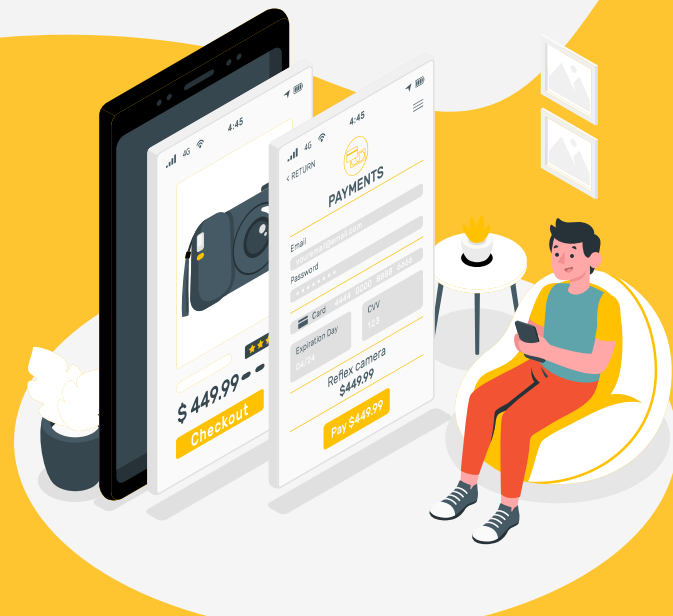


PROJET 5:

# Segmentation des clients d'un site e-commerce



Etudiant : Fatma Aidi  
Mentor : Kezhan Shi

Evaluateur : Mohammed Sedki  
Date : 15/04/2021

# TABLE

**01**

## **CONTEXT**

- présenter l'entreprise
- la mission

**04**

## **MODEL OPTICS**

- nombre de clusters
- stabilité
- interprétation

**02**

## **PREPARATION DES DONNEES**

- les bases de données
  - nettoyage
- analyse exploratoire

**05**

## **MODEL:CAH(Classification Ascendante Hiérarchique)**

- nombre de clusters
- stabilité
- interprétation

**03**

## **MODEL K-MEANS**

- nombre de clusters
- stabilité
- interprétation

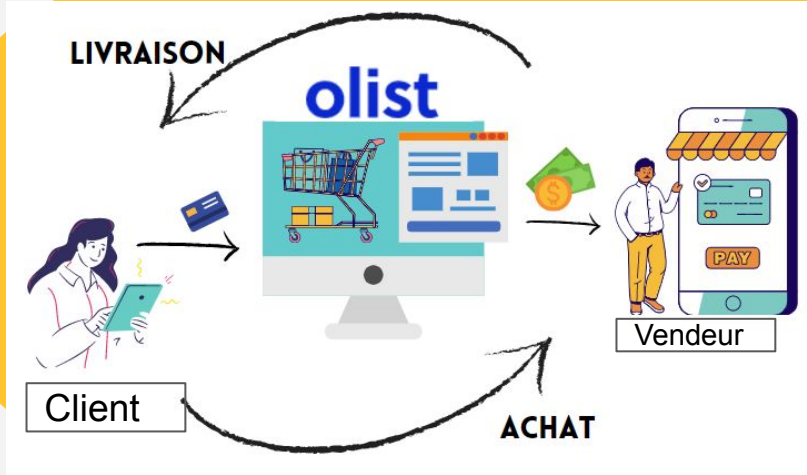
**06**

## **CONCLUSIONS**



olist

# -- L'entreprise



- **Market-place** en ligne au Brésil

- **Mission:** Segmentation des clients pour mener des actions marketing ciblées dans l'objectif de :

- Ramener plus de clients
- Fidéliser les clients
- Booster les ventes

- **Cahier des charges :**

- Résultat actionnable d'une segmentation,
- Proposition de contrat de maintenance basée sur une analyse de la stabilité des segments au cours du temps.





## 2-Préparation des données et analyse



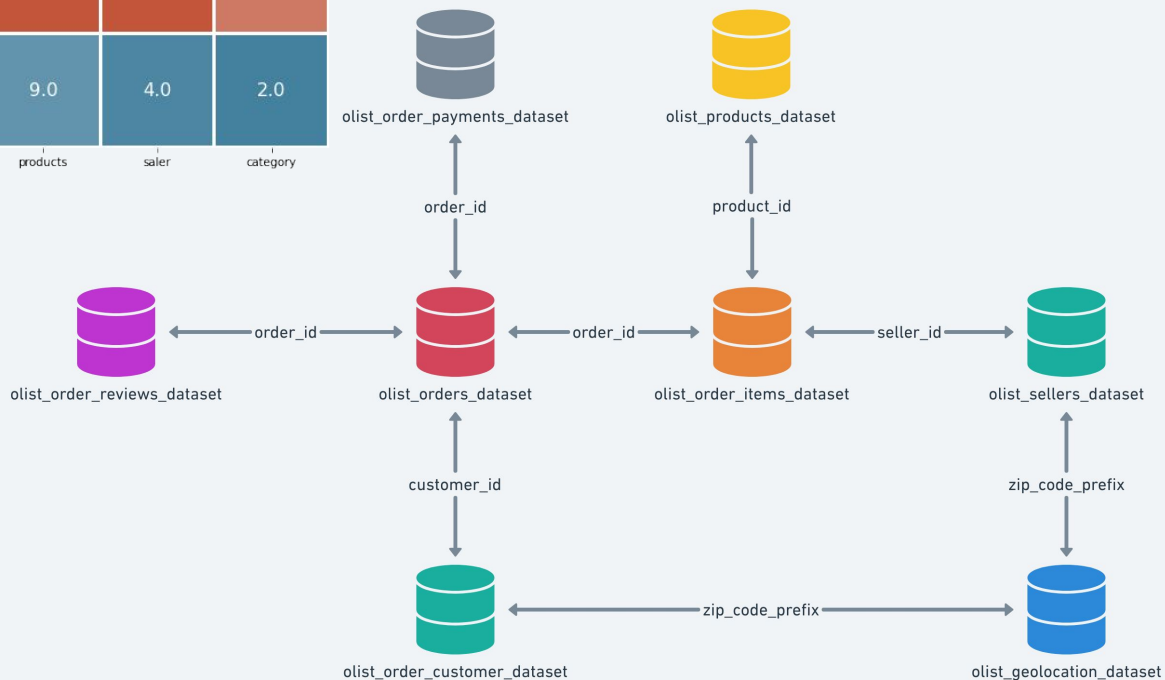
# Jeux de données

variables	customers	geolocation	items	payments	reviews	orders	products	saler	category
observations	99441.0	1000163.0	112650.0	103886.0	100000.0	99441.0	32951.0	3095.0	71.0
Data size	5.0	5.0	7.0	5.0	7.0	8.0	9.0	4.0	2.0

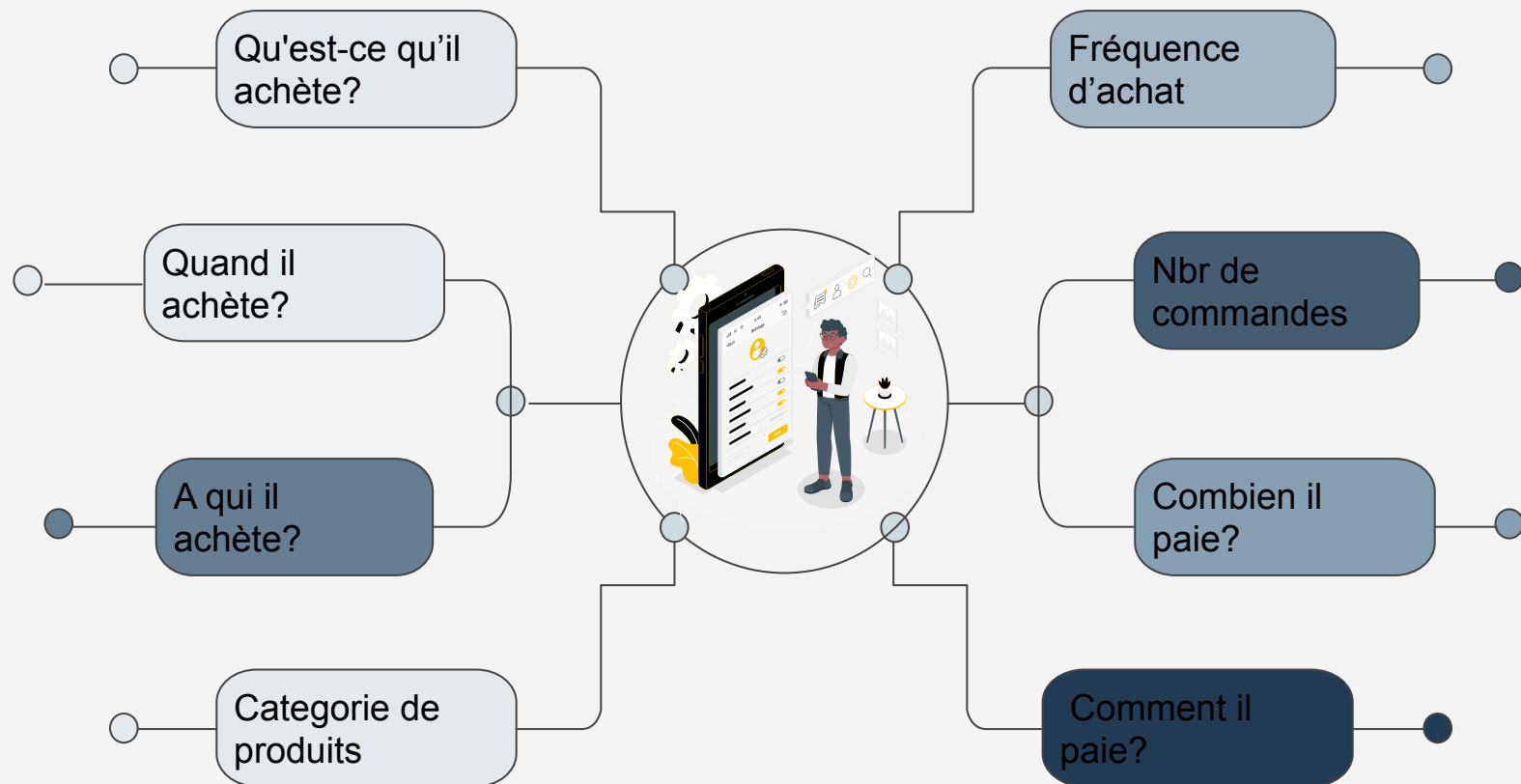
-8 bases de données

-les **tailles** varient entre :

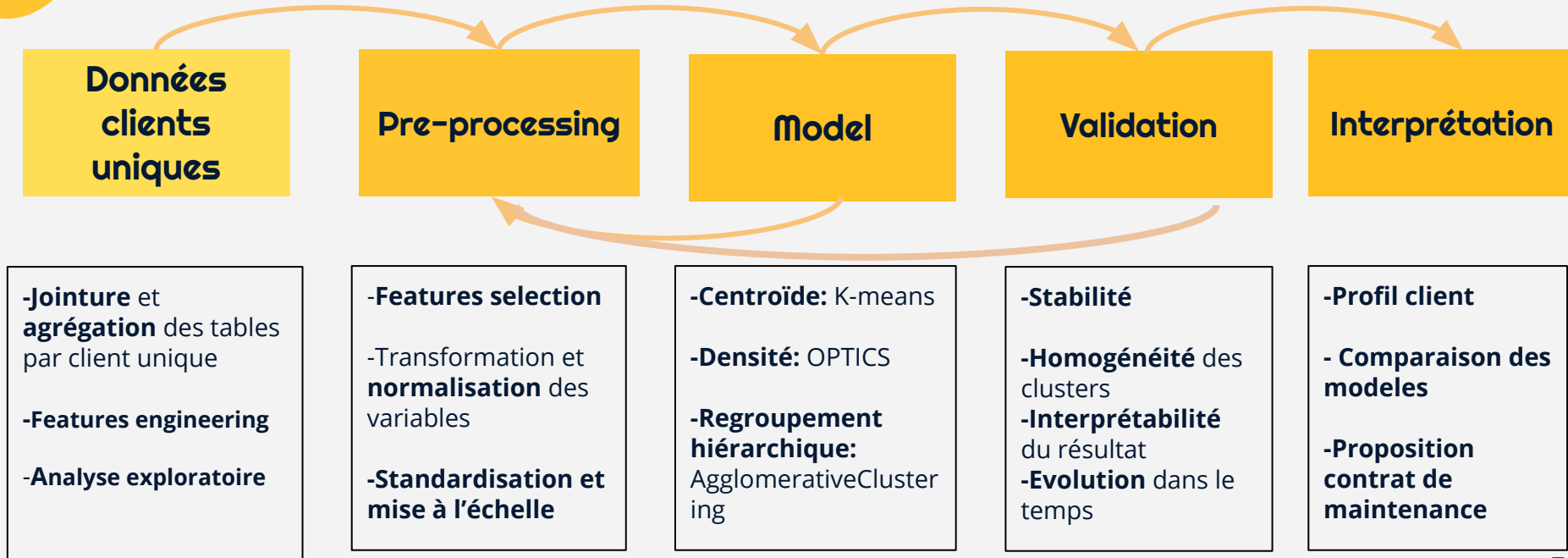
- 2 et 9 **variables**
- 71 et 112650 **observations**



