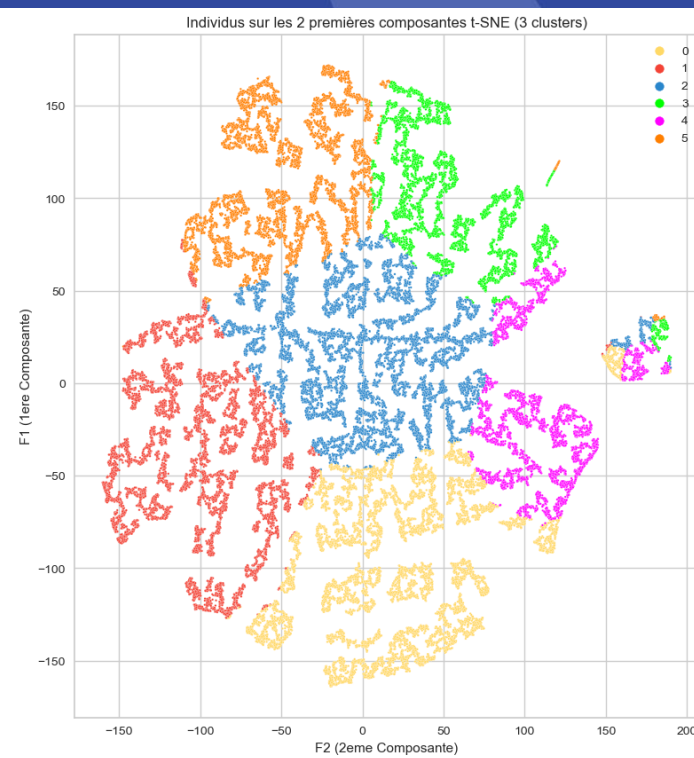
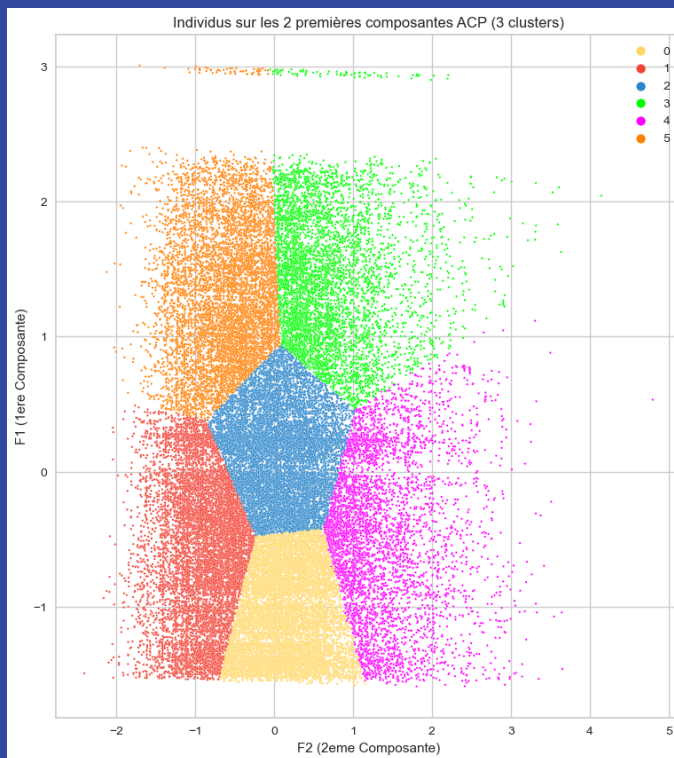


Clustering Hiérarchique

- 3 clusters:



- 6 clusters:

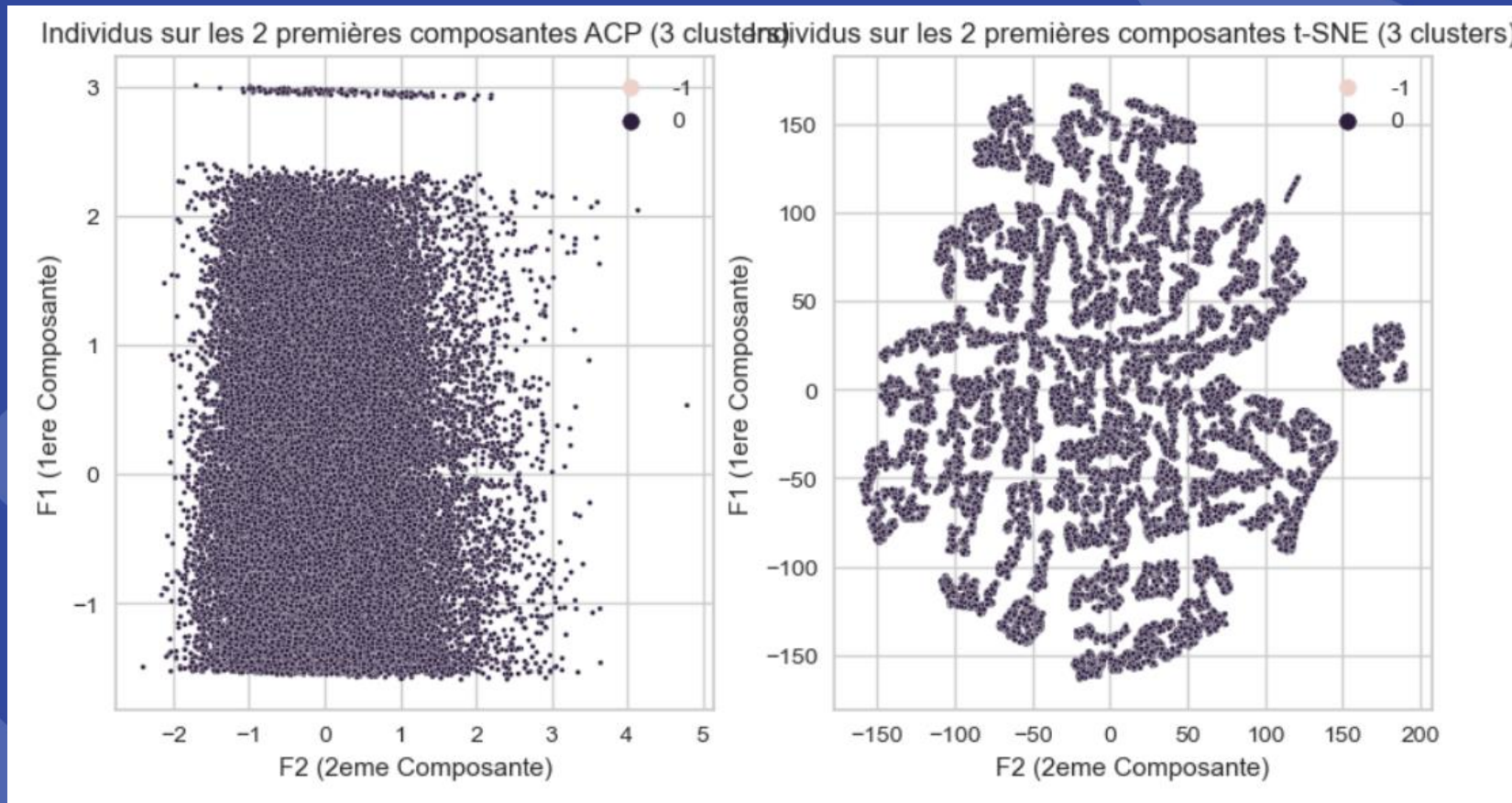


ARI entre les prédictions de k-Means et le clustering hiérarchique:

- Pour 3 clusters : $ARI = 0,28$
- Pour 6 clusters: $ARI = 1$

Pour des raisons matérielles, le modèle K-Means est préféré

Clustering via DBScan



DBScan n'est pas adapté, points trop proches les uns des autres. Quasiment tous réunis dans le même cluster ou tous considérés comme du bruit

