
Segmentation des clients du site e-commerce OLIST



Sommaire

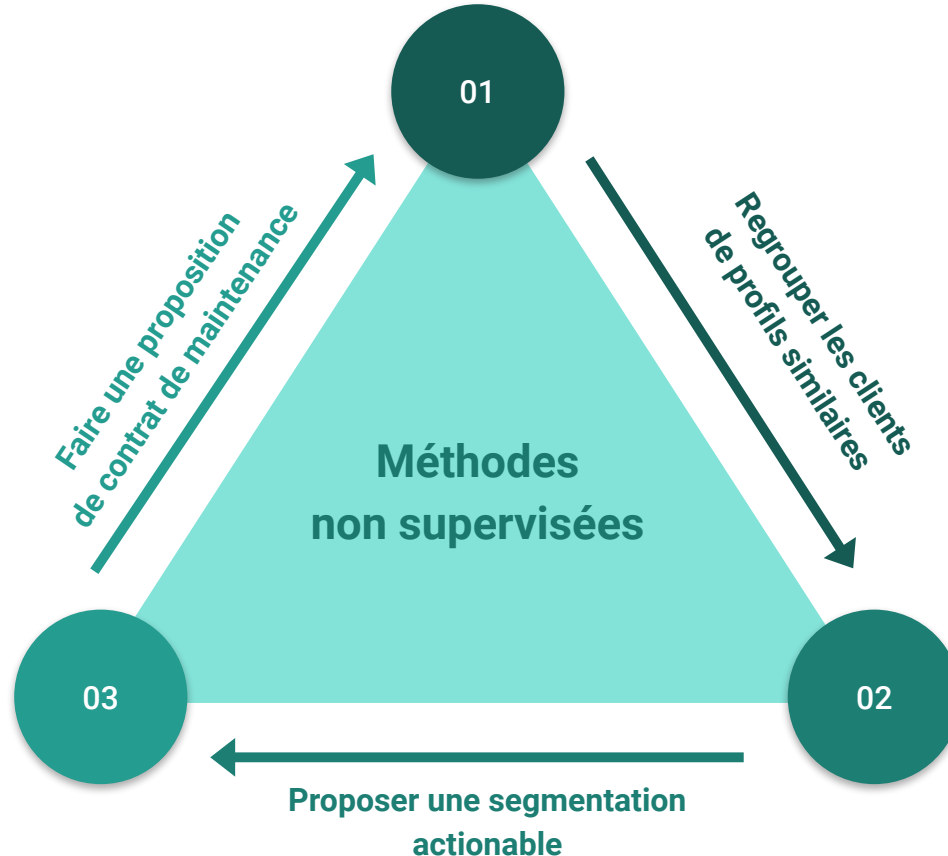
1. Enjeux
2. Analyse exploratoire
3. Pistes de modélisation
4. Etude de stabilité



Enjeux

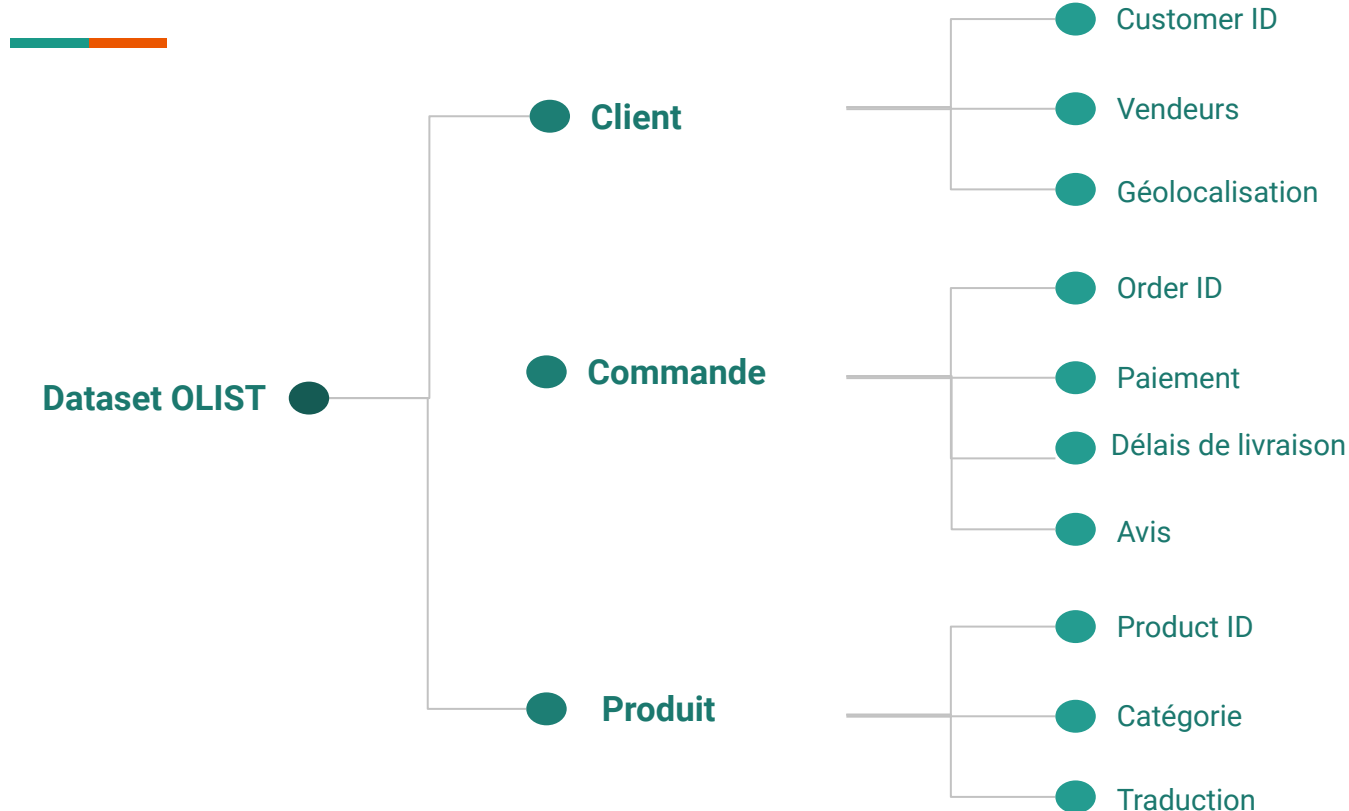


Fournir aux employés OLIST un outil de segmentation des clients



Analyse exploratoire

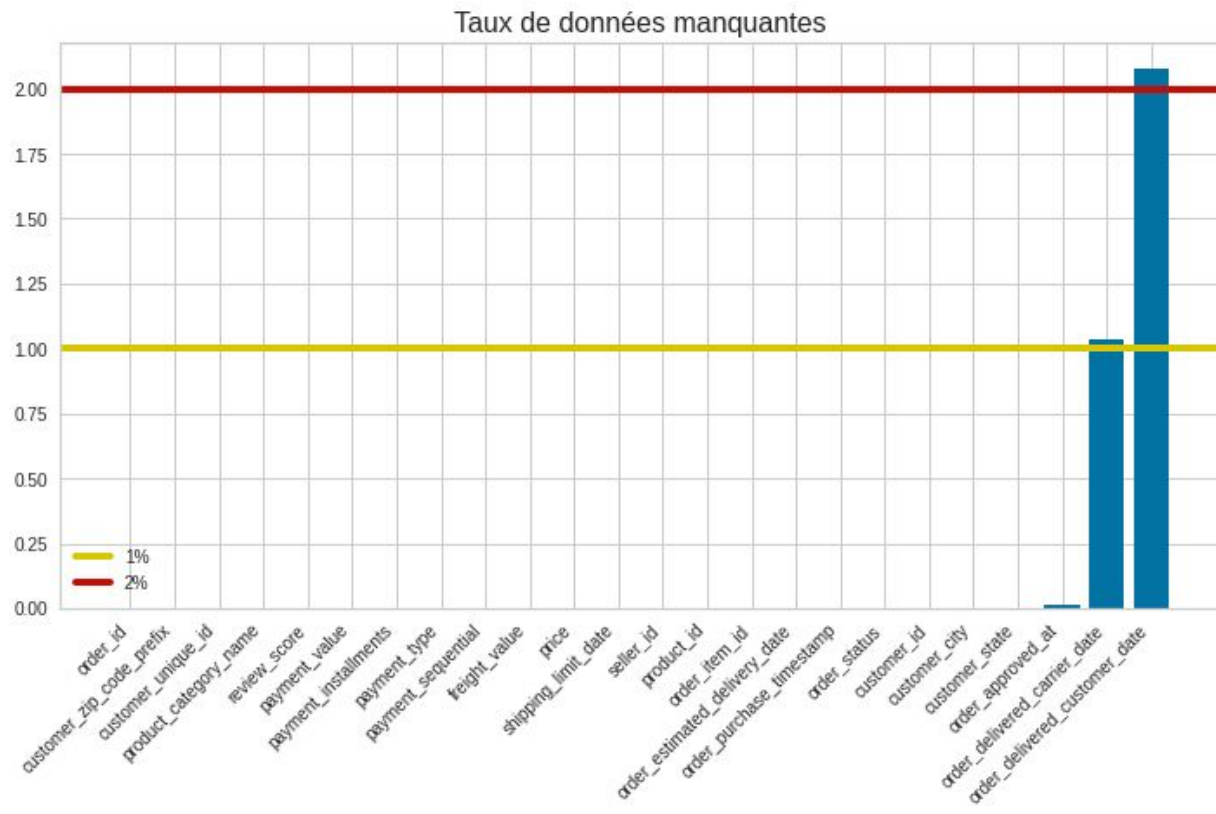
Base de données



Jeu de données final



Valeurs manquantes



Traitement des NaN

NaN avant

order_purchase_timestamp	0
order_approved_at	14
order_delivered_carrier_date	1195
order_delivered_customer_date	2400
order_estimated_delivery_date	0



NaN vs 'order_status'

shipped	1138
canceled	529
invoiced	358
processing	357
delivered	8
unavailable	7
approved	3



'order_status'

delivered	0.979
shipped	0.010
canceled	0.005
invoiced	0.003
processing	0.003
unavailable	0.000
approved	0.000

