

Projet 5 : Segmentation des clients d'un site e-commerce

Fanjamalala Rajaonalison
11/2021

Présentation de l'étude

A vertical line is positioned to the right of the title. In the bottom right corner of the slide, there is a large yellow right-angled triangle pointing towards the top-left.

Contexte

Client :



solution de vente sur les marketplaces en ligne

Mission : fournir aux équipes d'e-commerce une **segmentation des clients** qu'elles pourront utiliser au quotidien pour leurs campagnes de communication.

⇒ Base de données anonymisée comportant des informations sur l'historique de commandes, les produits achetés, les commentaires de satisfaction, et la localisation des clients.

Problématique

Problématique : Comprendre les différents types d'utilisateurs grâce à leur comportement d'achat et à leurs données personnelles



Objectifs : Segmentation des clients

- **Description actionable** de la segmentation et de sa logique sous-jacente pour une utilisation optimale,
- **Proposition de contrat de maintenance** basée sur une analyse de la stabilité des segments au cours du temps.

Méthodologie

Traitements et Analyses des données

- Mise en forme & Cleaning
- Exploration

Essais de différentes segmentation

- Segmentation globale
- Segmentation RFM
- Modèles : DBSCAN, K-MEANS

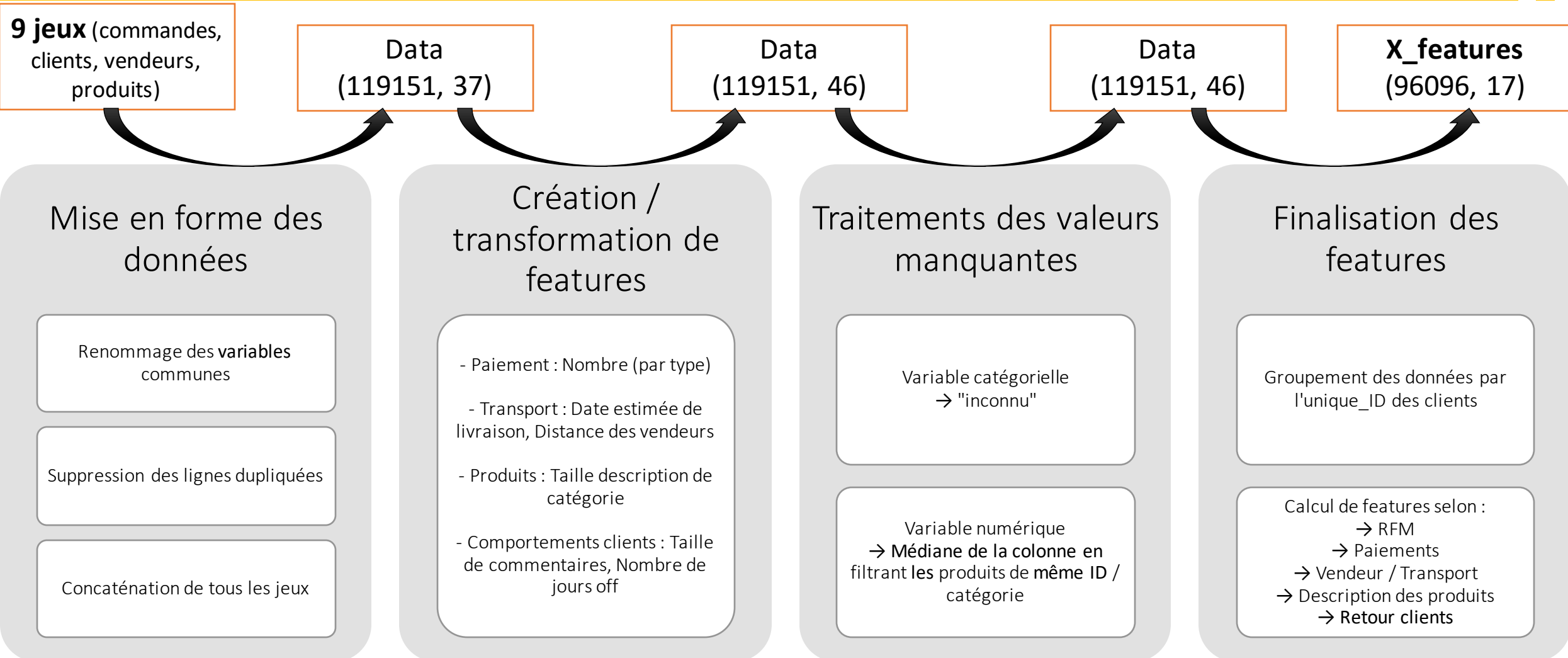
Finalisation de l'étude

- Sélection du modèle final
- Définition des groupes
- Maintenance du modèle

Traitements Analyses des données

A vertical line is positioned to the right of the text. A yellow triangle is located in the bottom right corner of the slide, pointing upwards and to the left.

Mise en forme & Cleaning



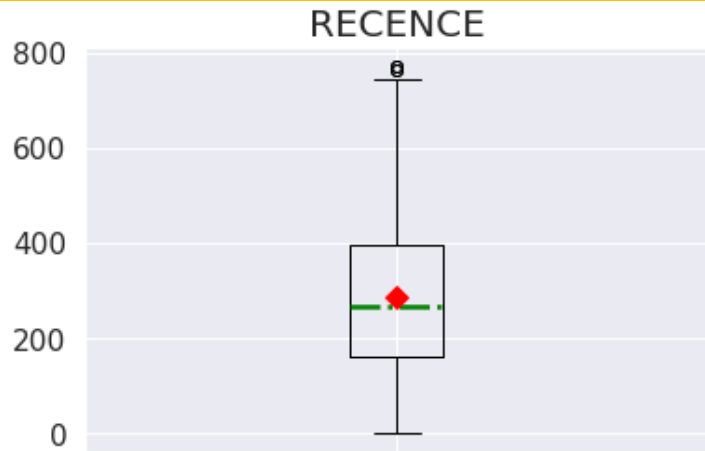
Mise en forme & Cleaning

#	Column	Non-Null Count	Dtype
0	recence	96096 non-null	int64
1	frequence	96096 non-null	int64
2	montant	96096 non-null	float64
3	n_payment	96096 non-null	float64
4	pay_credit_card	96096 non-null	uint8
5	pay_debit_card	96096 non-null	uint8
6	pay_boleto	96096 non-null	uint8
7	pay_voucher	96096 non-null	uint8
8	seller_distance	96096 non-null	float64
9	freight_value	96096 non-null	float64
10	estimated_delivery_days	96096 non-null	int64
11	product_description_lenght	96096 non-null	float64
12	product_category_lenght	96096 non-null	float64
13	product_photos_qty	96096 non-null	float64
14	day_off	96096 non-null	float64
15	review_score	96096 non-null	float64
16	review_comment_lenght	96096 non-null	float64

