

# Projet n° 4

Parcours Ingénieur  
Machine Learning

Présenté par  
Alfred Bazin

Mentor  
Amine Hadj-Youcef

25/09/2022

## Segmentez des clients d'un site e-commerce

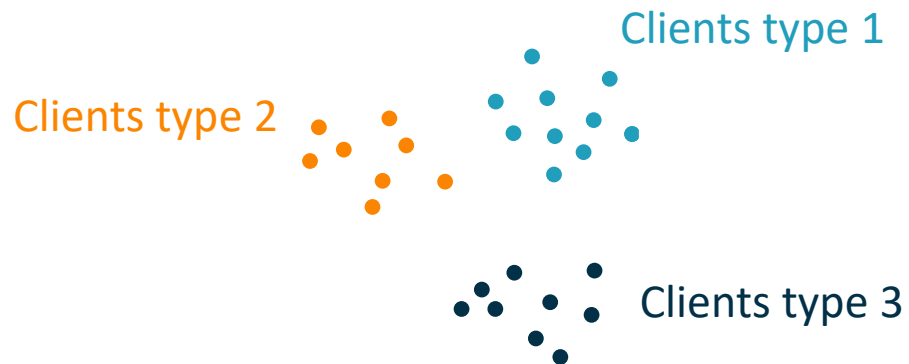
# Sommaire

- Présentation du projet
- Cleaning & Feature engineering
- Modélisation
- Simulation
- Conclusion

# Présentation du projet

« **Segmentation** des clients utilisable au quotidien pour les campagnes de communication »

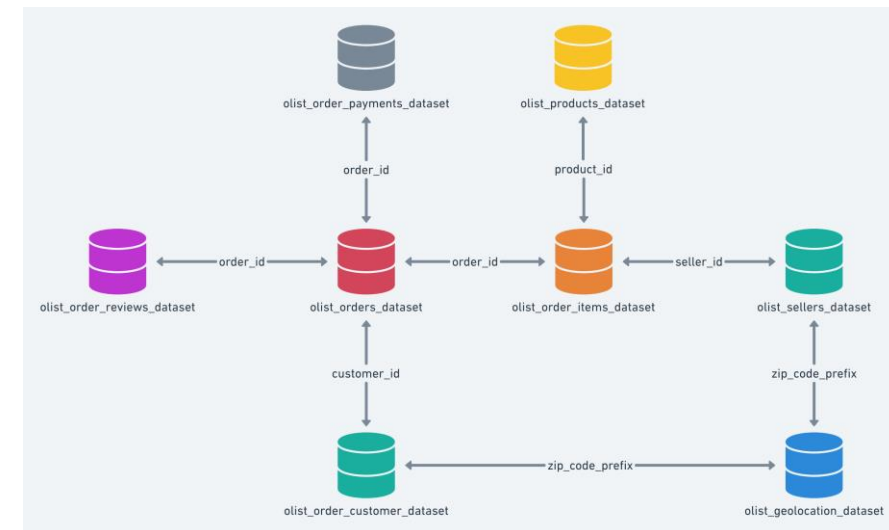
**olist**



Données sur Kaggle :

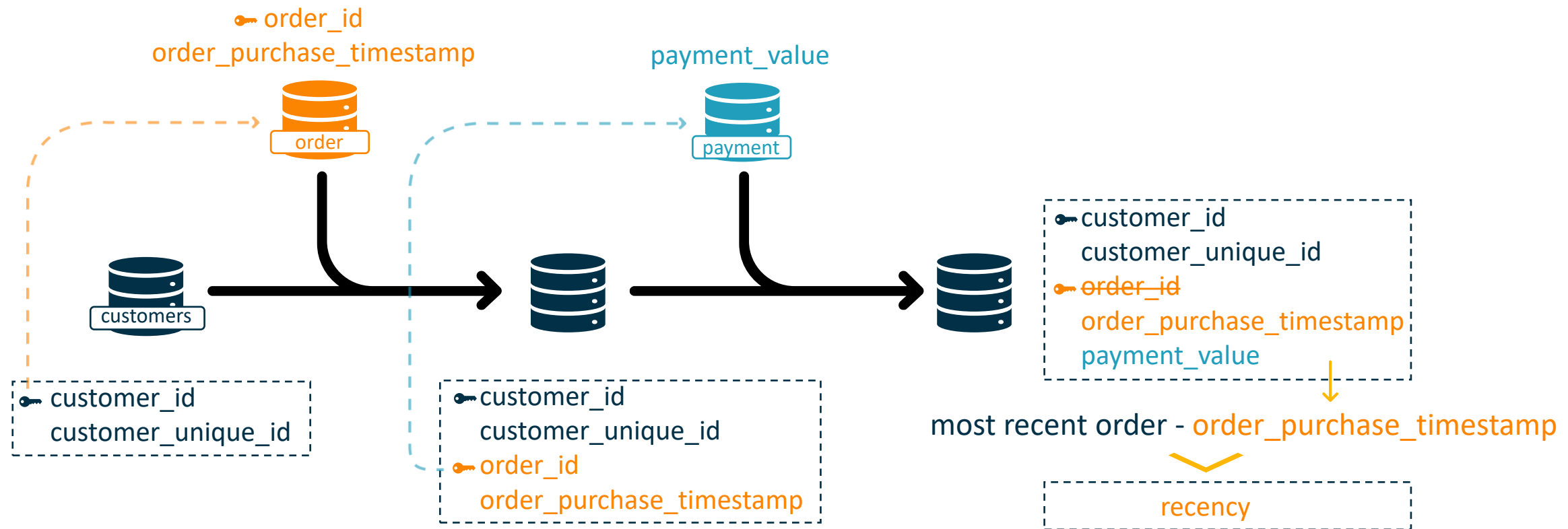
<https://www.kaggle.com/datasets/olistbr/brazilian-ecommerce>

- **8** bases de données
- **99441** commandes
- **96096** utilisateurs



# Nettoyage & Feature Engineering

Agrégation des bases de données (exemple RFM) :



# Nettoyage & Feature Engineering

Agrégation des bases de données (exemple RFM) :

Groupés par :

customer\_unique\_id

customer_id	count	Frequency
recency	min	Recency
payment_value	sum	Monetary Value

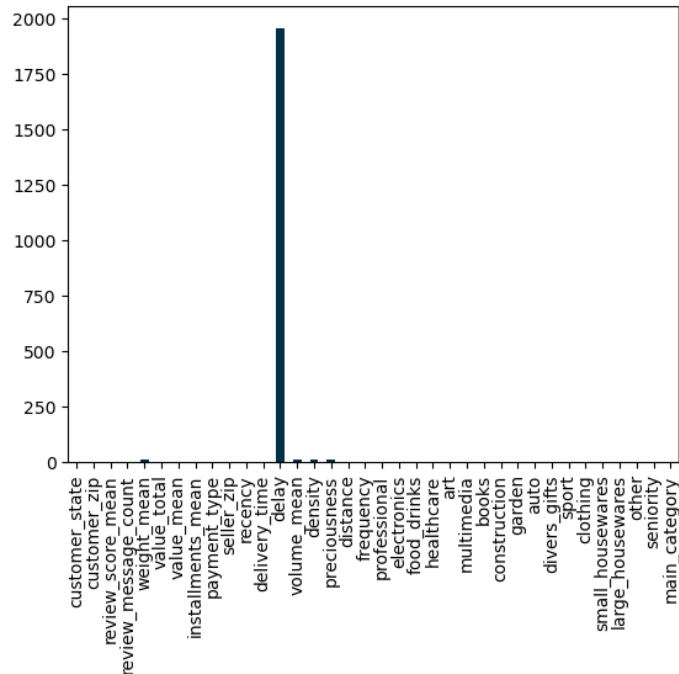
Autres features :

review\_score (mean)  
review\_message (count)  
Weight (mean)  
Volume (mean)  
Density (mean)  
Preciousness (mean)  
delivery\_time (mean)  
delay (mean)  
Distance (mean)  
Value (mean)  
Installments (mean)  
payment\_type (mode)  
customer\_state (mode)  
main\_category (Encoding...)