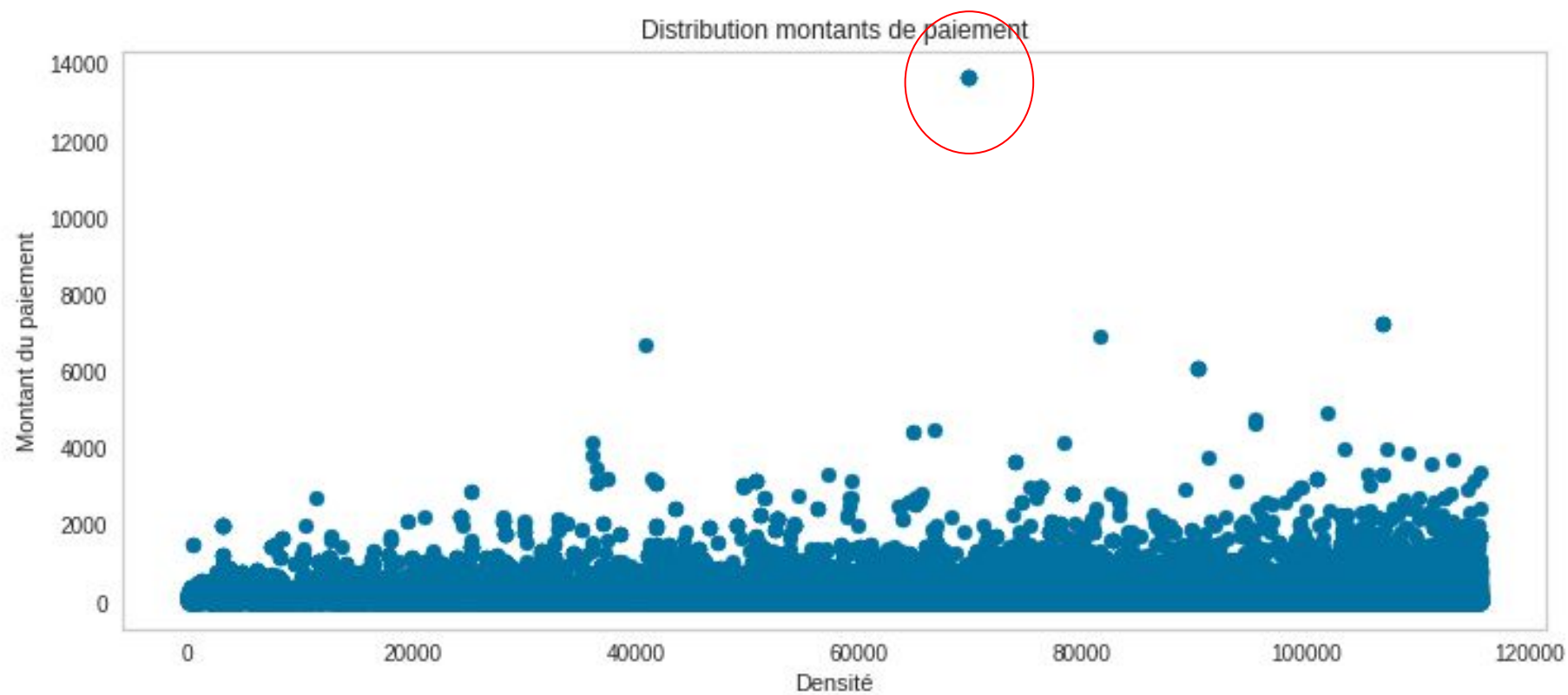
drop

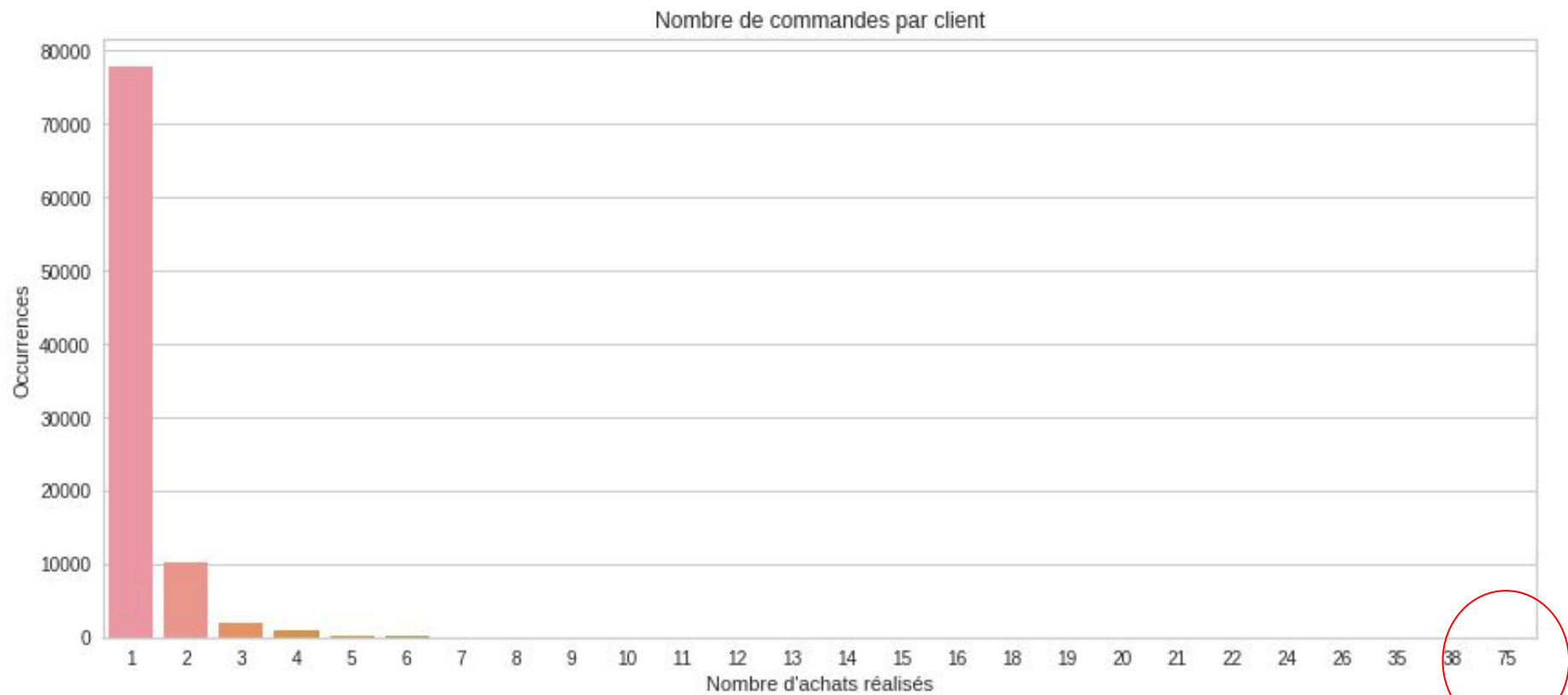


NaN après

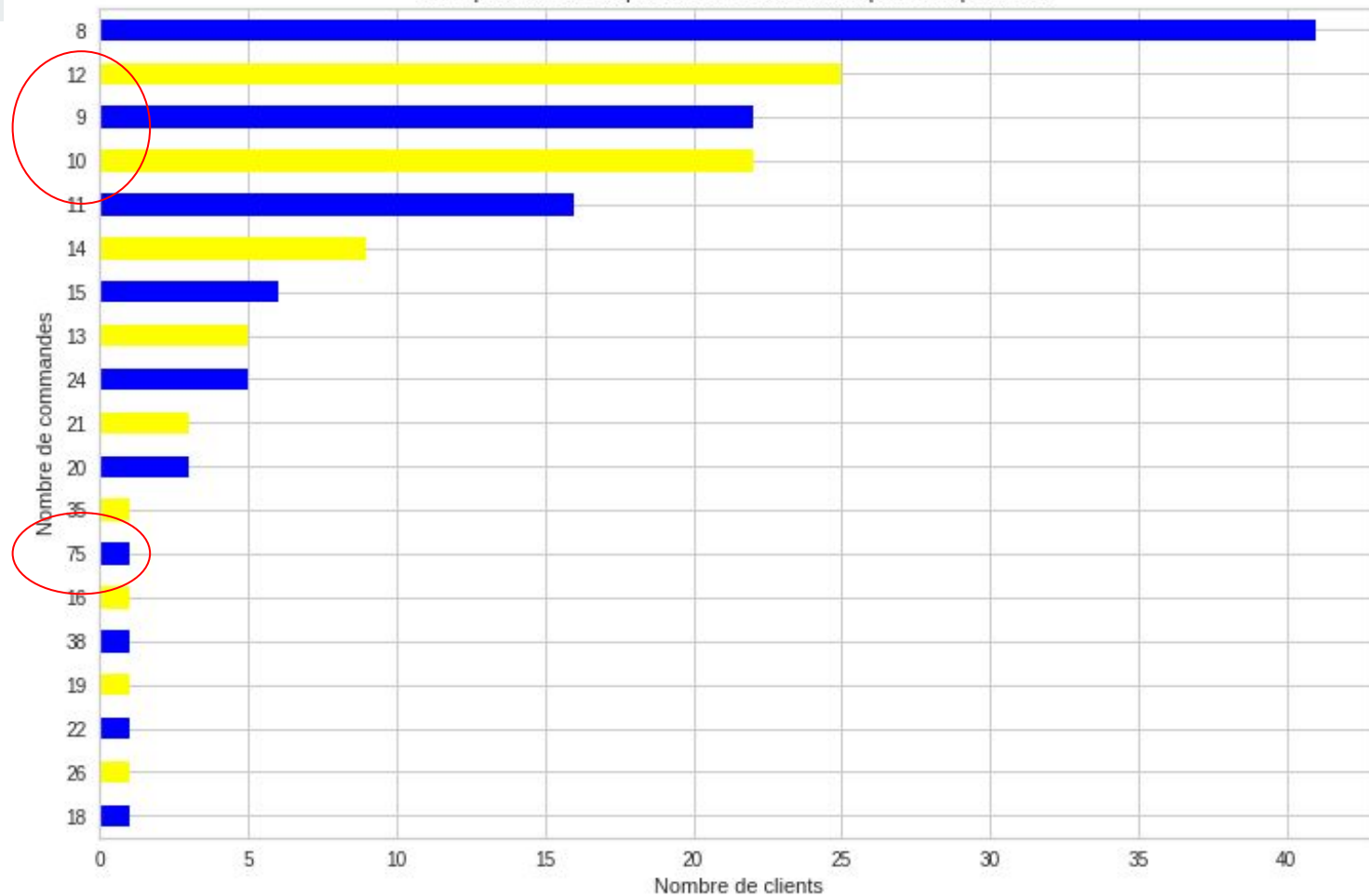
order_status	0
order_purchase_timestamp	0
order_approved_at	14
order_delivered_carrier_date	2
order_delivered_customer_date	8