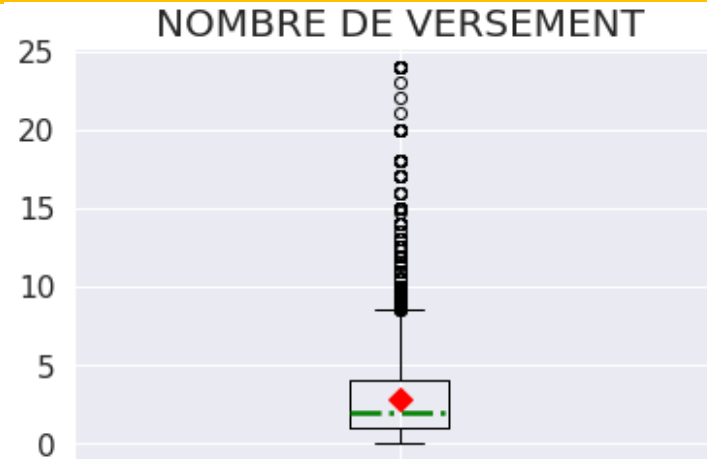
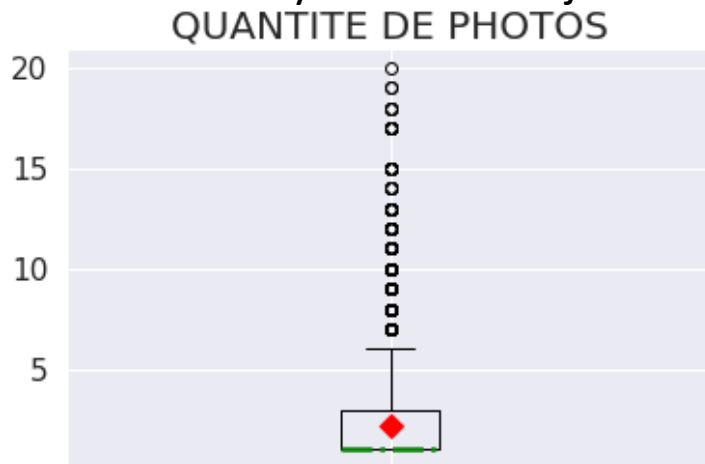
Données : (96096, 17)

- RFM
- Paiements
- Vendeur / Transport
- Description des produits
- Retour clients

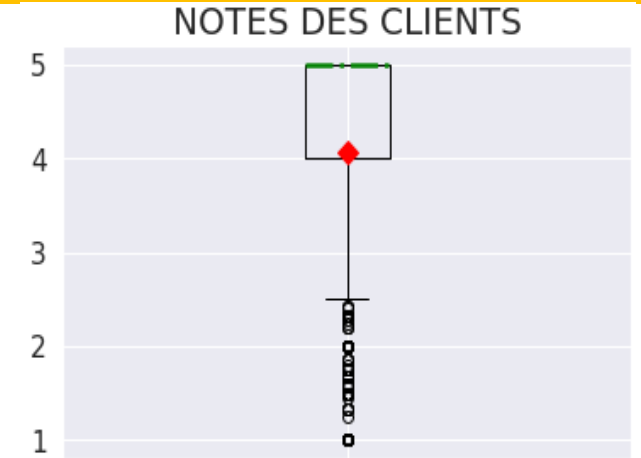
Exploration



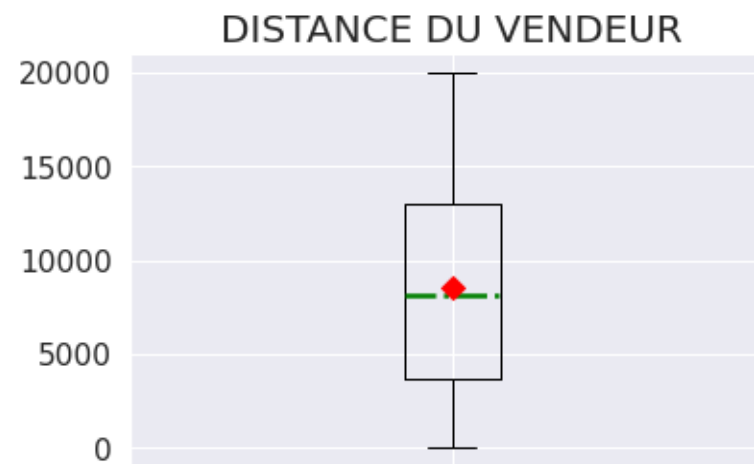
Le nombre de jours passés après la dernière commande varient fortement avec une moyenne de 287 jours.



Les clients paient généralement en 2, 3, 4 fois.



Les clients donnent en général une bonne note.



En général, les produits achetés :

- présentent peu de photos
- sont vendus par des vendeurs en moyenne à 10000km des clients

Exploration

	recence	frequence	montant	n_payment	pay_credit_card	pay_debit_card	pay_boleto	pay_voucher	seller_distance	freight_value	estimated_delivery_days	product_description_lenght	product_category_lenght	product_photos_qty	day_off	review_score	review_comment_lenght
recence	1	0.01	0.01	0.05	-0.01	-0.05	0.02	0.02	0.03	-0.05	0.23	-0.04	-0.04	-0	-0.06	-0.03	0.02
frequence	0.01	1	0.29	0.02	0.5	0.02	0.21	0.62	0.03	0.02	0.04	-0.05	0.04	-0.07	0.14	-0.08	0.05
montant	0.01	0.29	1	0.27	0.15	-0.01	0.02	0.21	0.04	0.38	0.08	0.16	0.04	0.02	0.05	-0.03	0.05
n_payment	0.05	0.02	0.27	1	0.32	-0.08	-0.29	-0.07	0.06	0.21	0.1	0.04	0.01	0	0.02	-0.03	0.05
pay_credit_card	-0.01	0.5	0.15	0.32	1	-0.14	-0.5	-0.02	0.02	0.03	0.03	-0.04	0.02	-0.06	0.13	-0.06	0.04
pay_debit_card	-0.05	0.02	-0.01	-0.08	-0.14	1	-0.04	-0.01	-0.01	-0.01	-0.03	0	0	-0.01	0.01	0.01	0.01
pay_boleto	0.02	0.21	0.02	-0.29	-0.5	-0.04	1	-0.05	0.01	-0	0.02	-0.01	0.03	-0.02	0.03	-0.02	0.01
pay_voucher	0.02	0.62	0.21	-0.07	-0.02	-0.01	-0.05	1	0.01	0	0.01	-0.01	-0	-0.01	0.01	-0.01	0.01
seller_distance	0.03	0.03	0.04	0.06	0.02	-0.01	0.01	0.01	1	0.17	0.26	0.01	0.02	-0.02	0.04	-0.02	0.01
freight_value	-0.05	0.02	0.38	0.21	0.03	-0.01	-0	0	0.17	1	0.29	0.09	0.06	0.02	0.02	-0.04	0.04
estimated_delivery_days	0.23	0.04	0.08	0.1	0.03	-0.03	0.02	0.01	0.26	0.29	1	-0	0.03	-0.04	0.06	-0.05	0.04
product_description_lenght	-0.04	-0.05	0.16	0.04	-0.04	0	-0.01	-0.01	0.01	0.09	-0	1	0.02	0.13	-0.05	0.02	-0
product_category_lenght	-0.04	0.04	0.04	0.01	0.02	0	0.03	-0	0.02	0.06	0.03	0.02	1	0.01	0	0.02	-0.01
product_photos_qty	-0	-0.07	0.02	0	-0.06	-0.01	-0.02	-0.01	-0.02	0.02	-0.04	0.13	0.01	1	-0.05	0.03	-0.01
day_off	-0.06	0.14	0.05	0.02	0.13	0.01	0.03	0.01	0.04	0.02	0.06	-0.05	0	-0.05	1	0.01	-0
review_score	-0.03	-0.08	-0.03	-0.03	-0.06	0.01	-0.02	-0.01	-0.02	-0.04	-0.05	0.02	0.02	0.03	0.01	1	-0.4
review_comment_lenght	0.02	0.05	0.05	0.05	0.04	0.01	0.01	0.01	0.01	0.04	0.04	-0	-0.01	-0.01	-0	-0.4	1
	recence	frequence	montant	n_payment	pay_credit_card	pay_debit_card	pay_boleto	pay_voucher	seller_distance	freight_value	estimated_delivery_days	product_description_lenght	product_category_lenght	product_photos_qty	day_off	review_score	review_comment_lenght

Les variables ne présentent pas de forte corrélation

Essais de différentes segmentations

A vertical line is positioned to the right of the text. A yellow triangle is located in the bottom right corner of the slide, pointing towards the top right.