Ajout de variables

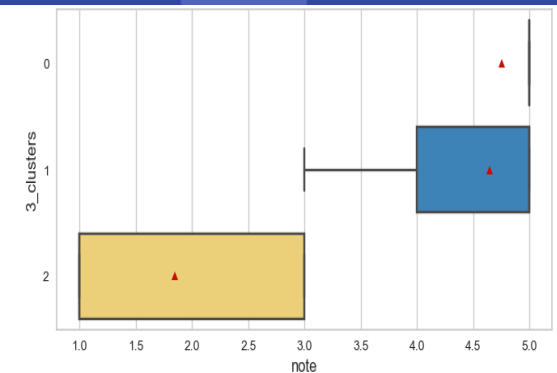
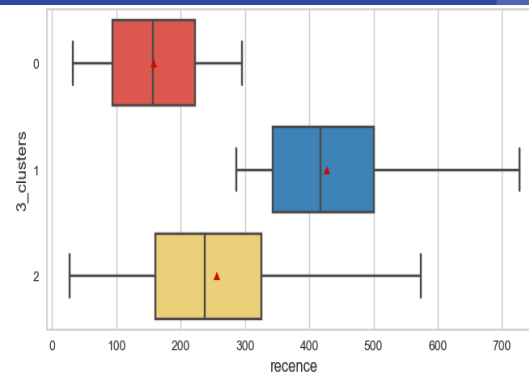
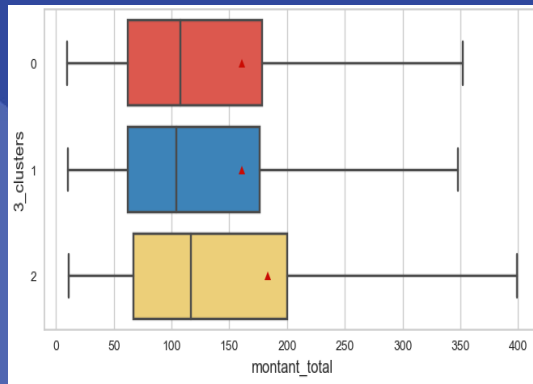
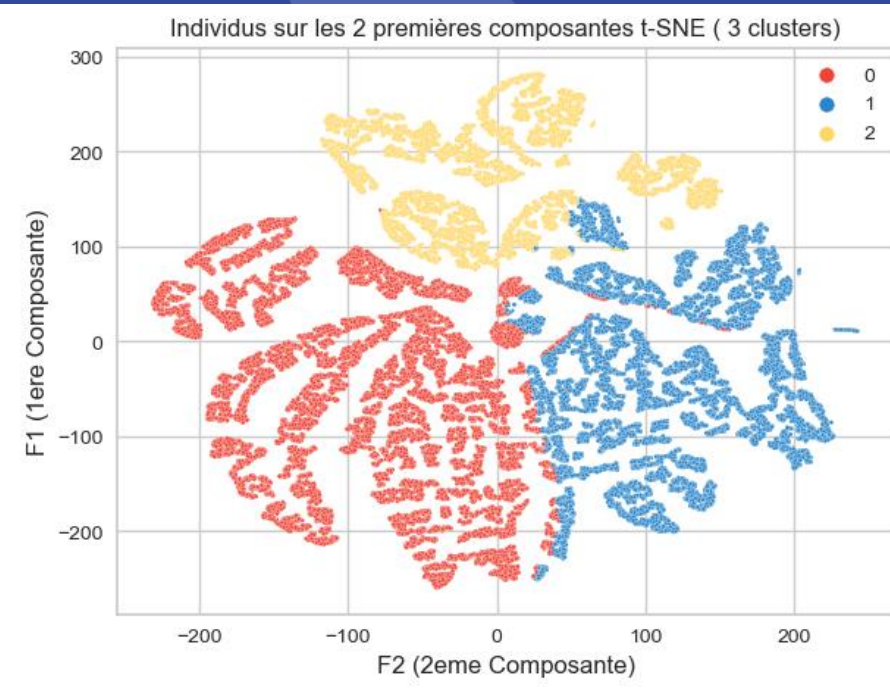
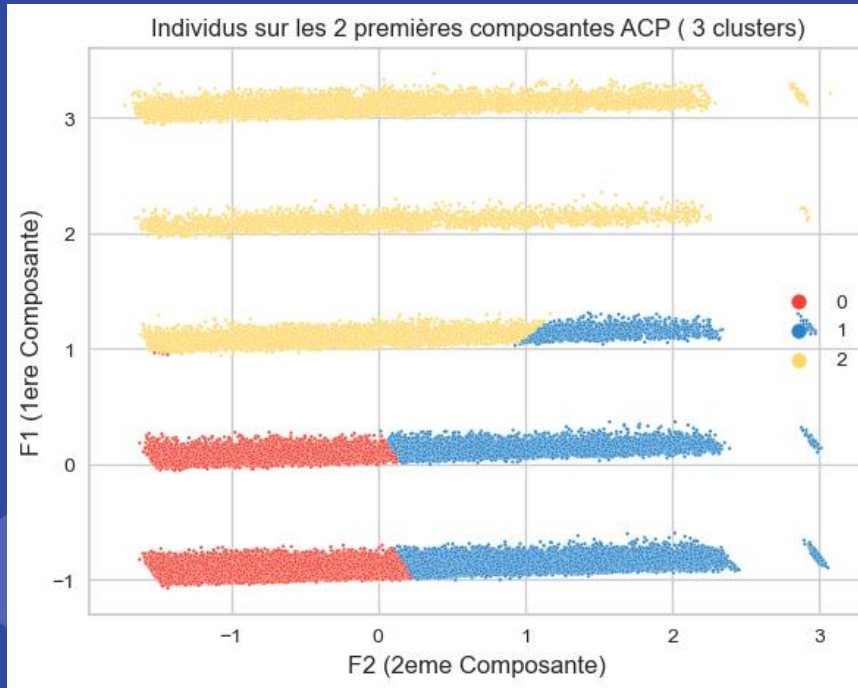
Le modèle utilisé sera k-Means, on regardera les 2 meilleurs choix de clustering, la répartition et la stabilité des prédictions sur 25 répétitions.

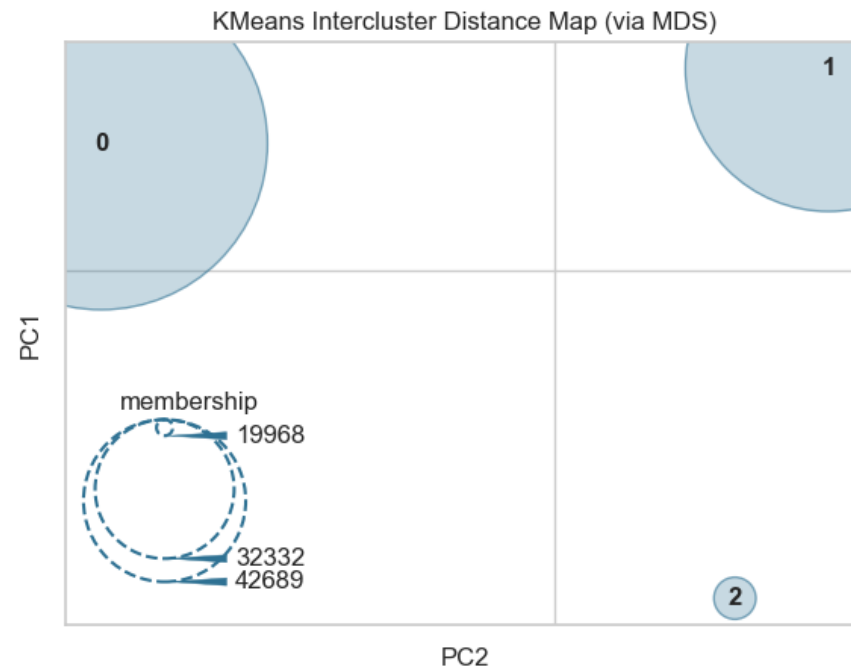
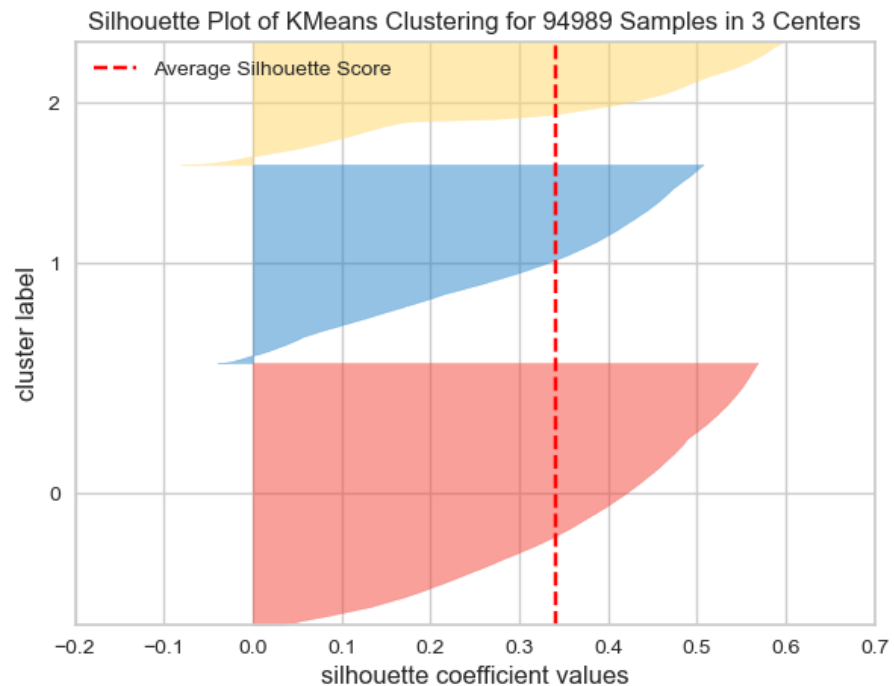
Ajout une à une de variable:

- Note , qui reflète la satisfaction de client
- Moyenne du nombre de produits, qui permet de voir si les clients font des grosses commandes
- Moyenne du nombre de paiement, qui permet de voir si l'acheteur est compulsif ou s'il fractionne ses paiements
- Groupe, qui permet de différencier la catégorie favori des clients

