



P5 : Segmentez des clients d'un site e-commerce

Defense

Formation : Data Scientist

Etudiante : Rocio Isorna

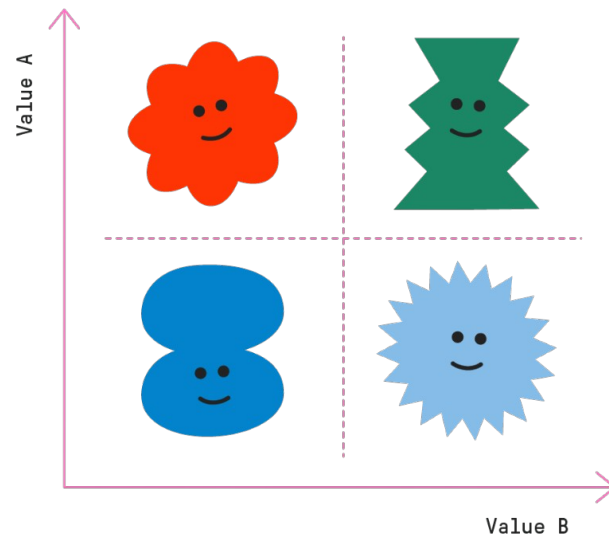
Evaluatrice : Souhail Toumndi

Client segmentation

ALL CUSTOMERS



SEGMENTED CUSTOMERS



@laptrinhx

Client segmentation

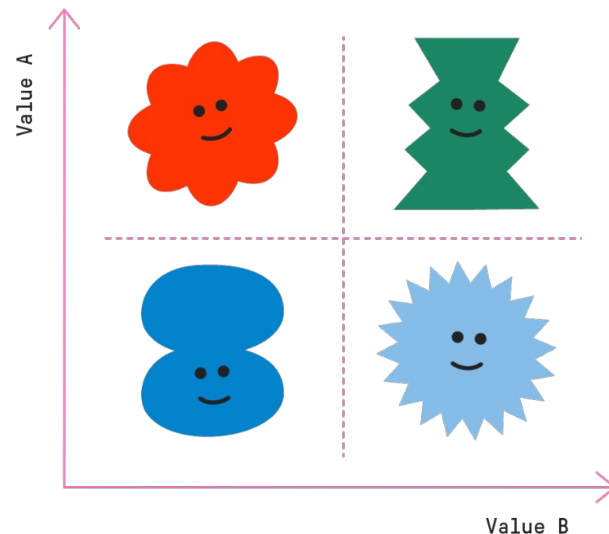


olist

ALL CUSTOMERS



SEGMENTED CUSTOMERS



@laptrinhx

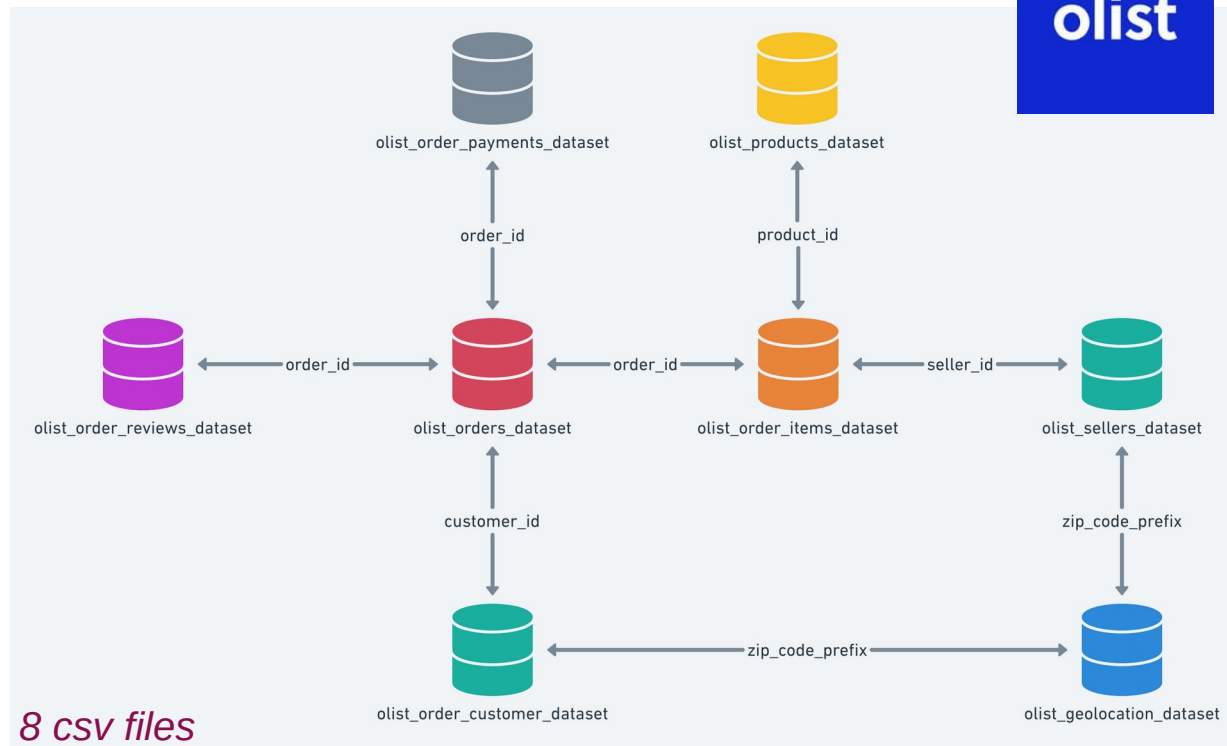
Comprendre les différents types des clients

Fournir à l'équipe marketing une description actionnable
Proposition de contrat de maintenance

Segmentation des clients

Data set : (99441 obs → achat)

- **Clients** : id et zip-code
- **Localisation**: villes
- **Items commandé** : id, item, vendeurs, prix (article et livraison)
- **Commandes** : id, statuts, date, livraison
- **Paielements** : type de paiement et montant
- **Reviews** : retour des clients
- **Vendeurs** : information des vendeurs