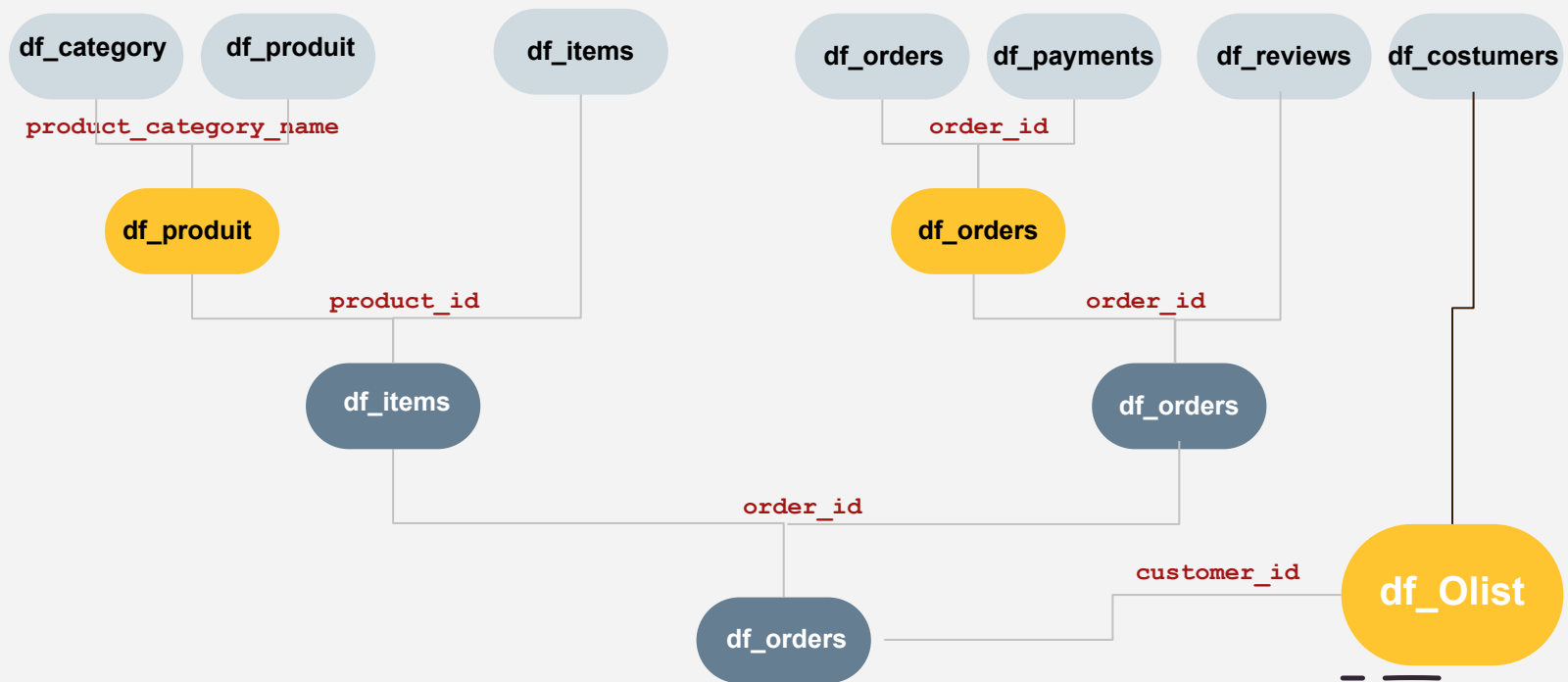
# Profil client recherché



# Démarche



# Jointure et agrégation des tables





# Features engineering

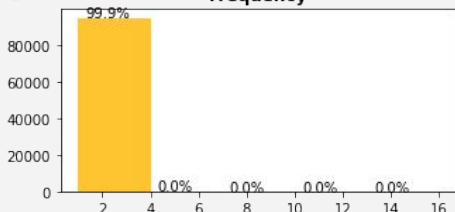
RFM	COMMANDE	PAIEMENT	PRODUIT	SATISFACTION
Recency	Product_nbr (nombre produit/ commande)	payment_sequential	density_product	review_len_message
Monetary	Nb_item (nombre d'article acheté)	payment_type_credit_card	product_weight_g	review_score_mean
Frequency	order_probl	payment_installments	Volume_product_mean	<b>Vendeur</b>
Frequency_age	price	nbr_payment_type	catg	seller_Nbr
age	finish_order			

# Analyse Exploratoire

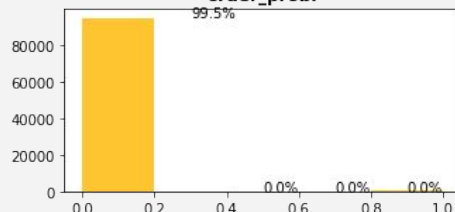
- Distributions positivement biaisées (non normales)
- 94%** ont commandé 1 fois
- 90%** ont acheté un seul article
- 70%** sont satisfaits
- 97%** ont passé une commande chez un seul vendeur

Dispersion of quantitative variables

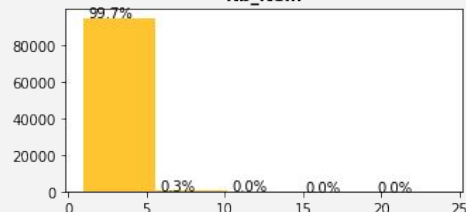
Frequency



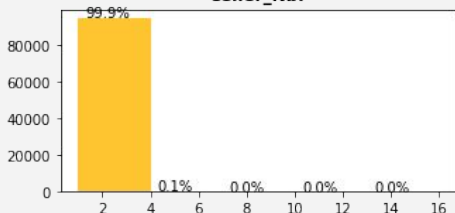
order\_probl



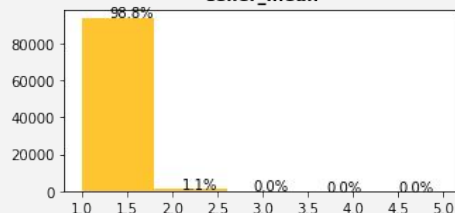
Nb\_item



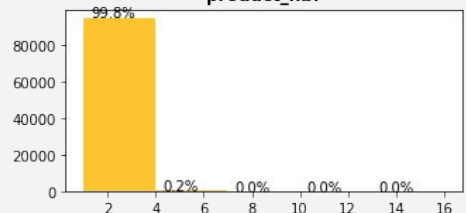
seller\_Nbr



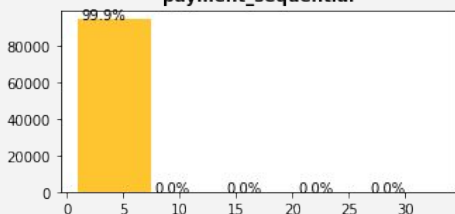
seller\_mean



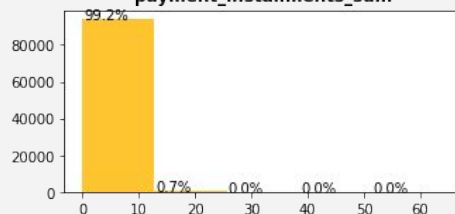
product\_nbr



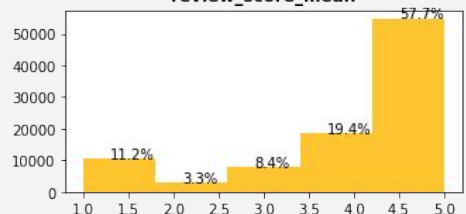
payment\_sequential



payment\_installments\_sum



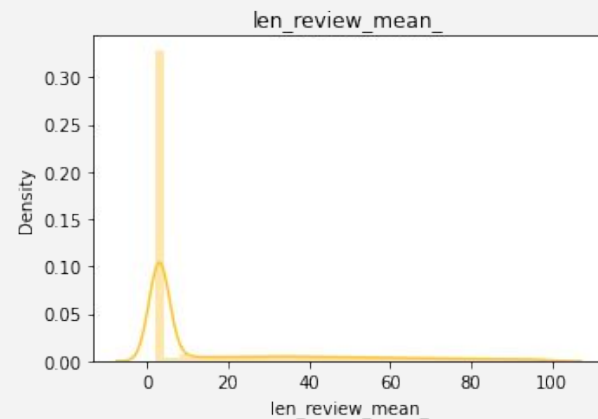
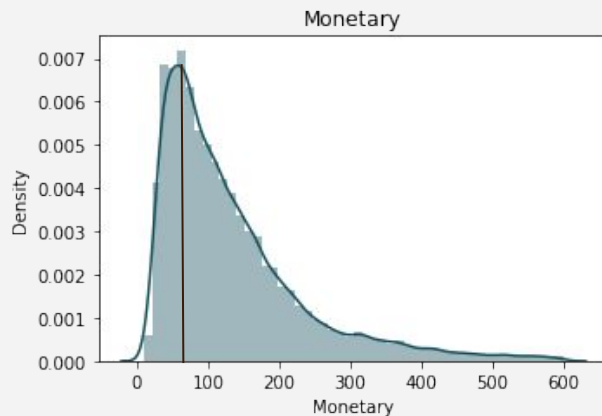
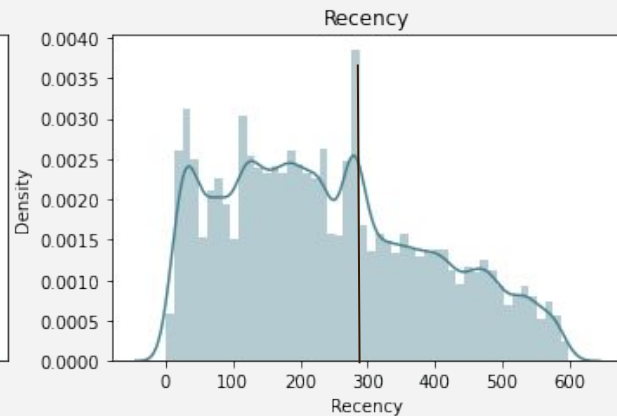
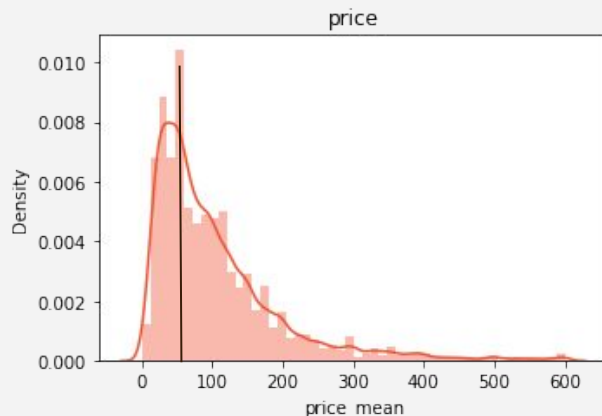
review\_score\_mean



# Analyse Exploratoire

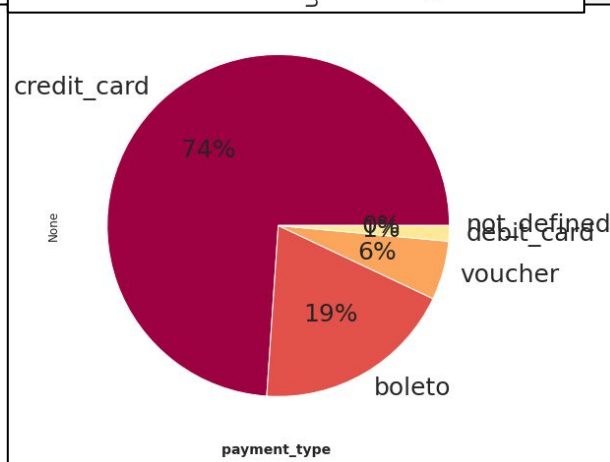
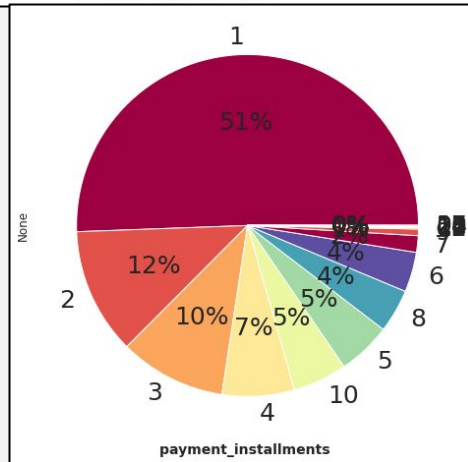
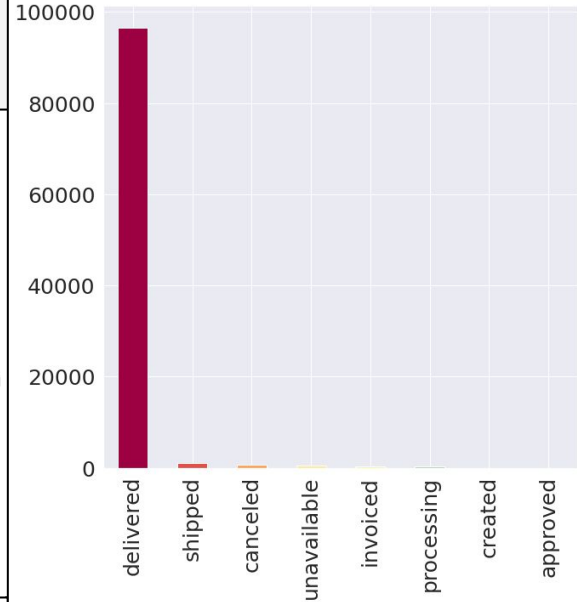
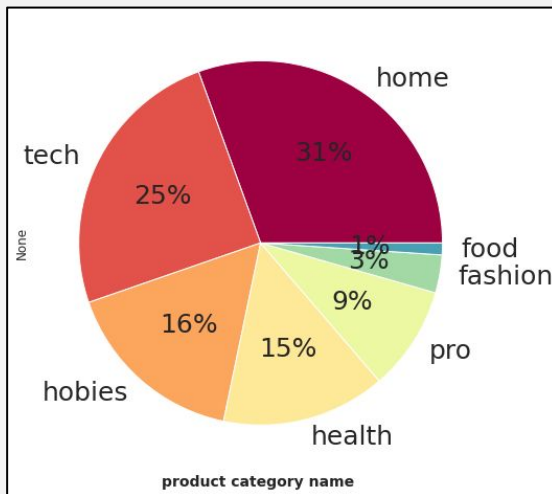
Distributions **positivement**  
**biaisées** avec **pic**

- **60** Réals prix moyen des articles achetés par client.
- **70** Réals dépense par client.
- **5** mots par commentaire
- **300** jours depuis la dernière commande



# Analyse Exploratoire

- Paiement par **carte bancaire** majoritaire (plus de **74 %**)
- Les **Catégories** maison et électronique représentent plus **50%** des articles vendus
- **95%** des commandes sont **sans problèmes de livraison**
- **50%** des clients paient en une seule fois.