Segmentation globale

Principe :

Segmentation sur l'ensemble des features (17 features)

Essais de 2 modèles :

- DBSCAN
- K-MEANS

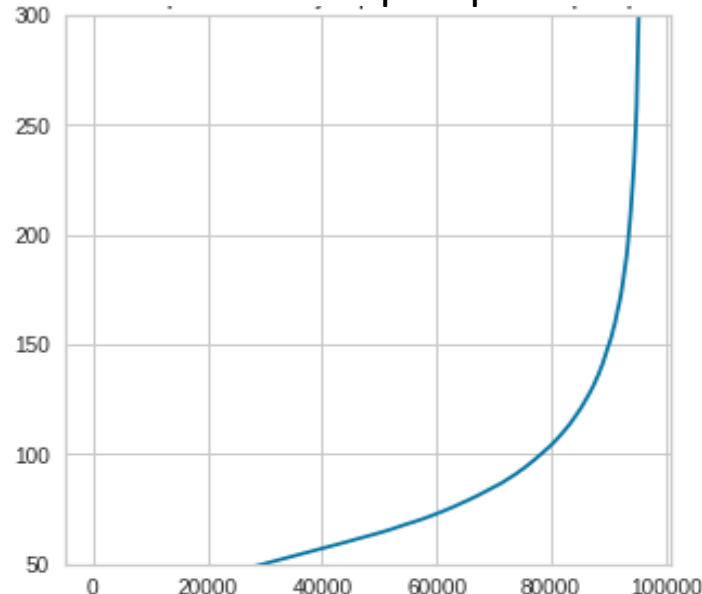
Segmentation globale

1. DBSCAN

a. Estimation des paramètres : eps & min_sample

- min_sample : ~2 fois la taille du jeu $\Rightarrow 34$

- eps : Points triés par distance jusqu'au 34e voisin le plus proche



$\Rightarrow 125$

Segmentation globale

1. DBSCAN

b. Implémentation du modèle

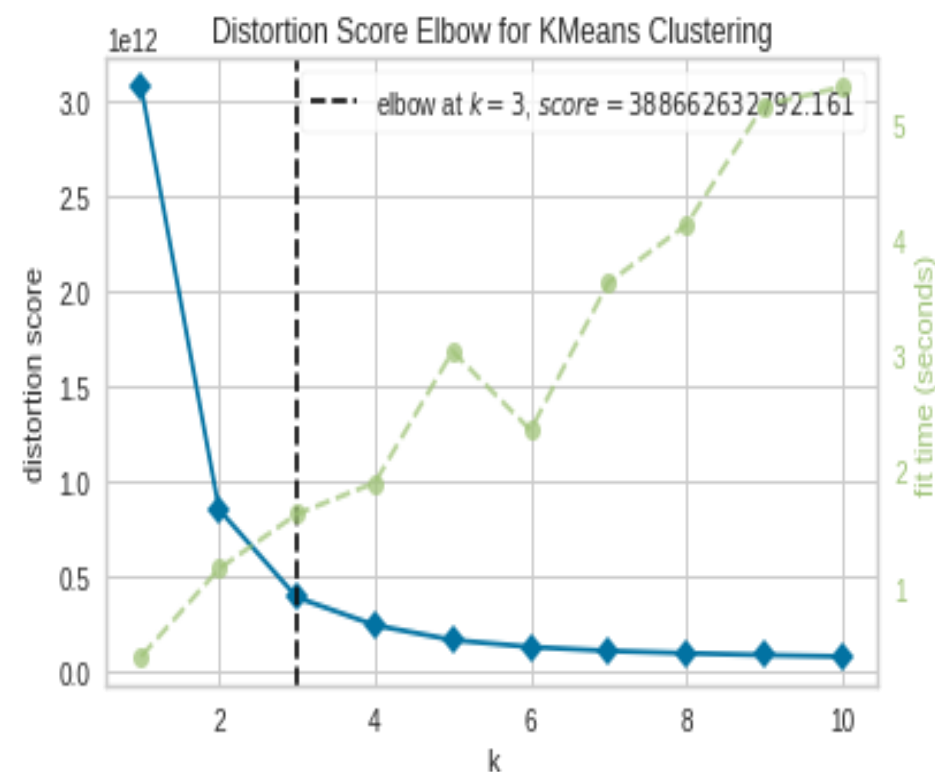
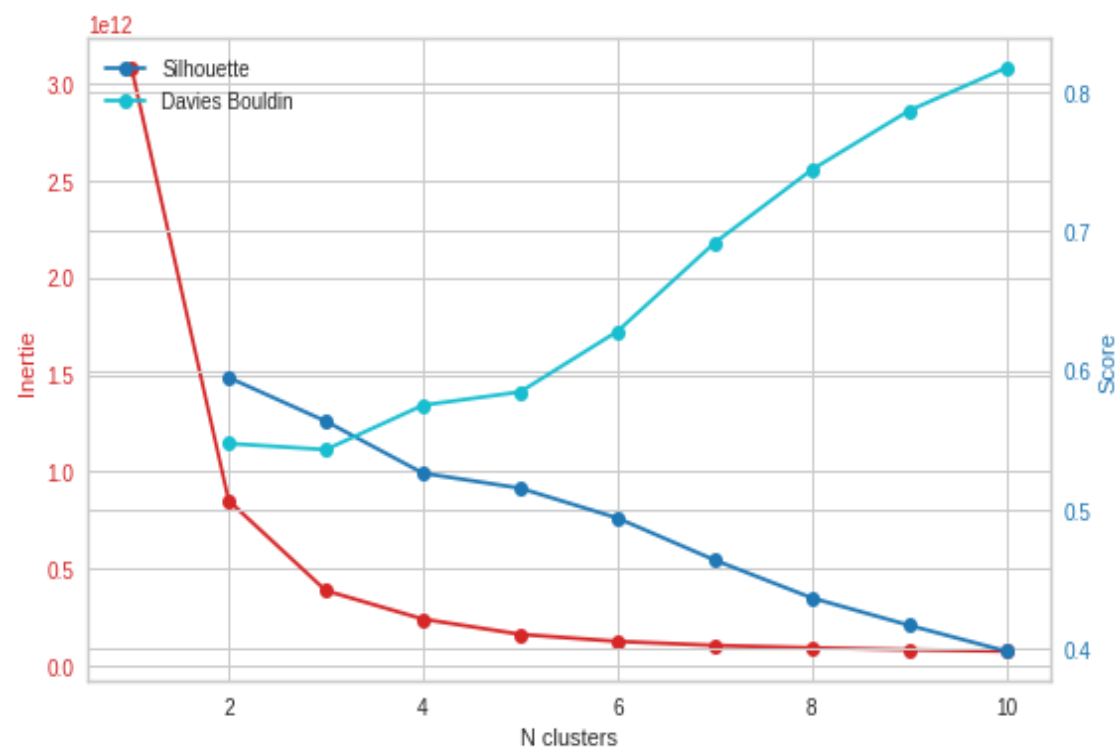
- Nombre de clusters : 14 \Rightarrow Assez grand mais acceptable
 - Silhouette Coefficient: -0.752
 - Score Davies Bouldin: 2.426
 -
- } \Rightarrow Scores pas bons