# Nettoyage & Feature Engineering

## Valeurs manquantes et aberrantes

### Valeurs manquantes



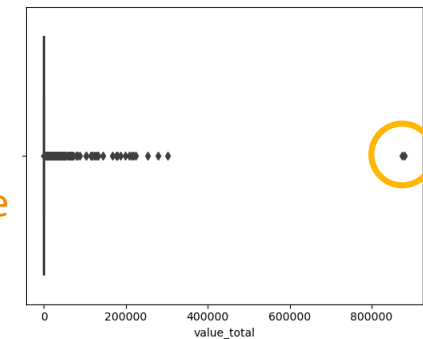
Weight  
Volume  
Density  
preciousness

} Supprimées

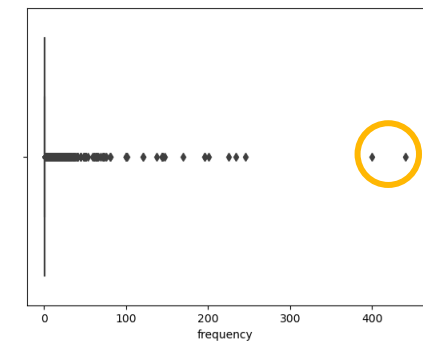
Delay → Imputés avec  
delivery time

### Valeurs aberrantes et atypiques

Valeur  
monétaire totale



Fréquence

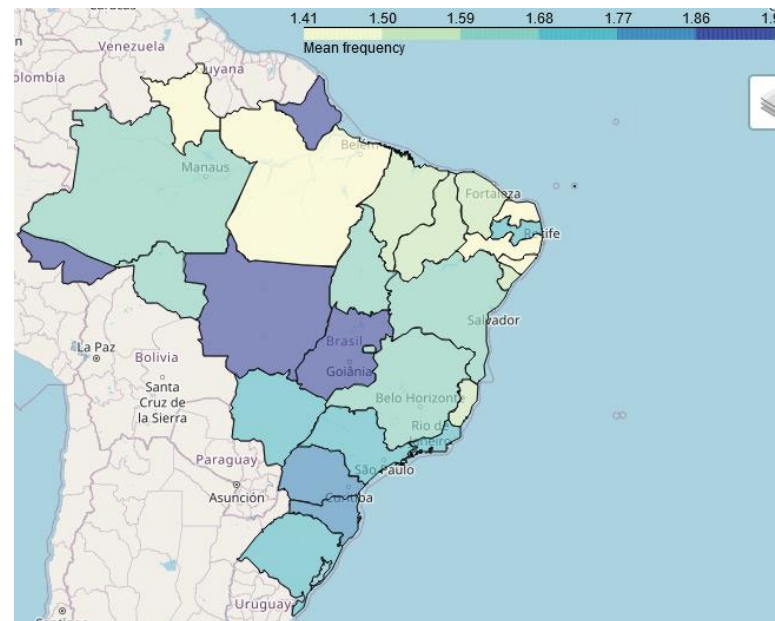


# Nettoyage & Feature Engineering

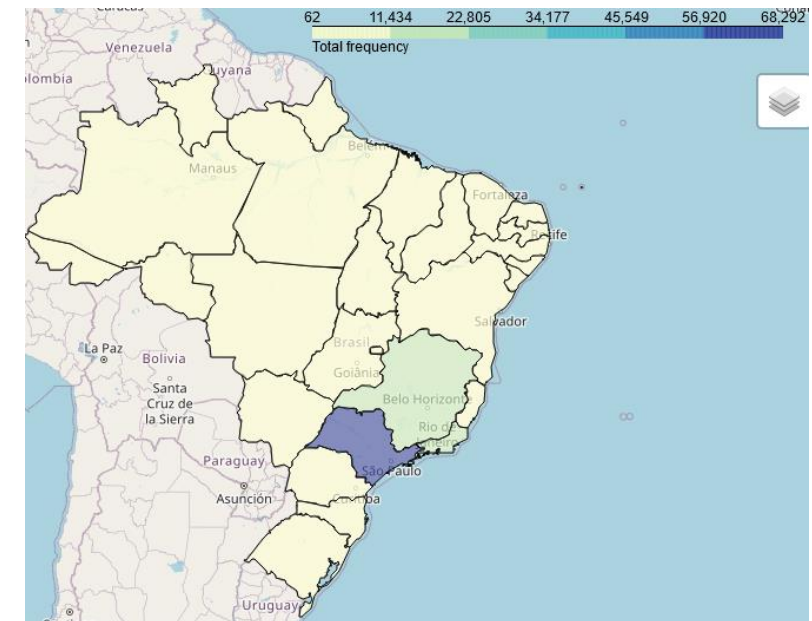
## Analyse exploratoire

- La **fréquence moyenne** est équilibrée entre les états
- Les **états** plus **peuplés** représentent une **large majorité des commandes**

Fréquence **moyenne** par état

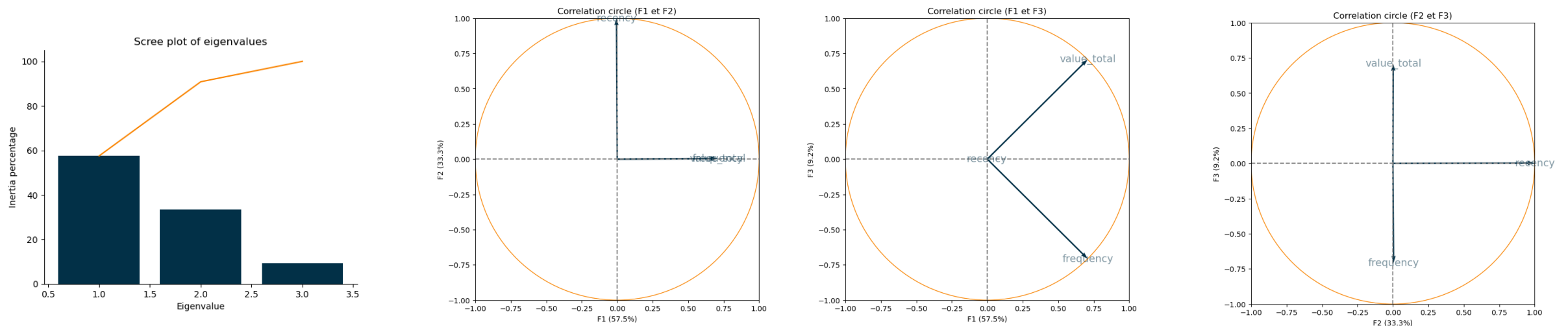


Fréquence **totale** par état



# Nettoyage & Feature Engineering

## Analyse exploratoire



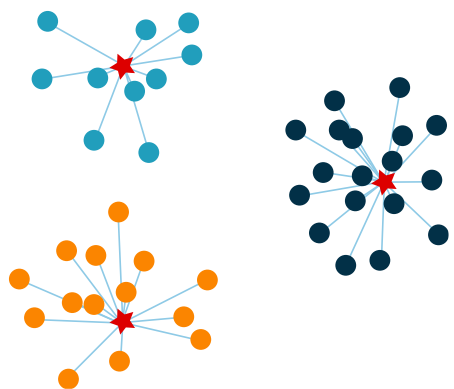
- F1 : Quantité d'achats (*value\_total* + *frequency*)
- F2 : La récence des achats (*recency*)
- F3 : Le prix des achats (*frequency* vs *value\_total*)