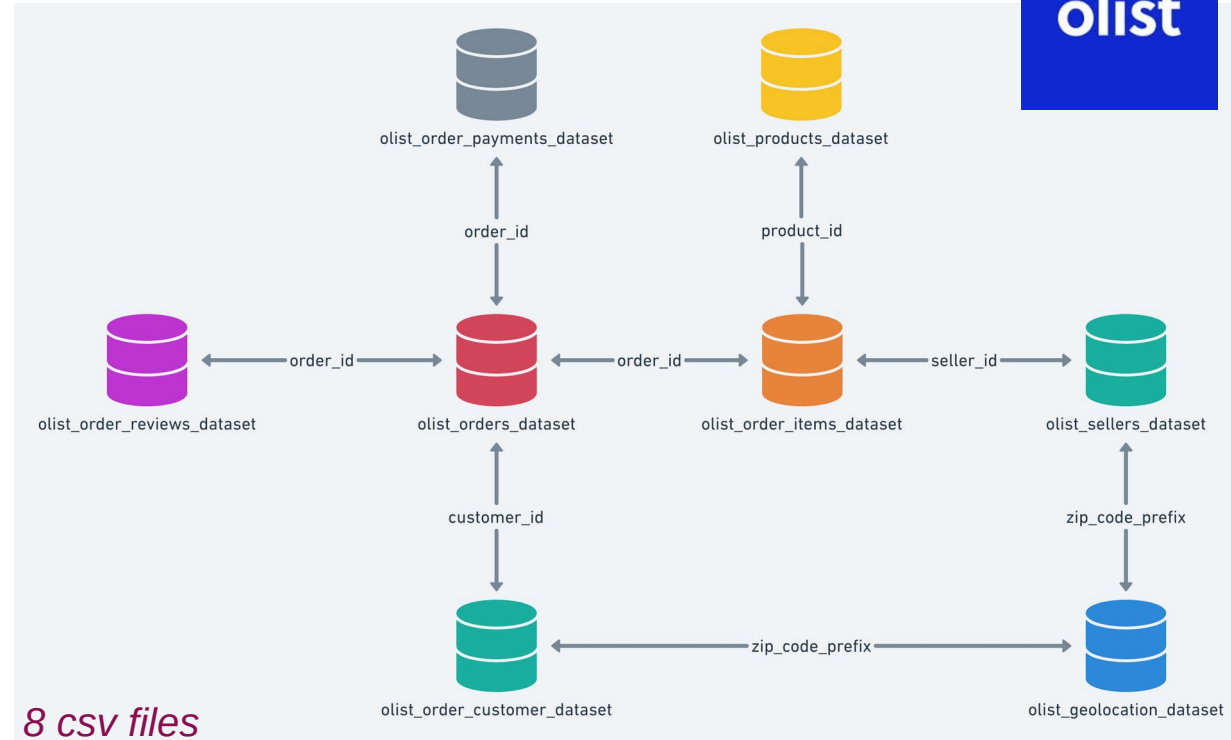


Segmentation des clients



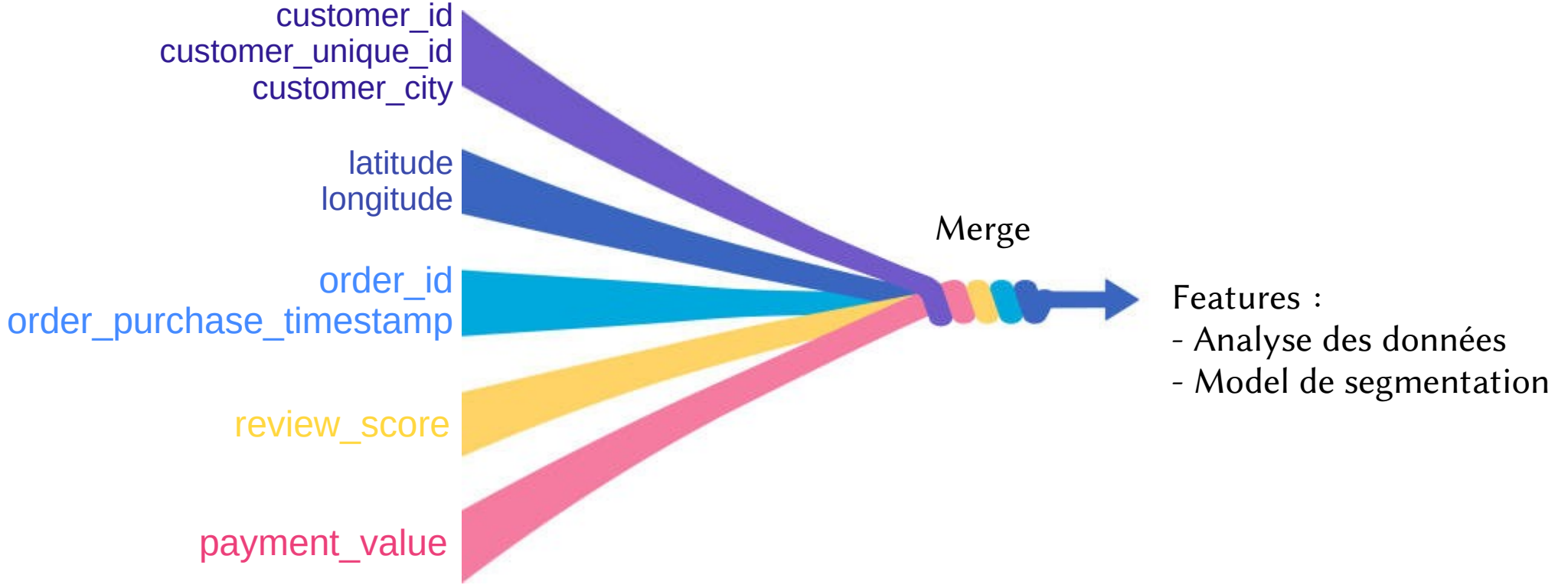
Data set : (99441 obs → achat)

- **Clients** : id et zip-code
- **Localisation** : villes
- **Items commandé** : id, item, vendeurs, prix (article et livraison)
- **Commandes** : id, statuts, date, livraison
- **Paielements** : type de paiement et montant
- **Reviews** : retour des clients
- **Vendeurs** : information des vendeurs