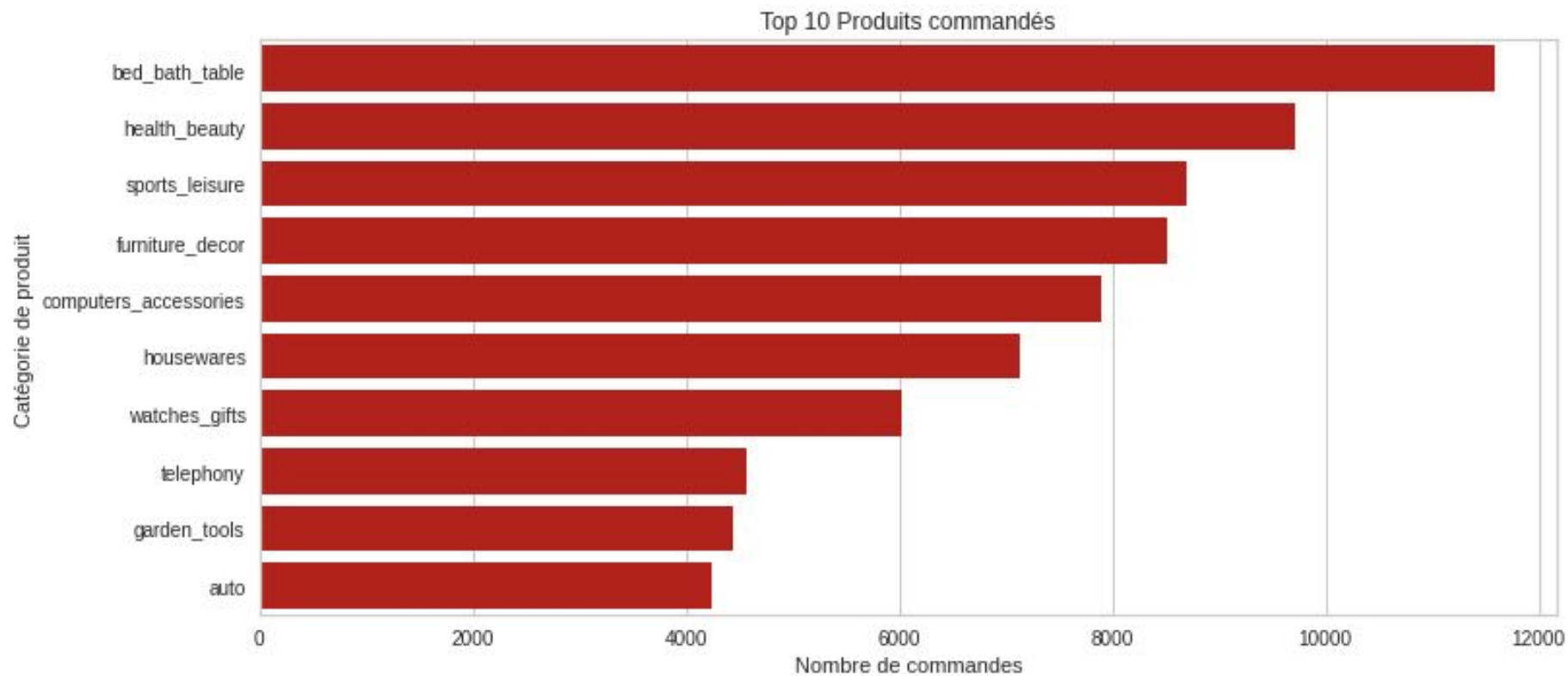
Exploration des variables

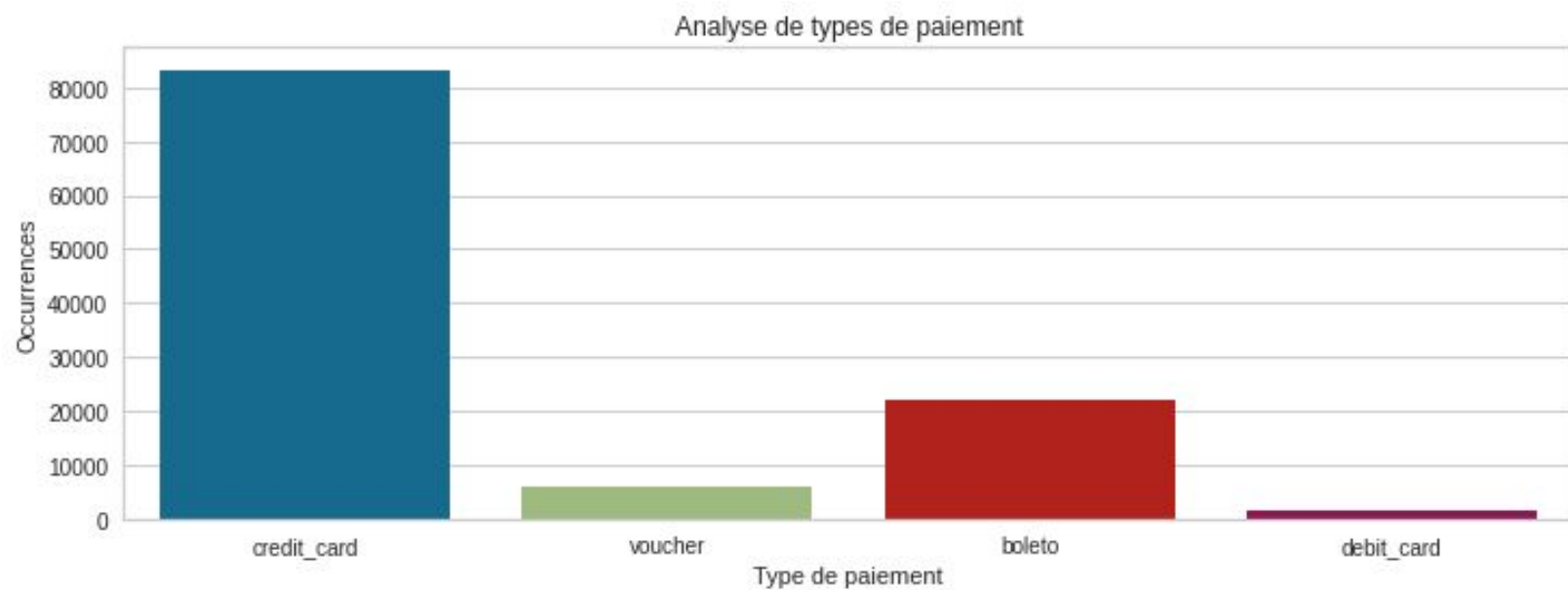


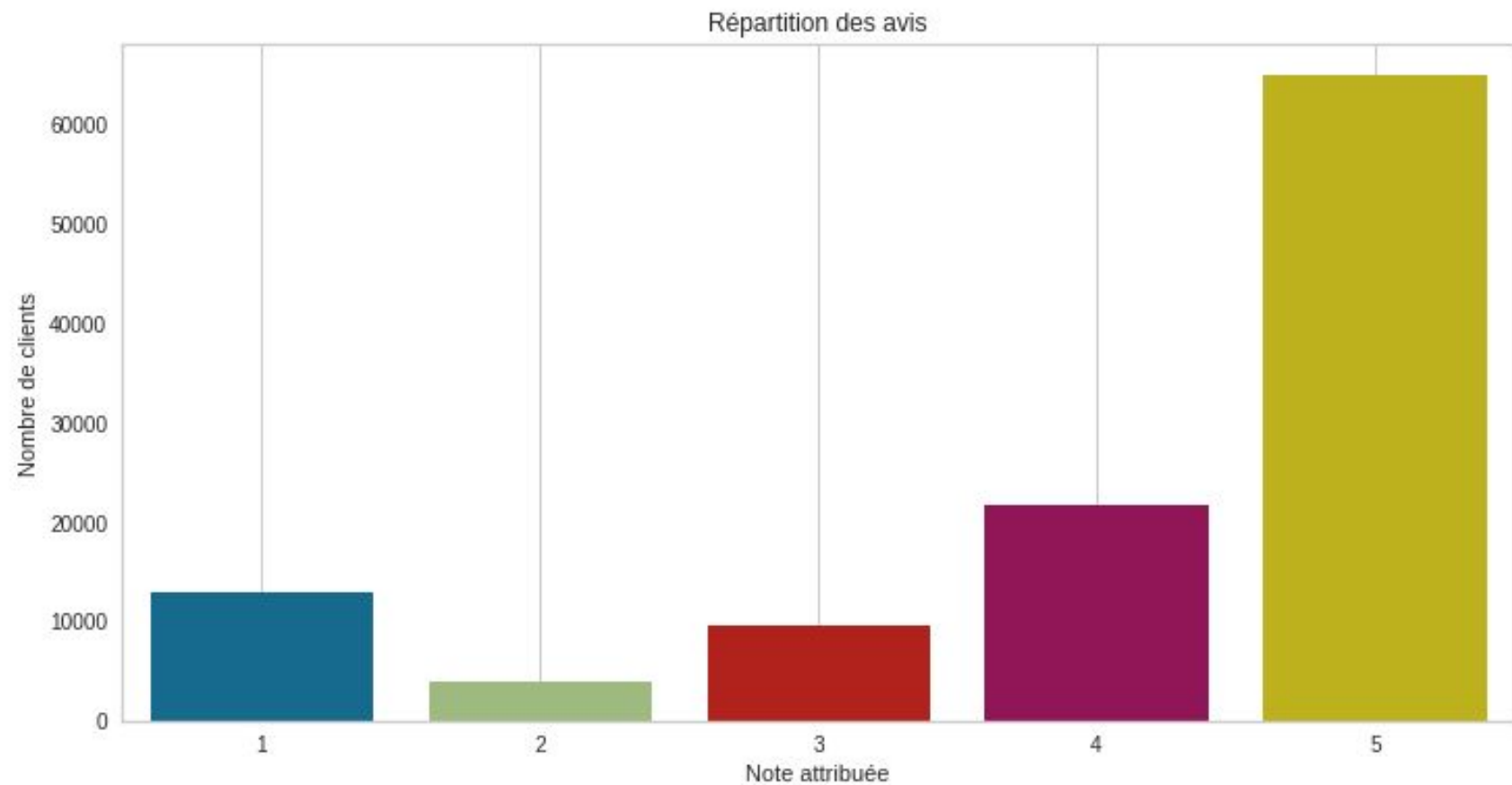


La répartition de la quantité de commandes passées par client

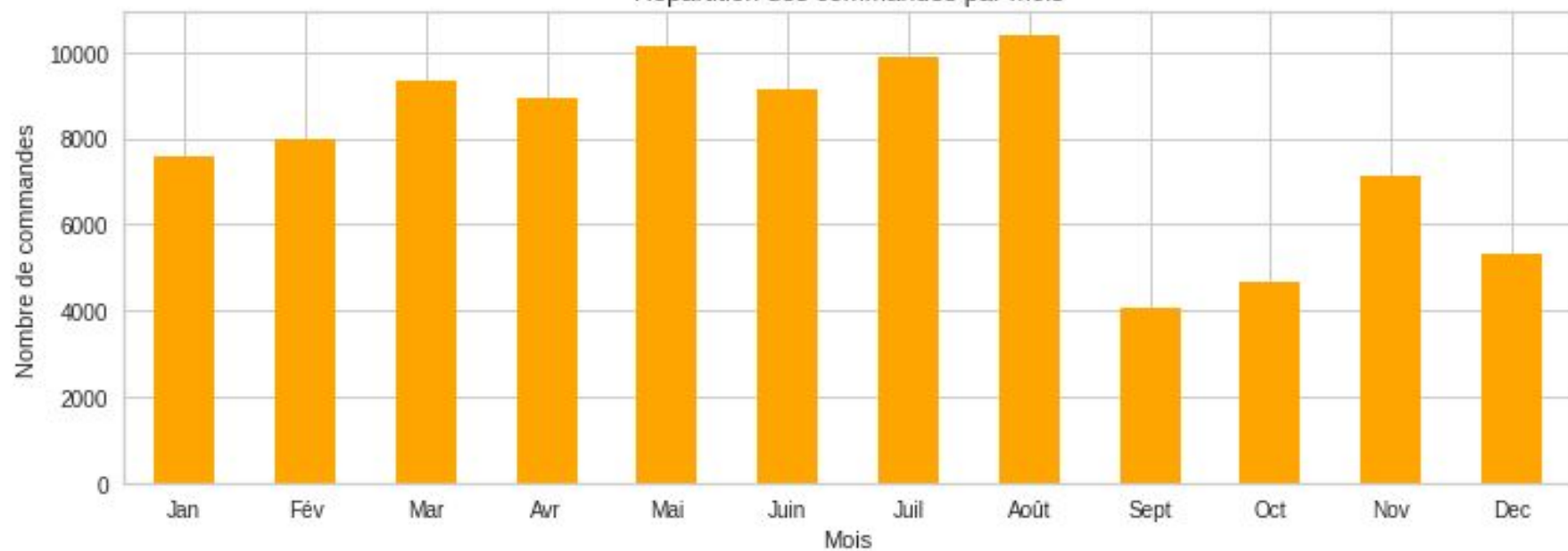




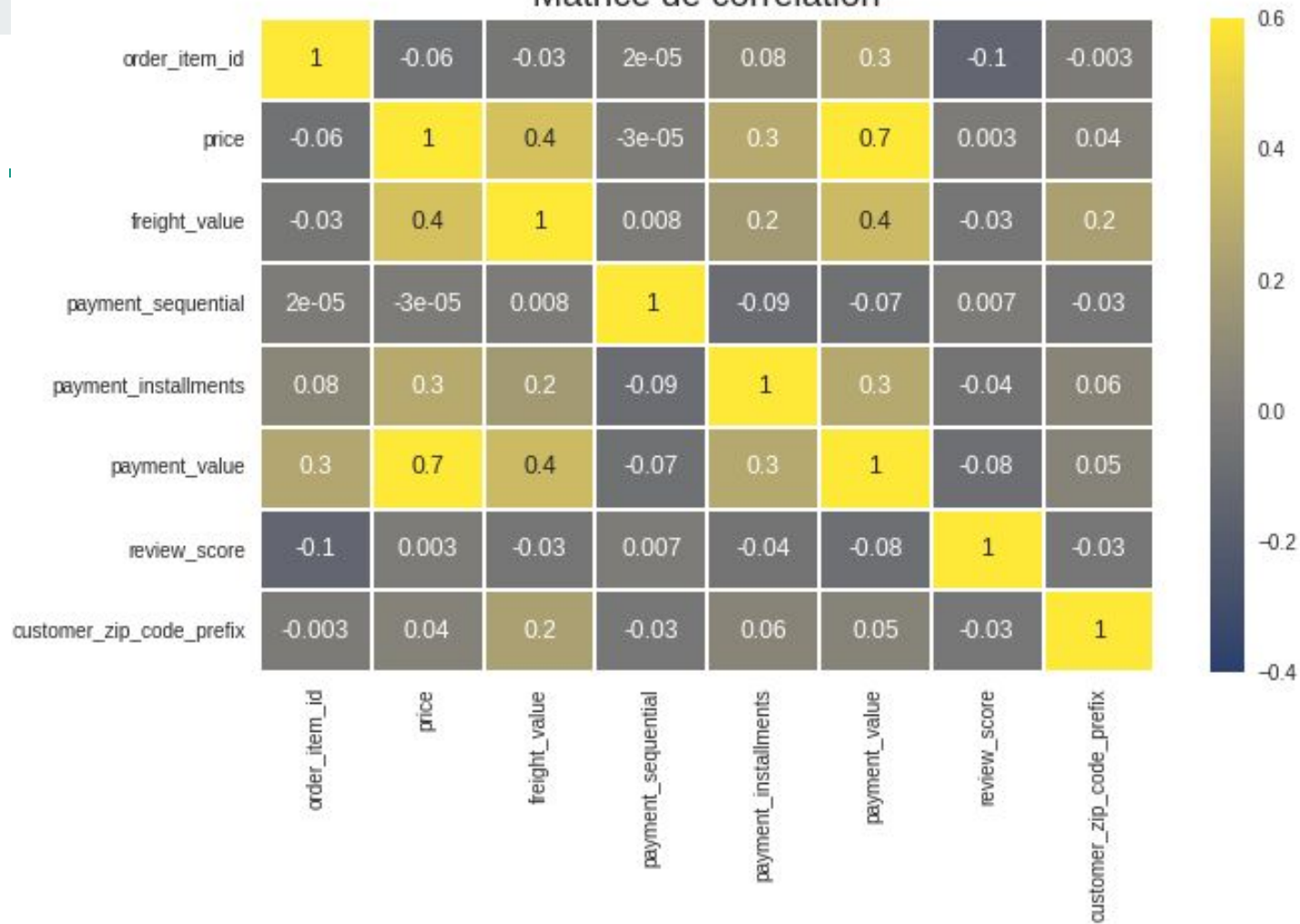




Répartition des commandes par mois



Matrice de corrélation



Pistes de modélisation

Segmentation RFM



Récence

Fréquence

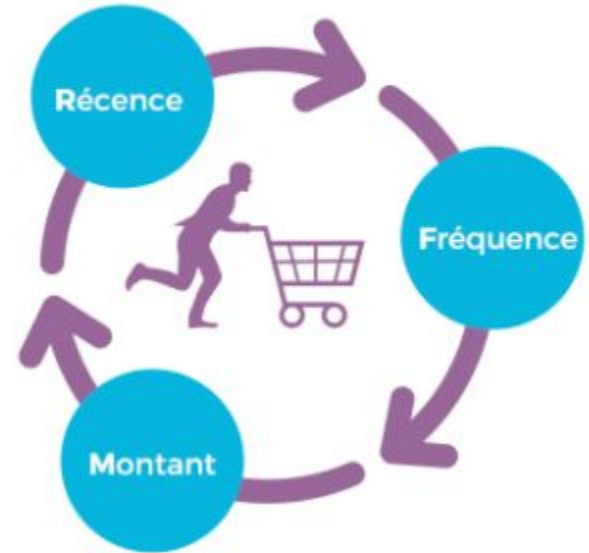
Montant

Ceux qui ont acheté le plus récemment

Ceux qui ont acheté le plus fréquemment