# Modélisation

## RFM - KMeans

( Sur un échantillon de  
~10 000 individus )

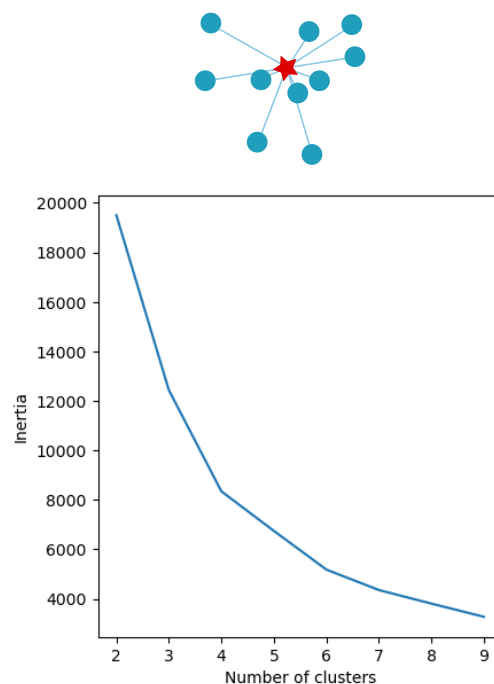


### Paramètres :

- n\_clusters
- init
- n\_init

## Inertie

Somme de distance des points au centroïde :

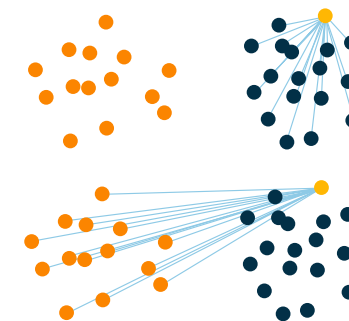


$$\sum_{i=0}^n \min_{\mu_j \in C} (\|x_i - \mu_j\|^2)$$

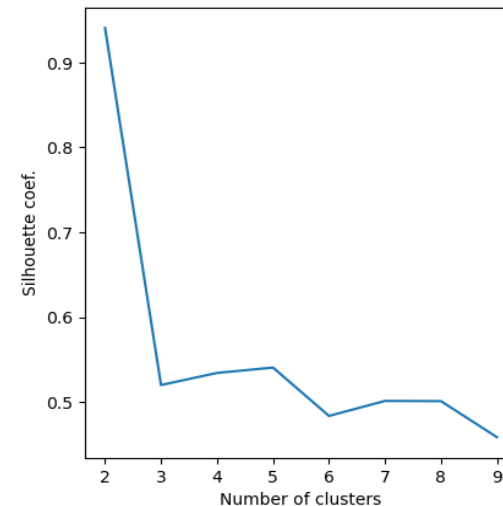
## Métriques :

### Coefficient de silhouette

- a : Distance moyenne entre un **point** et ses **voisins** du **même cluster**
- b : Distance moyenne entre un **point** et ses **voisins** du **cluster le plus proche**



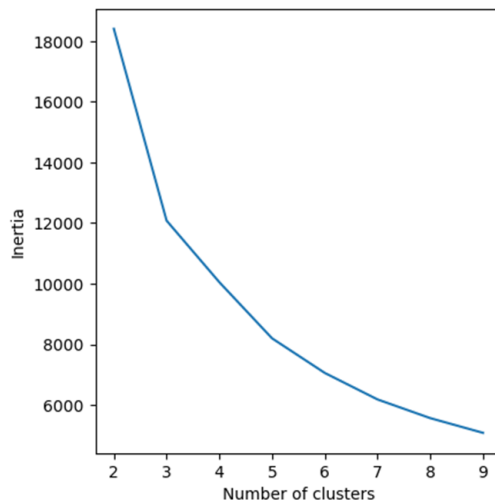
$$s = \frac{b - a}{\max(a, b)}$$



- 1** : Parfait
- 0** : Les clusters se recouvrent
- 1** : Mauvais clustering

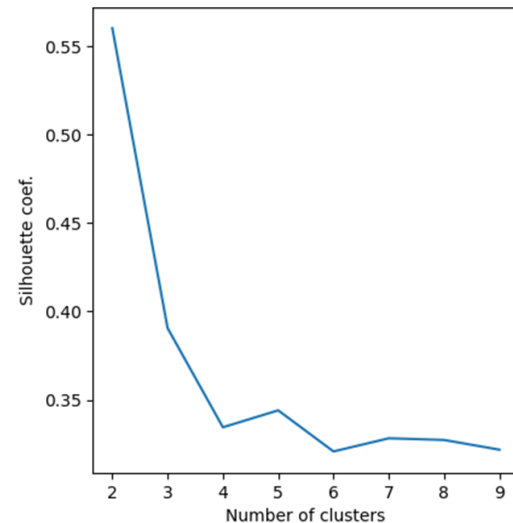
# Modélisation

## RFM (log) - KMeans



Inertie

- Coude entre **3** et **5** clusters

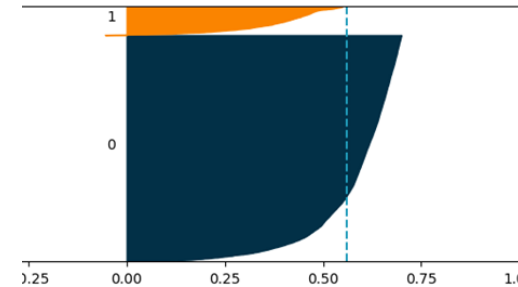


Coef. de silhouette

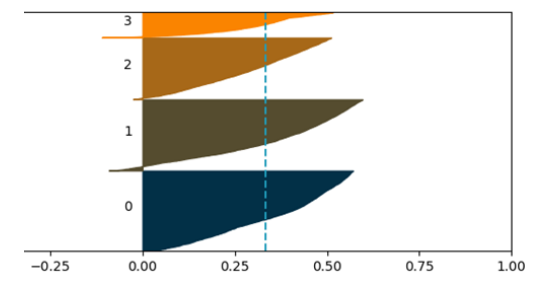
- Max à **2** clusters
- Chute avec un max à **5** clusters