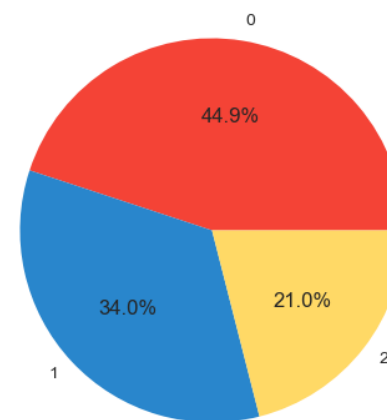
Variable Note

3 Clusters:

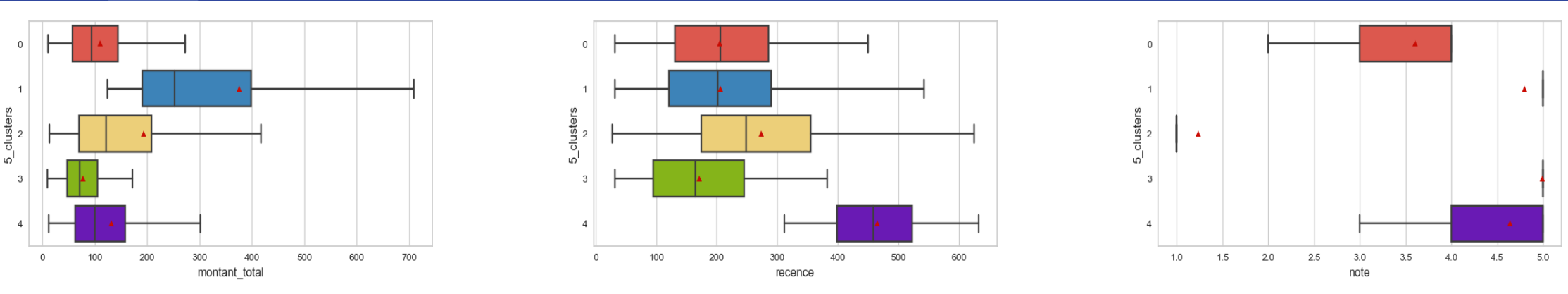
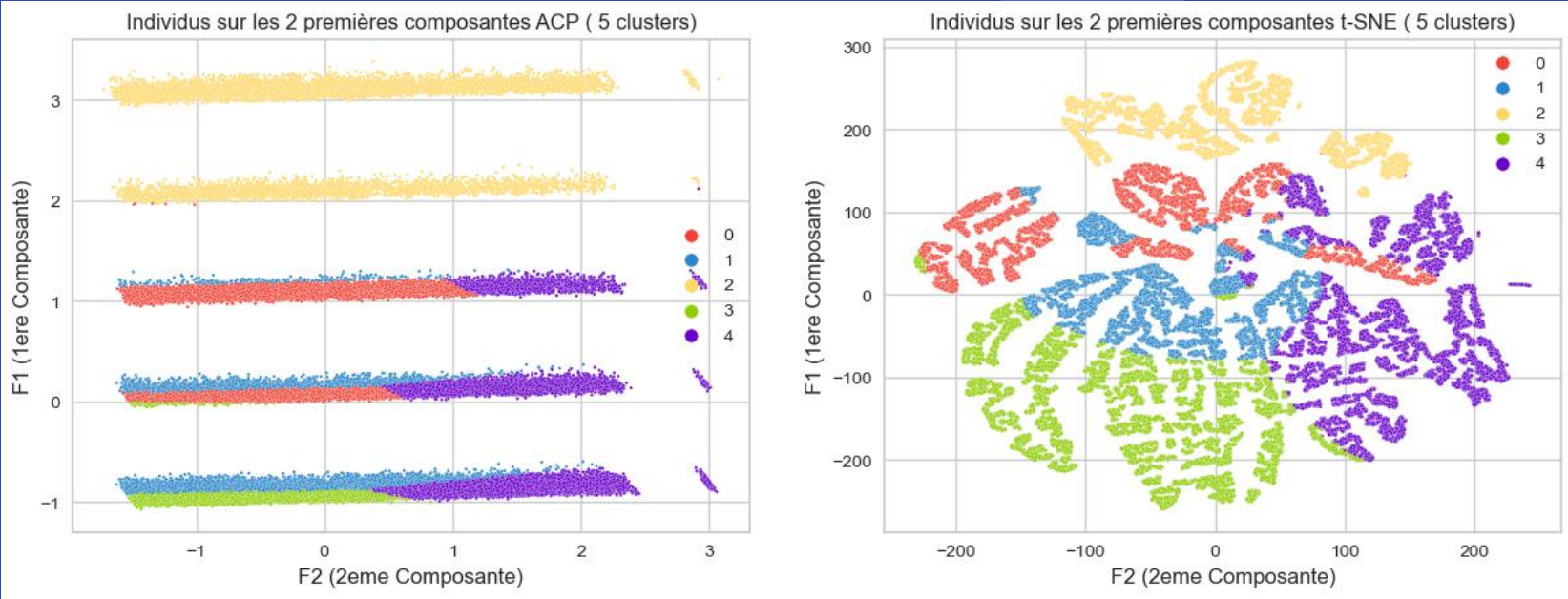


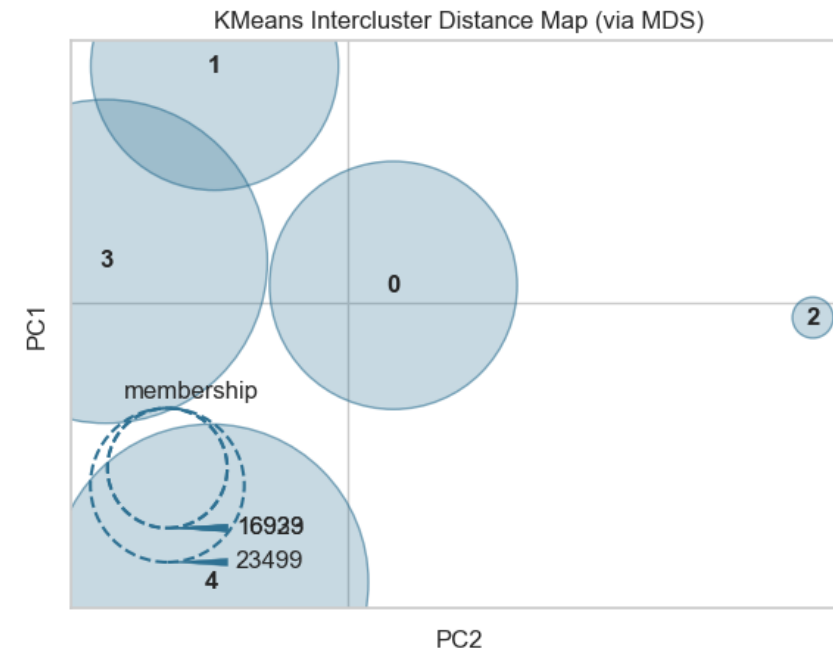
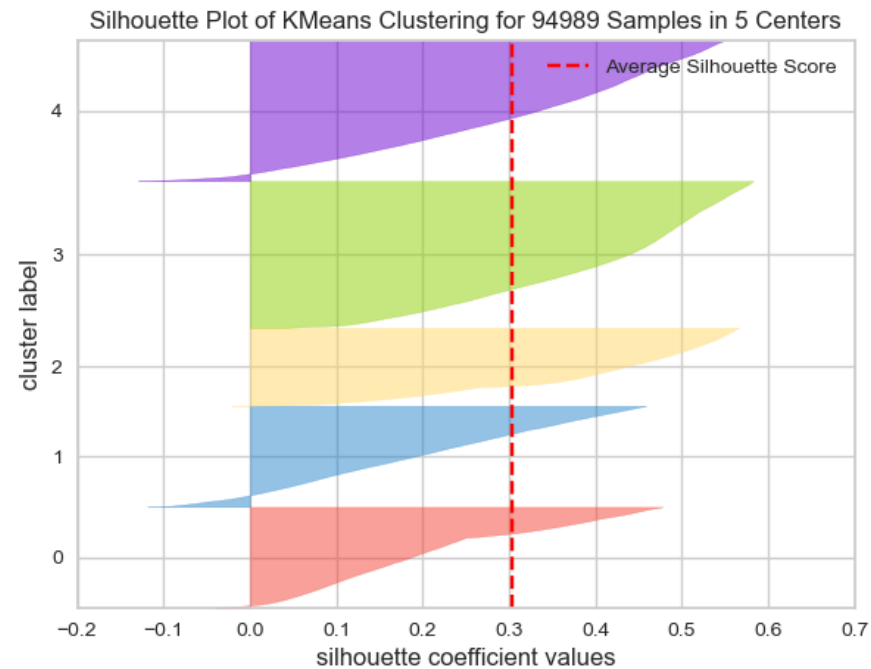