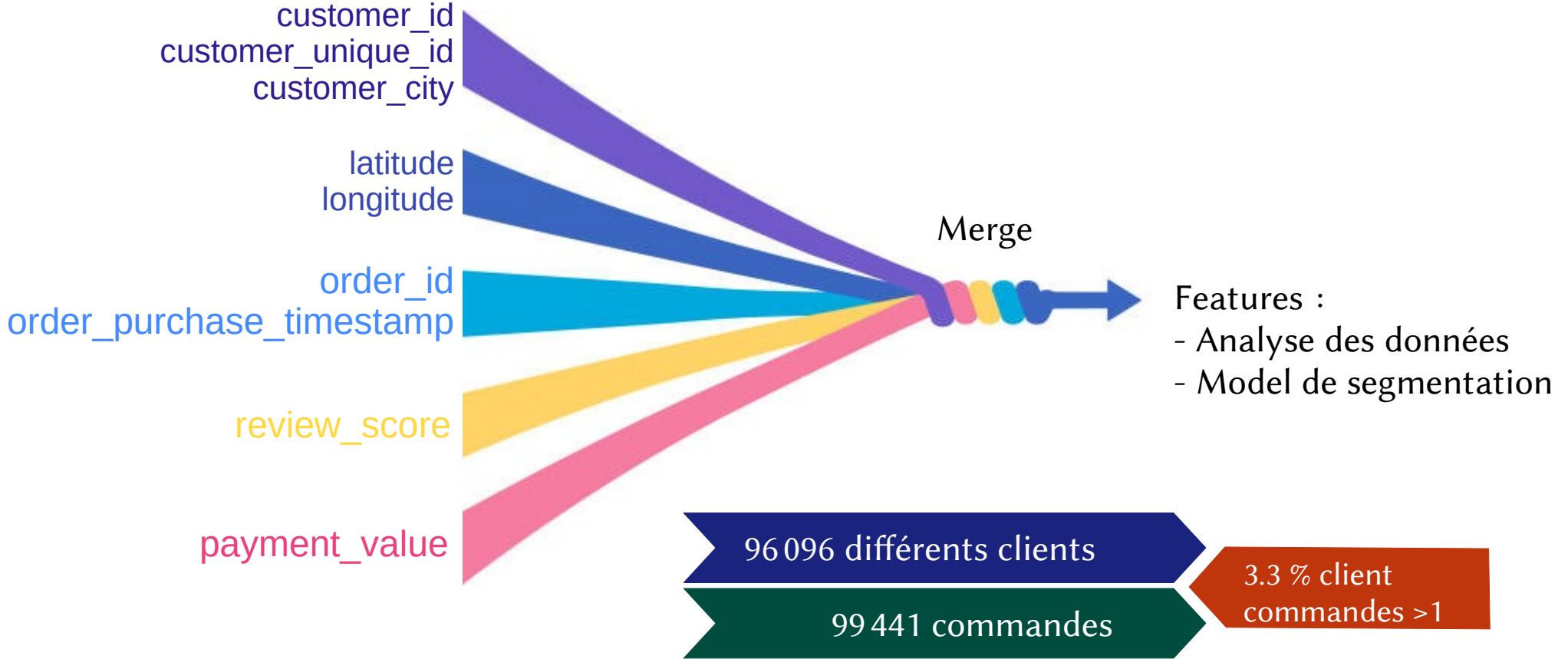


Customer_id vs Customer_unique_id

Sélection des features à partir