## Profil de silhouette

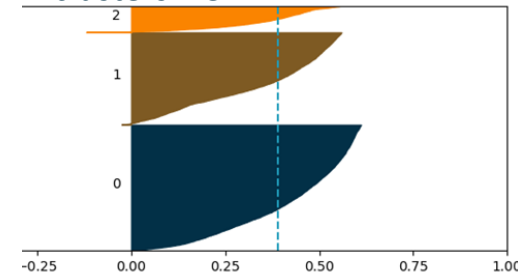
n clusters = 2



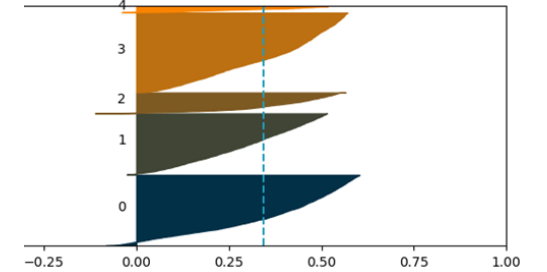
n clusters = 4



n clusters = 3



n clusters = 5



Profil de silhouette

- Clusters plus **équilibrés**
- Toujours **chevauchement**



Résultats **similaires** avec **ACP**

# Modélisation

RFM (log) - KMeans

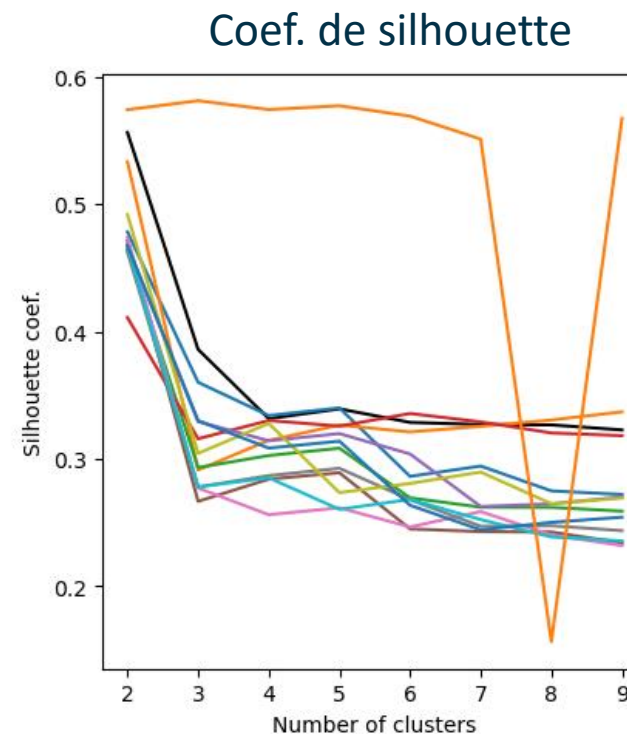
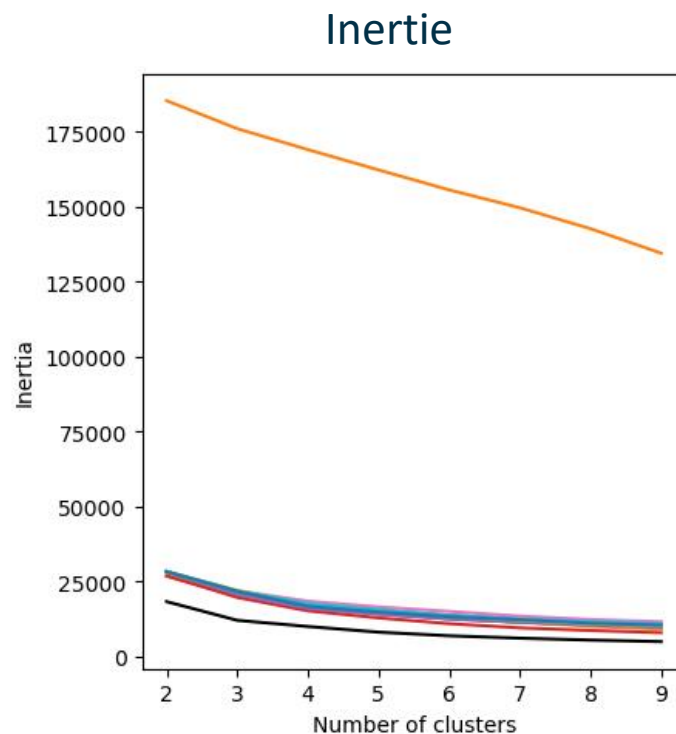
Ajout de features

- Ajout de features à RFM :

Inertie ↗ Silhouette ↘

- Exceptions :
  - review\_score\_mean
  - review\_message\_count
  - value\_mean
  - Cat. Features

↘ Explosion de l'inertie



Features

- RFM
- review\_score\_mean
- review\_message\_count
- weight\_mean
- value\_mean
- installments\_mean
- delivery\_time
- delay
- volume\_mean
- density
- preciousness
- distance
- Cat. Features

# Modélisation

RFM (log) - KMeans

Ajout de features catégoriques

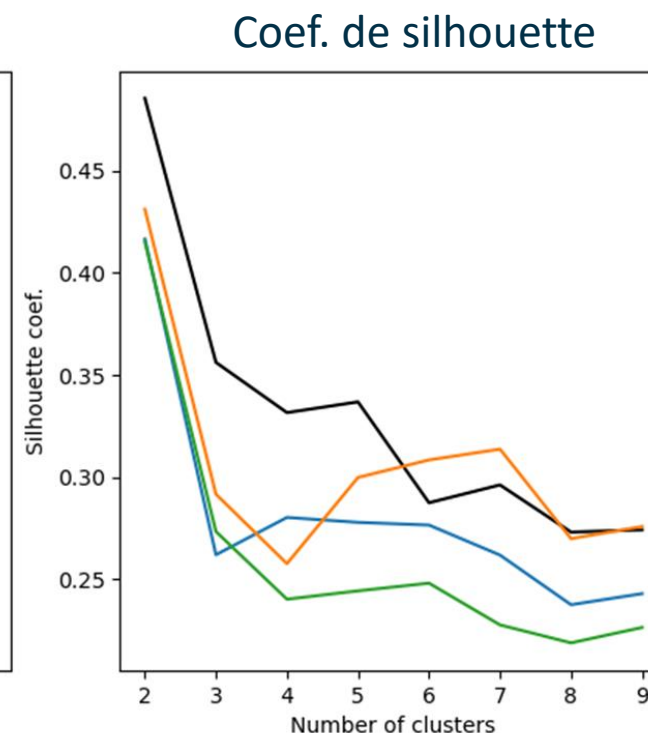
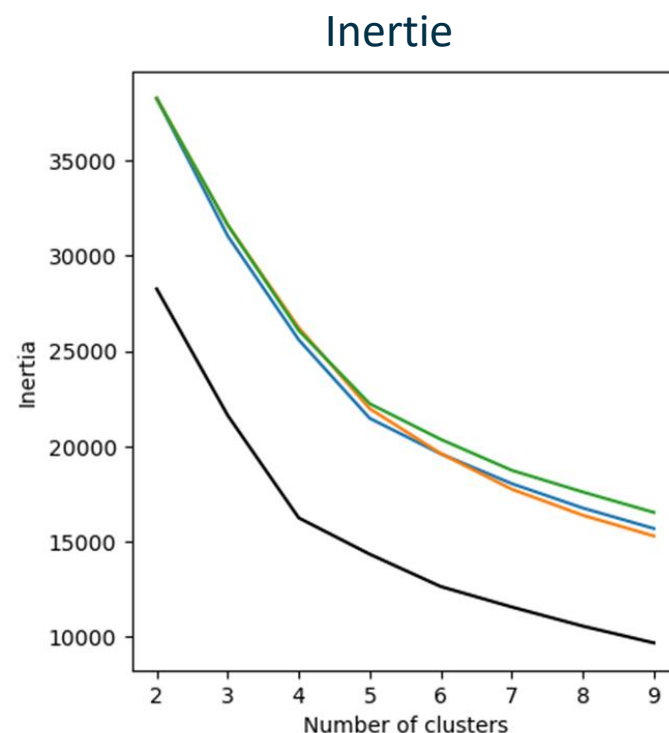
- Ajout de features à RFM :