\Rightarrow **Modèle pas pertinent pour les données**

Segmentation globale

2. K-MEANS

a. Estimation du nombre de cluster



⇒ N_clusters = 3

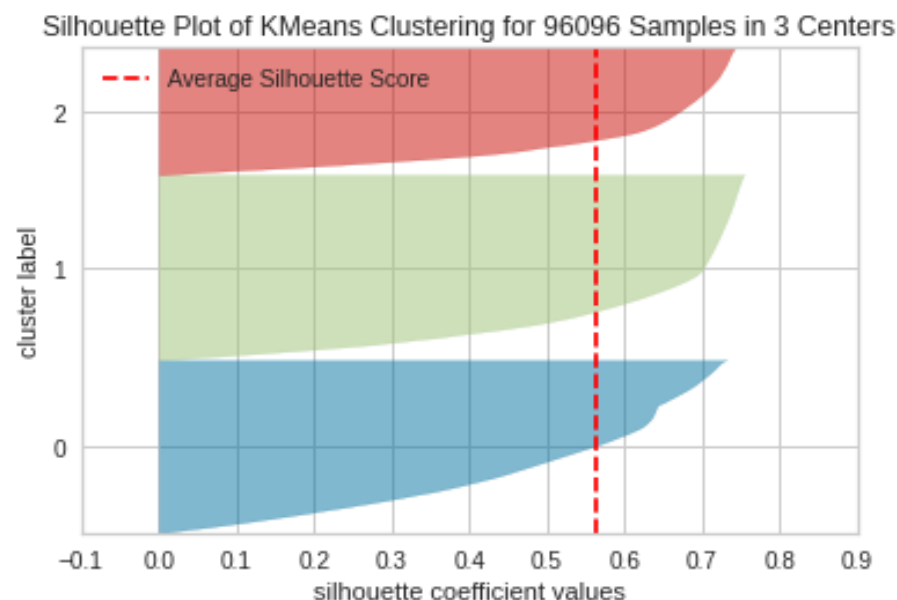
Segmentation globale

2. K-MEANS

b. Implémentation du modèle

- Nombre de clusters : 3
- Silhouette Coefficient: 0.564

⇒ score assez bon



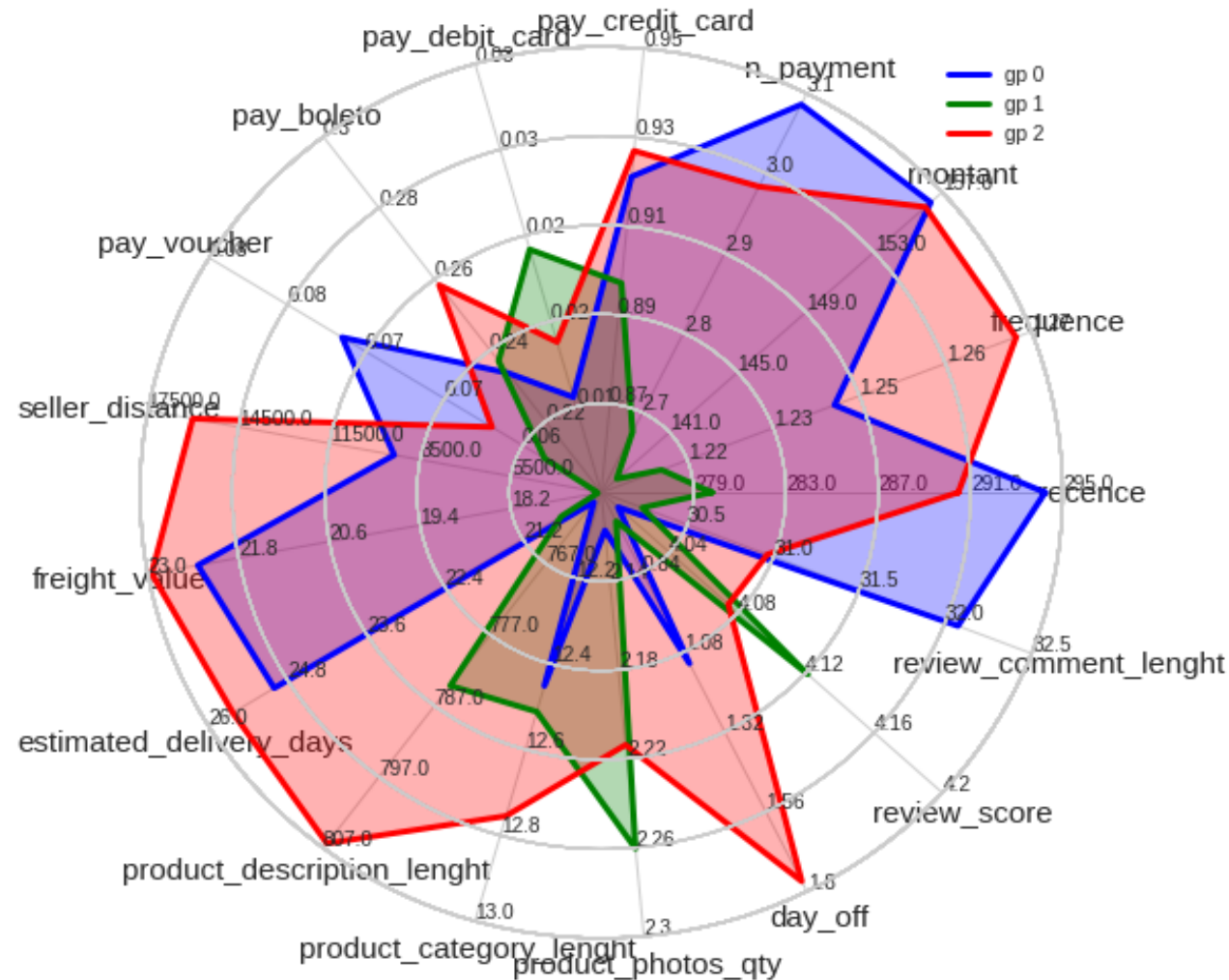
- clusters avec un silhouette score > silhouette moyen
⇒ chaque groupe bien éloigné des autres
- épaisseurs du tracé des silhouette des clusters similaires
⇒ taille des clusters assez similaires entre eux.