ARI moy = 1,0 (arrondi)
Stable mais trop généraliste

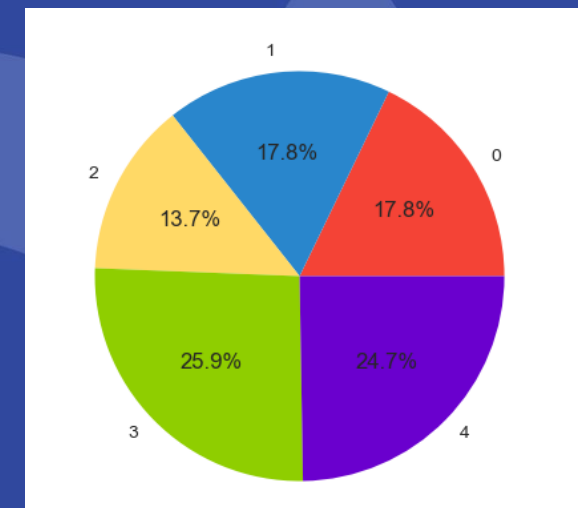


5 Clusters:

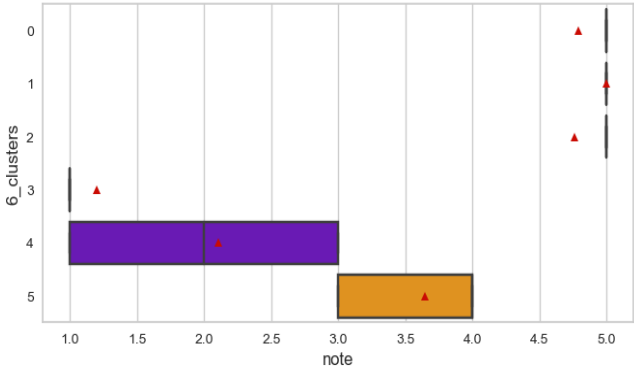
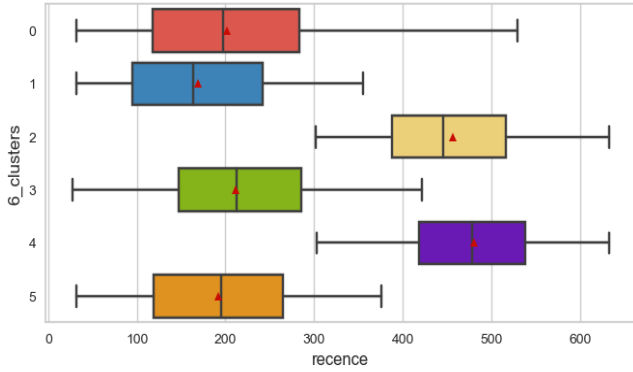
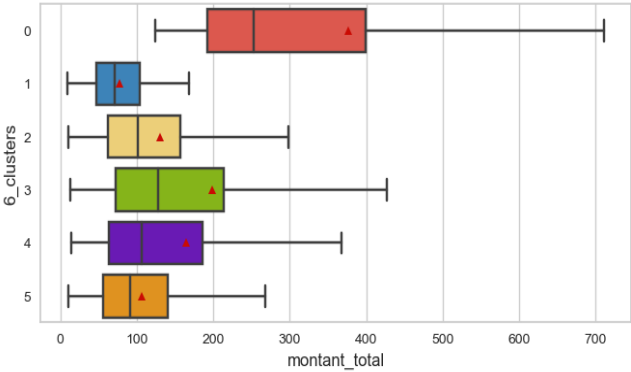
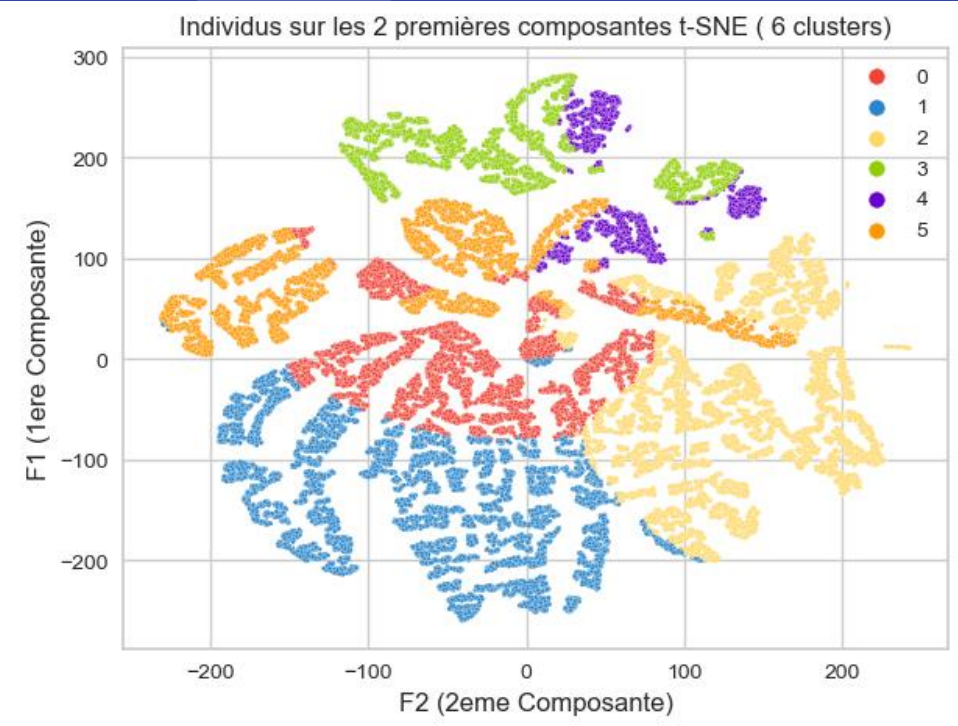
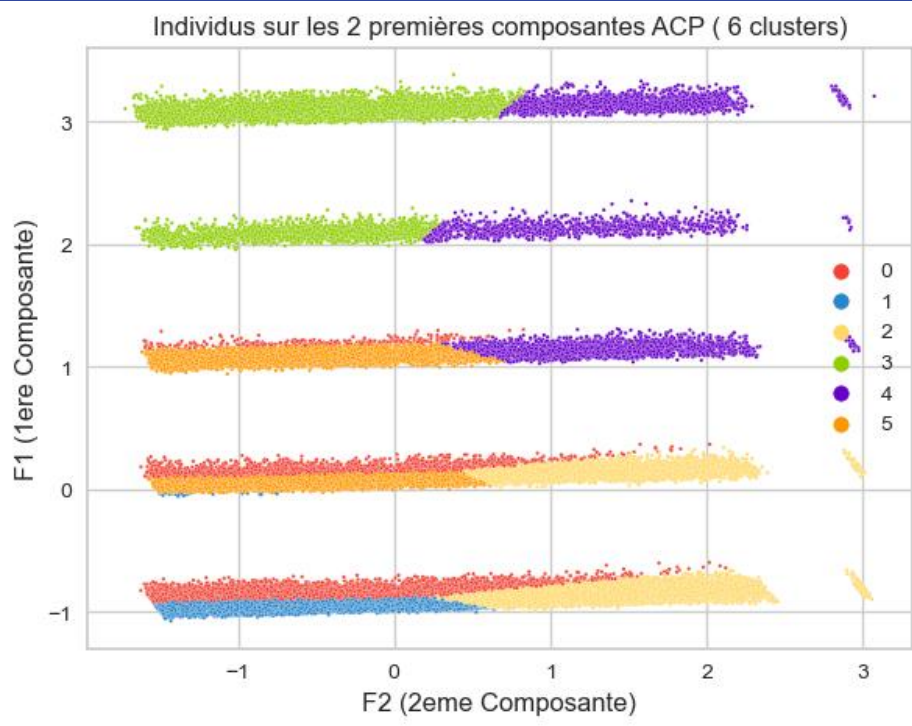


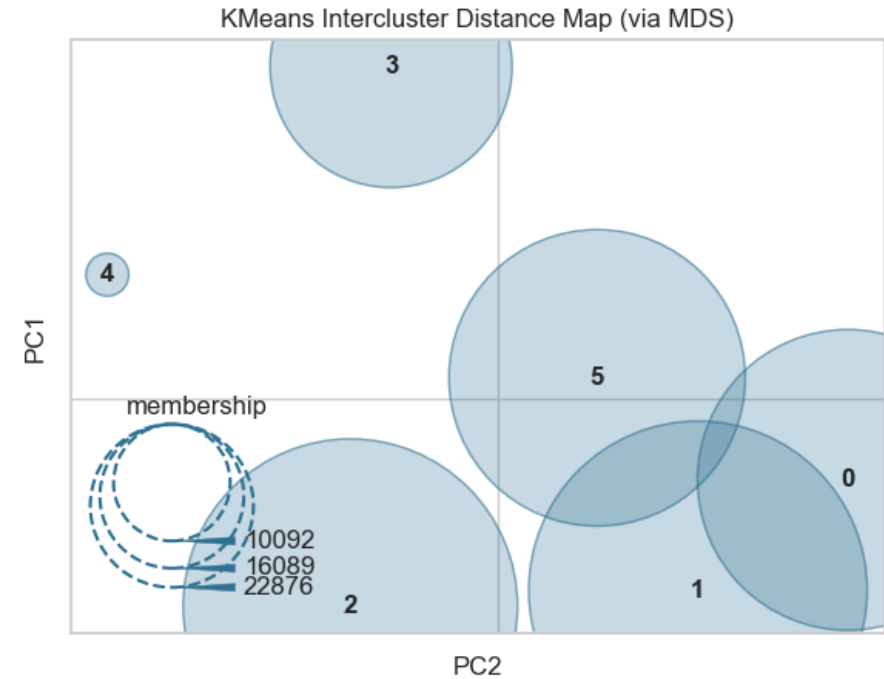
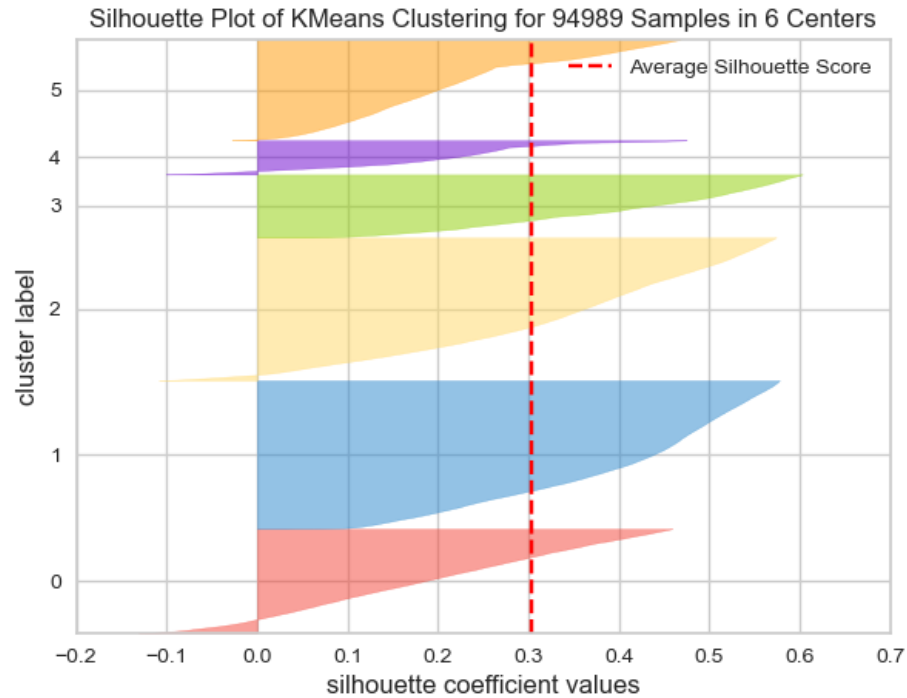