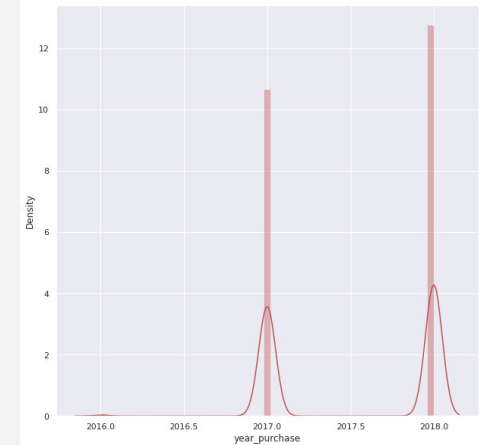
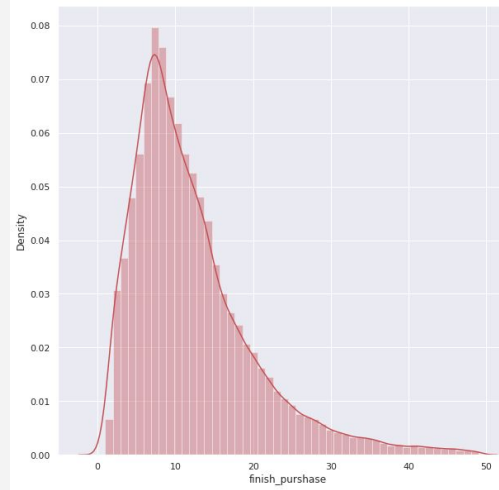
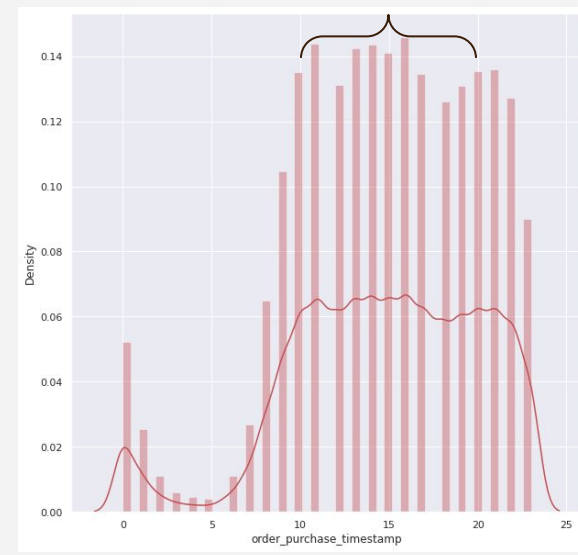
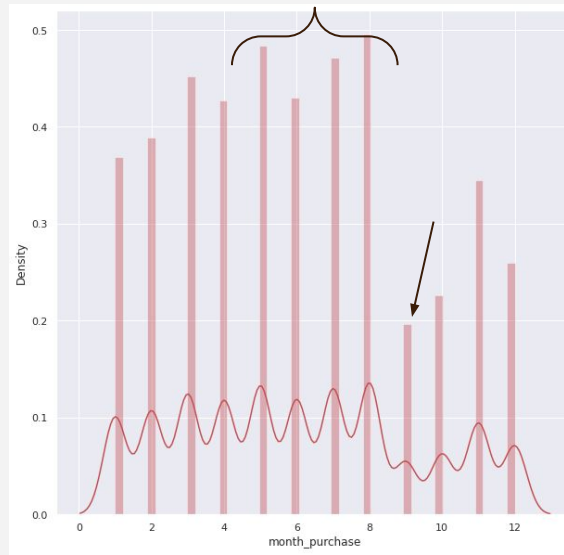
# Analyse Exploratoire

-**Opération d'achat** dure en moyenne 10 jrs

-Augmentation de nombre de commandes en **2018**

-Les meilleurs **mois** sont de Mai à Août

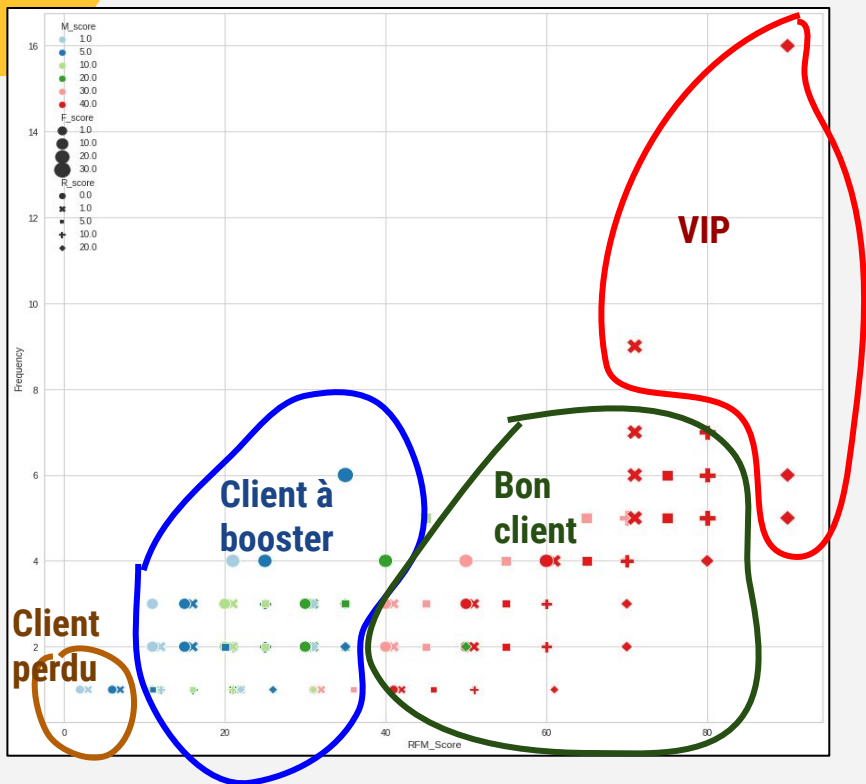
-La majorité des ventes se font entre 12h et 19h



## Récapitulatif

- **96% des commandes contiennent un seul article**
- **2%** des clients ont acheté chez plusieurs vendeurs
- **94%** de nos clients ont acheté sur le site une **seule fois**.
- **95%** ont payé avec un **seule mode de paiement**
- **50%** ont choisi un paiement en **plusieur fois**.

# Segmentation RFM



## RÈGLE pour RFM SCORE

### Recency/JR

0-60 : 20pt  
 60-120 : 10pt  
 120-180 : 5pt  
 180-240 : 1pt  
 >=240 : 0pt

### FREQUENCE

1 : 1pt  
 2-3 : 10pt  
 4 : 20pt  
 >=5 : 30pt

### MONTANT

0-100 : 1pt  
 100-200 : 5pt  
 200-300 : 10pt  
 300-400 : 20pt  
 400-500 : 30 pt  
 >=500 : 40 pt

- Client **Haut\_VIP** RFMS>80
- Client **bon** 40<RFMS<80
- Client à **booster** 10<RFMS<40
- Client **perdu** RFMS<10

Pour la suite on isole les **VIP** que l'équipe marketing traite à part. Ils sont **fidèles**, achètent au moins **5 fois** et sont dépensiers: **panier moyen élevé**

RFM signifie Récence, Fréquence et Valeur monétaire, chacun correspondant à une caractéristique d'un client. Ces mesures sont des indicateurs importants du comportement d'un client:

- **Récence**: à quand remonte la dernière fois qu'ils ont acheté?
- **Fréquence**: à quelle fréquence et pendant combien de temps ont-ils acheté?
- **Valeur monétaire / ventes**: combien ont-ils acheté

# 3-MODEL NON SUPERVISE: K-MEANS



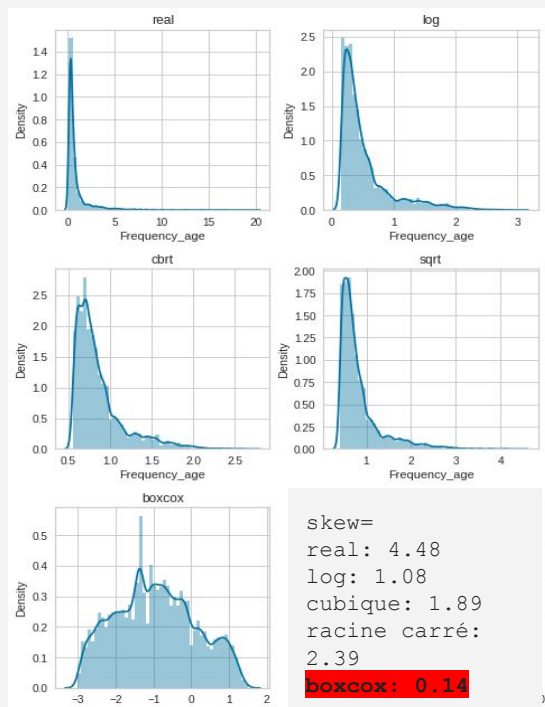


# K-MEANS : Pre-processing

**-Choix K-means:** C'est l'un des algorithmes de clustering les plus répandus. Il s'appuie sur la mesure de distance pour créer les clusters en minimisant la somme des carrés des distance entre un point et la moyenne des points de son cluster (centroïde).

**-PowerTransformer** pour stabiliser la variance et rendre les distributions plus normalisées. le résultat de calcul d'asymétrie des données est presque nul.

**-Transformation et Mise à l'échelle**  
:MinMaxScaler ,StandardScaler,  
QuantileTransformer



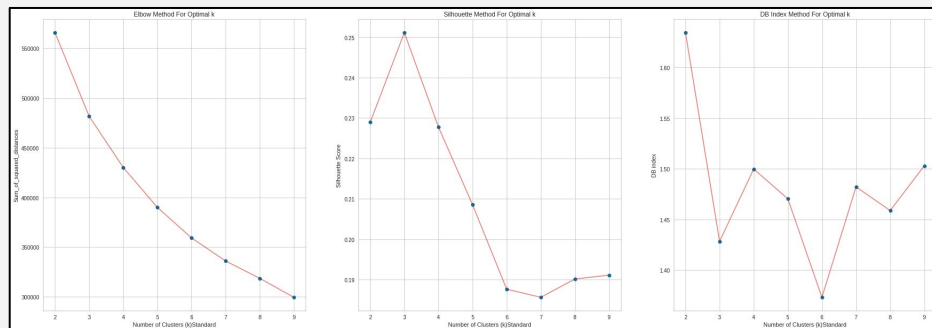
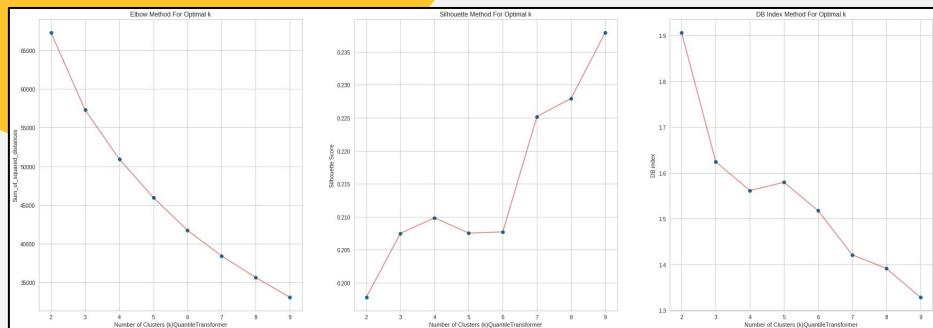
## Features selections: 8

<b>Recency</b>	La période depuis la dernière commande
<b>Frequency_age</b>	Taux d'achat par jour
<b>Monetary</b>	Le montant d'achat
<b>age</b>	La période depuis la première commande
<b>payment_installments_sum</b>	nombre de tranches de paiement
<b>review_score_mean</b>	Le score moyen donné par le client
<b>len_review_mean</b>	La longueur moyenne du commentaire
<b>Nb_item</b>	Nombre d'article acheté

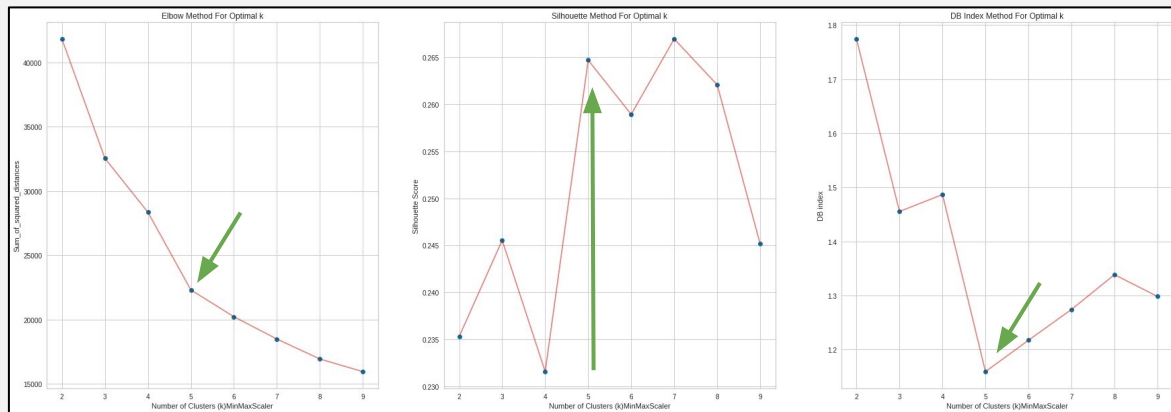
# K-MEANS :Nombre de clusters

QuantileTransformer nbr:7,9,5

StandardScaler:6,3,5



MinMaxScaler nbr:5



-**Indice Silhouette** (à maximiser): mesure la **cohésion** et la **séparation**: calcul de la différence entre la distance «intra-classe» et la distance au centroïde le plus proche d'un autre cluster.