- Score Davies Bouldin: 0.543
- ⇒ score assez bon

⇒ **Modèle pertinent pour les données**

Segmentation globale

2. K-MEANS c. Interpretation



Comportements différents pour les 3 groupes obtenus

⇒ Bonne segmentation des clients

Segmentation RFM

Principe :

Segmentation sur 3 features :

- Récence : date de la dernière commande (en nombre de jours)
- Fréquence : fréquence de commandes
- Montant : montant total de tous les commandes passées

Essais de 2 modèles :

- DBSCAN
- K-MEANS

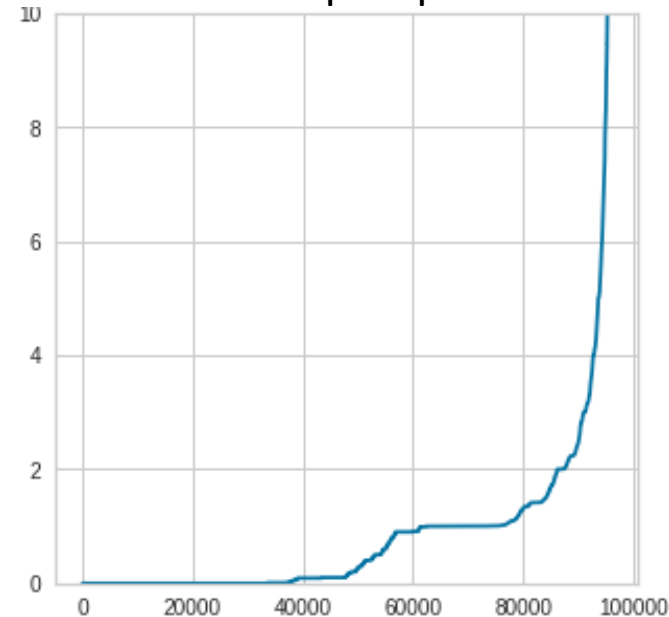
Segmentation RFM

1. DBSCAN

a. Estimation des paramètres : eps & min_sample

- min_sample : ~2 fois la taille du jeu $\Rightarrow 6$

- eps : Points triés par distance jusqu'au 6e voisin le plus proche



$\Rightarrow 1.25$

Segmentation RFM

1. DBSCAN

b. Implémentation du modèle

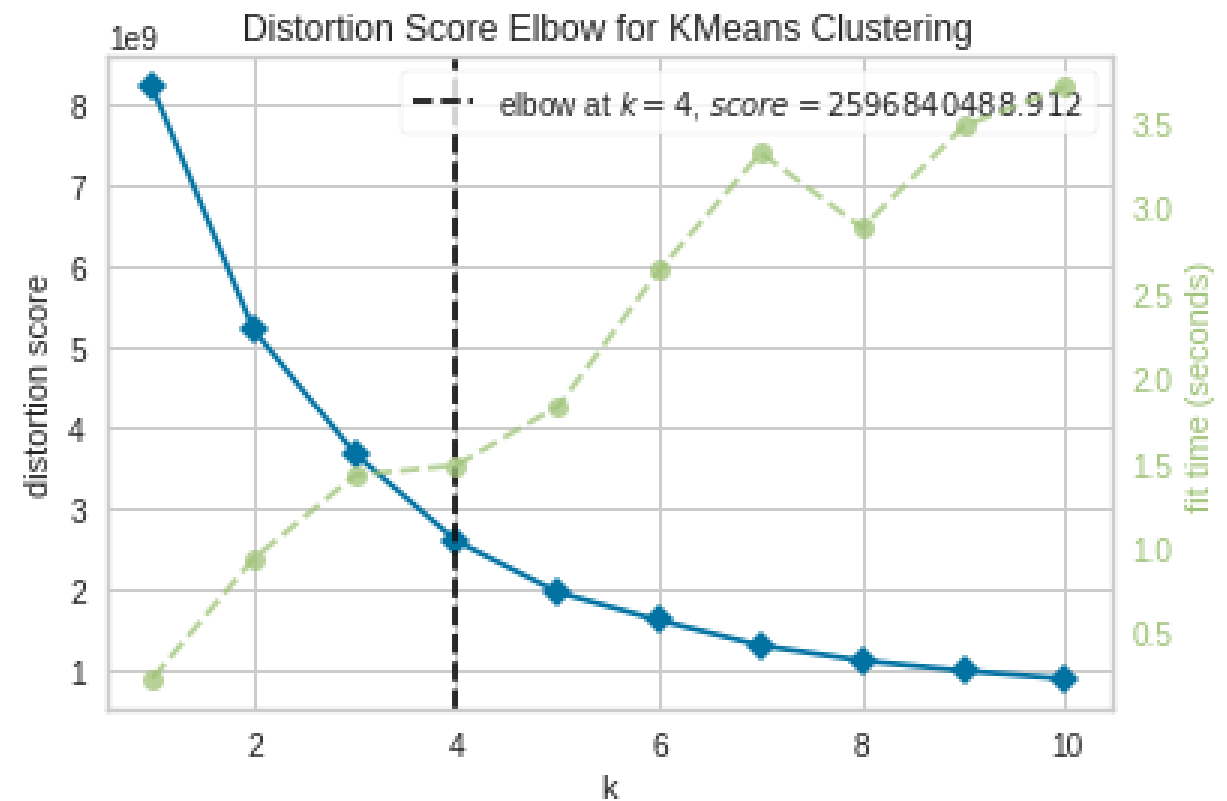
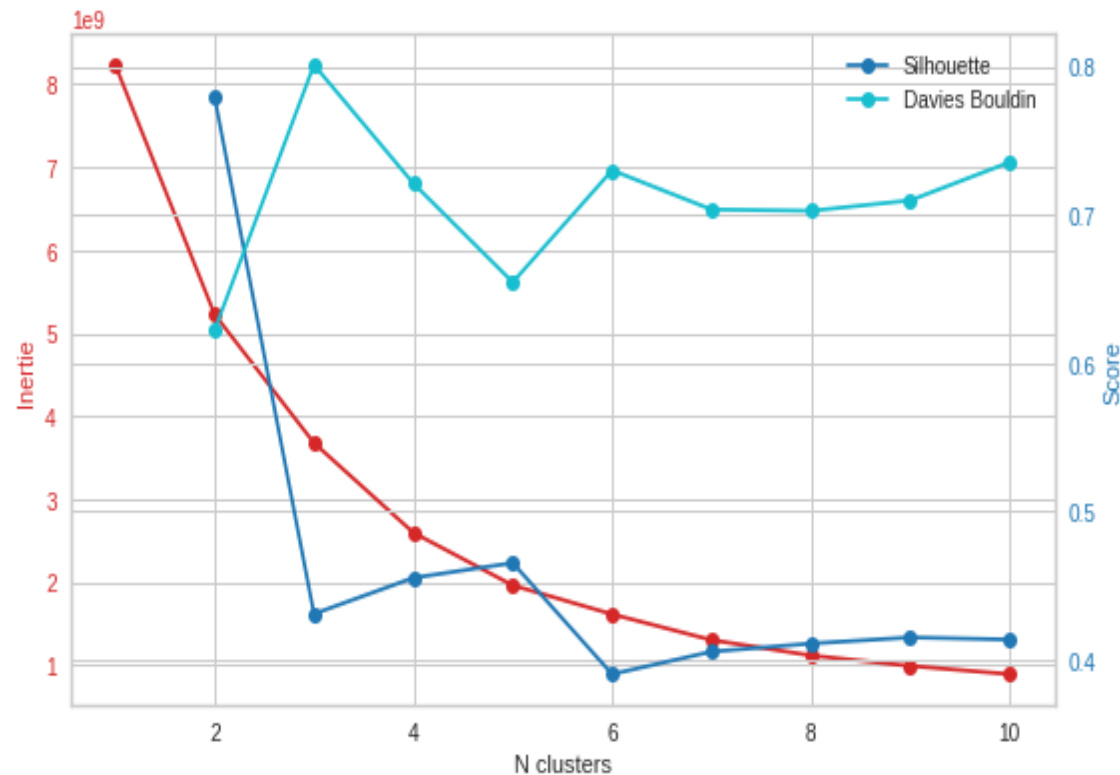
- Nombre de clusters : 110 \Rightarrow Trop grand
 - Silhouette Coefficient: -0.660
 - Score Davies Bouldin: 1.784
- \Rightarrow Scores pas bon