ARI moy = 0,94
ARI Min = 0,74
Perd en stabilité

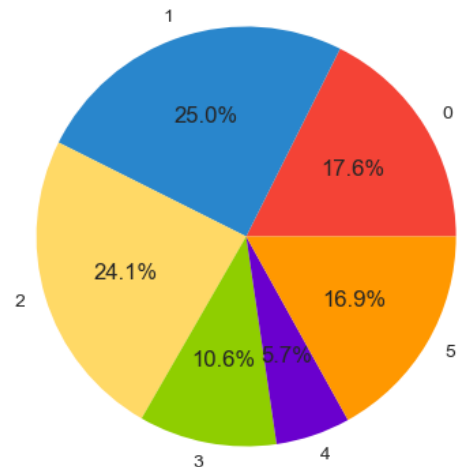


6 Clusters:



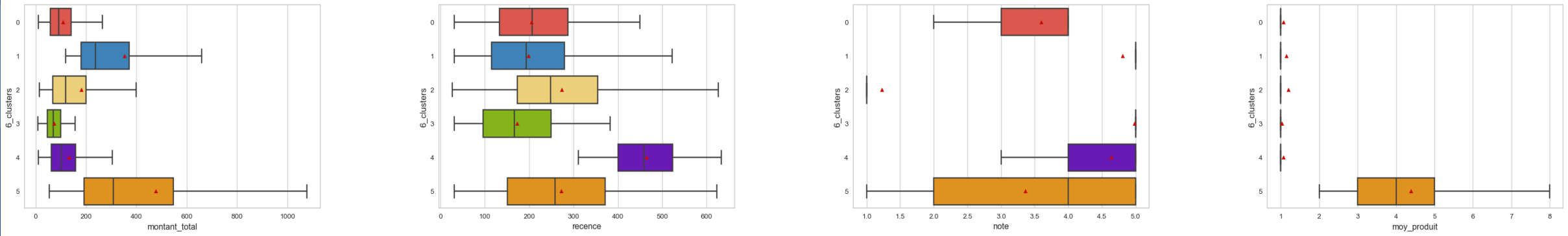
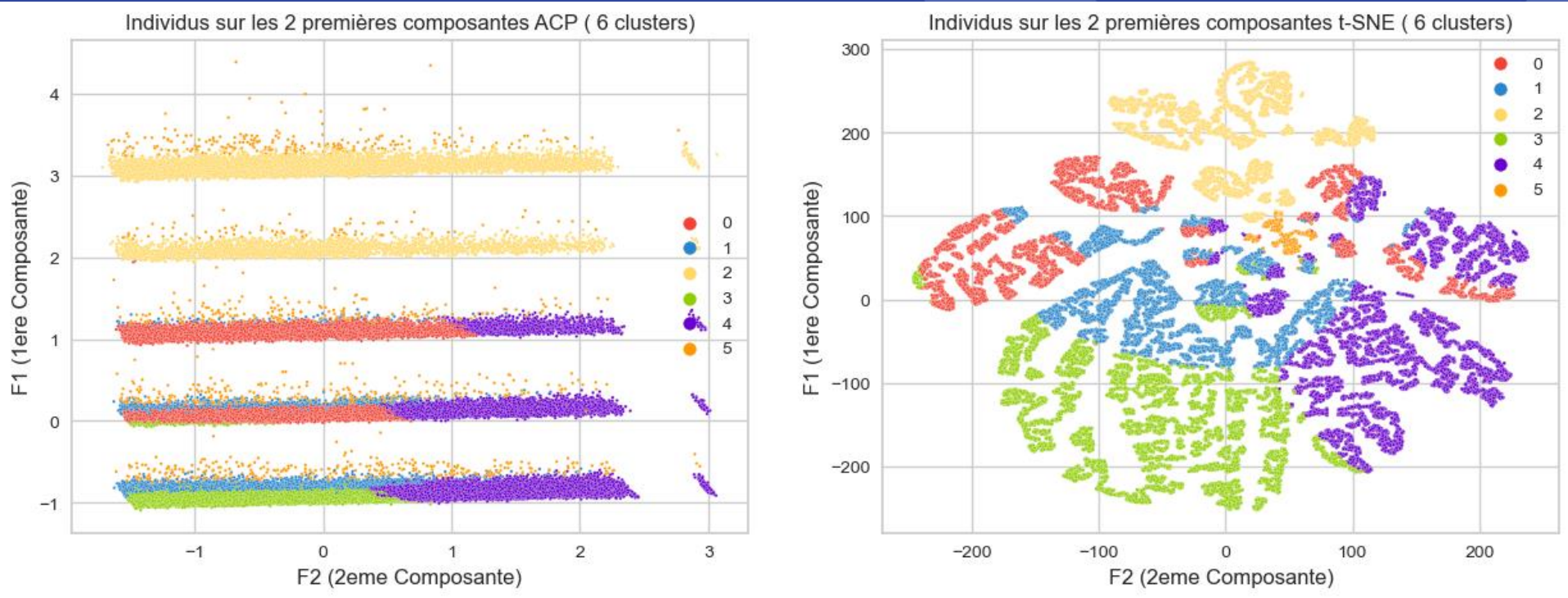


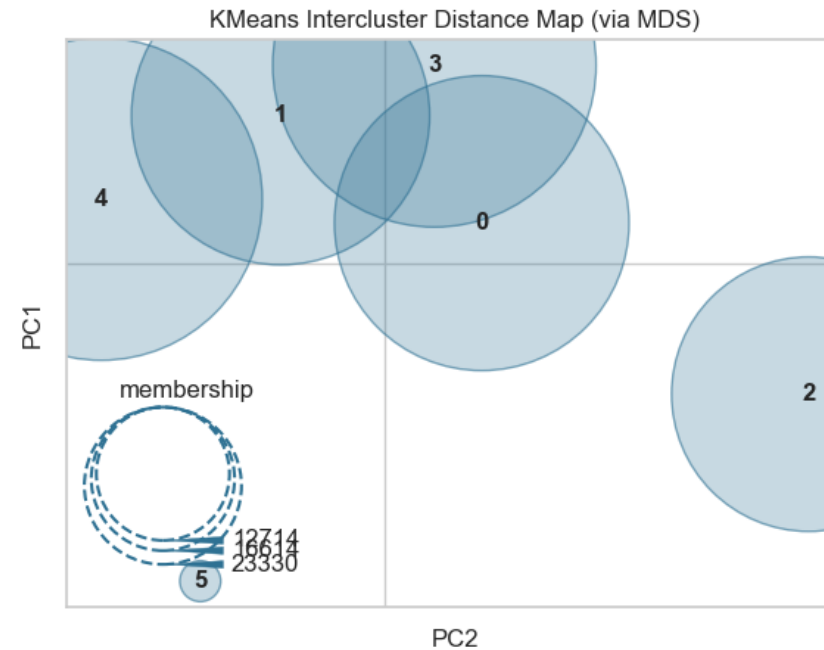
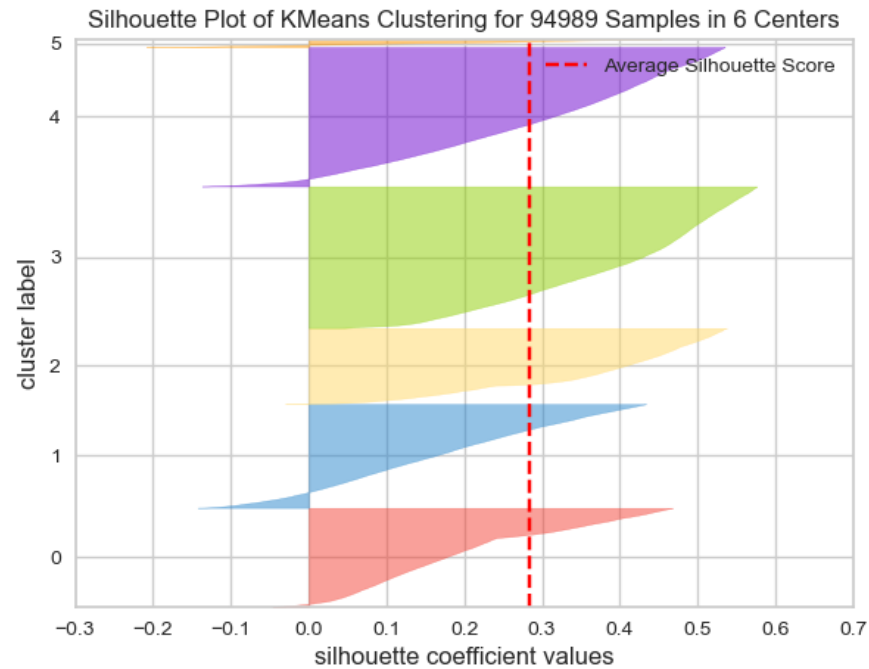
ARI moy = 0,96
ARI Min = 0,56
Perd encore en stabilité



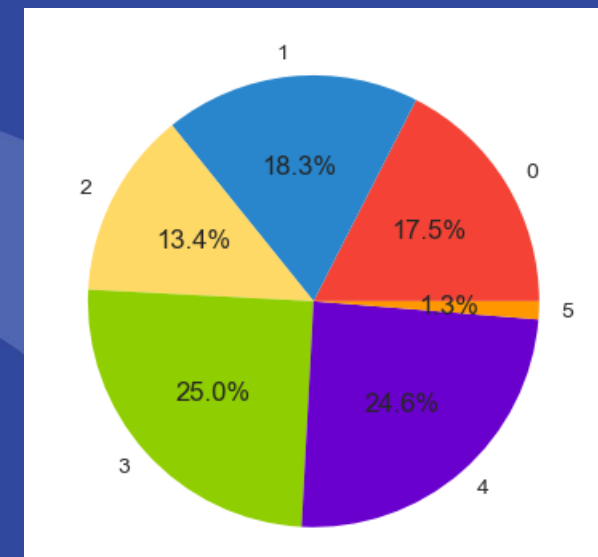
Variable nombre de produit moyen

6 Clusters:



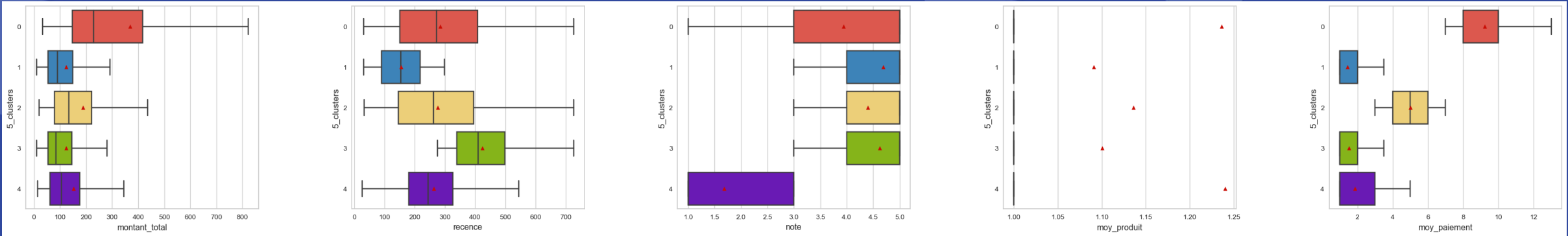
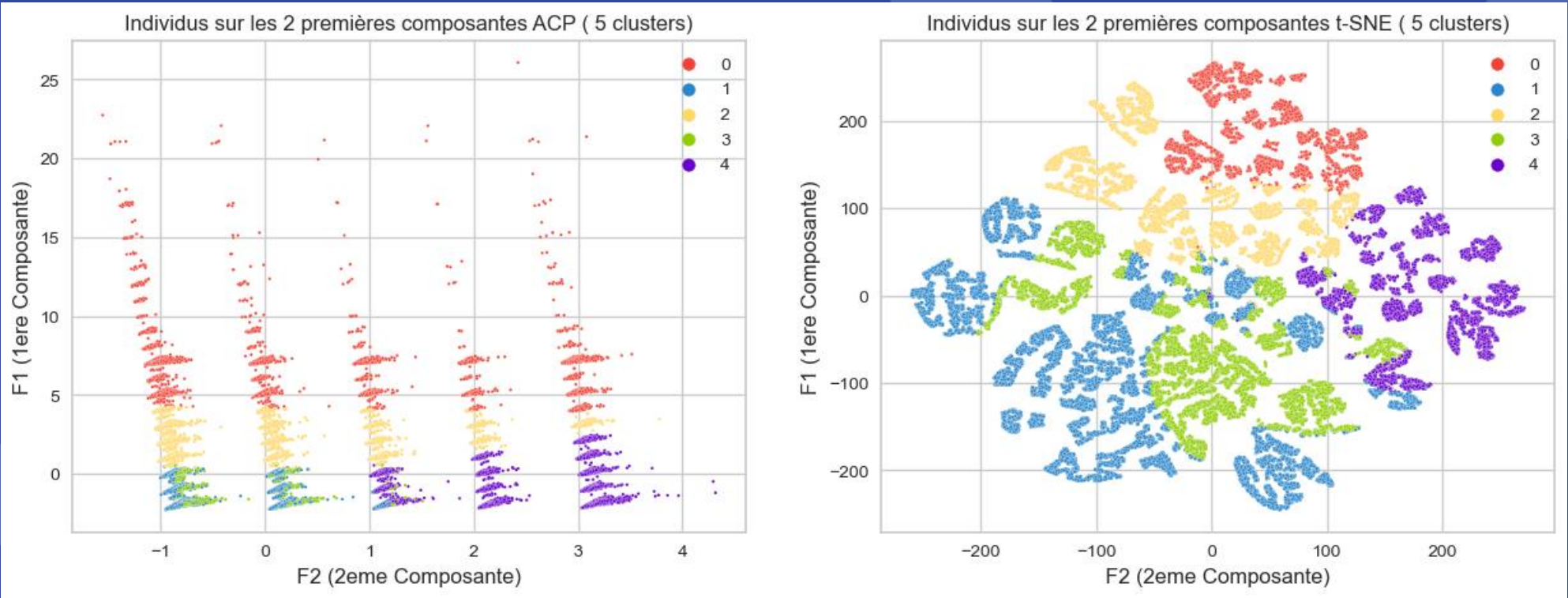