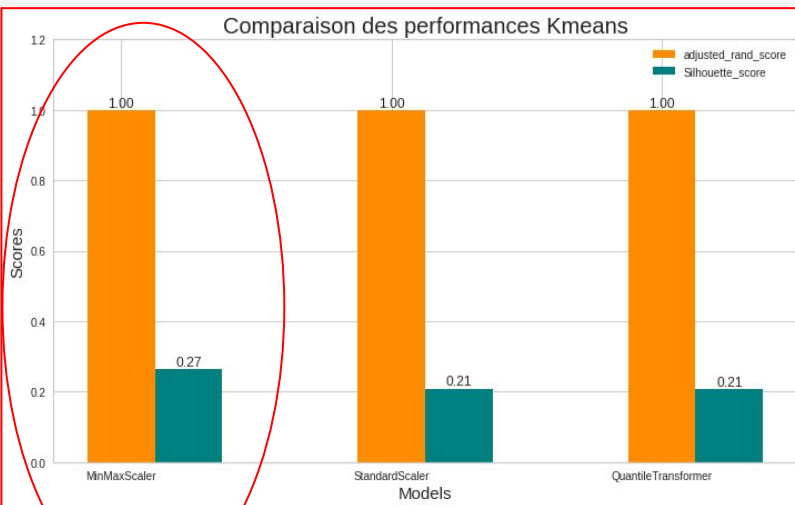
-**Distortion (coude)**: mesure la **variance** «intra-classe»: Moyenne des sommes des distances quadratiques des points au centroïde le plus proche.

-**Davies\_Bouldin** (à minimiser) il mesure la **l'homogénéité** et la **séparation**: est basé sur le rapport entre les distances «intra-classe» et «inter-classe»

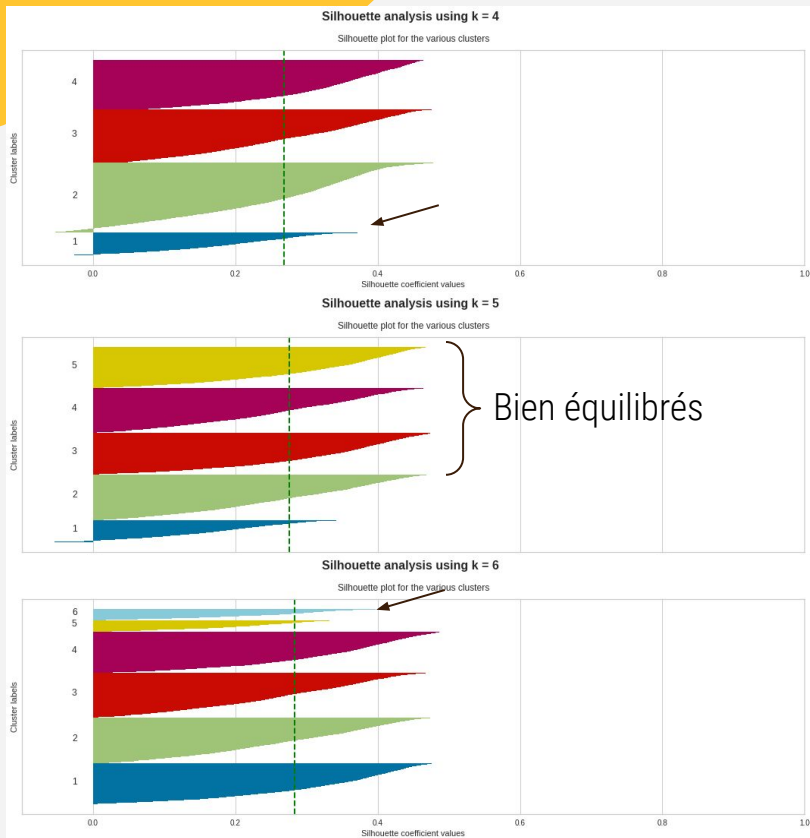
# Stabilité des clusters

- **ARI: adjusted rand score:** est une mesure de similarité entre deux classifications, un indice qui utilise les étiquettes externes (réelles). Dans notre cas, on l'utilise pour voir si, à l'initialisation, les clients changent de cluster ou non.
- les 3 modèles sont stables mais on retient le **MinMaxScaler** car il donne un meilleur score de silhouette à 0.27.

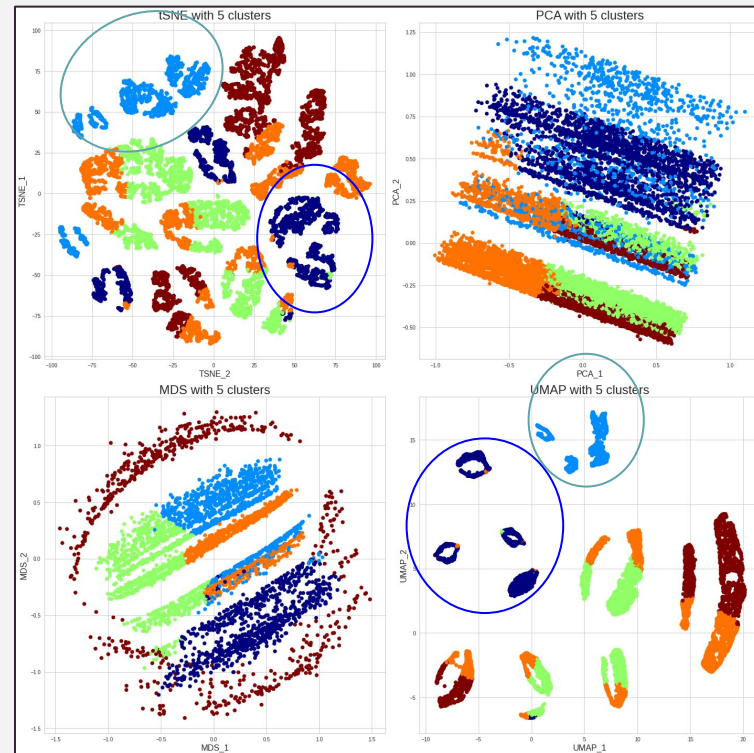
```
1 def score_ari(df,model, nb_itr=6):  
2     # Fitting the model  
3     model.fit(df)  
4     labels_true=model.labels_  
5     # Calculate the ARI scores  
6     ARI_scores = []  
7     # Iterating  
8     for i in range(0,nb_itr):  
9         # Fitting the model  
10        model.fit(df)  
11        labels_predict= model.labels_  
12        # Compute the ARI score with labels_true  
13        ARI_score= adjusted_rand_score(labels_predict,labels_true)  
14        ARI_scores.append(ARI_score)  
15    #Compute the mean of ARI scores  
16    ARI_mean = statistics.mean(ARI_scores)  
17    return ARI_mean
```



# K-MEANS :Visualisation



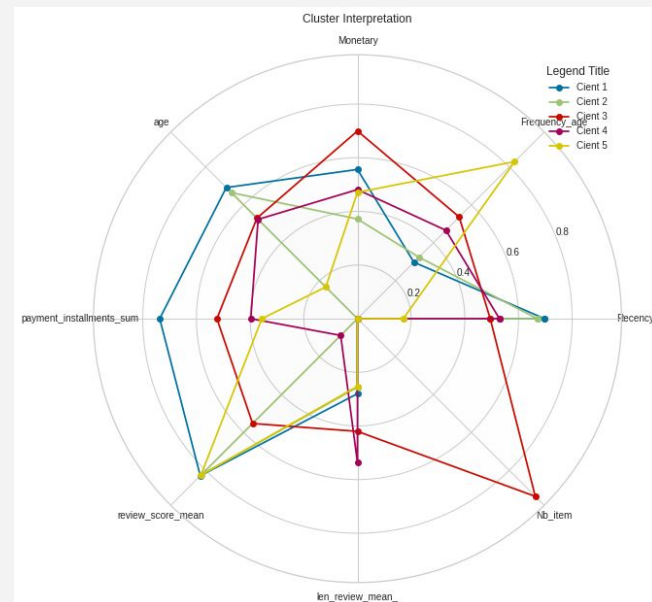
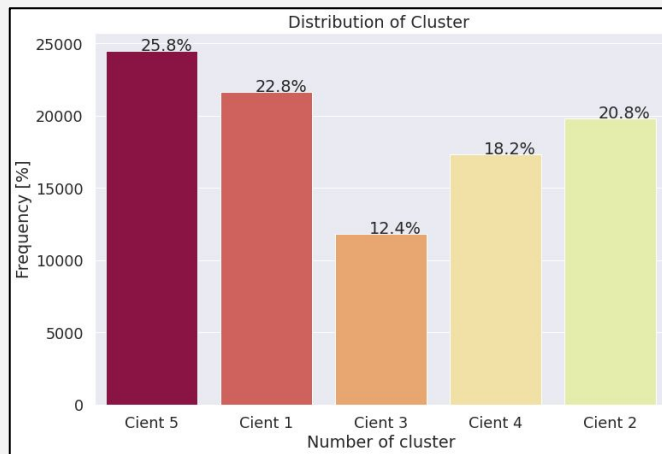
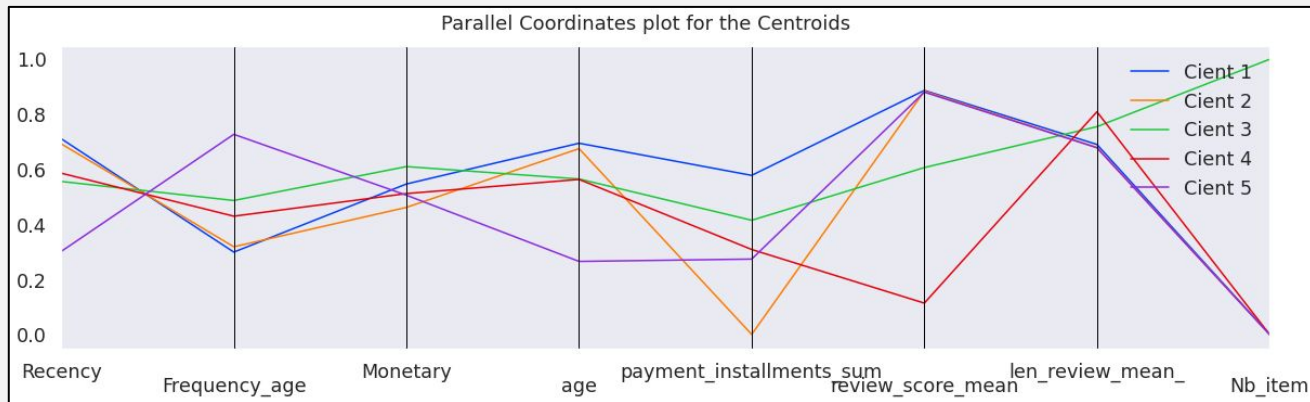
- La projection 2D de tSNE , MDS ,PCA et UMAP montrent que au moins 2 clusters sont bien regroupés
- 5 est le nombre de cluster optimal



# INTERPRETATION

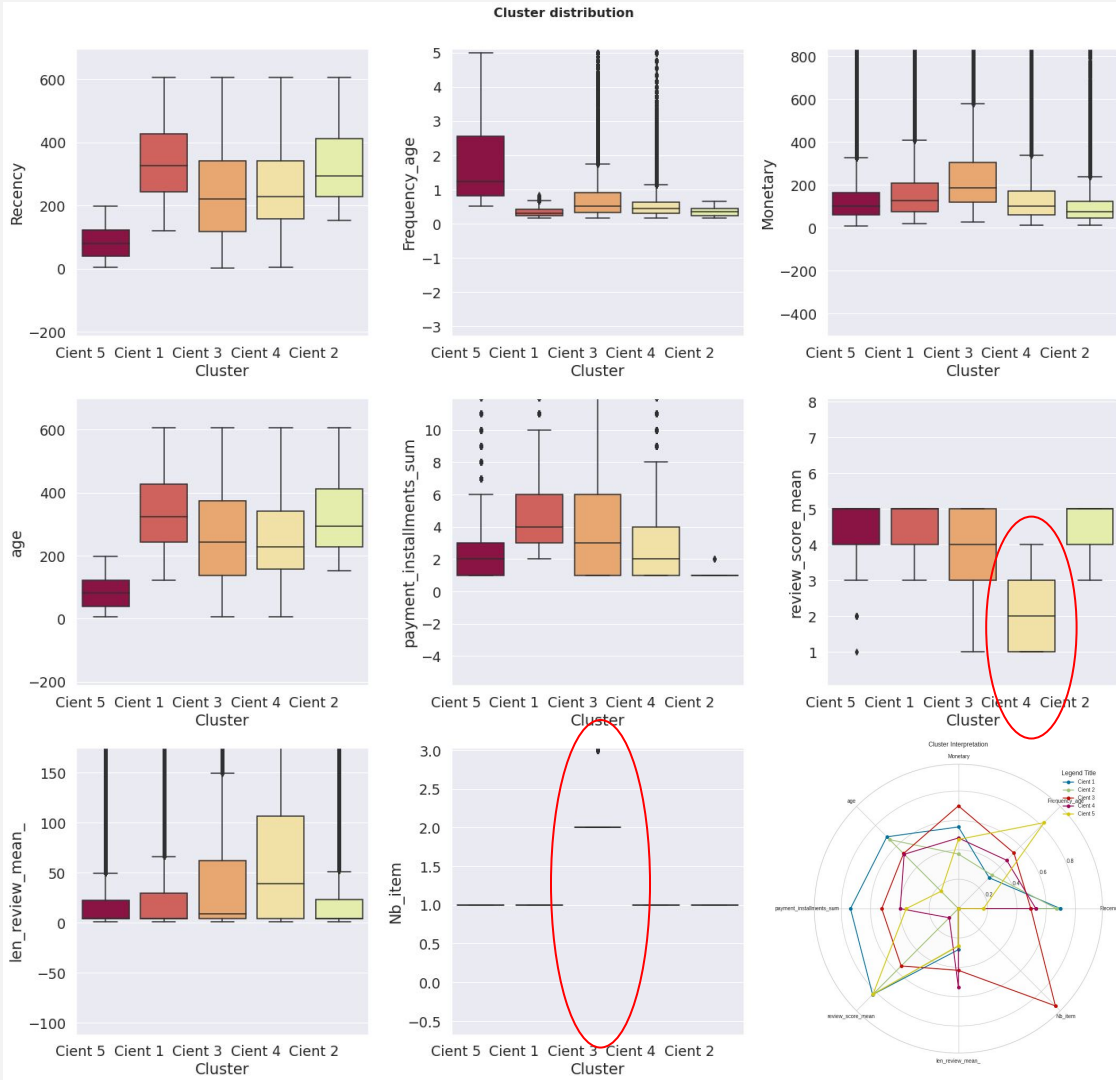
## - Clusters sont :