Inertie ↗ Silhouette ↘

- Exceptions :

- Payment\_type

↘ Haut nombre de  
clusters avec  
recouvrement



Features



# Modélisation

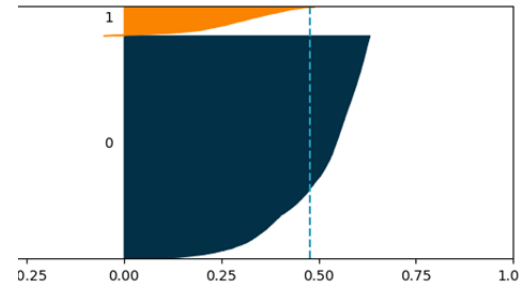
RFM+ (log) - KMeans

Features :

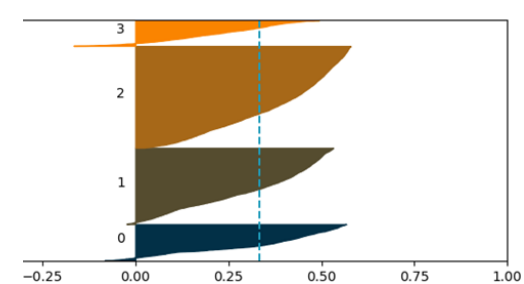
- recency
- frequency
- monetary value
- review\_score\_mean

Profil de silhouette

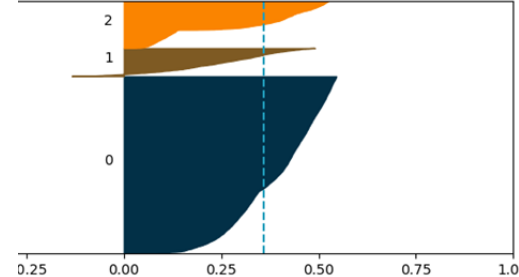
n clusters = 2



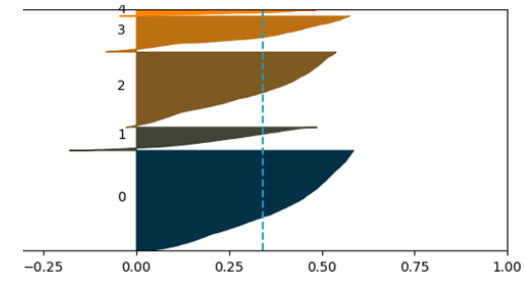
n clusters = 4



n clusters = 3



n clusters = 5



Profil de silhouette

- Toujours **chevauchement**



Résultats **similaires** avec ACP

# Modélisation

RFM+ (log) - DBSCAN

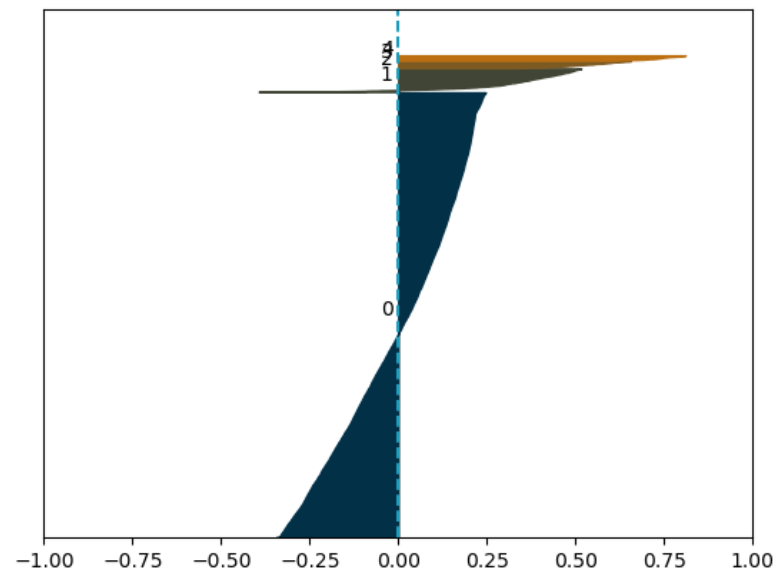


Paramètres :