\Rightarrow **Modèle pas pertinent pour les données**

Segmentation RFM

2. K-MEANS

a. Estimation du nombre de cluster



⇒ N_clusters = 5 (meilleurs scores)

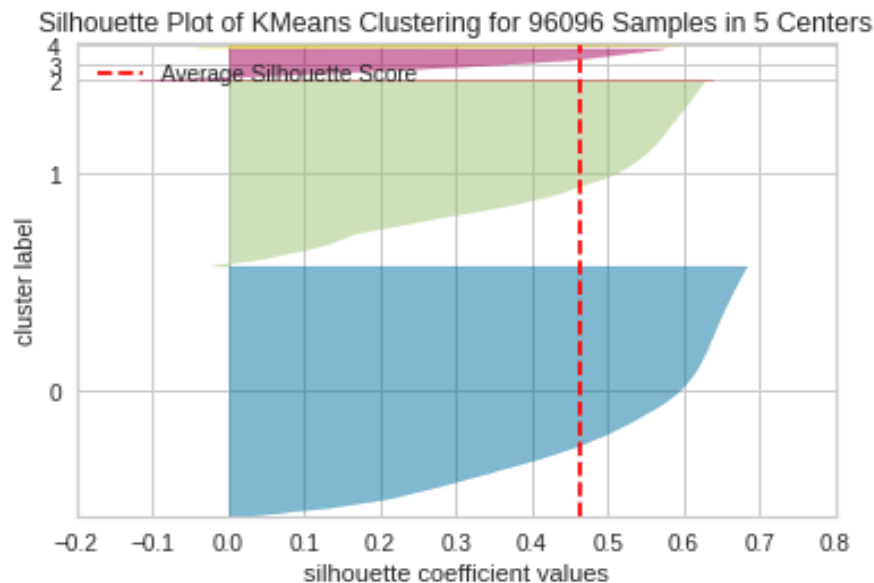
Segmentation RFM

2. K-MEANS

b. Implémentation du modèle

- Nombre de clusters : 5
- Silhouette Coefficient: 0.465

⇒ score moyennement bon



- clusters avec un silhouette score > silhouette moyen
⇒ chaque groupe bien éloigné des autres
- épaisseurs du tracé des silhouette des clusters non similaires
⇒ tailles des clusters pas similaire entre eux.

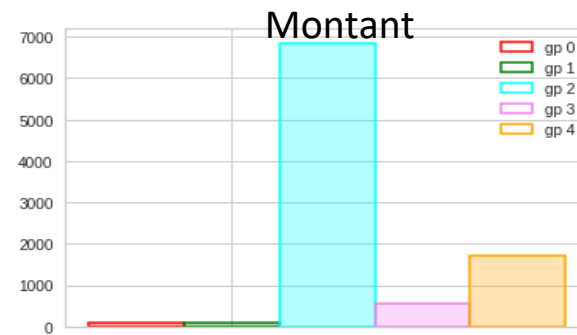
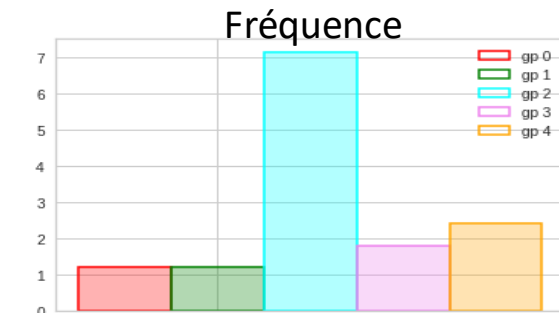
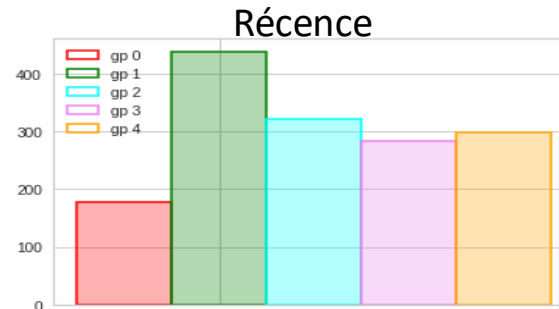
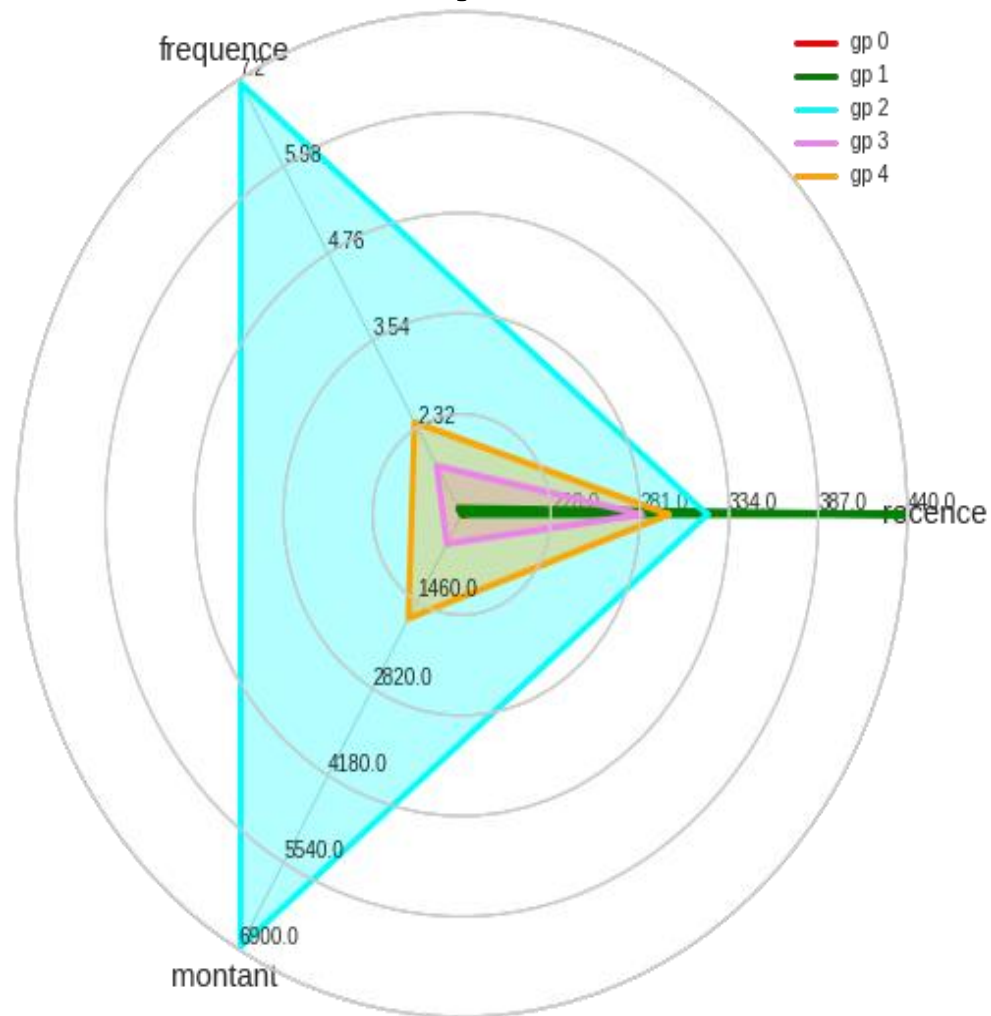
- Score Davies Bouldin: 0.655 ⇒ score moyennement bon