- Homogènes
- Équilibrés
- Lisibles
- Interprétables



# INTERPRETATION

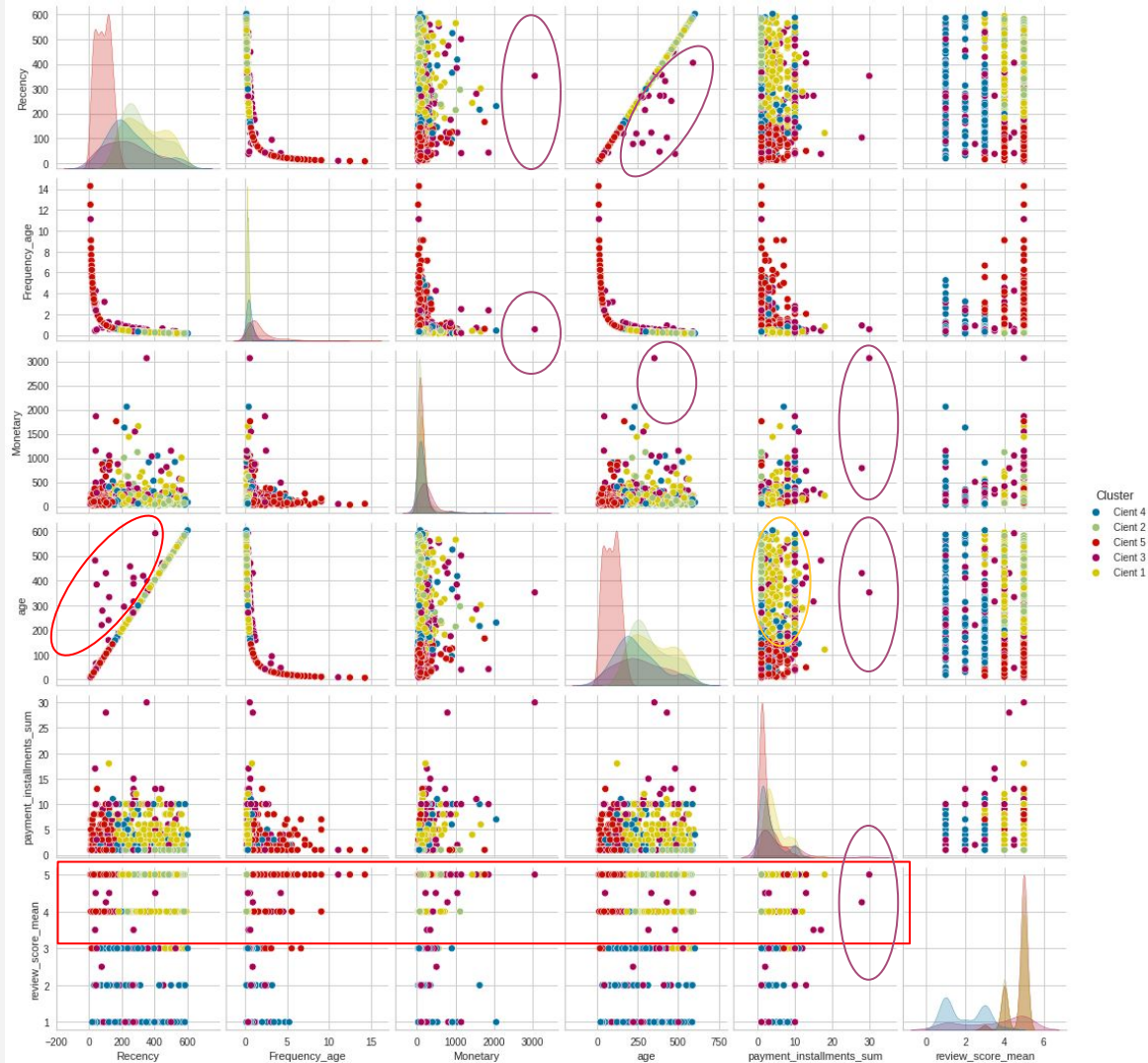
- **Client 2, 1 et 5** sont satisfaits
- **Client 5** est nouveau.
- **Client 2** paye en seule fois.
- **Client 4 et 3** sont les moins satisfaits.
- **Client 3** est le plus dépensier et qui a acheté au moins 2 articles.



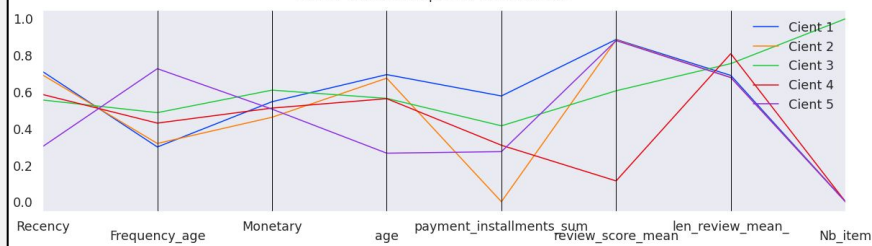


# INTERPRETATION

- **Recency /age** on arrive à séparer les clients qui ont acheté plusieurs fois (**6%** du nombre total) quand  $\text{age} = \text{Recency}$ .
- Plus Monetary est grand plus le paiement échelonné.
- **Client 3** représente le client outlier .
- **Client 1** est l'ancien ou il n'a pas acheté depuis longtemps.



Parallel Coordinates plot for the Centroids



## PROFIL CLIENT K-méans

23%



### Ancien et aléatoire

Les clients **anciens** qui n'ont pas fait d'achat **depuis longtemps**, ils ont fait des achats pour des montants **moyens (200 Real)**.

21%



### Ancien économe

Les clients **anciens** qui n'ont pas fait d'achats **depuis longtemps**, font des achats pour des montants **faibles** et il paient en **une seule fois**.

13%



### Dépensier

Les clients plutôt **anciens** :qui achètent **différents** articles, pour des montants **élevés** (>500 real)et ont besoin d'un paiement **échelonné** (>5 fois).

18%



### Insatisfait

Les clients qui ont fait des achats **récemment**, ils achètent **rarement**, font des achats pour des montants **faibles** à **moyens**. Ils sont **mécontents** (ils laissent des commentaires assez longs)

26%



### Nouveau et satisfait

Les clients **nouveaux** : font des achats **récemment** pour des montants **faibles** à **moyens** et ils sont satisfaits de leur achat.



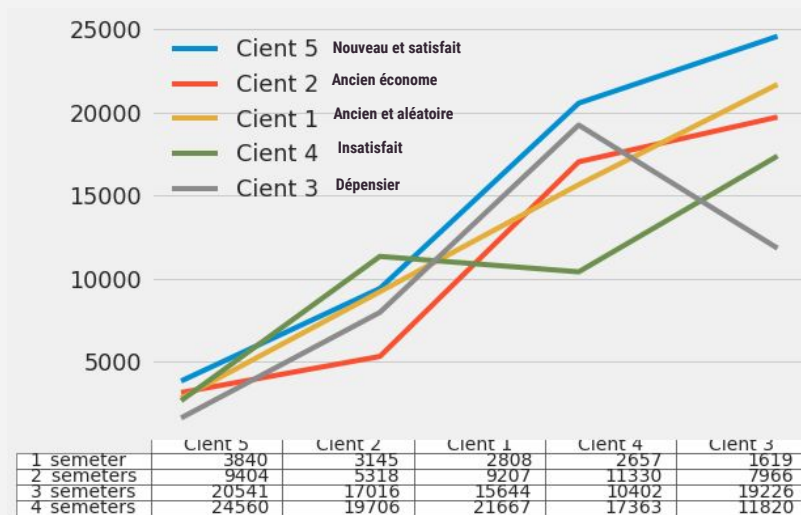
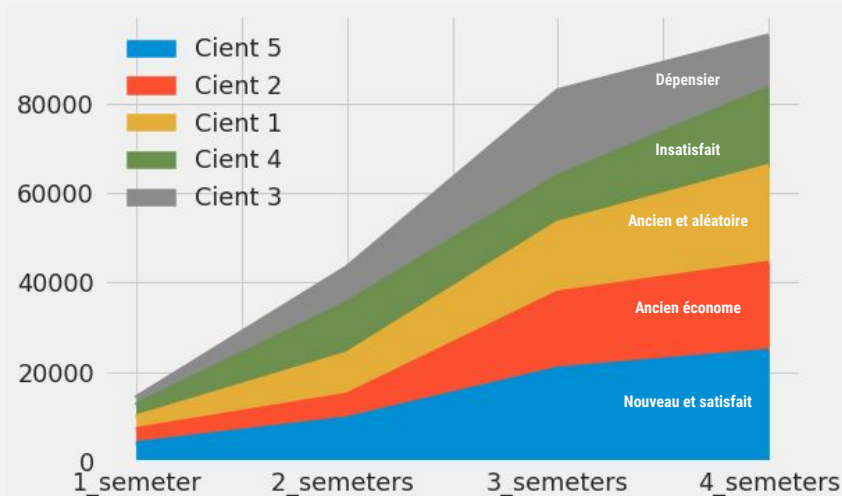
# Stabilité des clusters dans le temps

- Diviser la data en 4 temps et faire la classification à chaque fois et recalculer les noms clusters (donner aléatoirement par K-means).

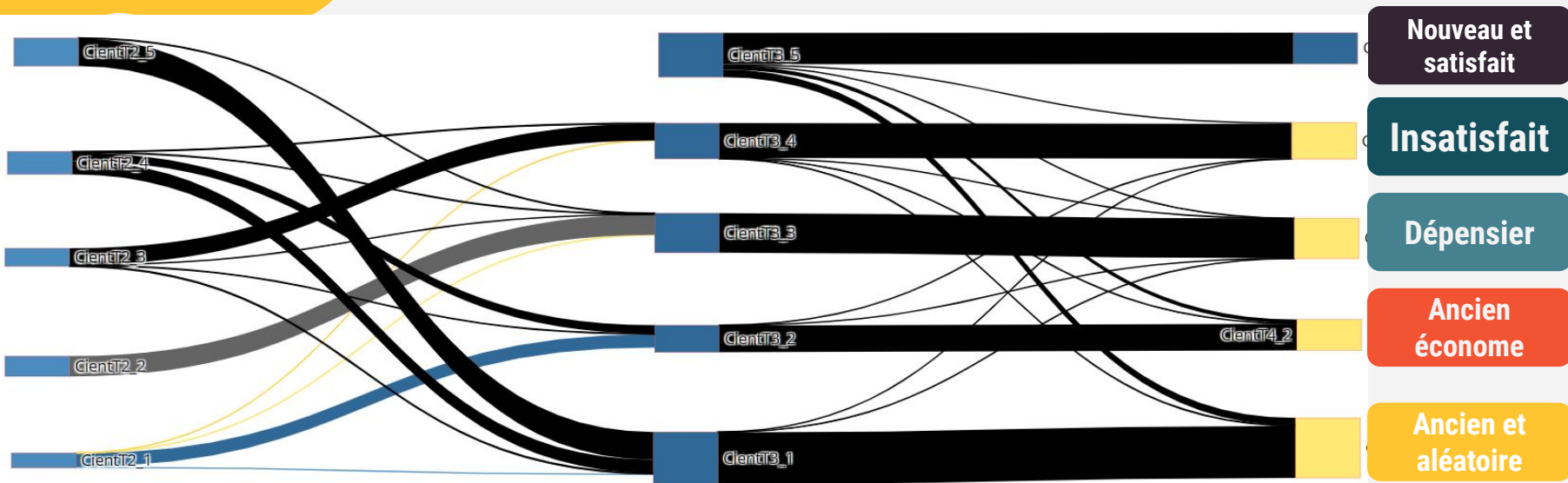
**-Stabilité** dans le temps du **silhouette\_score**

-Baisse du **nombre de clients dépensiers** (client 3) sur le 4ème semestre.

- Hausse du **nombre de nouveaux clients** (client 5) sur le 4ème semestre.



# Contrat de maintenance



-On fait la classification sur chaque base de données de chaque période (condition du futur contrat de maintenance)

-Diagramme de **Sankly**:

- Les clients **migrent** d'un groupe à un autre chaque semestre.
- un **contrat de maintenance** est nécessaire tous les **6 mois**

# 4- OPTICS: Ordering Points To Identify the Clustering Structure



# OPTICS : Pre-processing

-**Choix OPTICS**: s'appuie sur la densité estimée des clusters pour faire la classification. On travaille sur un échantillon de 30 000 clients.

-**PowerTransformer** pour normaliser des distributions.

-**Transformation et Mise à l'échelle**:MinMaxScaler ,StandardScaler, QuantileTransformer.