- eps
- min\_sample

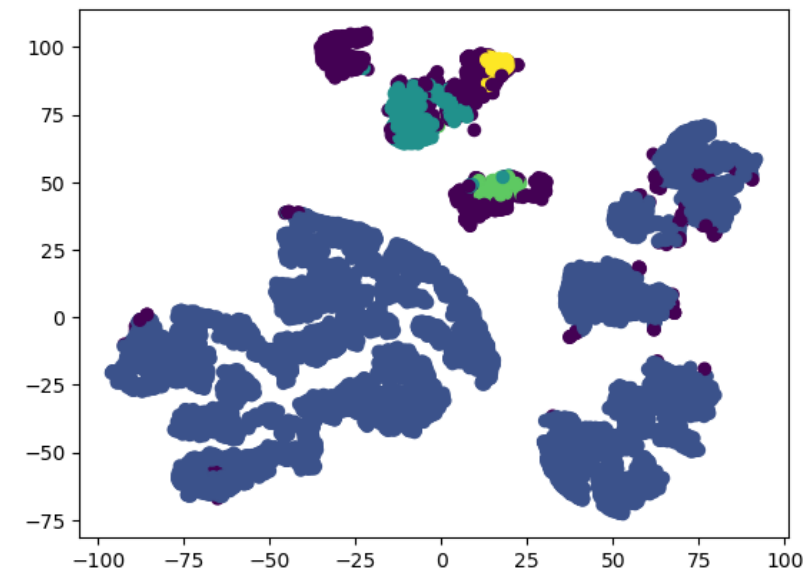
Profil de silhouette

n\_clusters = 5



- Obtention de clusters **déséquilibrés** et se **chevauchant** fortement

t-SNE



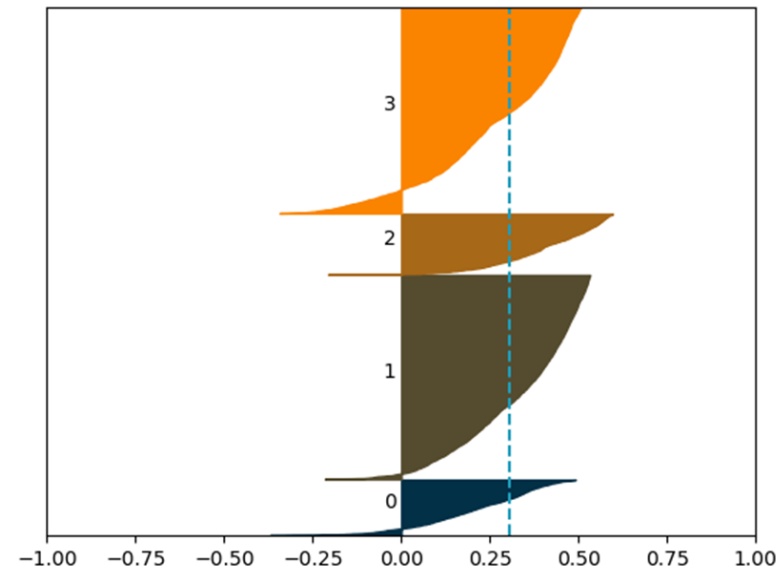
# Modélisation

RFM+ (log) – Clustering hiérarchique

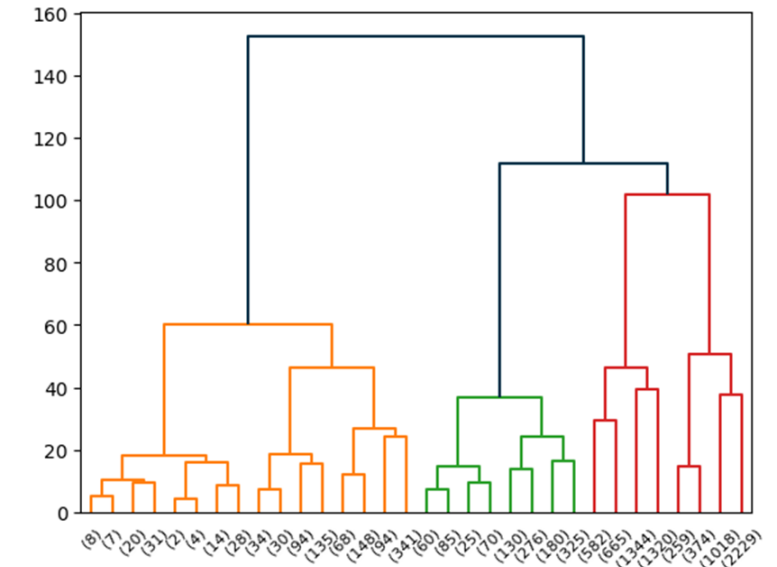


Profil de silhouette

n\_clusters = 4



Dendrogramme



- Résultats très proches du **KMeans** mais avec plus de **chevauchement**

# Modélisation

RFM+ (log) – KMeans

Interprétation des clusters (n\_clusters = 4)

## Cluster 1:

Clients qui ne sont **pas satisfaits** de leur commande.

## Cluster 2 :

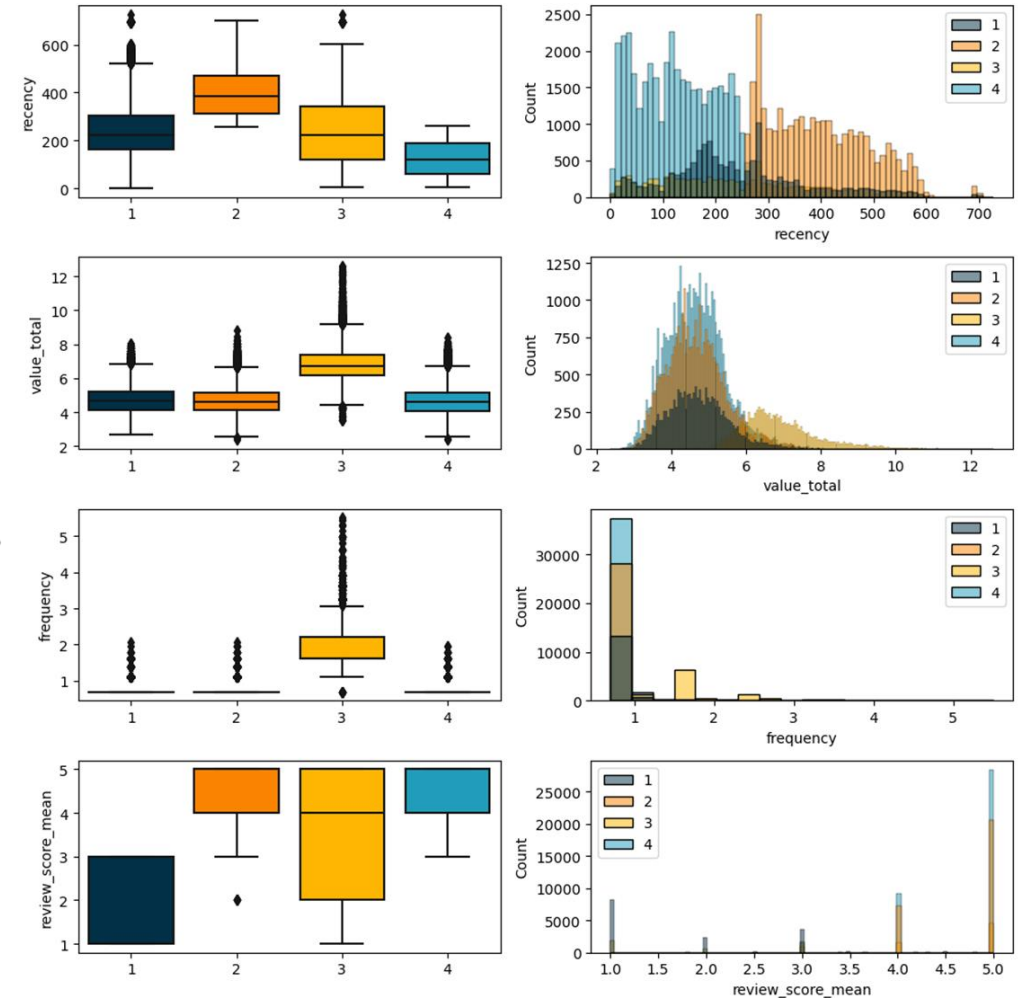
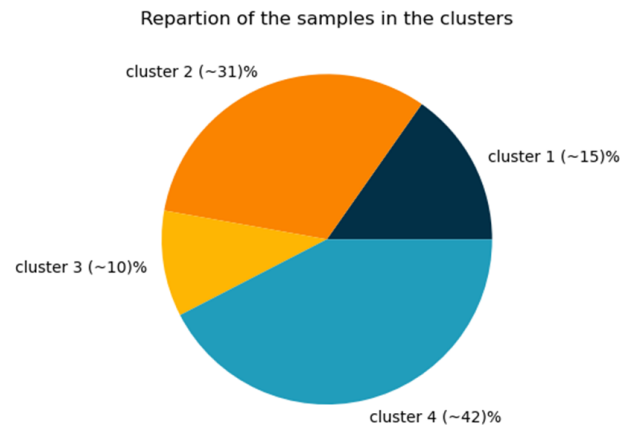
Clients qui ont commandé peu de fois **il y a longtemps**.

## Cluster 3 :

Clients qui ont commandés **plusieurs fois** sur la plateforme.

## Cluster 4 :

Clients **récents** qui n'ont pas commandé beaucoup et sont content.