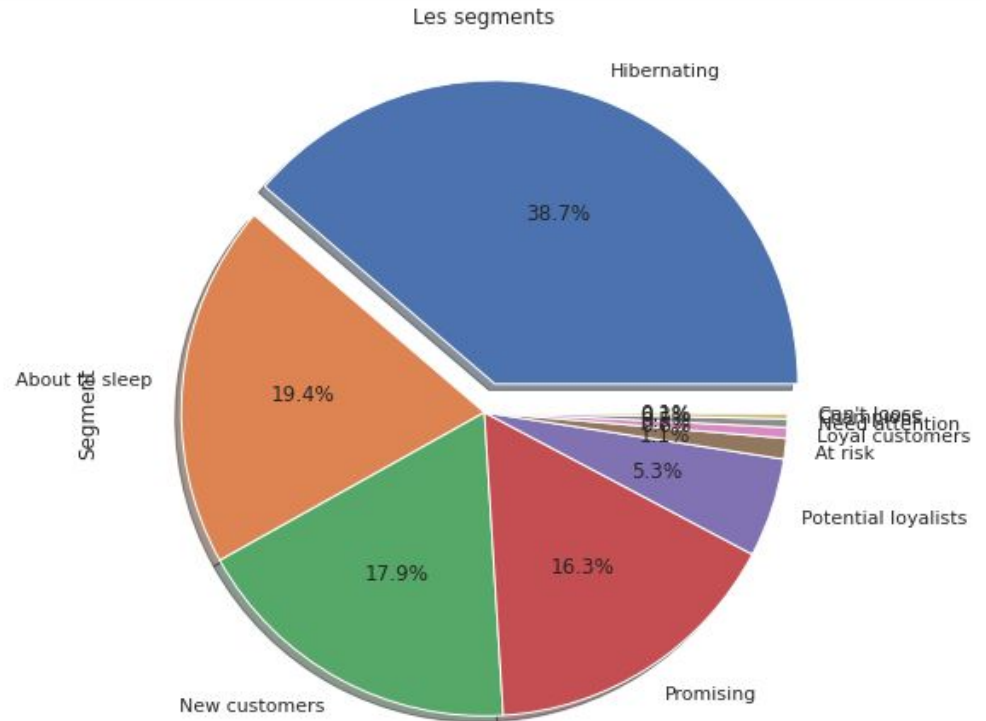
Ceux qui dépensent le plus

Segmentation RFM



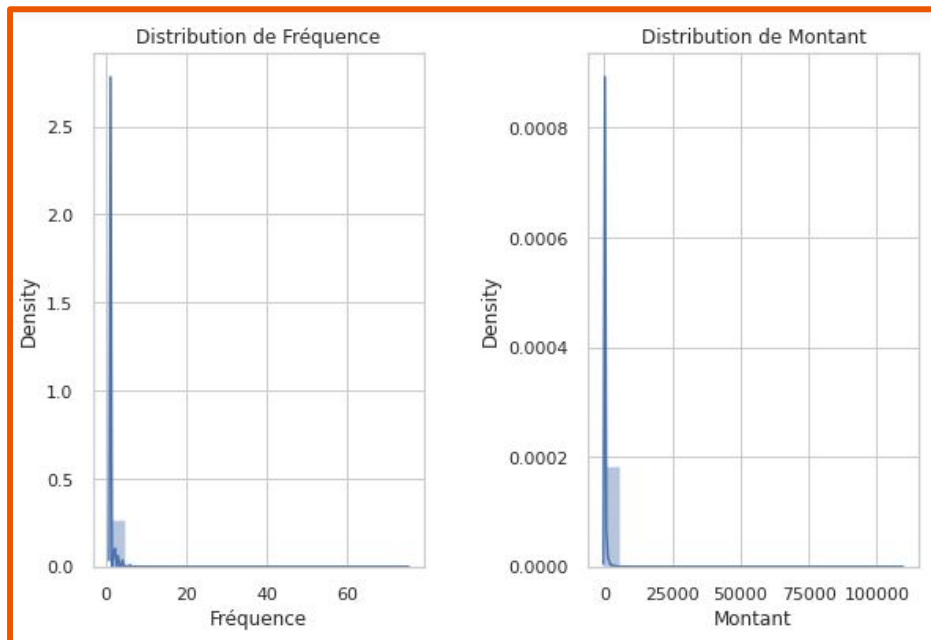
Profils clients

Segment	Récence	Fréquence	Montant	
	mean	mean	mean	count
About to sleep	122.500	1.100	165.100	7865
At risk	237.900	3.700	501.200	447
Can't loose	224.400	8.500	413.200	21
Champions	26.100	4.800	567.500	117
Hibernating	236.900	1.100	159.600	15723
Loyal customers	91.800	4.700	596.900	232
Need attention	122.300	3.000	496.200	183
New customers	24.100	1.000	137.200	7261
Potential loyalists	43.700	2.200	342.100	2144
Promising	62.600	1.000	149.400	6604