⇒ **Modèle pertinent pour les données**

Segmentation RFM

2. K-MEANS

c. Interpretation



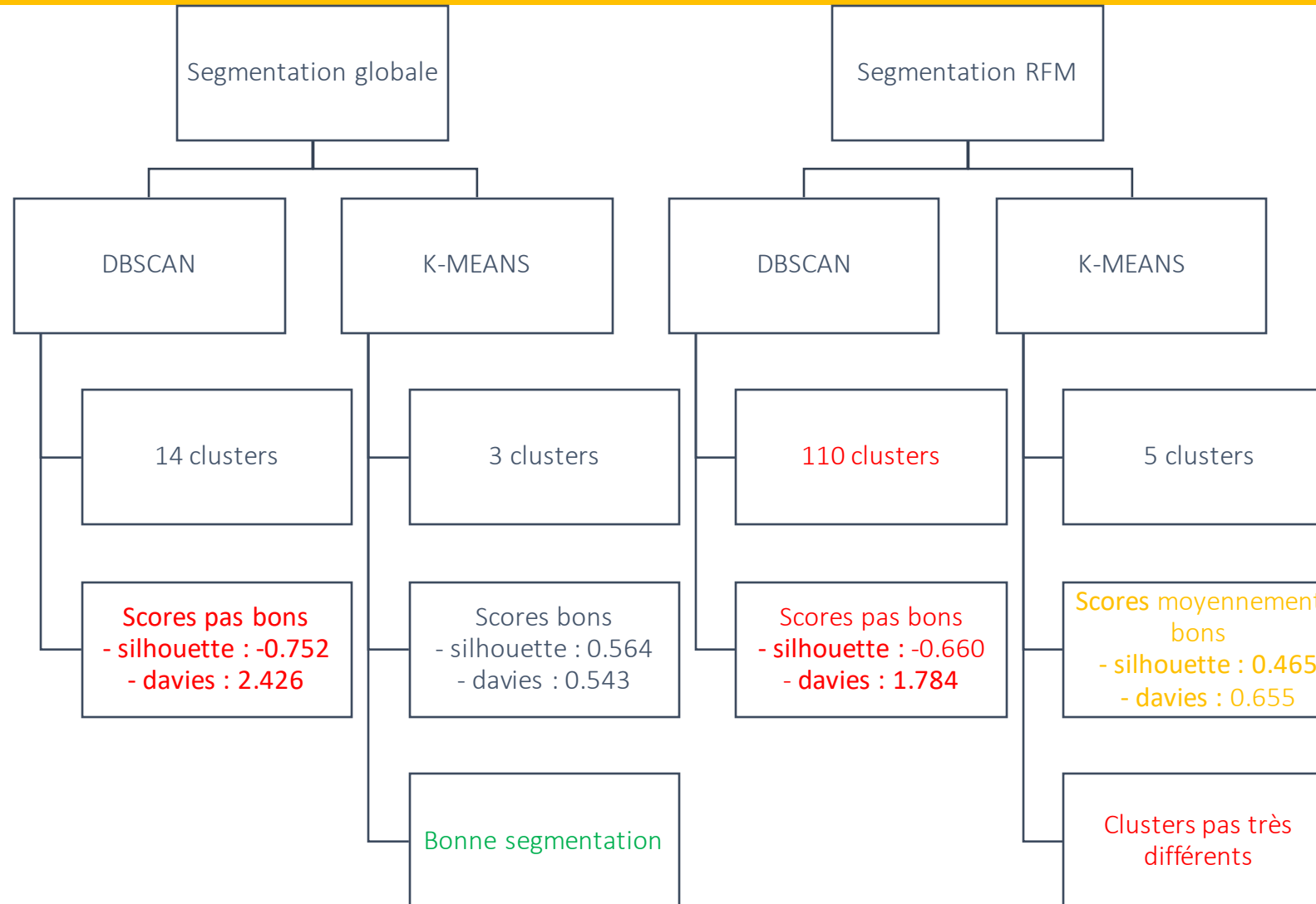
- Comportements assez similaires pour les groupes 0 et 1 (seule différence la date de dernière commande)
- Comportements assez similaires pour les groupes 3 et 4 (seule différence la date de dernière commande)

⇒ **Segmentation des clients moyennement bonne**

Finalisation de l'étude

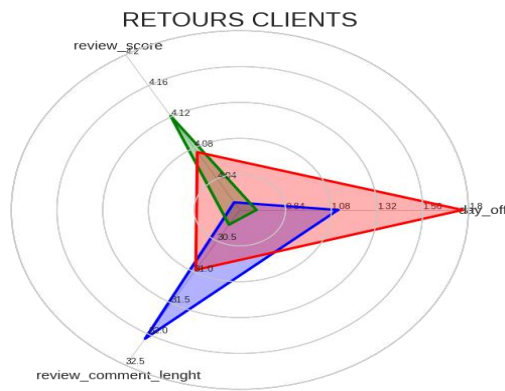
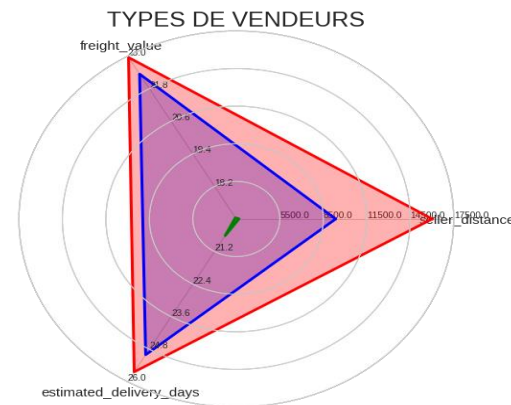
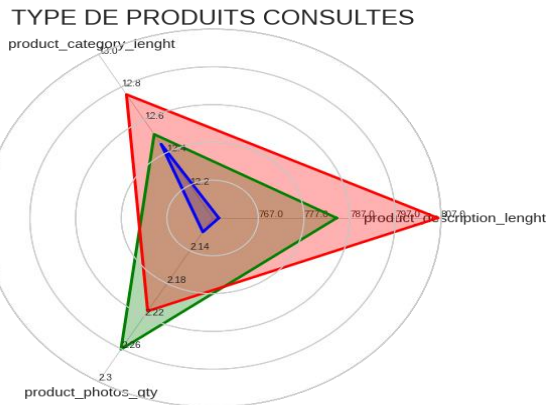
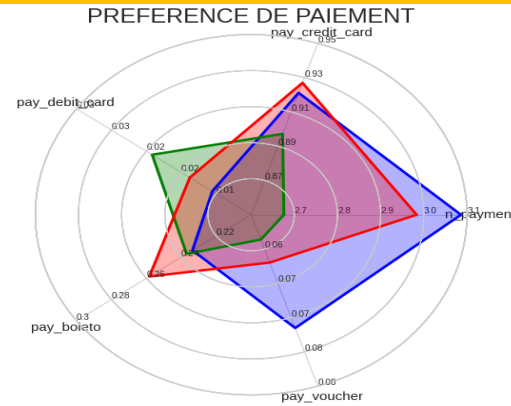
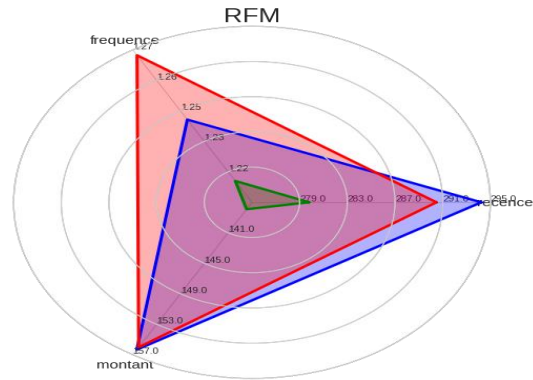
A vertical line is positioned to the right of the text. In the bottom right corner of the slide, there is a large yellow right-angled triangle pointing towards the top-left.