## Features selections: 8

<b>price_mean</b>	Le prix moyen des produits
<b>order_item_mean</b>	Nombre moyen des produit par commande
<b>product_nbr</b>	Le nombre de différent produit par commande
<b>seller_Nbr</b>	La période depuis la première commande
<b>Frequency_age</b>	Taux d'achat par jour
<b>payment_installments_sum</b>	nombre de tranches de paiement
<b>age</b>	La période depuis la première commande
<b>item_Monetary</b>	Panier moyen par client(montant/nbr de commande)

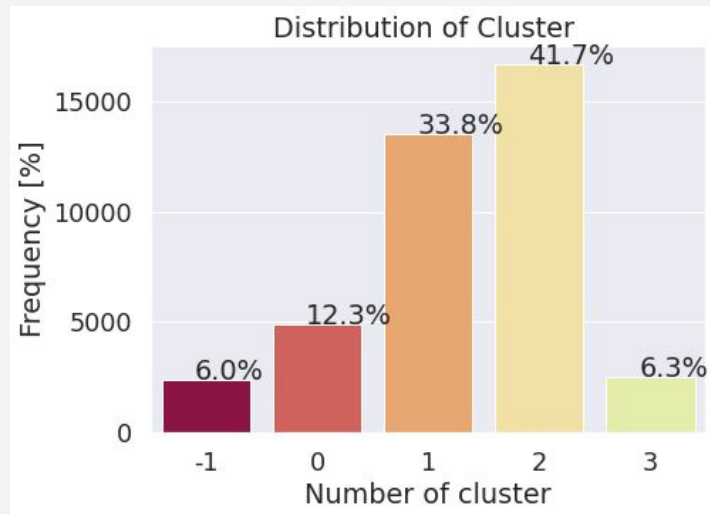
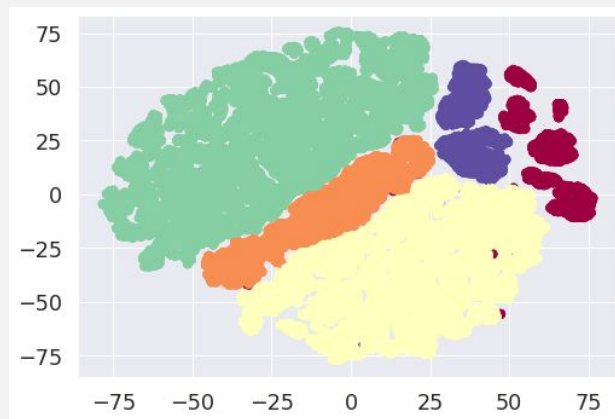
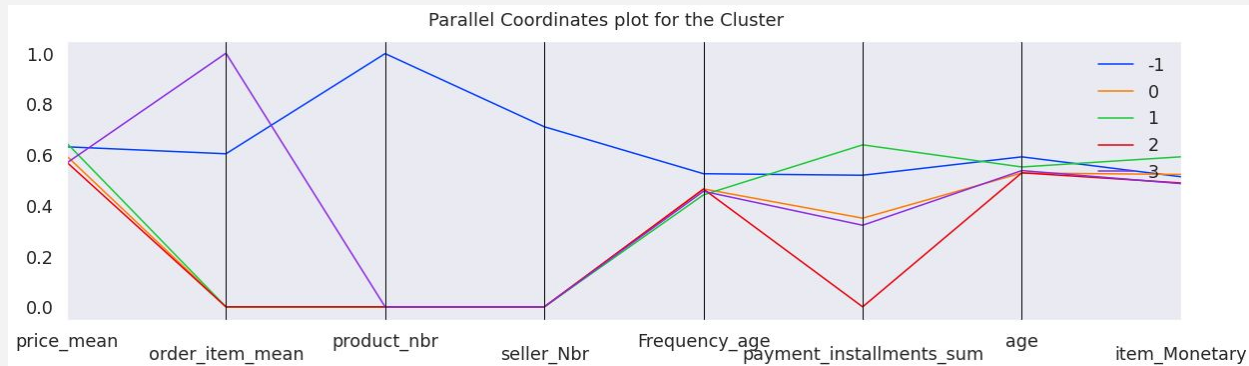
# OPTIC+ MinMaxScaler

**-Nombre de clusters :5**

**-Client 2** représente les clients qui payent en une seul fois.

**-Client -1** représente les clients qui achètent plusieurs produits différents et chez des vendeurs différents.

**-On a 6% de bruit:**-Clusters déséquilibrés mais interprétables.

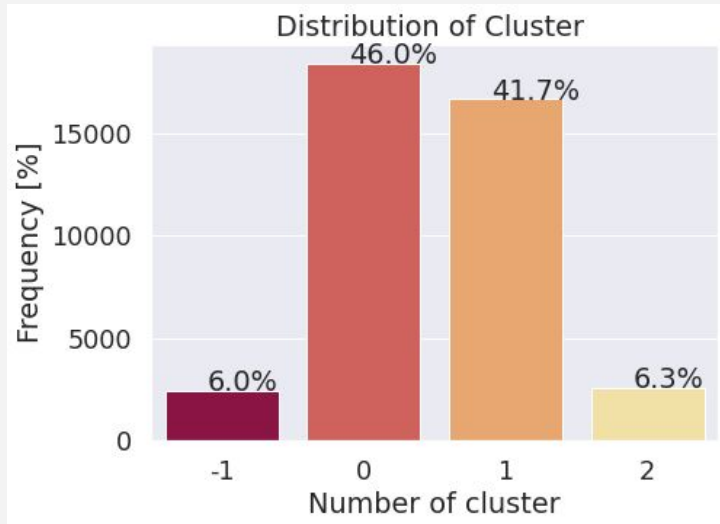
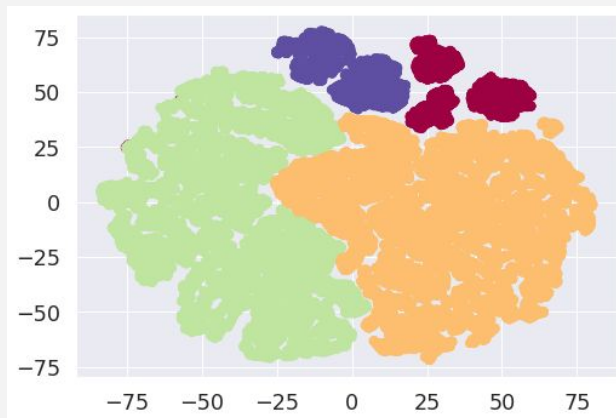
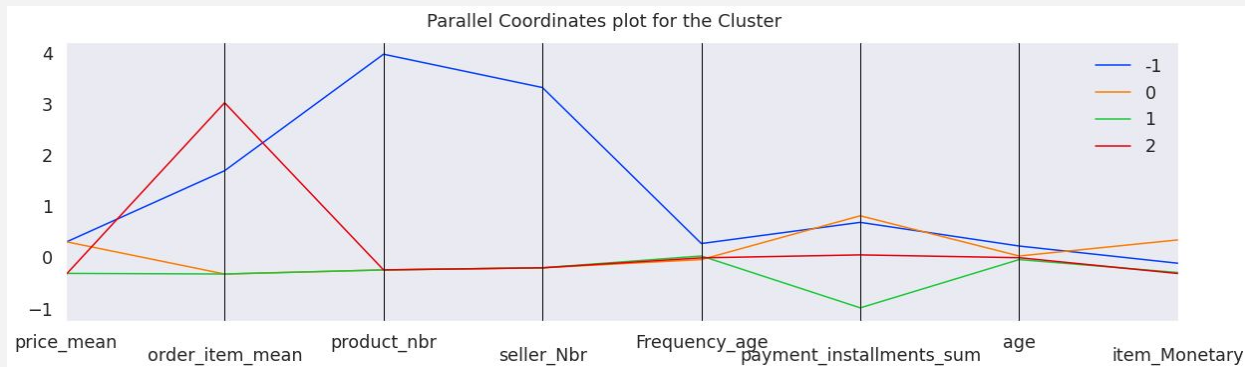


# OPTIC+ StandardScaler

**-Nombre de cluster :4**

**-Client (-1)** représente les clients qui achètent chez plusieurs vendeurs.

**-On a 6% de bruit:**-Clusters déséquilibrés mais interprétables.

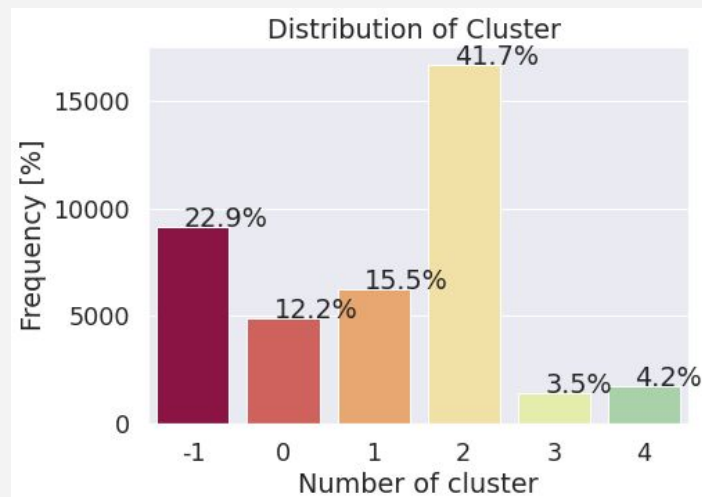
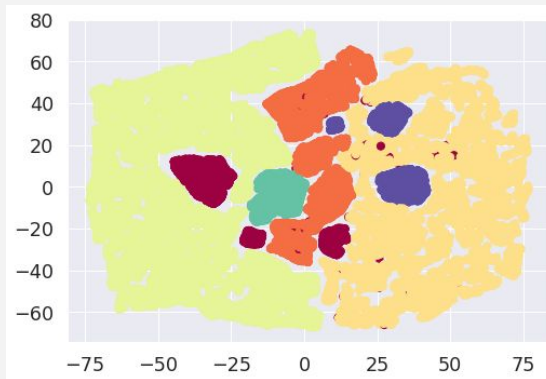
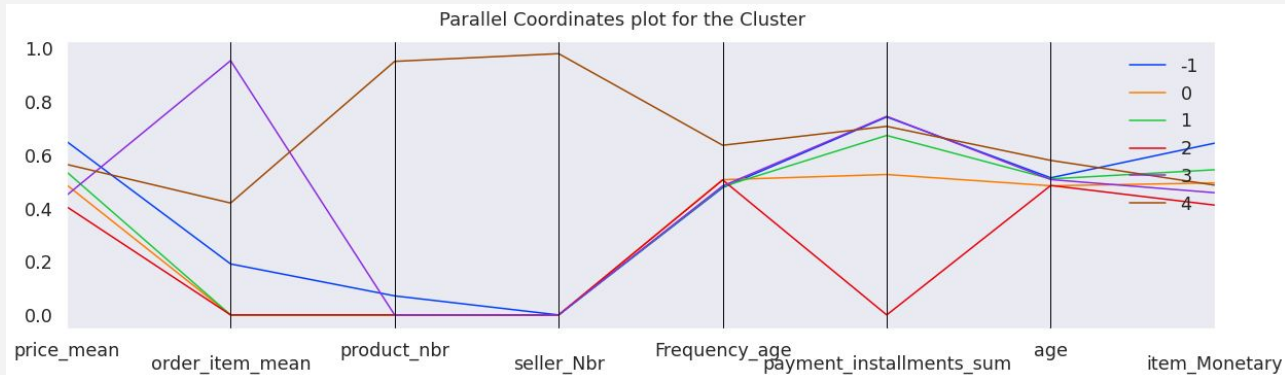


# OPTICS+ QuantileTransformer

-Nombre de cluster :6

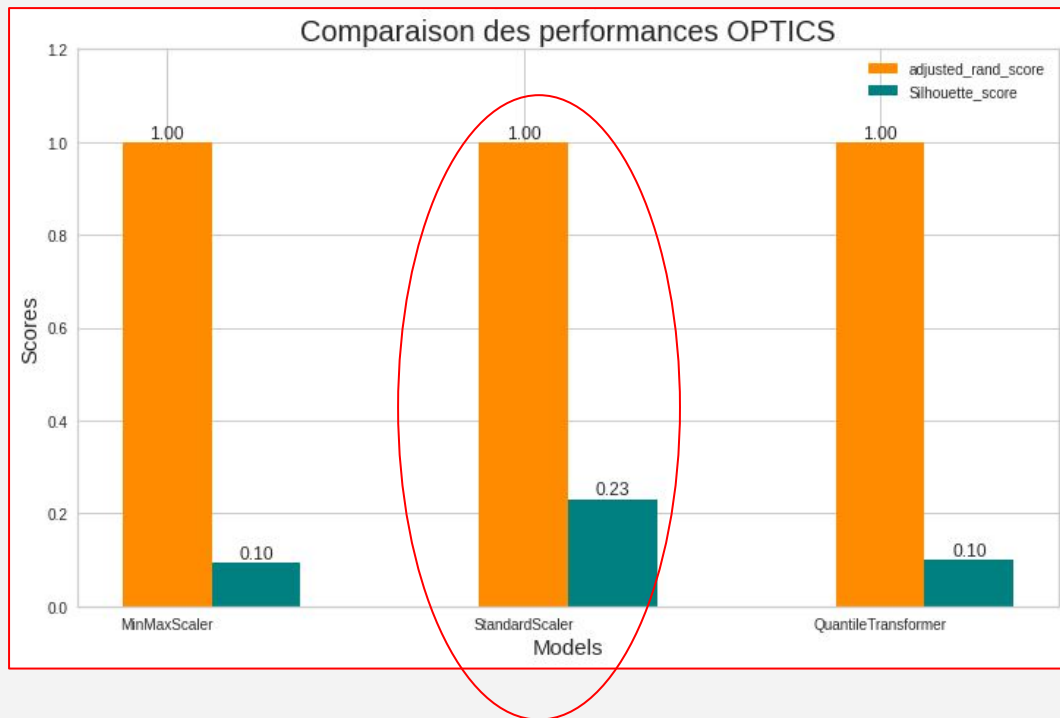
-Client 4 représente les clients qui achètent chez plusieurs vendeurs.

-On a 23% des données sont des bruits:-Clusters déséquilibrés mais interprétables.