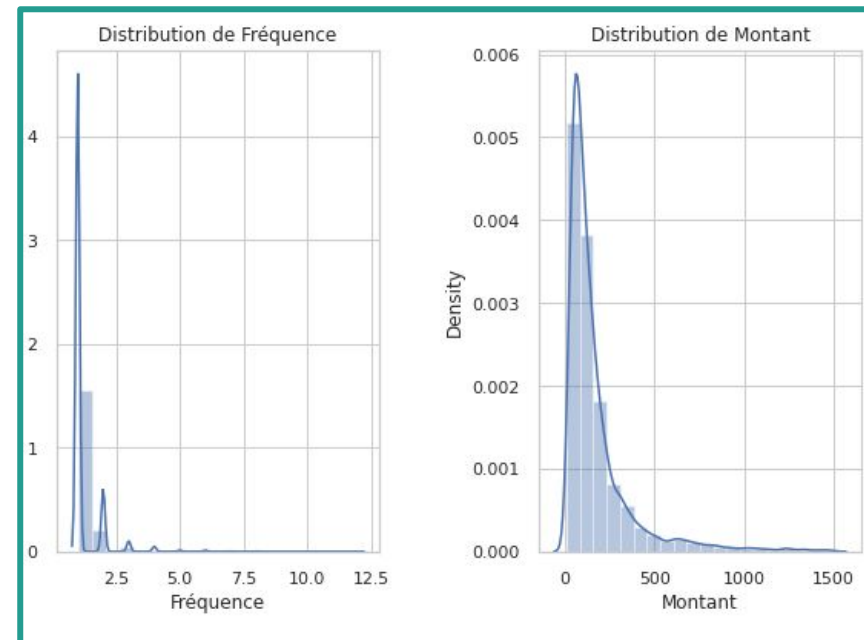


Outliers

AVANT

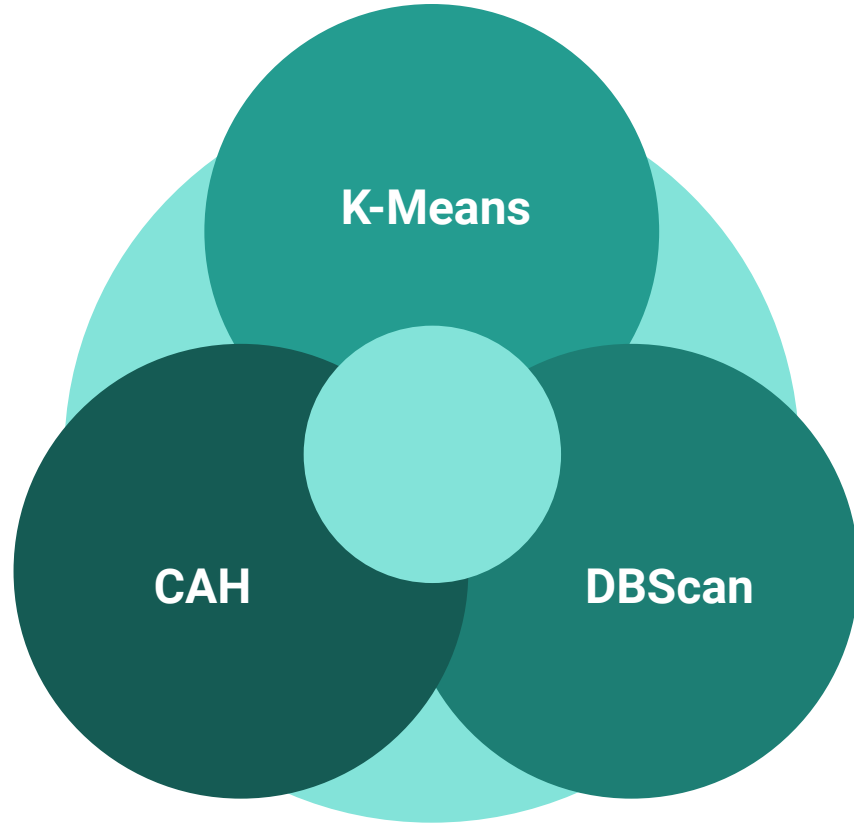


APRES

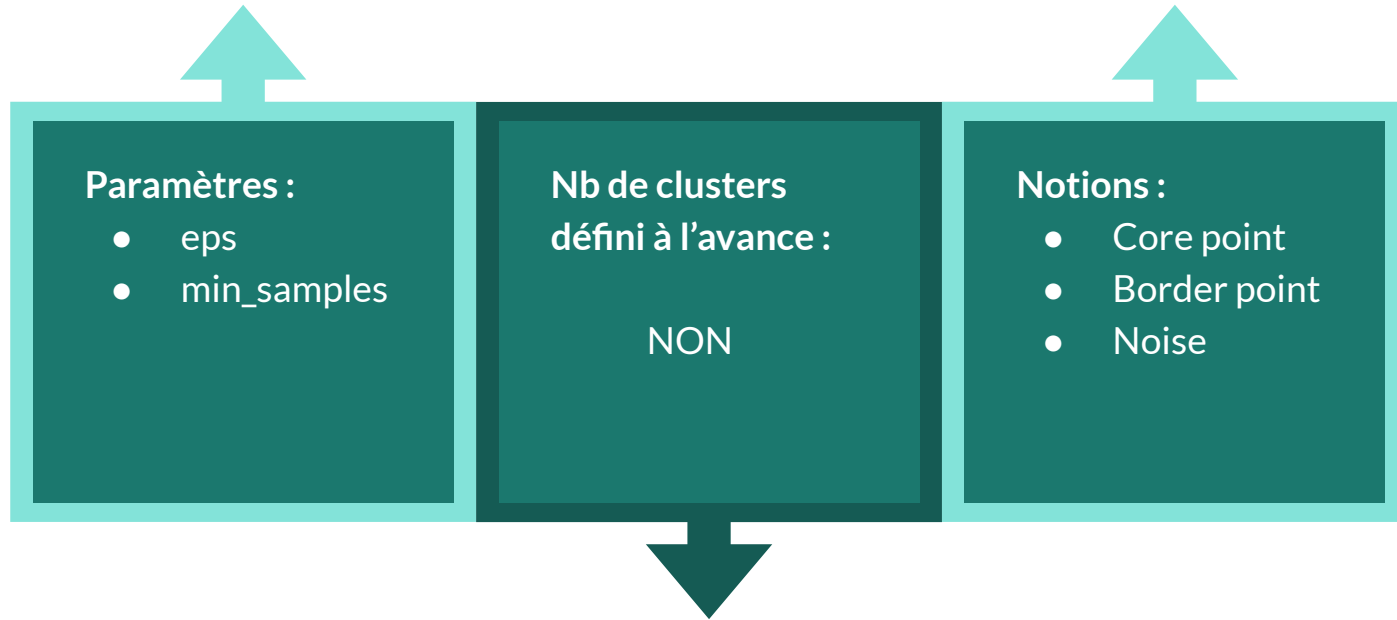


Clustering

Algorithmes utilisés

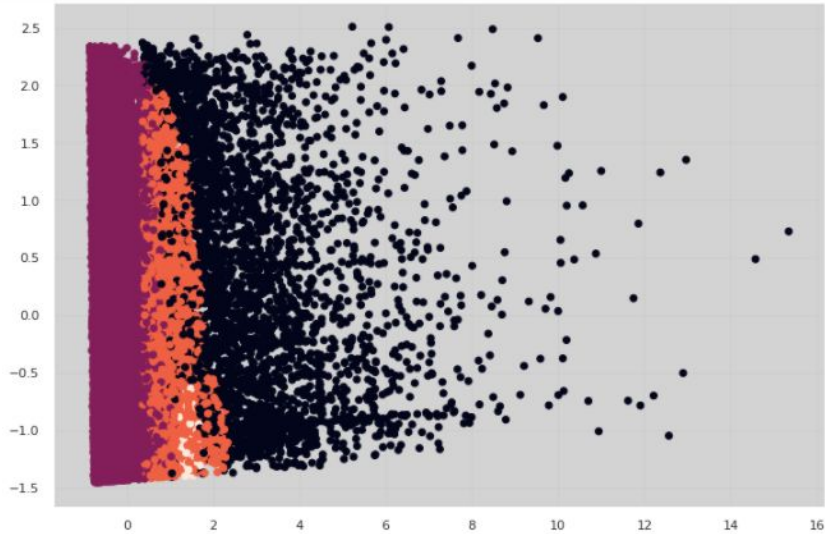


DBScan



Essai n° 1 :

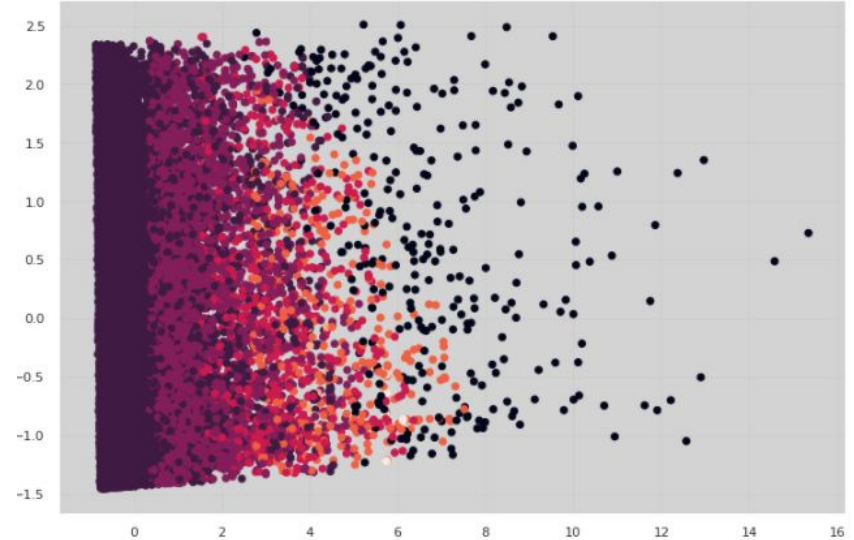
- min_samples = 100
- eps = 0.3, 0.4, 0.5



- eps=0.3 : 3 clusters et 3549 anomalies
- eps=0.4 : 2 clusters et 2759 anomalies
- eps=0.5 : 2 clusters et 2311 anomalies

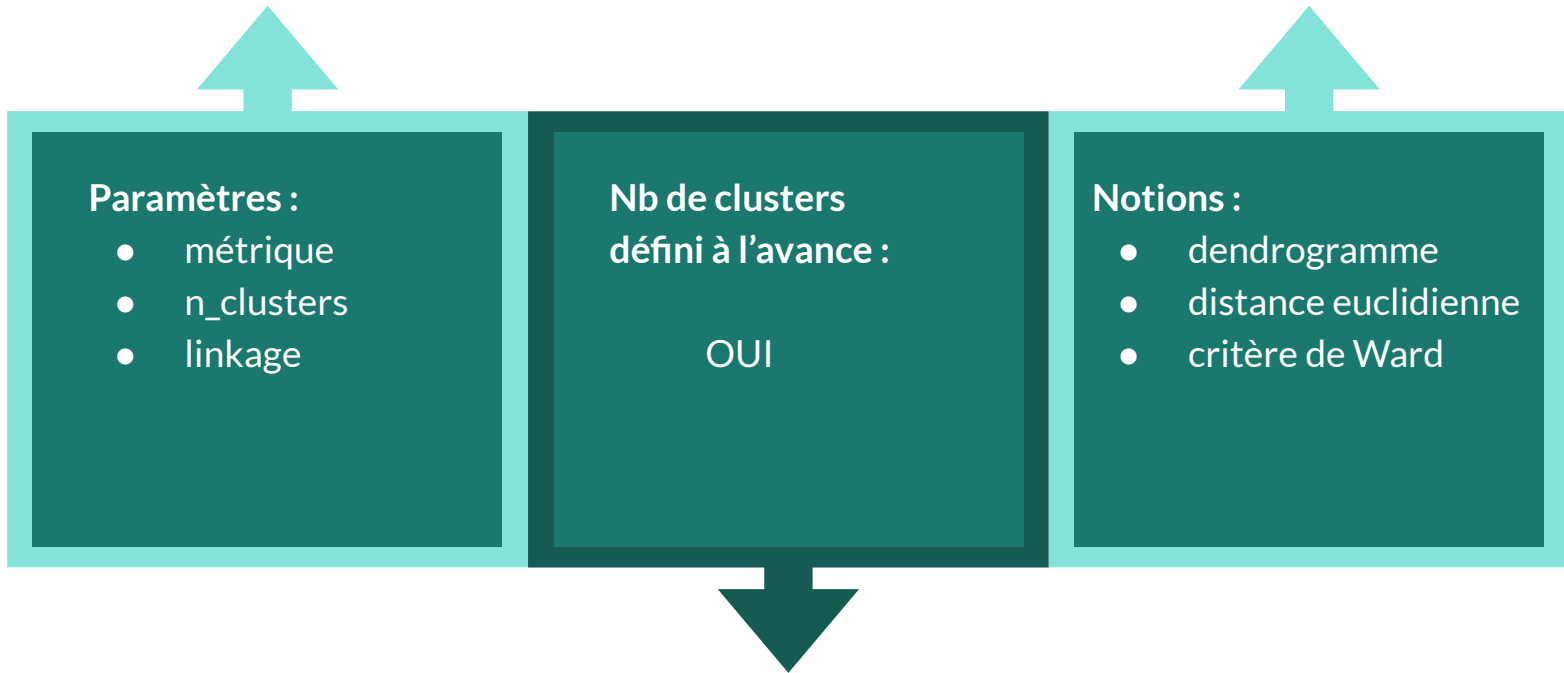
Essai n° 2 :