du Data-set



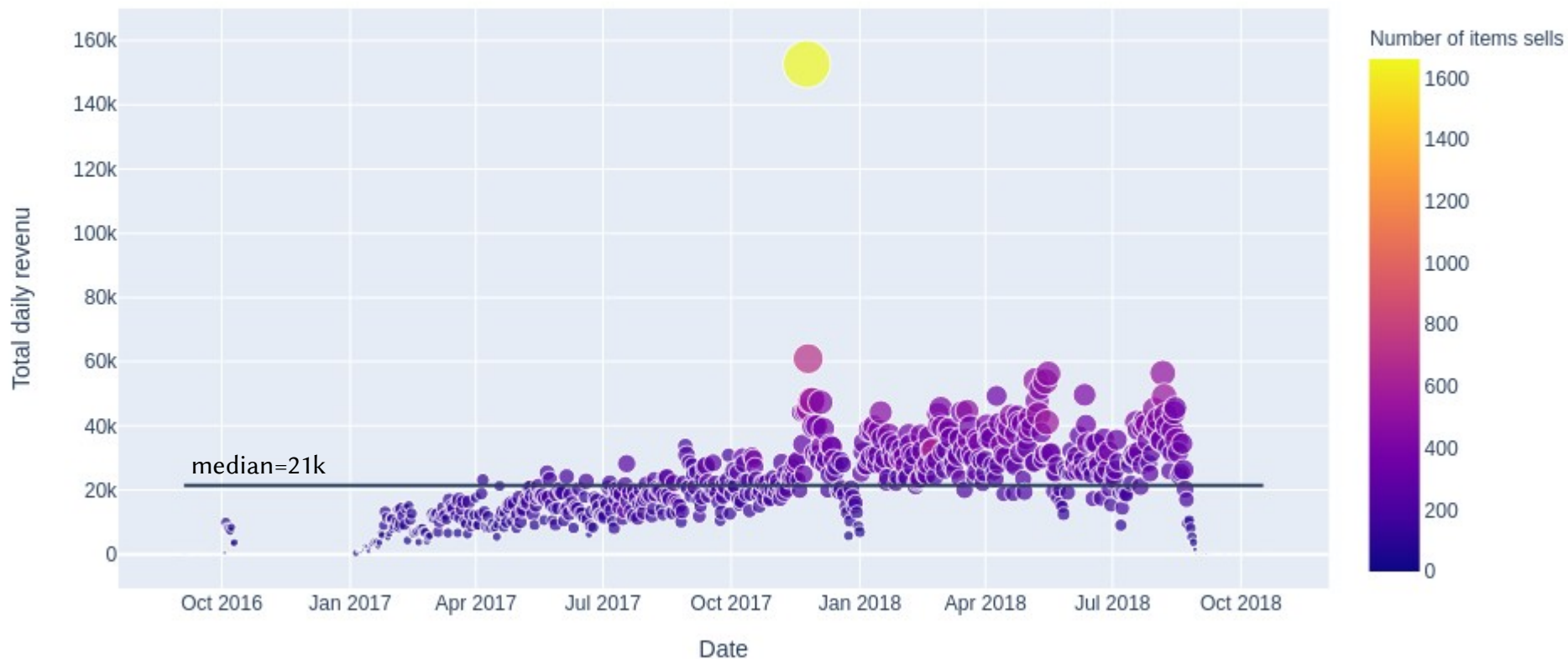
Sélection des features à partir du Data-set