# Stabilité des clusters

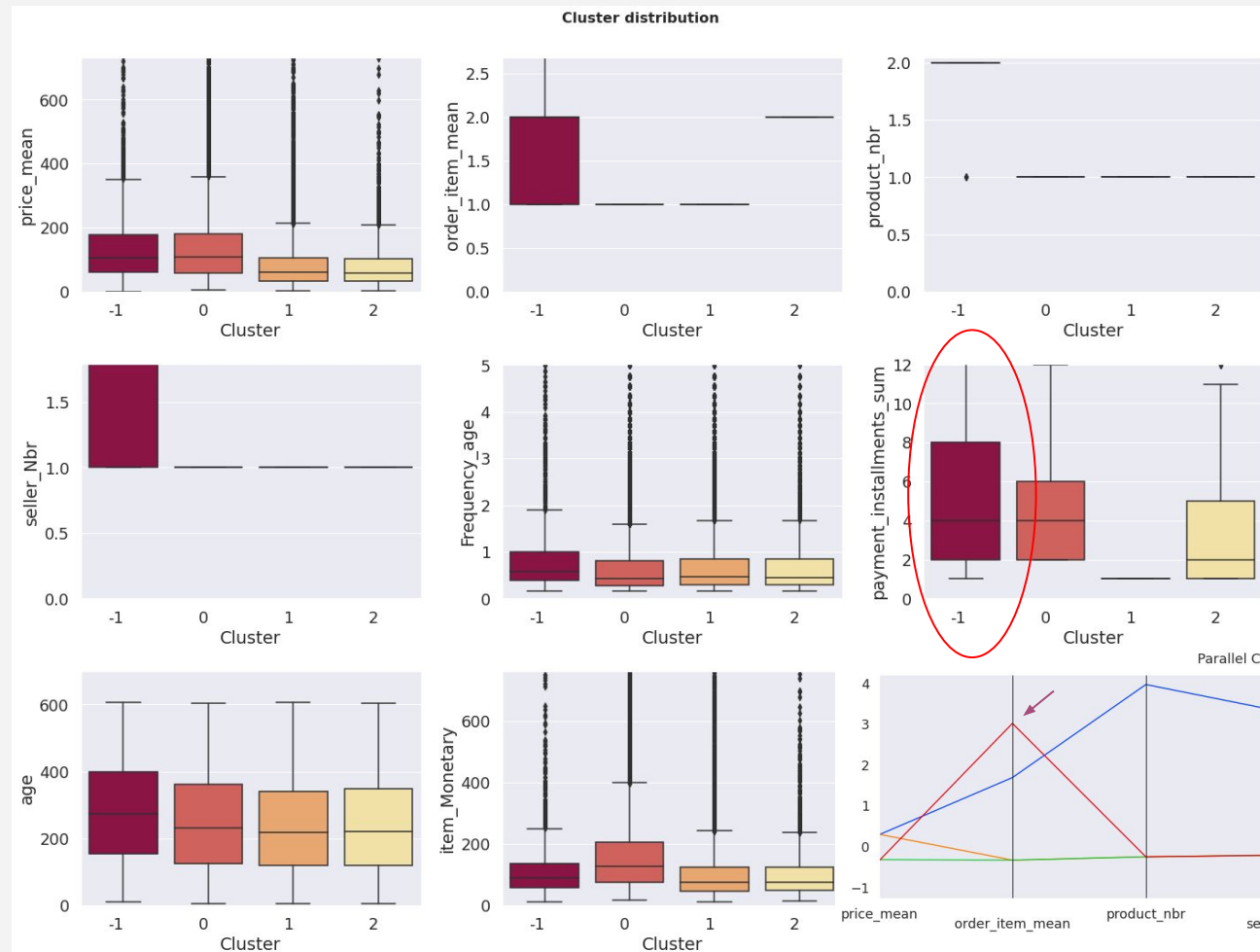
- les 3 modèles sont **stables** mais on retient le **StandardScaler** car il donne un meilleur score de silhouette à 0.23.



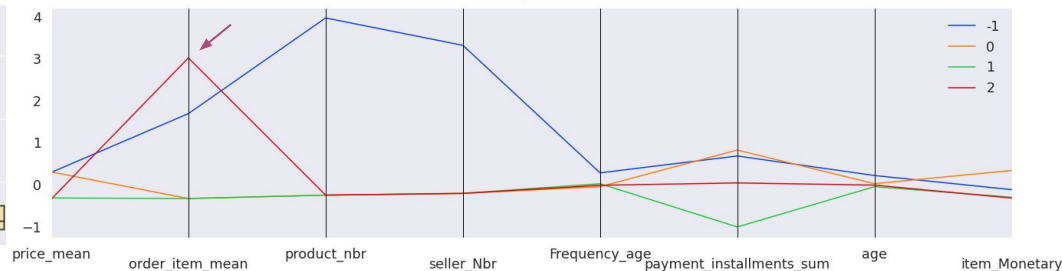


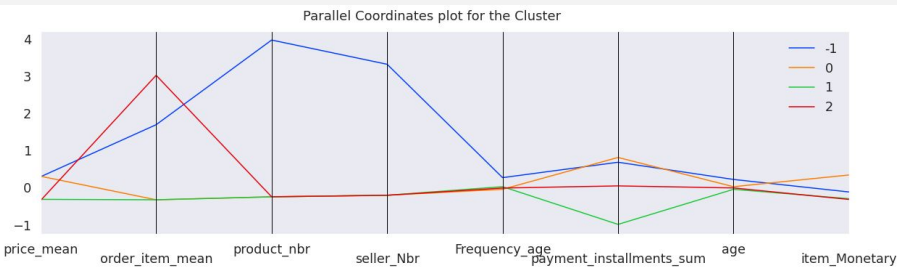
# INTERPRETATION

- **Le client -1** achètent chez plusieurs vendeurs et des articles différents
- **Le client 2** achètent le même article en quantité



Parallel Coordinates plot for the Cluster





# PROFIL CLIENT OPTICS



**Fidèle**

Les clients **anciens** qui achètent **régulièrement**, chez **plusieurs vendeurs**, et plusieurs produits **différents (product\_nbr)** avec une valeur de **panier moyennement élevée** et ont besoin d'un paiement **échelonné**.



**Panier moyen**

Les clients qui font un seul achat d'un seul produit avec un prix relativement **élevé (100 Réal)** et ils payent en plusieurs fois.



**Econome**

Les clients qui font des achats d'articles **pas chers** et payent en une **seule fois**.



**Fidèle à un article**

Les clients qui achètent le même article, en quantité, pour des montants relativement **faibles à moyens**. Ils ont besoin de payer en **plusieur fois**

# 4-MODEL: Hierarchical Clustering



# Algorithmes de clustering hiérarchique : Pre-processing

-**Choix AgglomerativeClustering**: Approche ascendante qui commence par de nombreux petits clusters et les fusionne pour créer de plus grands clusters.

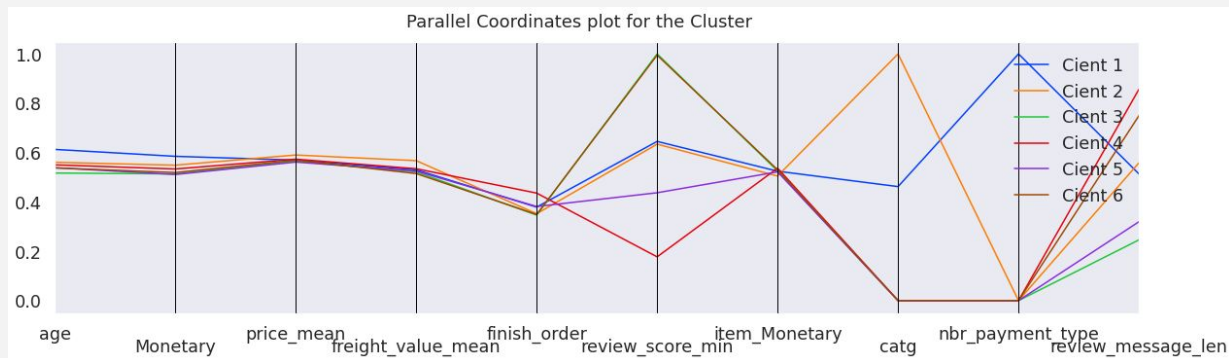
-**PowerTransformer** pour normaliser des distributions.

-**Transformation et Mise à l'échelle** :MinMaxScaler ,StandardScaler, QuantileTransformer

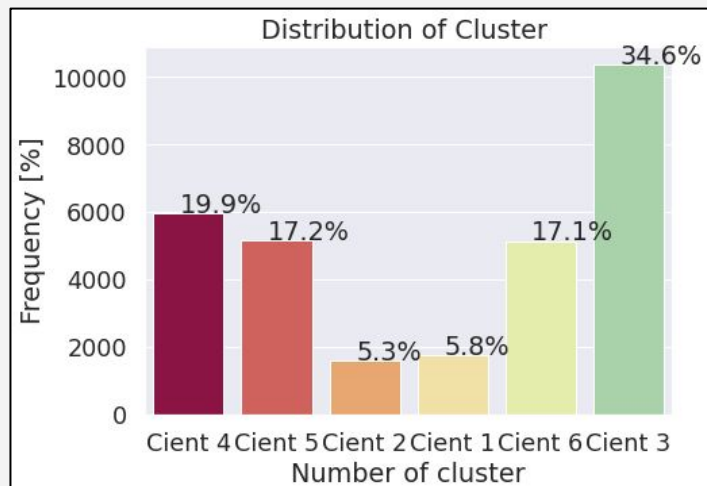
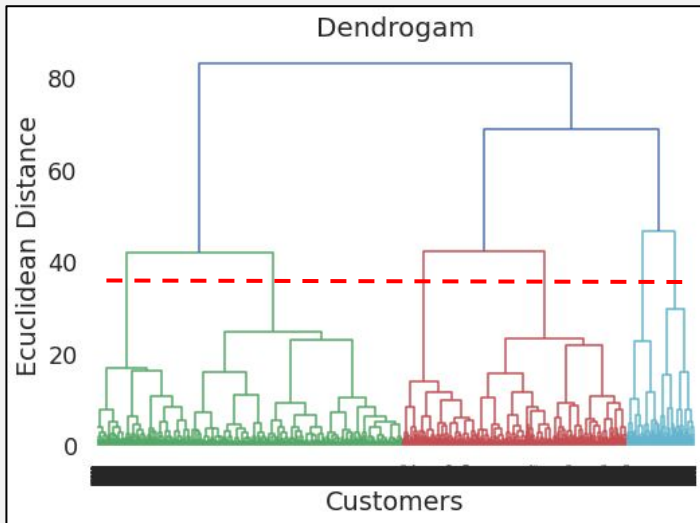
## Features selections: 10

<b>age</b>	La période depuis la première commande
<b>Monetary</b>	Le montant d'achat
<b>price_mean</b>	Le prix moyen des produits
<b>freight_value_mean</b>	Poids moyen des produits achetés
<b>finish_order</b>	La durée pour clôturer une commande
<b>review_score_min</b>	Le score minimum donné par le client
<b>item_Monetary</b>	Panier moyen par client(montant/nbr de commande)
<b>catg</b>	Le nombre des catégories des produits commandés
<b>nbr_payment_type</b>	Le nombre de types de paiements utilisés par le client
<b>review_message_len</b>	La longueur moyenne du commentaire /message

# CAH+ MinMaxScaler

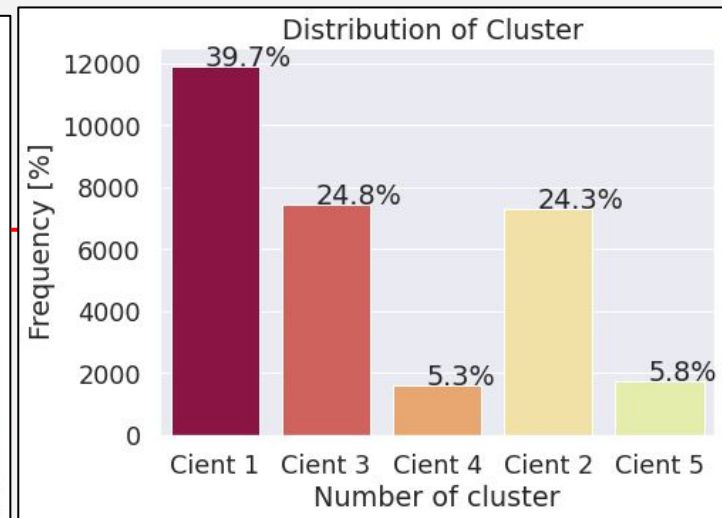
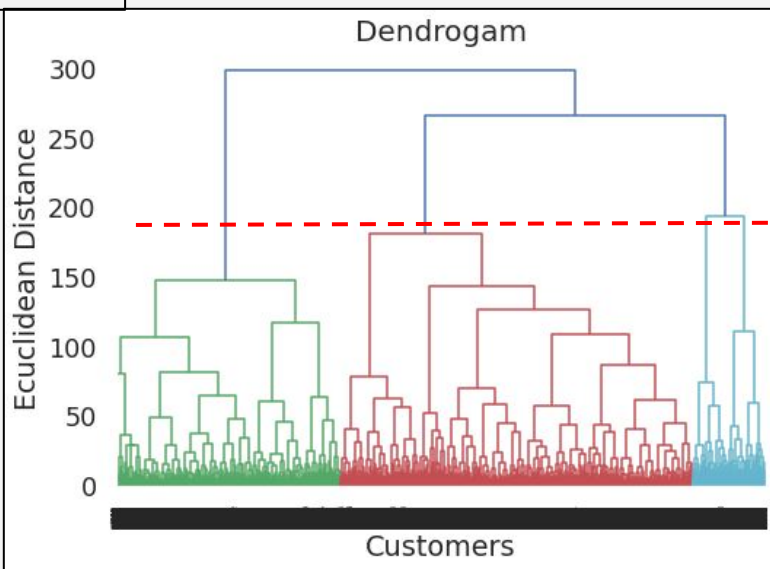
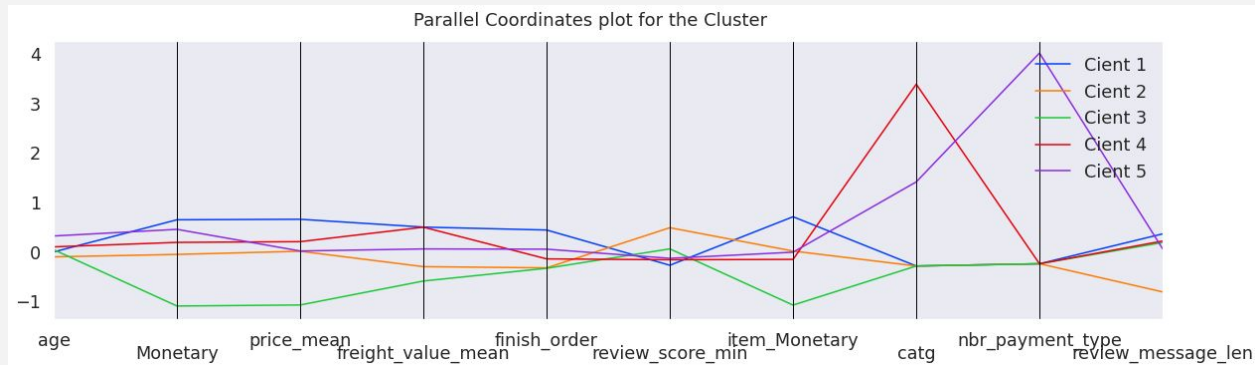