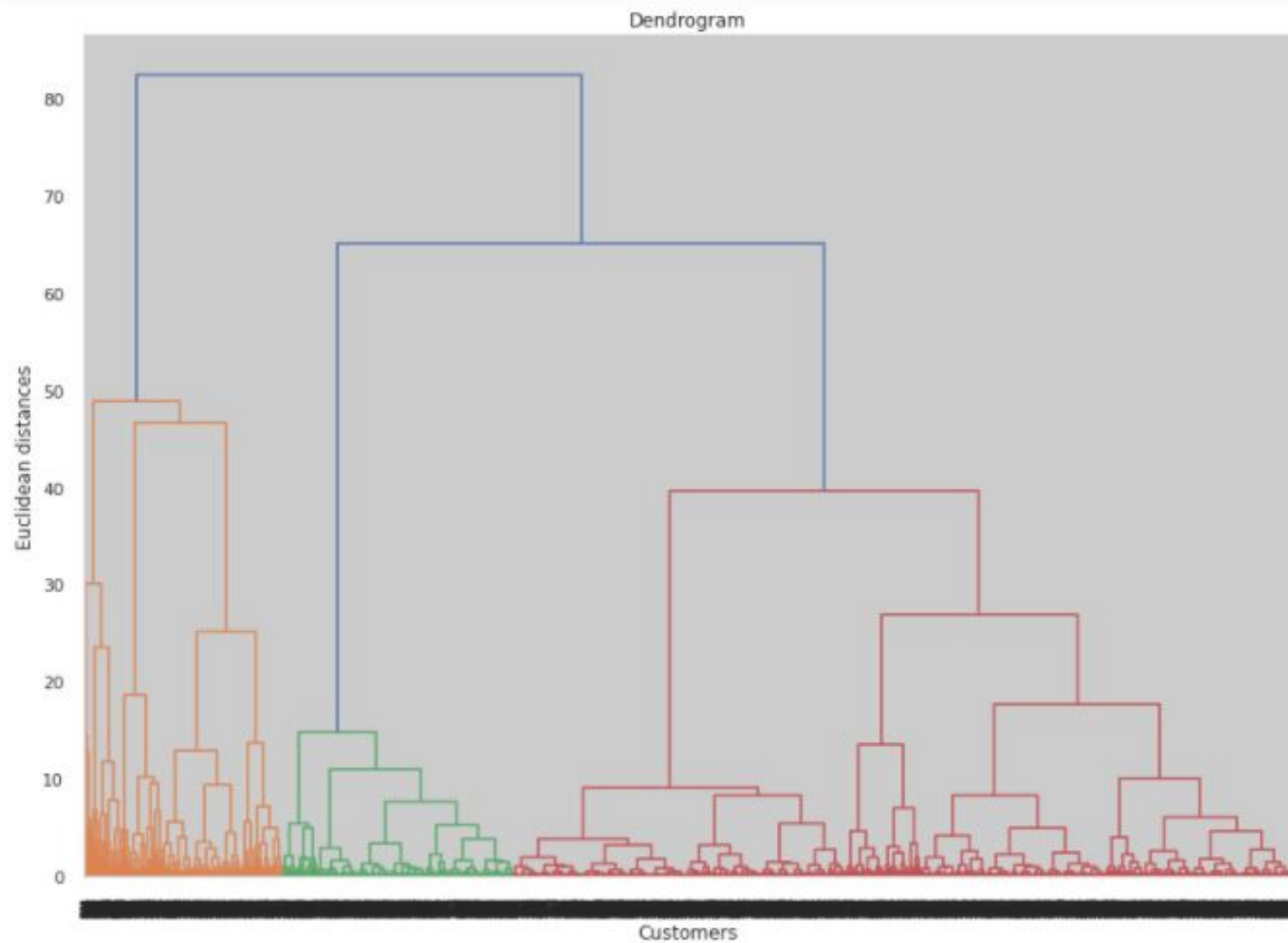
- min_samples = 10
- eps = 0.3, 0.4, 0.5

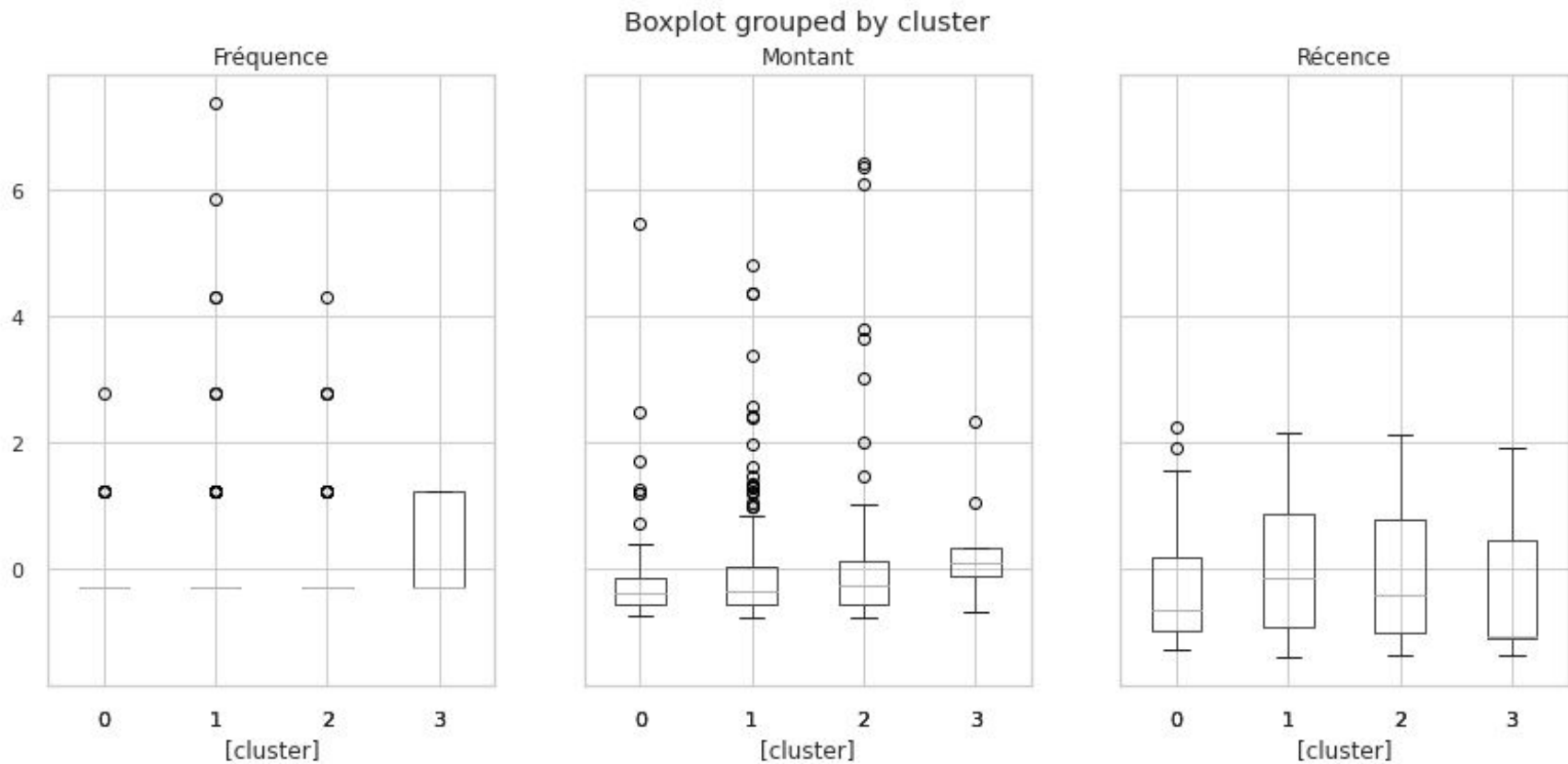


- eps=0.3 : 17 clusters et 854 anomalies
- eps=0.4 : 9 clusters et 450 anomalies
- eps=0.5 : 6 clusters et 298 anomalies

Classification Ascendante Hiérarchique

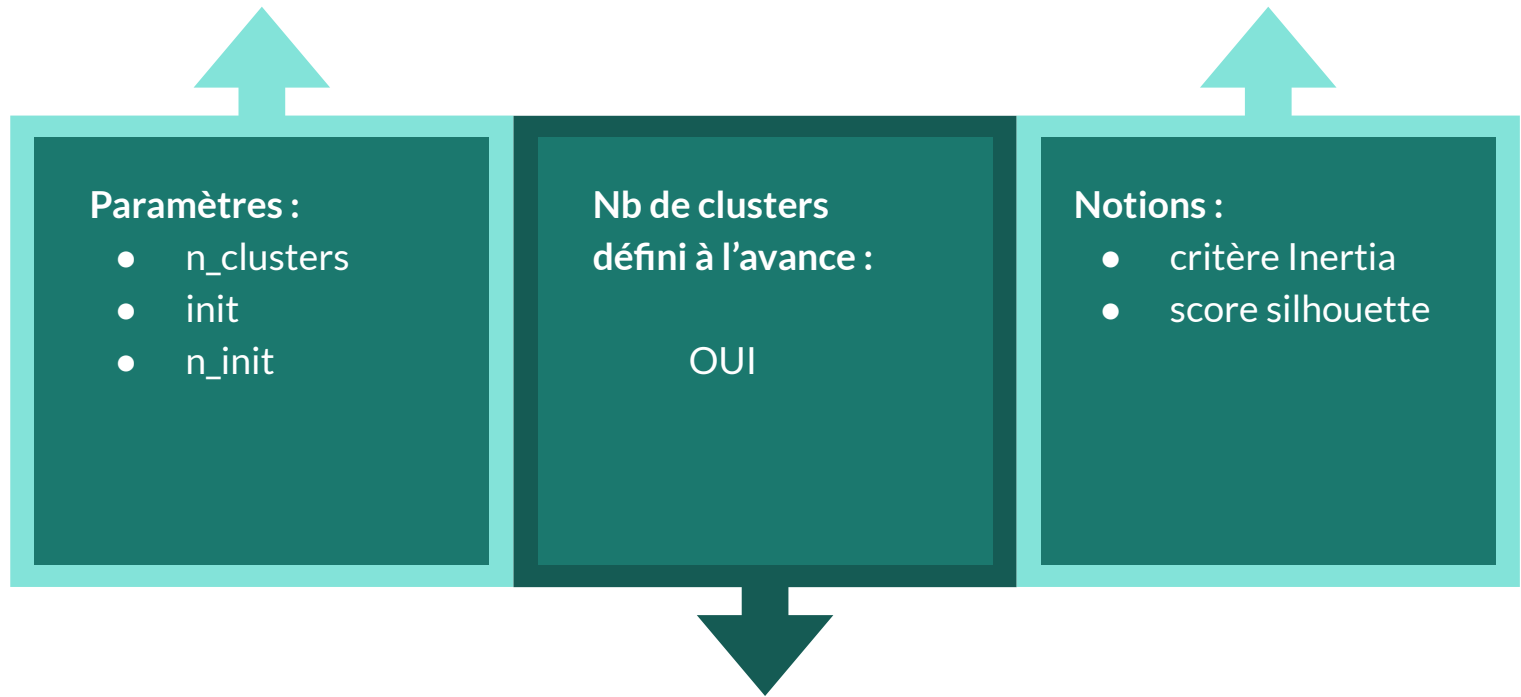






Résultats satisfaisants MAIS lents à l'exécution...