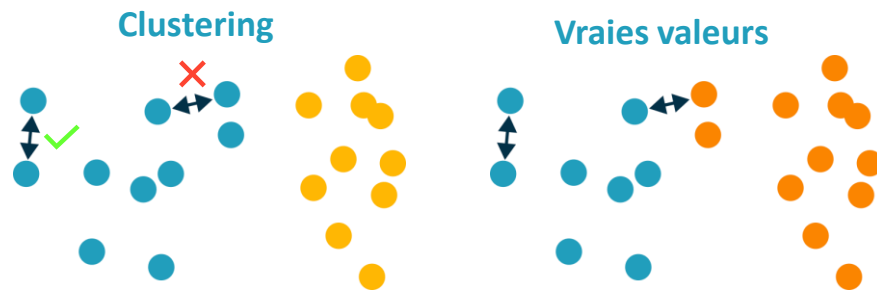




# Simulation

ARI

Index de rand (RI)



Index de rand ajusté (ARI)

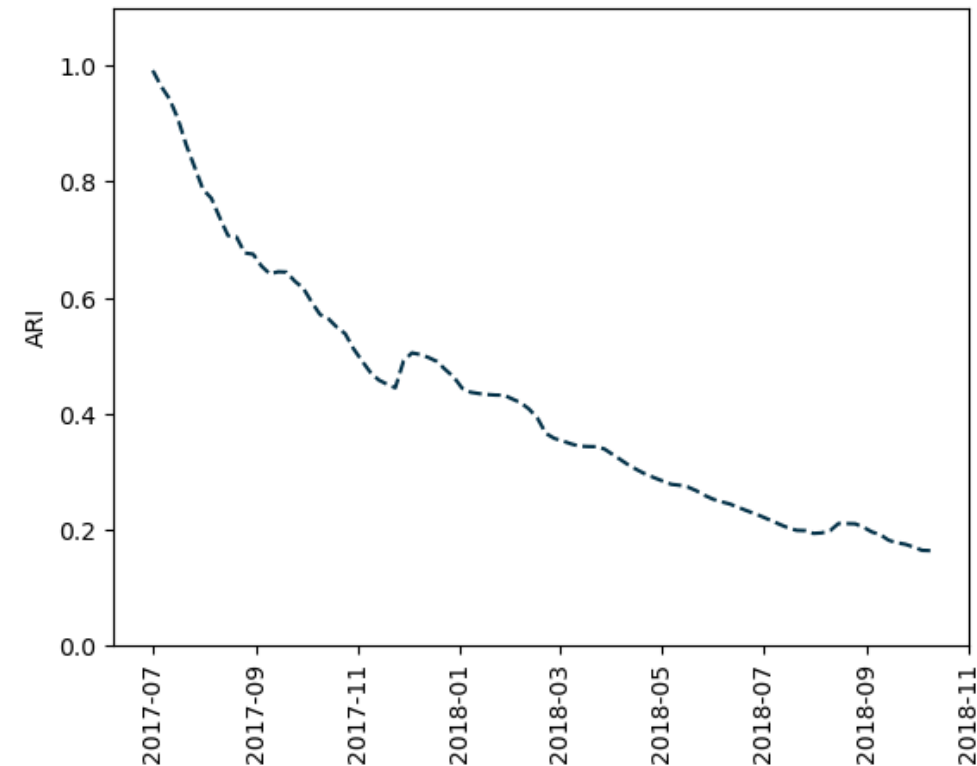
$$ARI = \frac{RI - \mathbb{E}(RI)}{1 - \mathbb{E}(RI)}$$

# Simulation

(début après un semestre de données)

- Sans maintenance
  - Diminution jusqu'à un score de  $\sim 0,2$

ARI sans maintenance

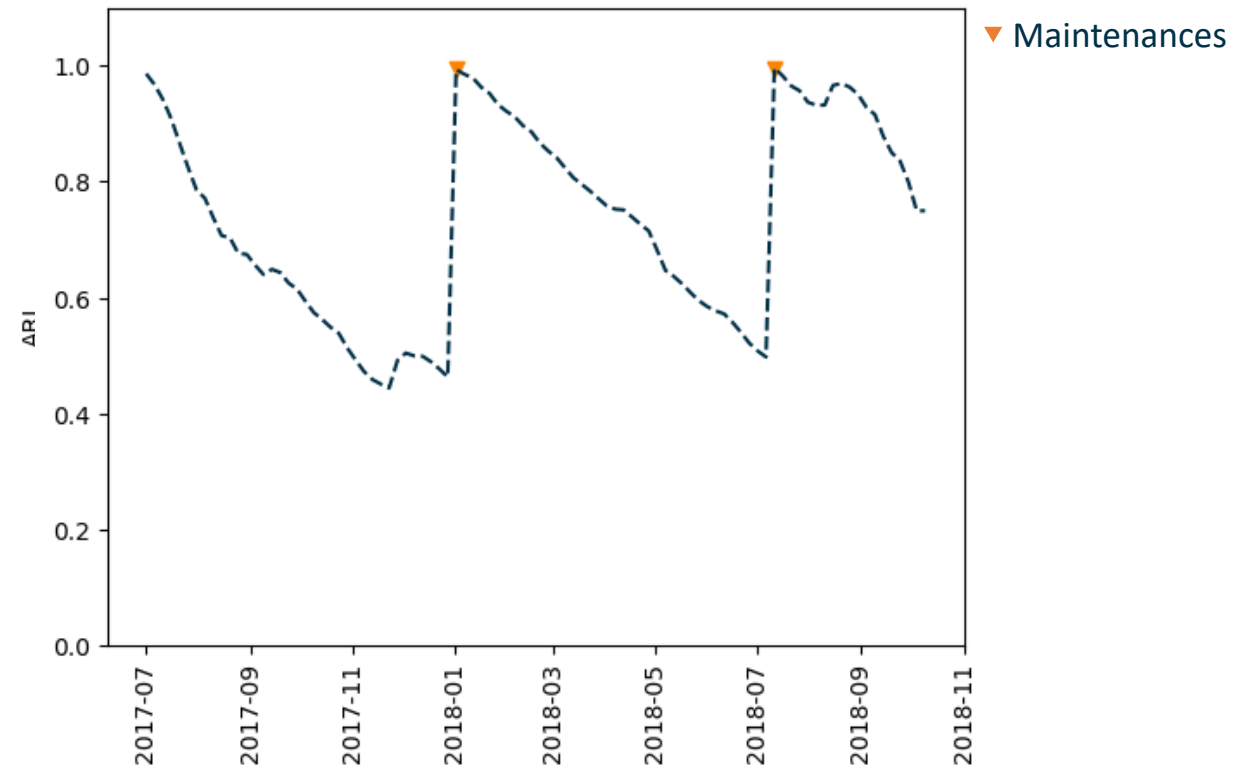


# Simulation

(début après un semestre de données)

- Sans maintenance
  - Diminution jusqu'à un score de  $\sim 0,2$
- 6 mois
  - Diminution jusqu'à un score de  $\sim 0,4$