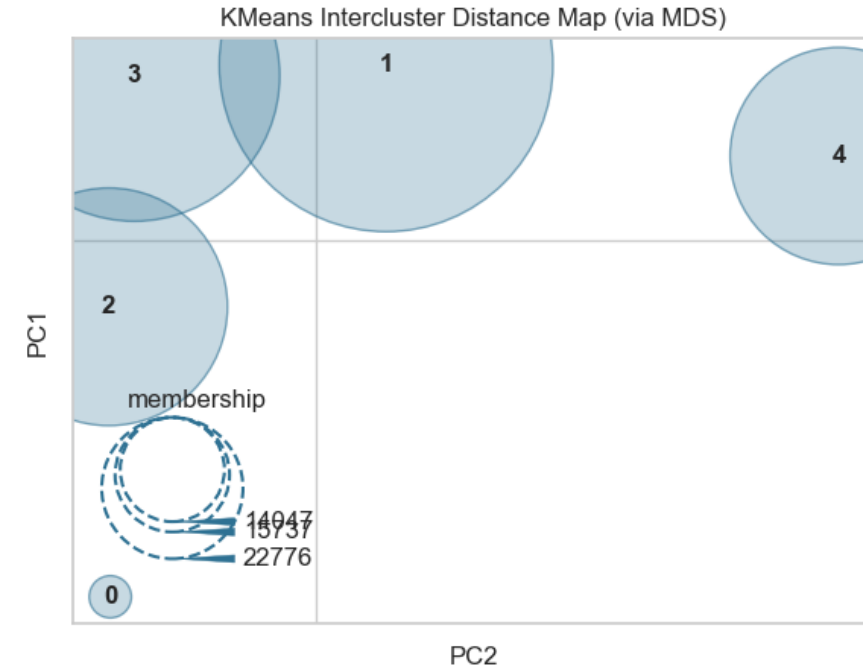
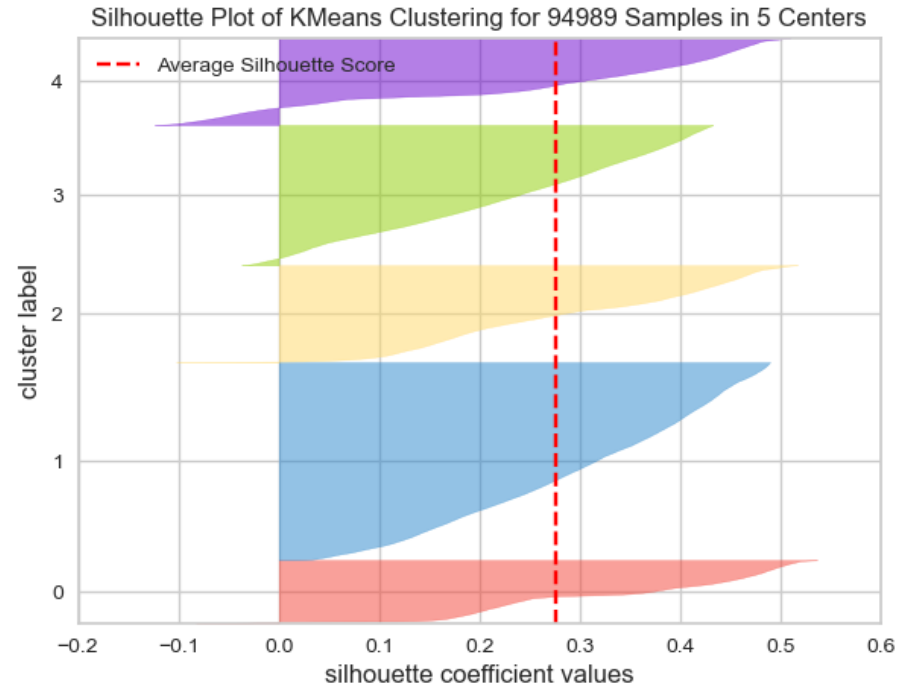
ARI moy = 0,9
ARI Min = 0,62
Pas très stable



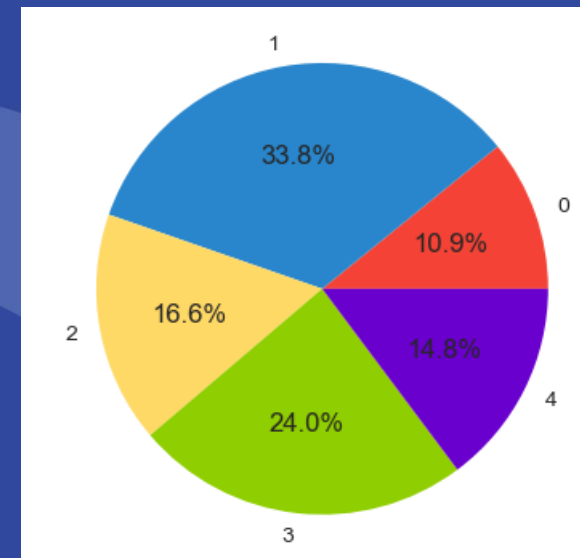
Variable nombre de paiement moyen

5 Clusters:



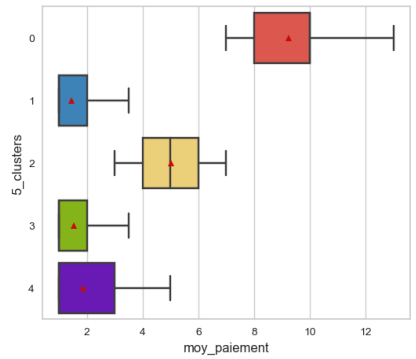
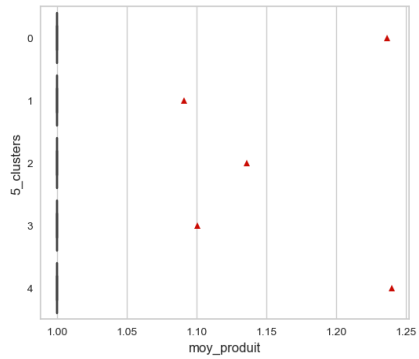
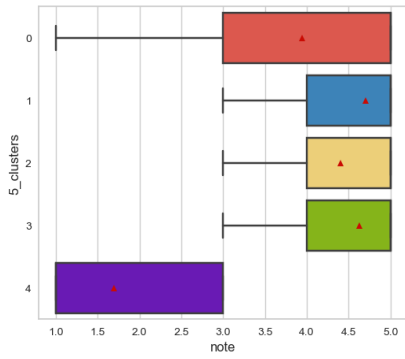
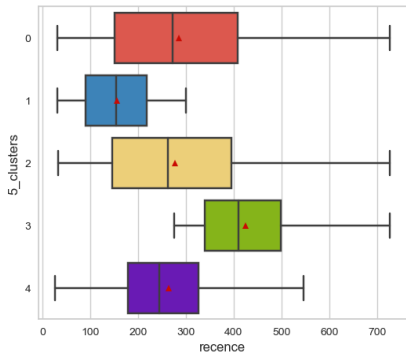
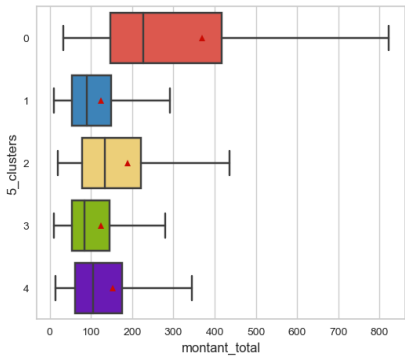
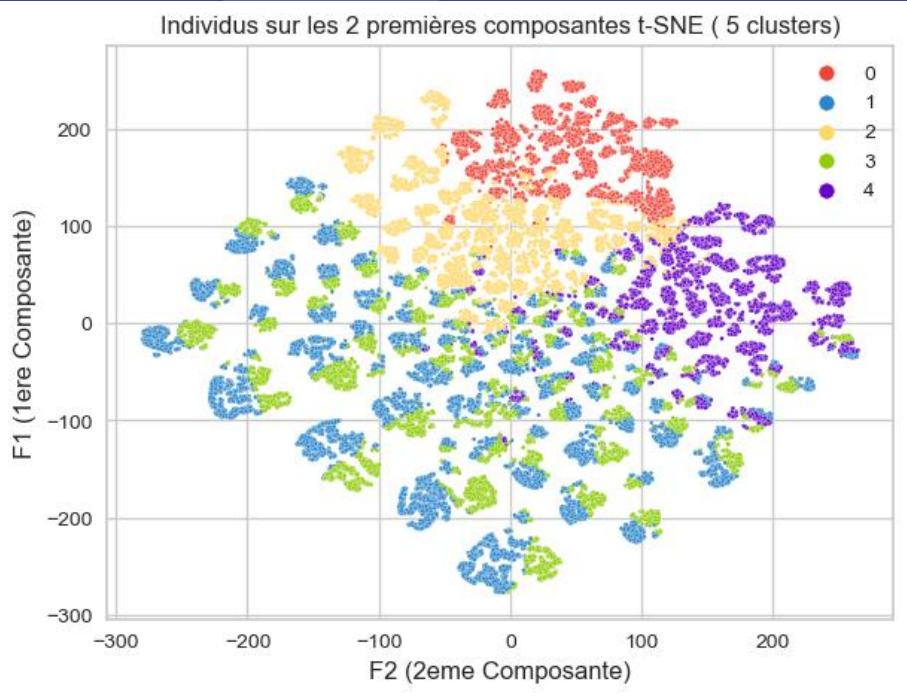
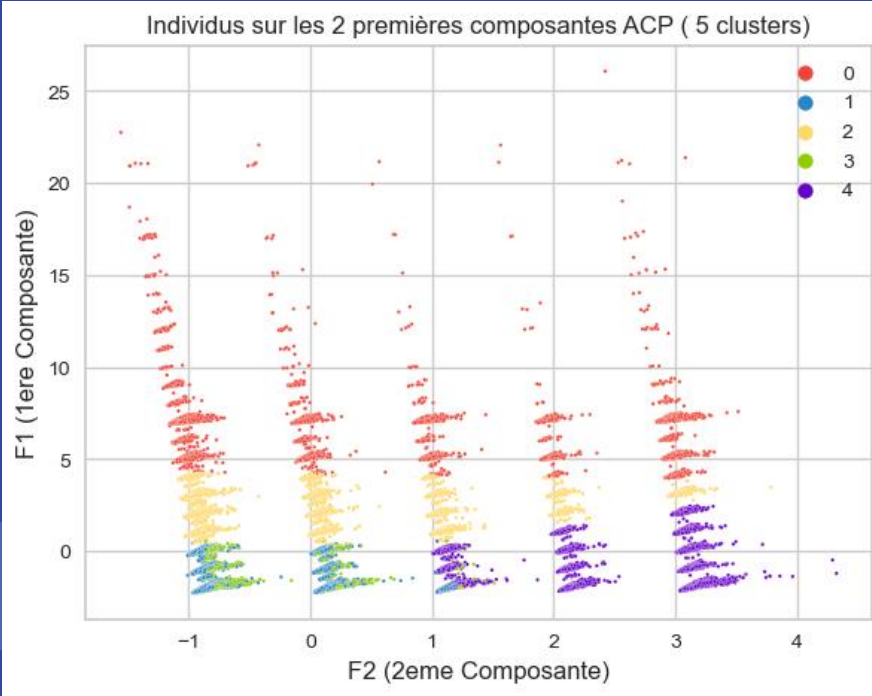