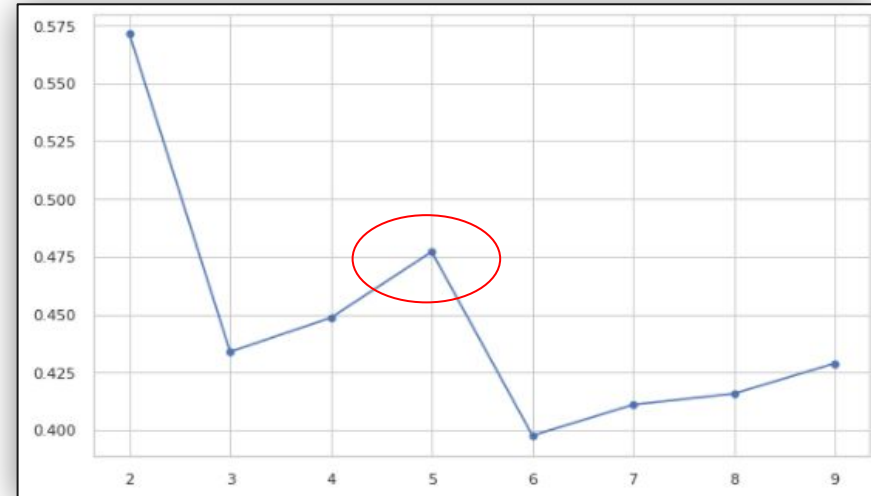
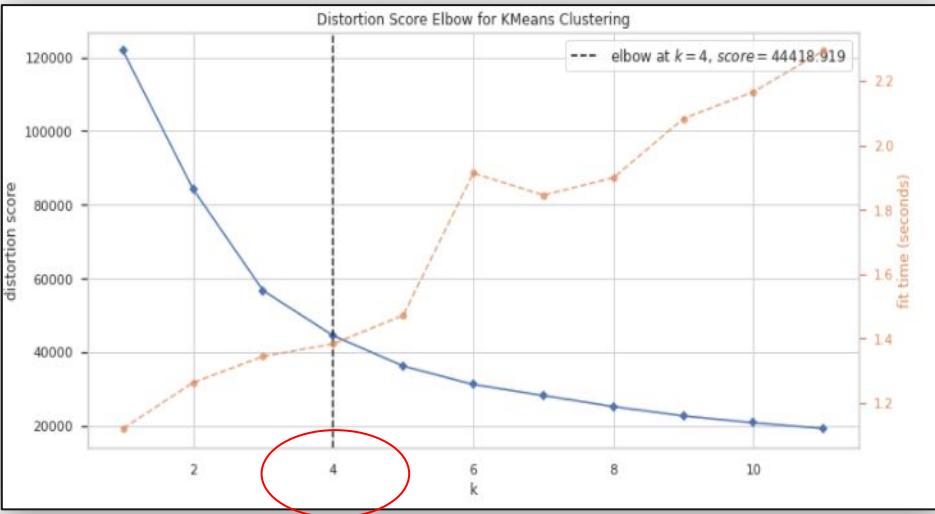
K-Means

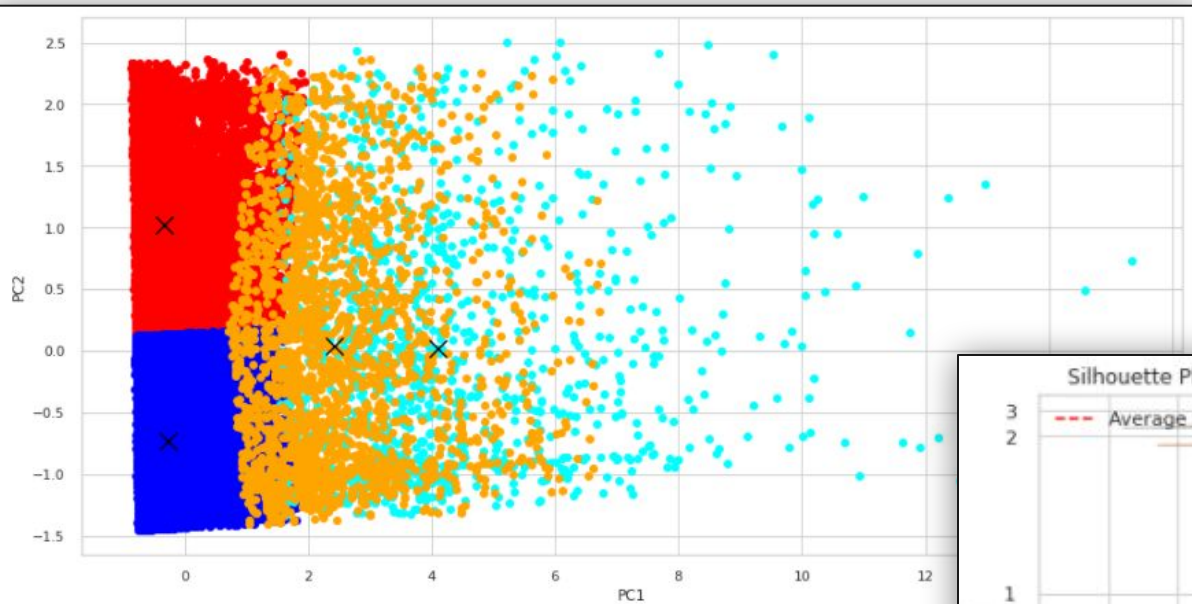


K-Means sur 3 variables : RFM

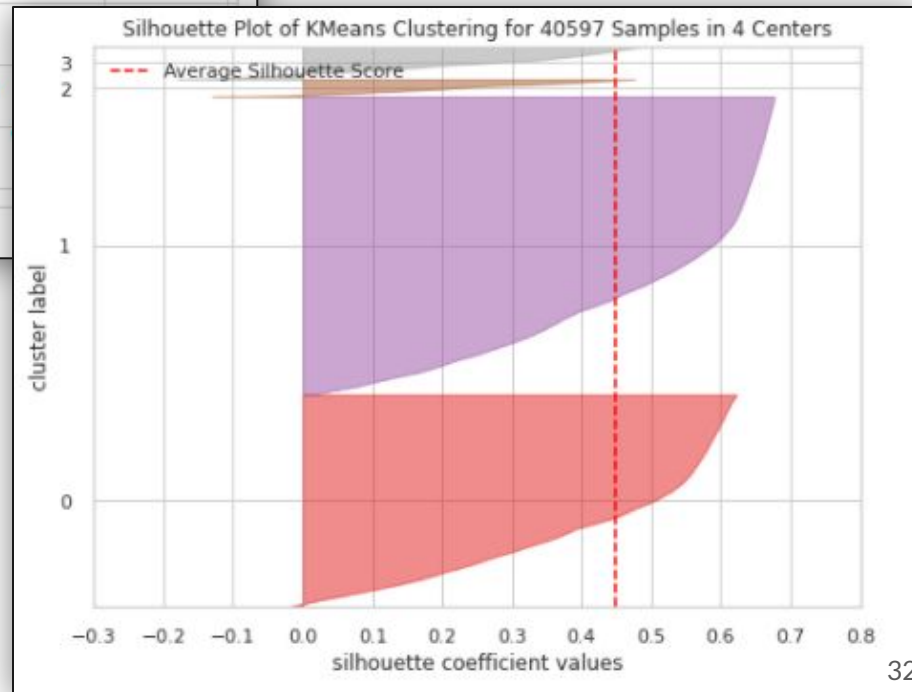
Méthode du coude

Score Silhouette

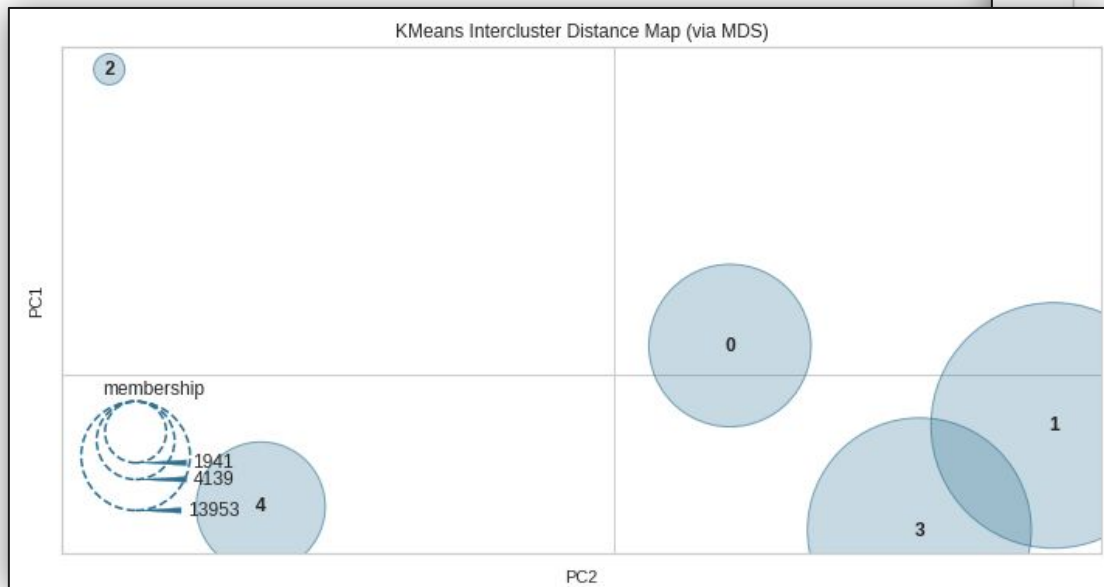
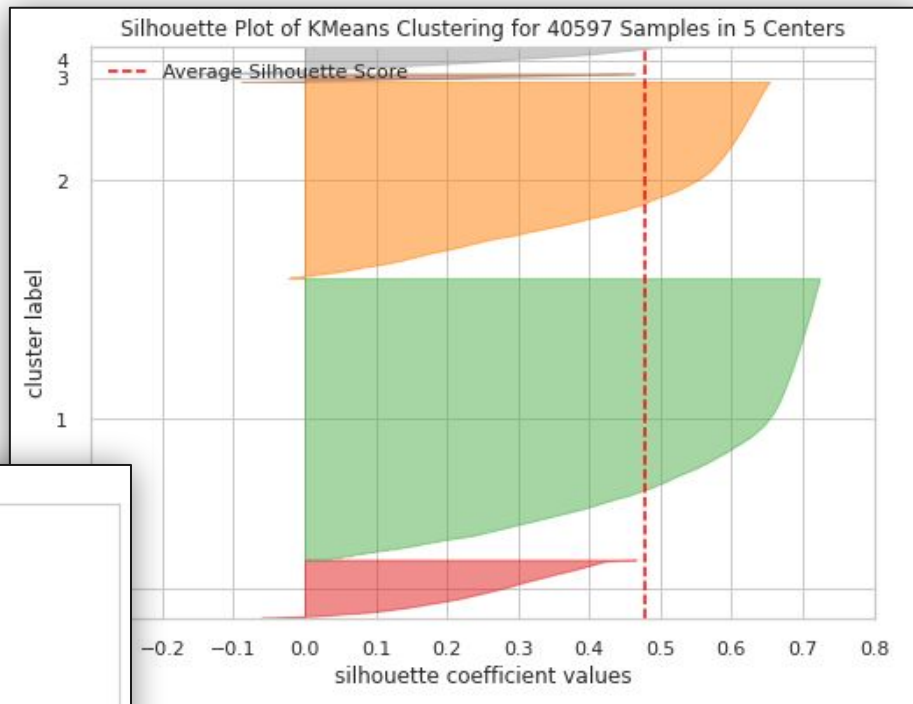