Maintenance tous les 6 mois

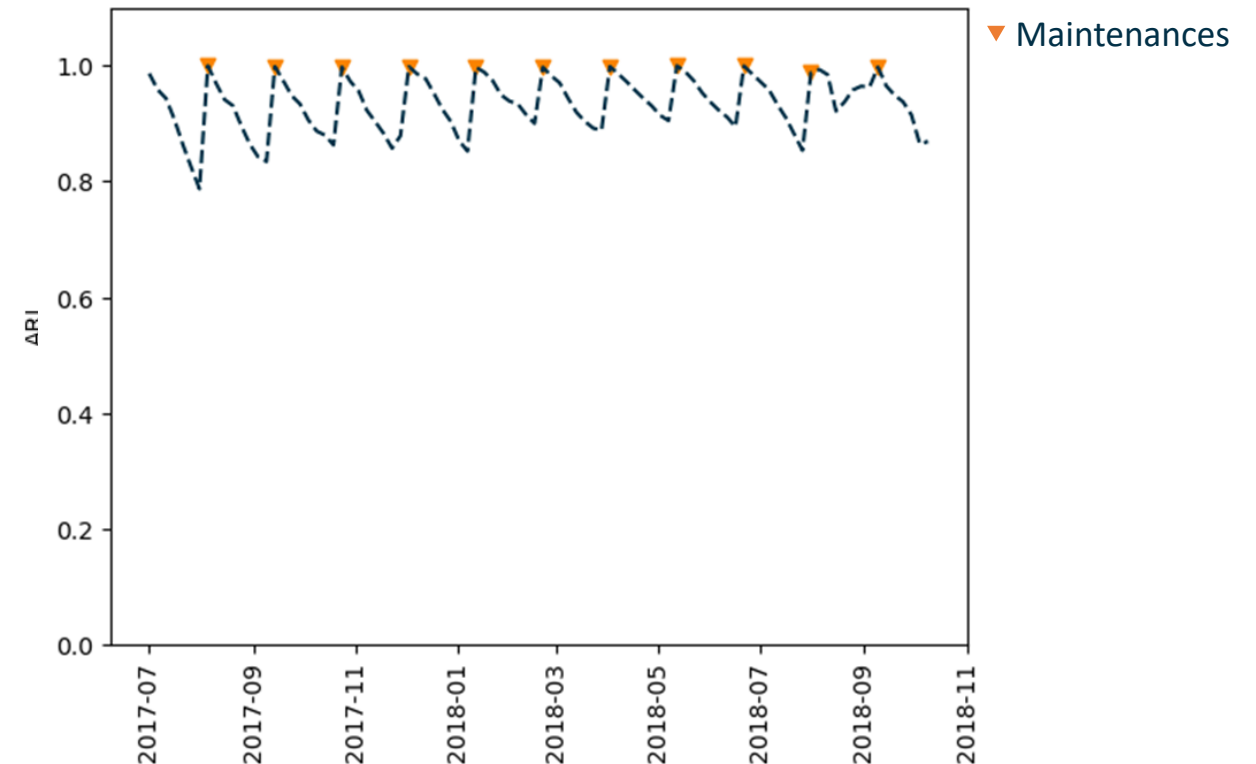


# Simulation

(début après un semestre de données)

- Sans maintenance
  - Diminution jusqu'à un score de  $\sim 0,2$
- 6 mois
  - Diminution jusqu'à un score de  $\sim 0,4$
- 1 mois
  - Diminution jusqu'à un score de  $\sim 0,8$

Maintenance tous les mois

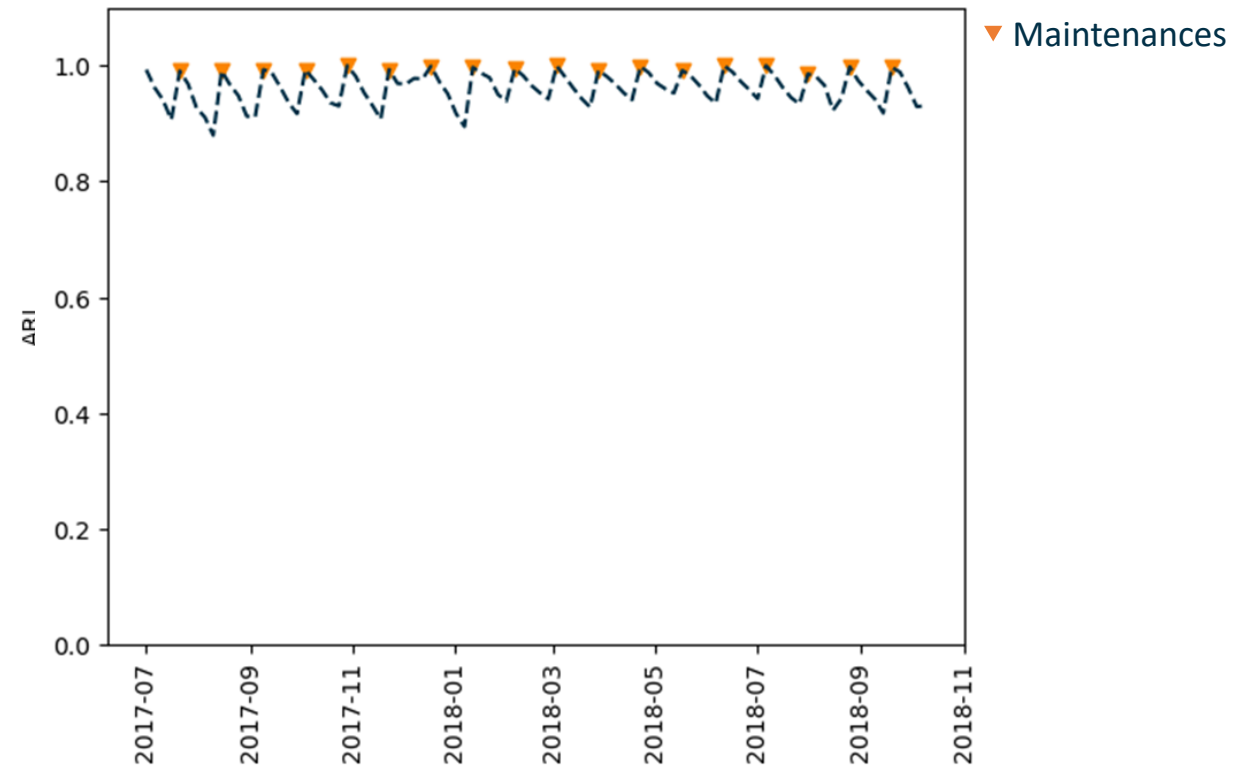


# Simulation

(début après un semestre de données)

- Sans maintenance
  - Diminution jusqu'à un score de  $\sim 0,2$
- 6 mois
  - Diminution jusqu'à un score de  $\sim 0,4$
- 1 mois
  - Diminution jusqu'à un score de  $\sim 0,8$
- 15 jours
  - Diminution jusqu'à un score de  $\sim 0,85$

Maintenance tous les mois



# Conclusion

- Clustering avec features **RFM** et **score moyen** :

Rency      Frequency      Monetary value      Mean review score

- Parmi plusieurs modèles (**DBSCAN**, **Clustering hiérarchique...**), **KMeans** est le **plus performant**.
- **4 Clusters** formés :
  - Clients **récents**
  - Clients qui ne **reviennent pas**
  - Clients **fréquents**
  - Clients **mécontents**
- Maintenance de **15 jours** recommandée