- Nombre de cluster **6**
- **Non** interprétable
- Non équilibrés



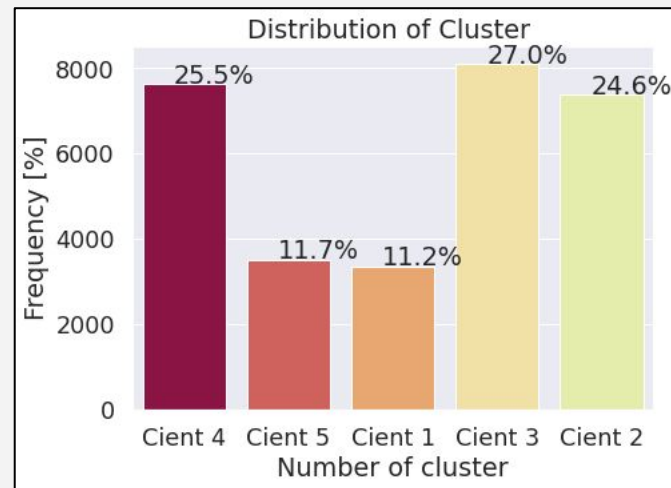
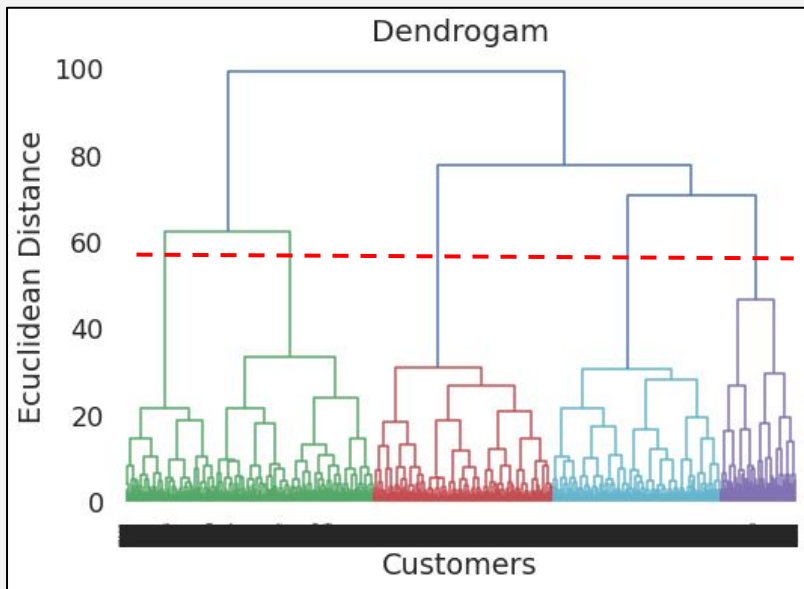
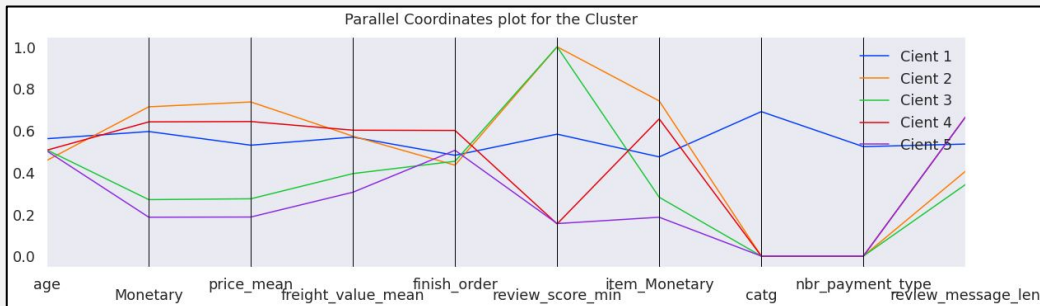
# CAH+ StandardScaler

- Nombre de cluster **5**
- interprétable
- Non équilibré en nombre de client par cluster



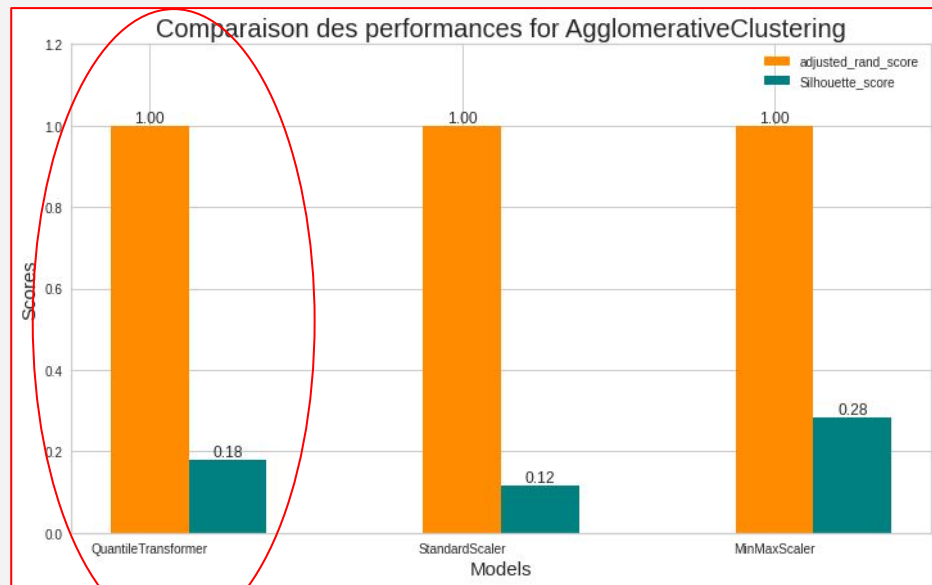
# CAH+ QuantileTransformer

- Nombre de cluster **5**
- **interprétable**
- équilibré



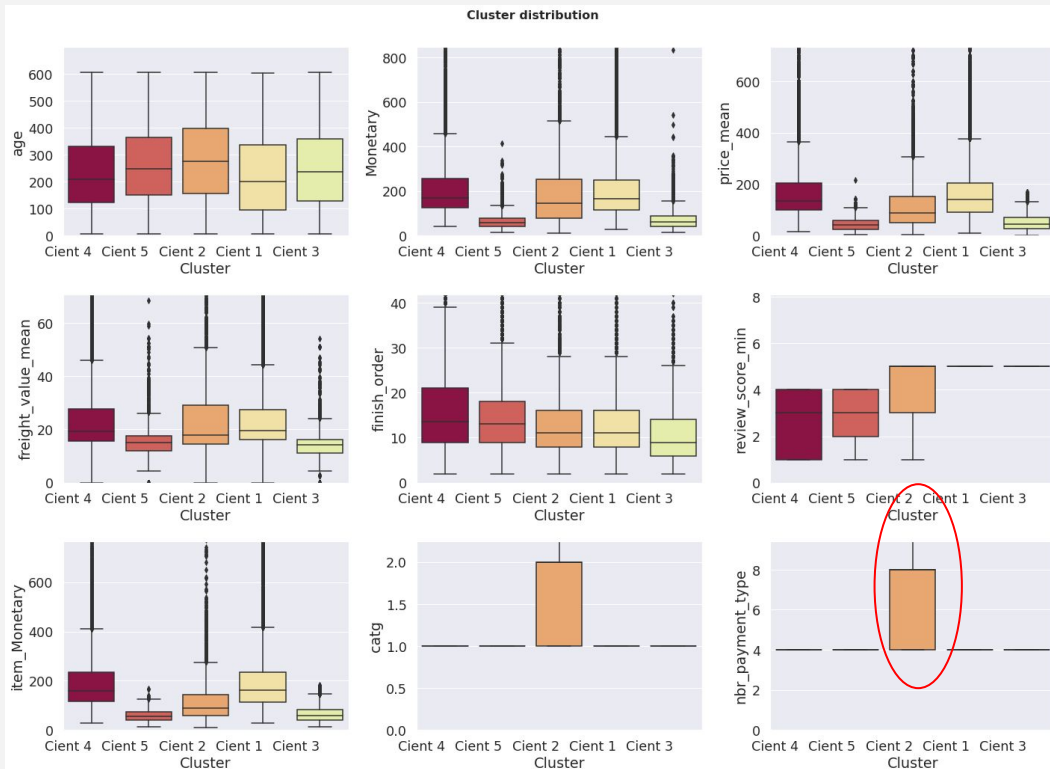
# Stabilité des clusters

- les 3 modèles sont **stables** mais on retient le **QuantileTransformer** car il donne un meilleur score de silhouette à 0.18 et le mieu interprétable





# INTERPRETATION



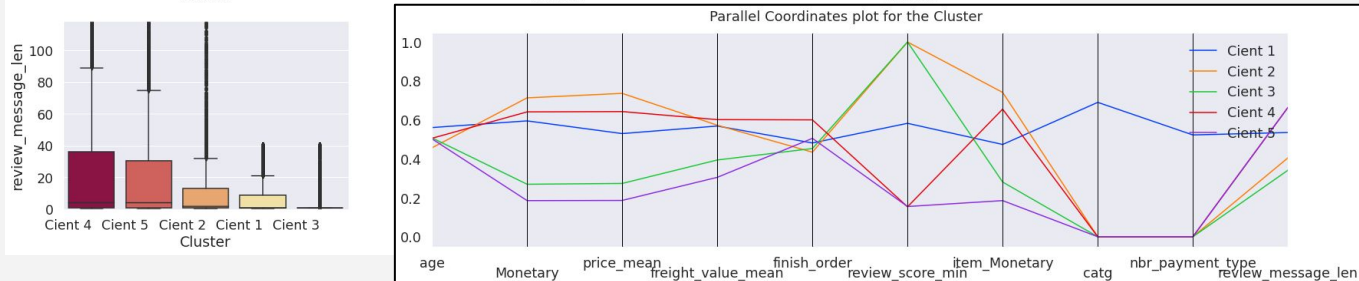
-Le client 4: sa commande a pris beaucoup de temps à être livré.

- Le client 5 et 3: les plus économes.

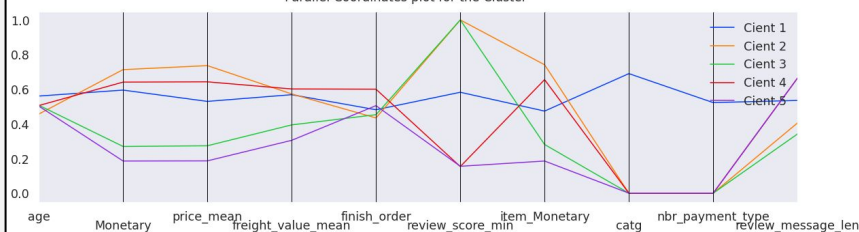
-Le client 3: achète plusieurs catg d'articles

-Le client 2 le plus anciens

-Le client 4 et 1 sont les plus dépensiers



Parallel Coordinates plot for the Cluster



# PROFIL CLIENT CAH



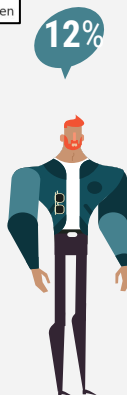
## Fidèle

Les clients **anciens et réguliers** qui utilisent plusieurs **types de paiement** et achètent différentes **catégories** de produits.



## Nouveau dépensier

Les clients **nouveaux** : font des achats pour des montants **élevés** et ils sont **contents**.



## Econome Satisfait

Les clients plutôt **anciens** qui font des achats pour des montants **faibles** à **moyens** et ils sont **contents**.



## Dépensier Insatisfait

Les clients qui sont plutôt **anciens**, font des achats pour des montants relativement **élevés**, achètent des produits assez volumineux et ils sont **mécontents** car le délais de la **livraison** était **long**.



## Econome Insatisfait

Les clients plutôt **anciens**, font des achats pour des **montants faibles**, ils sont **insatisfaits** de leurs achats et ils ont laissé un commentaire assez **long**.

# Comparaison

	Kmeans	OPTICS	CAH
Nombre de clusters	5	4	5
Taille de la dataframe	full	40 000	30 000
Stabilité	Bon	Bon	Bon
Interpretabilité	oui	oui	oui
Temps de calcul	rapide	long	long
Equilibre des clusters	Oui	Non	Non

# Conclusions

Contrat de maintenance à **renouveler** tous les 6 mois .

Revoir les résultats avec l'équipe Marketing pour choisir le modèle qui répond aux mieux à leurs attentes.

## **Pistes d'amélioration:**

Manque de précision de la part de l'équipe de marketing pour nous guider sur la partie de features selection .

Manque de plus de données sur les clients comme l'age , sexe..

Essayer d'autres modèles de classification