Choix optimal de clusters : 4



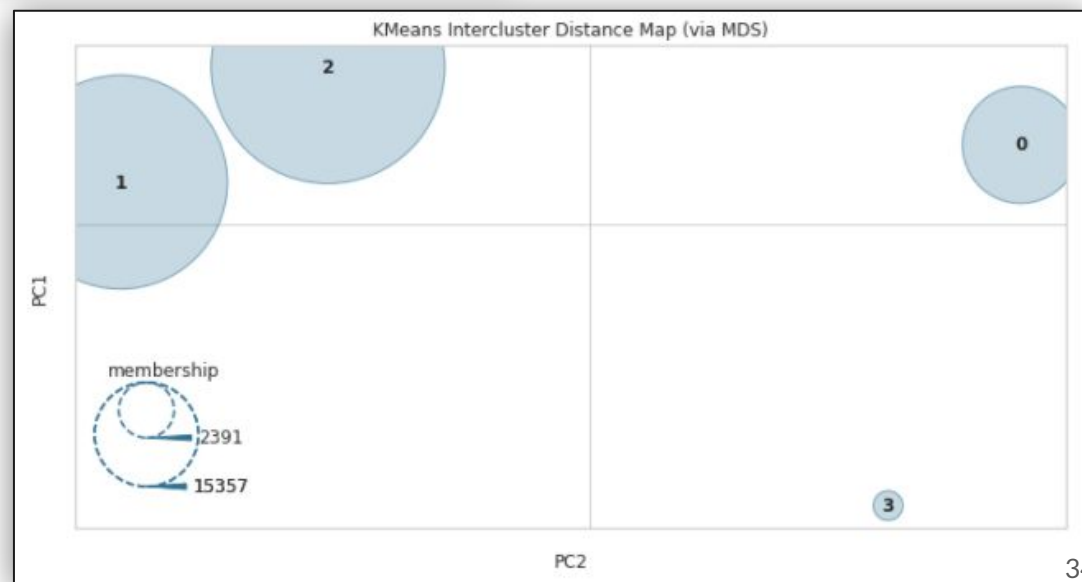
Essai avec 5 clusters

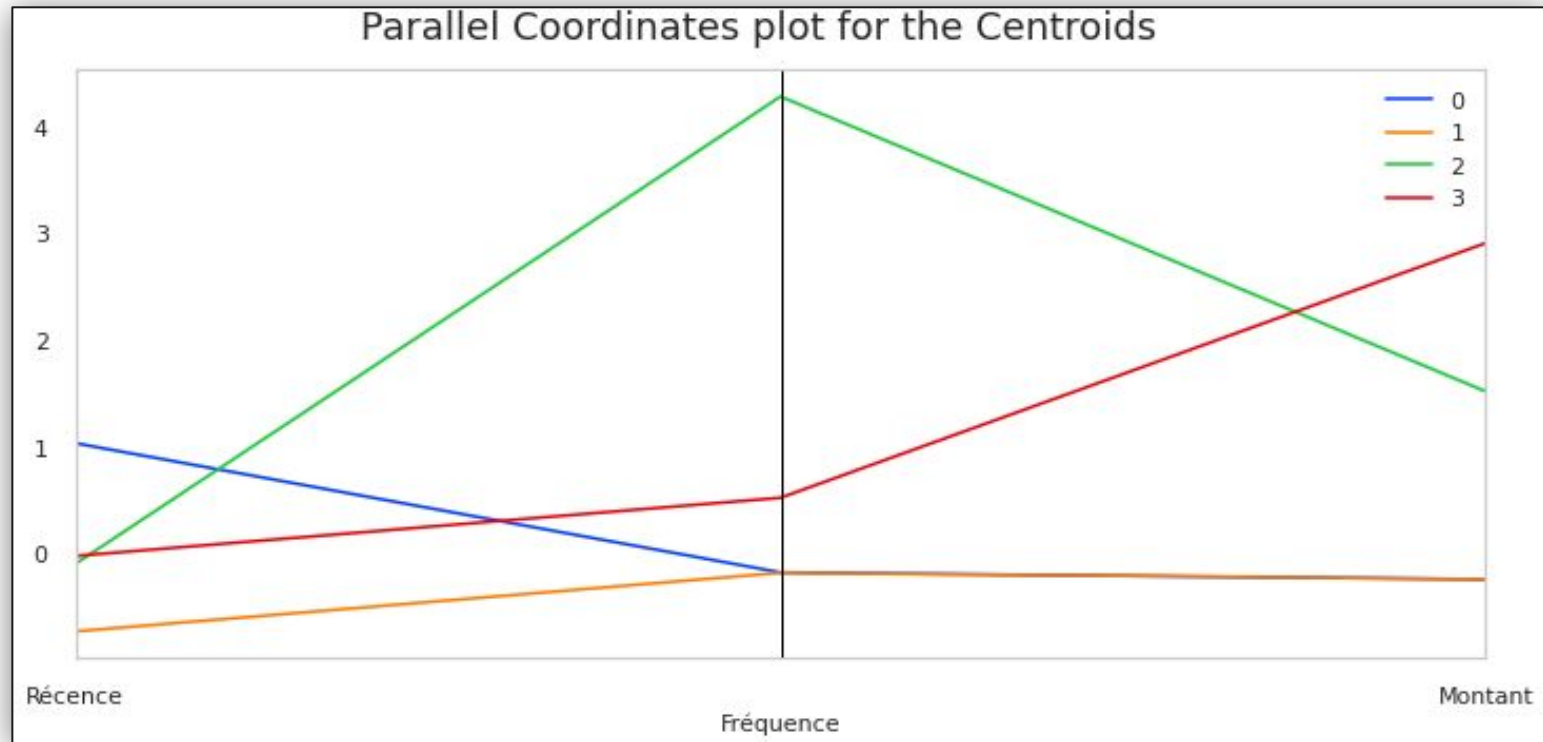


Cluster	Récence	Fréquence	Montant	Segment	
	mean	mean	mean	unique	count
0	234.000	1.000	127.000	[Hibernating, About to sleep, At risk]	15357
1	67.000	1.000	127.000	[Promising, About to sleep, New customers, Pot...]	21612
2	127.000	4.000	473.000	[At risk, Potential loyalists, Loyal customers...]	1237
3	134.000	2.000	746.000	[About to sleep, Potential loyalists, Hibernat...]	2391

Critères :

- taille des silhouettes
- volumétrie par cluster
- présence d'erreurs





- **Cluster 0** : les clients dont la Récence est élevée (le dernier achat date d'il y a longtemps) et qui ont dépensé peu d'argent
- **Cluster 1** : les clients dont la Récence est basse (le dernier achat est récent) et qui ont dépensé peu d'argent
- **Cluster 2** : les clients qui achètent souvent et dont le panier moyen est important
- **Cluster 3** : les clients qui n'achètent pas souvent et dont le panier moyen est très élevé

Etude de stabilité

Objectifs et méthodologie



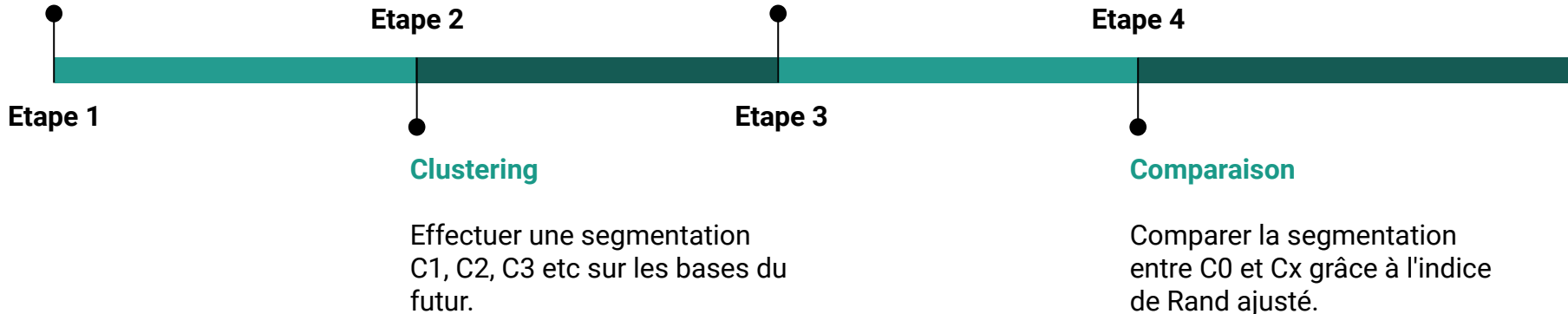
Tester la stabilité de l'algorithme **dans le temps** pour savoir à quel moment **les clients changent de cluster**.

Jeux de données du futur

Créer plusieurs datasets : B0 et les B1, B2, B3 etc du futur.

Prédictions

Générer des labels B1, B2, B3 etc.





Indice de Rand (ajusté)

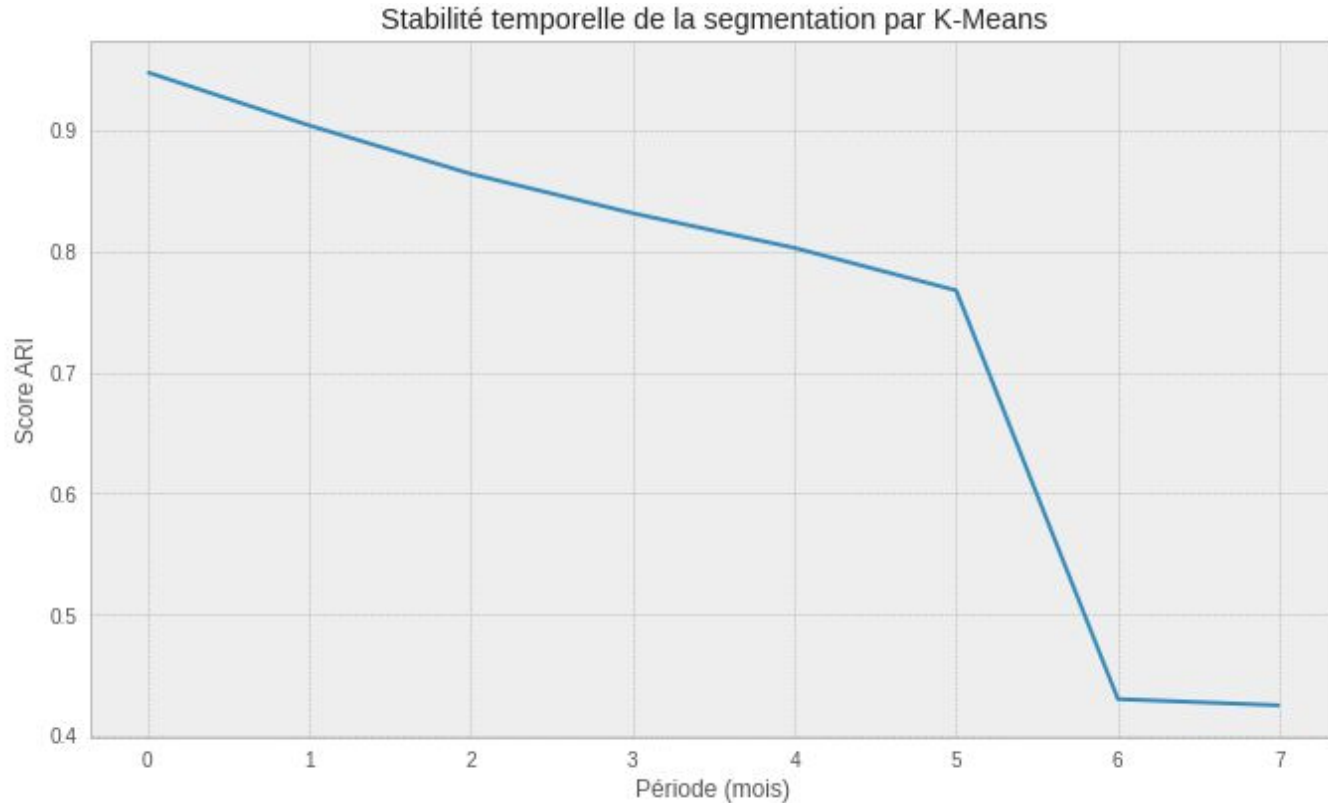
Indice de Rand :

utilisé pour mesurer la **similarité des points de données** présents dans les clusters.

ARI :

- proche de 0 pour un clustering aléatoire
- égal à 1 uniquement quand le clustering correspond exactement à la partition initiale

Calcul du score ARI



Conclusion



- Parmi toutes les méthodes de clustering appliquées, **c'est le K-Means qui semble être le plus efficace.**
- Le CAH donne des résultats similaires, en revanche **les temps de performances laissent à désirer.**
- Les résultats de K-Means sont parfaitement **interprétables et utilisables** par l'équipe marketing.
- La segmentation RFM semble fournir **une analyse des profils client tout aussi pertinente** (note : rétrécir la liste des comportements d'achat en fusionnant certains profils - ceux similaires ou ceux les moins nombreux).

Avez-vous des questions ?