Choix de la Segmentation



⇒ Segmentation globale avec K-MEANS

Segmentation globale avec K-MEANS



PETITS CLIENTS

- **Comportement RFM** : peu d'achats, petits montants, dernière commande pas trop longtemps.
- **Préférence de paiement** : en peu de fois et plutôt par carte de crédit et/ou carte de débit.
- **Type de produits consultés** : Qté d'informations (description, catégorie) moyennes et beaucoup de photos.
- **Type de vendeurs fréquentés** : vendeurs proches, avec très peu ou pas de frais de livraison et peu de temps de livraison.
- **Retour** : bonnes notes, très peu de commentaires.

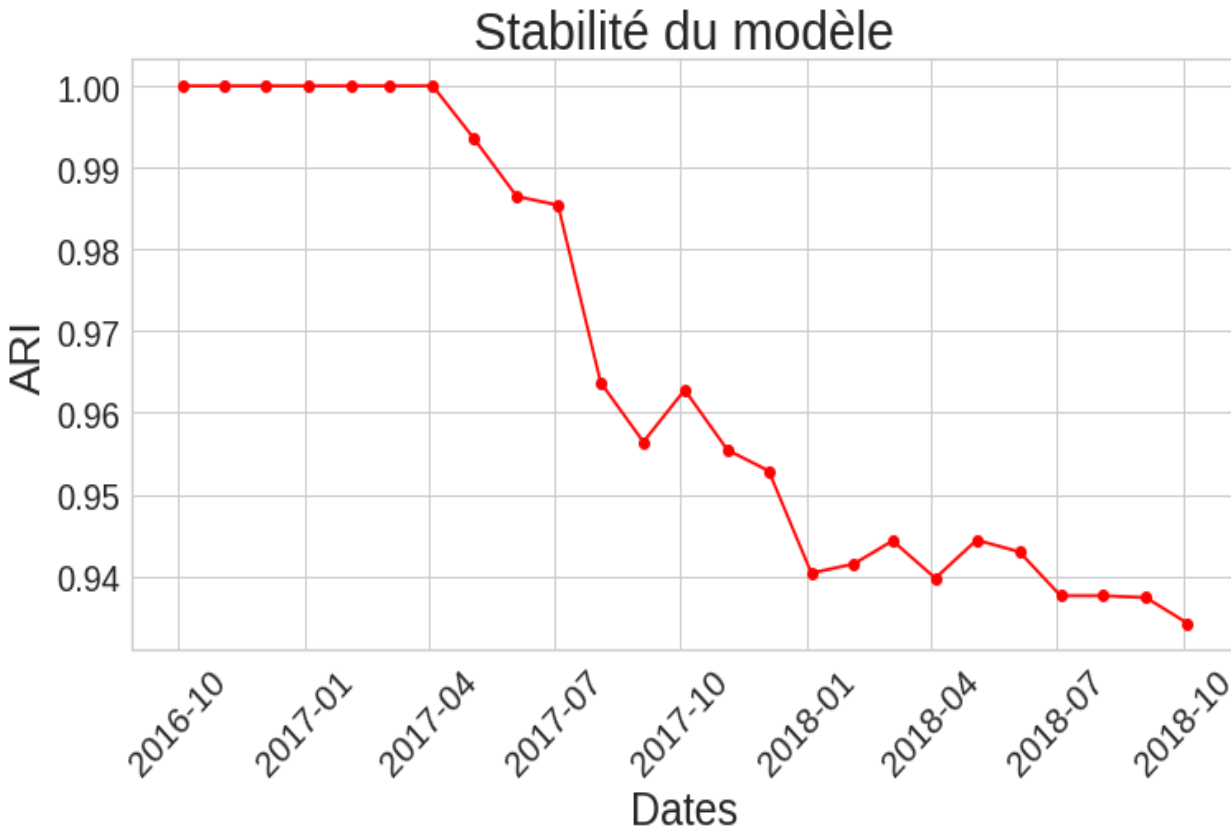
MOYENS CLIENTS

- **Comportement RFM** : achats à fréquence moyenne, montants élevés, dernière commande depuis un grand moment.
- **Préférence de paiement** : en plusieurs fois et par carte de crédit et/ou par voucher.
- **Type de produits consultés** : peu de description, autres d'informations (catégorie, photo) moyennes.
- **Type de vendeurs fréquentés** : vendeurs moyennement loins, avec des frais et temps de livraison assez chers.
- **Retour** : notes moyennes, beaucoup de commentaires.

GROS CLIENTS

- **Comportement RFM** : achats très fréquents, montants élevés, dernière commande depuis un certain temps.
- **Préférence de paiement** : versements moins nombreux et paiement par carte de crédit et/ou par billets.
- **Type de produits consultés** : bonne description, catégories bien renseignés, avec Qté moyenne de photos.
- **Type de vendeurs fréquentés** : vendeurs plus loins, donc avec des frais et temps de livraison plus conséquents.
- **Retour** : notes moyennes, peu de commentaires.

Maintenance du modèle



- Modèle de base : K-MEANS (n=3) fitté sur des données de référence (6 premiers mois)
- Ajout de données par mois
- Modèle : K-MEANS (n=3) fitté sur les nouvelles données
- Evaluation :
⇒ ARI (groupes base % groupes nouv.)

⇒ Modèle stable sur la durée des données

⇒ Tendance négative su fur et à mesure du temps

Améliorations

- Essais d'autres features (selon les besoins, en concertation avec eq. marketing)
- Essais de segmentation avec d'autres features (*ex* : RFM + notes)
- Essais d'autres modèles de segmentation
- Essais d'autres évaluateurs (*ex* : Calinski_Harabasz score)
- Correction aux cas où modèle n'est plus stable :
 - ⇒ vérifier minutieusement les données à partir de la date,
 - ⇒ réaliser un autre modèle à partir de cette date.

MERCI

Questions



Incompréhensions