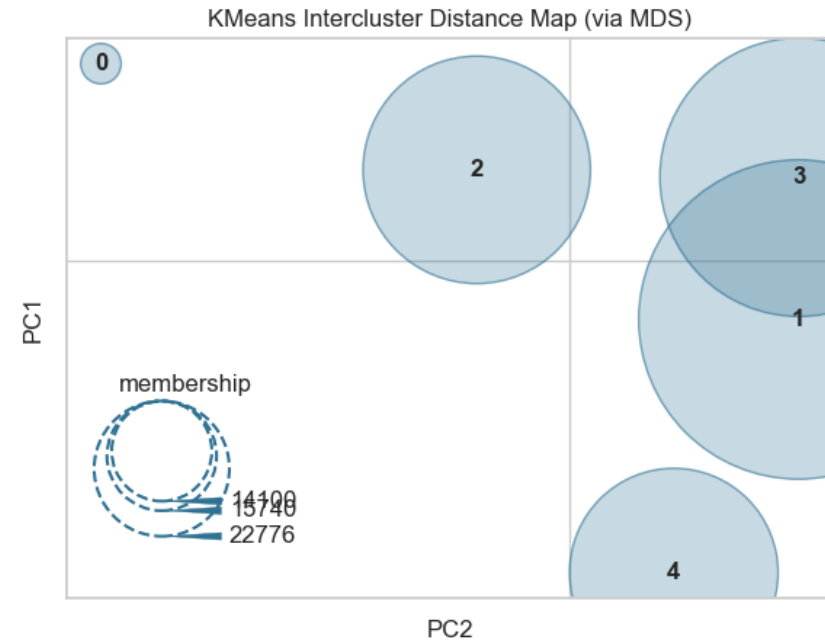
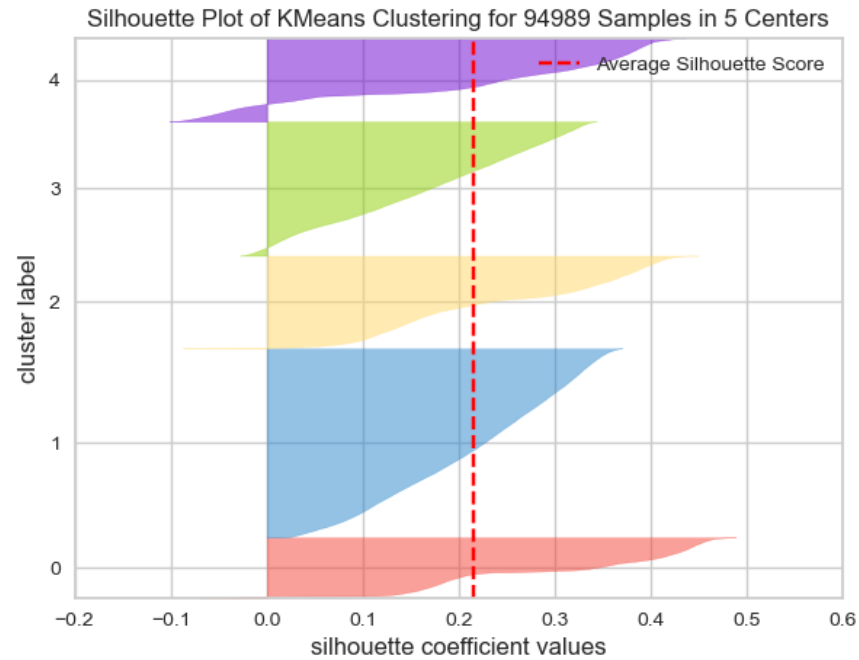
ARI moy = 0,99
ARI Min = 0,94
Très stable



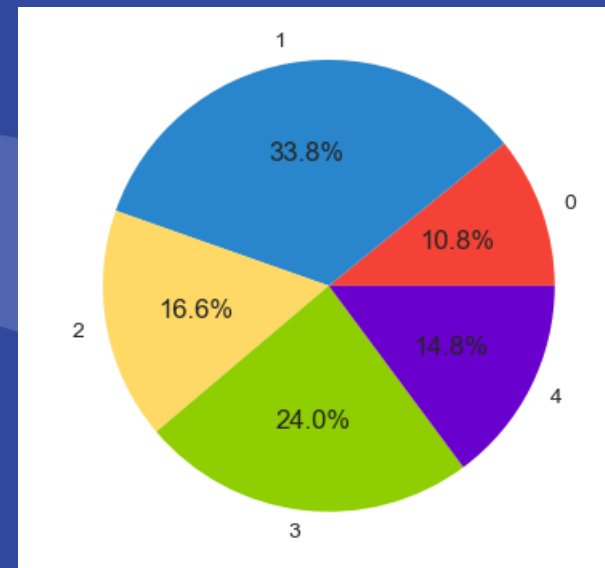
Variable groupe

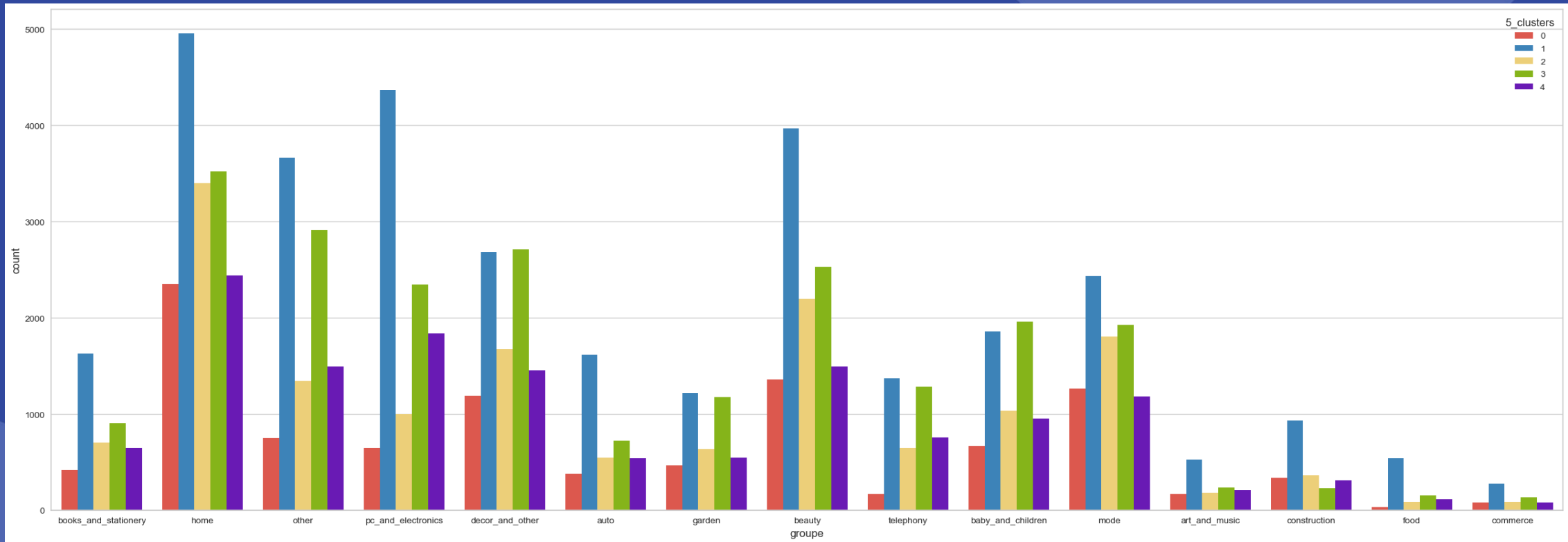
5 Clusters:





ARI moy = 1
ARI Min = 0,99
Très stable





Clustering sélectionné (k-Means)

Le nouveau client (régulier à en devenir):

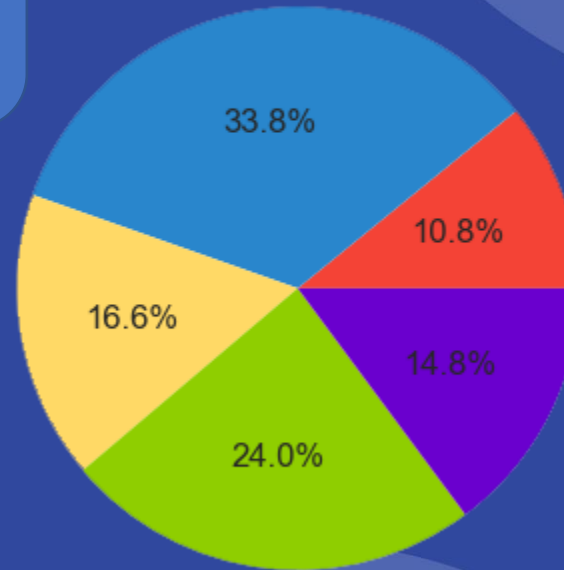
- Dépense peu
- À commander dans les 6-7 derniers mois
- Satisfait des produits

Commande ce dont il a besoin:

- Dépense raisonnablement
- Étale les paiements sur peu de mensualités

Client à récupérer:

- N'as pas commandé durant la dernière année
- Dépense peu mais, est satisfait des produits
- N'étale pas les paiements



Acheteur Compulsif:

- Dépense beaucoup
- Étale les paiements sur un grand nombre de mensualités

Client à séduire:

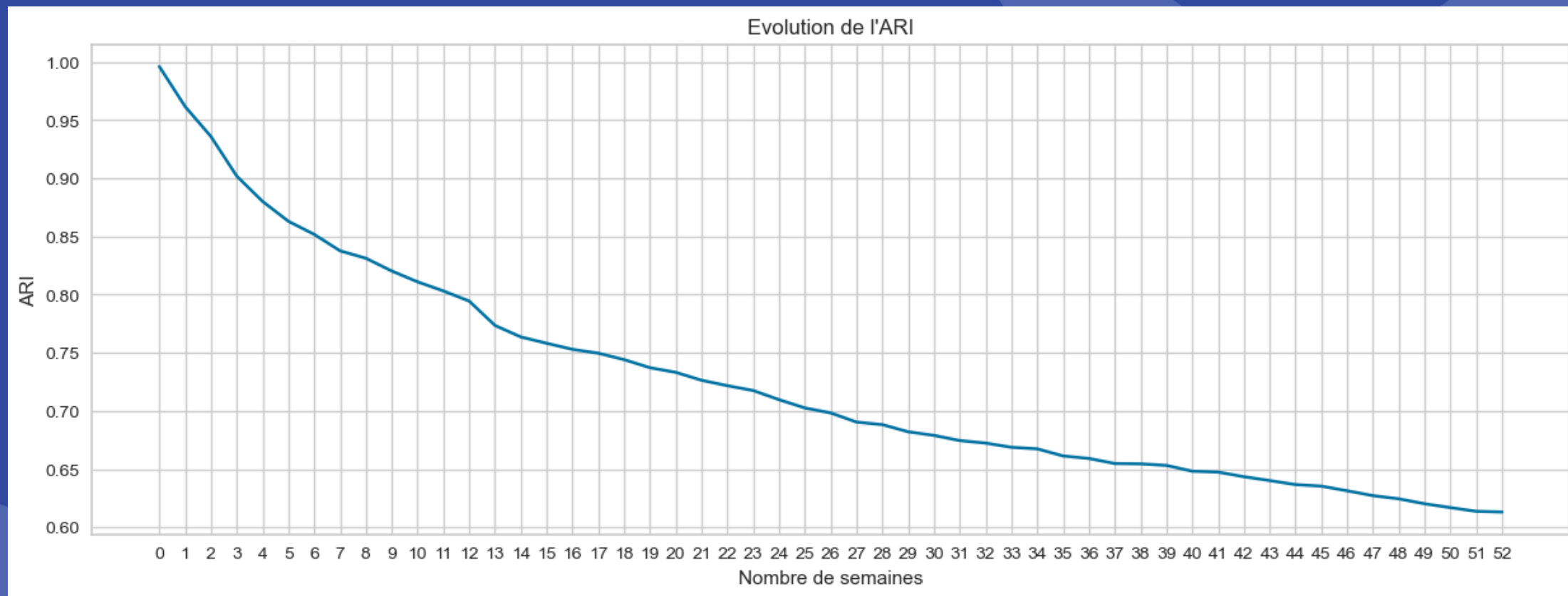
- Peu satisfait des produits

Maintenance du modèle

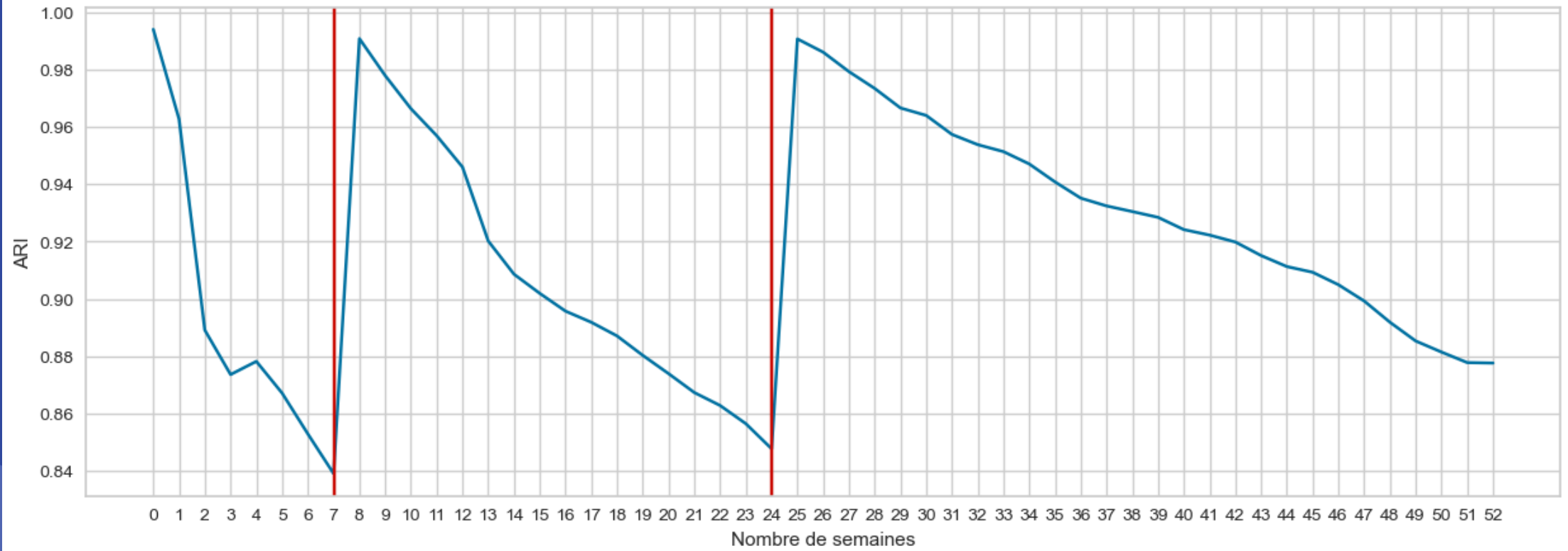
Procédure:

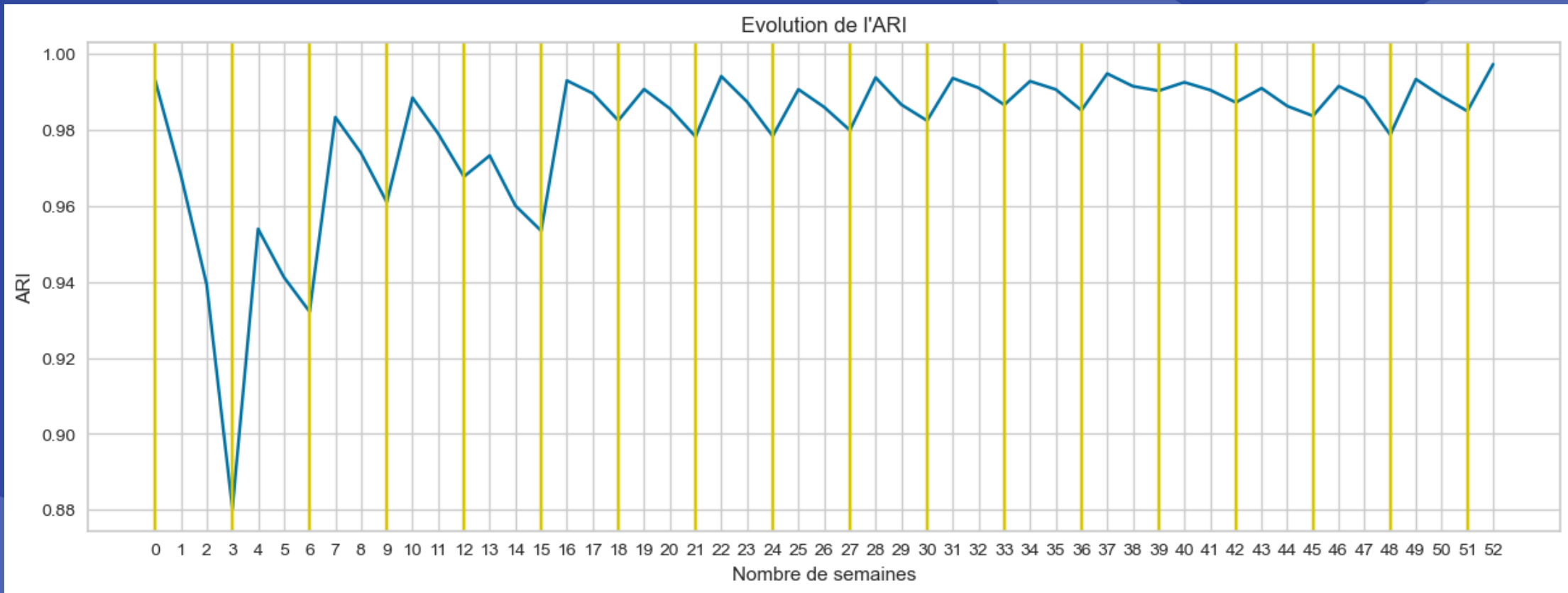
- On commence avec les commandes entre septembre 2016 et septembre 2017 ce qui représentent 23% de l'effectif total des clients.
- Puis on ajoute les commandes semaine par semaine sur 53 semaines, on observe l'évolution de l'ARI.
- On fixe un seuil à 0,85 pour l'ARI.
- On observe la fréquence à laquelle l'ARI atteint ce seuil.
- On définit une fréquence de mise à jour du cluster.

Evolution de l'ARI sur 53 semaines



Evolution de l'ARI





Une mise à jour avec une fréquence de 3 semaines semble bien adaptée,
ce qui conserve l'ARI entre 1 et 0,88



MERCI DE VOTRE ATTENTION