Analyse exploratoire

Sells per day





Analyse exploratoire

Localisation des acheteurs

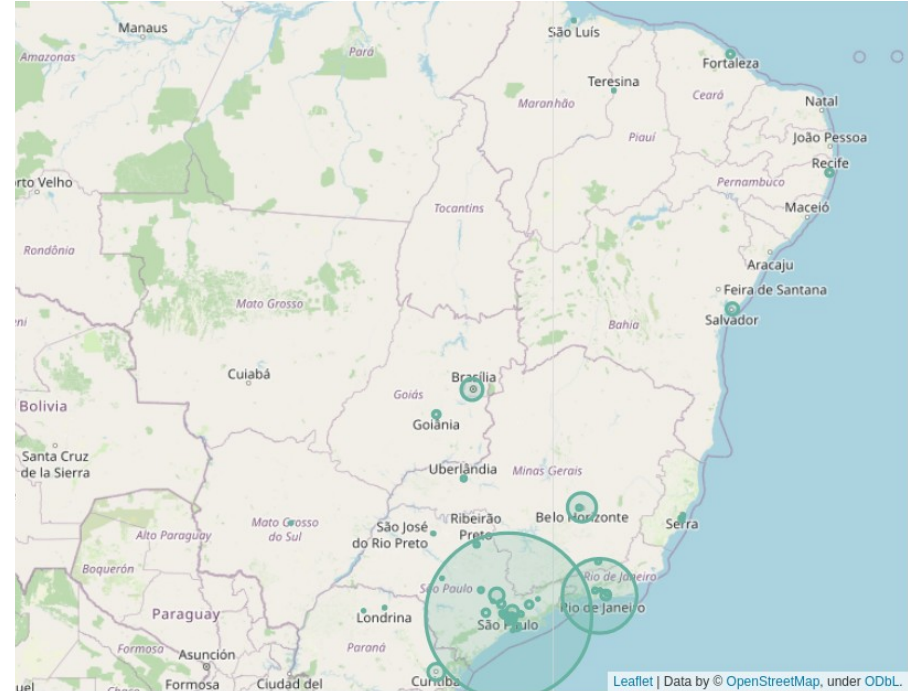


4069 villes

50 % < 3 commandes

06/01/2022 Soutenance P5

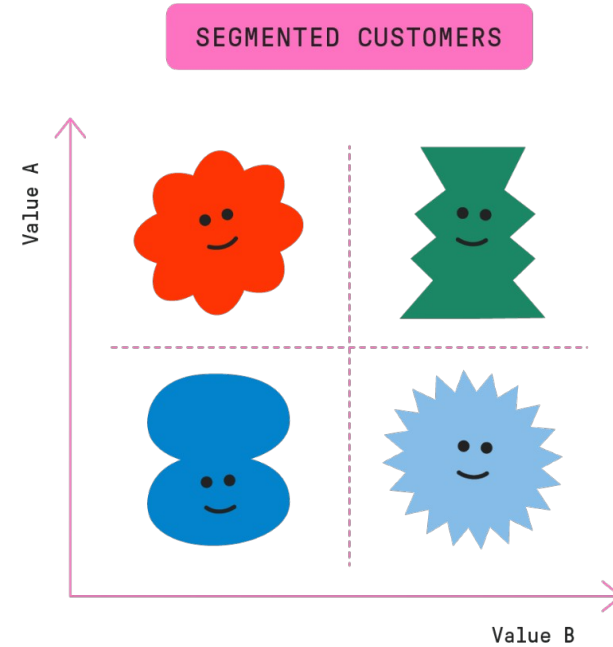
50 principales villes



Sao Paulo → 15k

Rio de Janeiro → 7k

Rocio Isorna



Évaluation des différentes méthodes de clustering

Méthode analytique

R F M Score
e r o
c e n
e q e
n u t
c e a
y n r
c y
y



Méthode analytique

R F M Score
e r o
c e n
e q e
n u t
c e a
y n r
c y
y



Méthodes de apprentissage
automatique **non-supervisées**

- K-Mean
- DBSCAN
- Agglomerative approach

Variables de classification

Méthode analytique

R F M Score
e r o
c e n
e q e
n u t
c e a
y n r
c y
y



Méthodes de apprentissage
automatique **non-supervisées**

- K-Mean
- DBSCAN
- Agglomerative approach

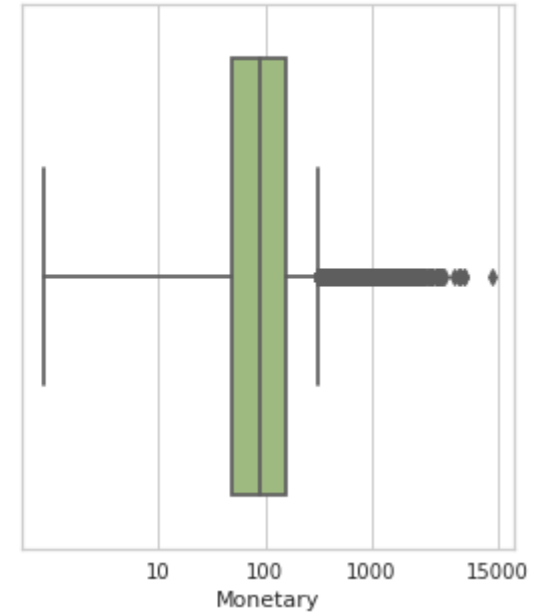
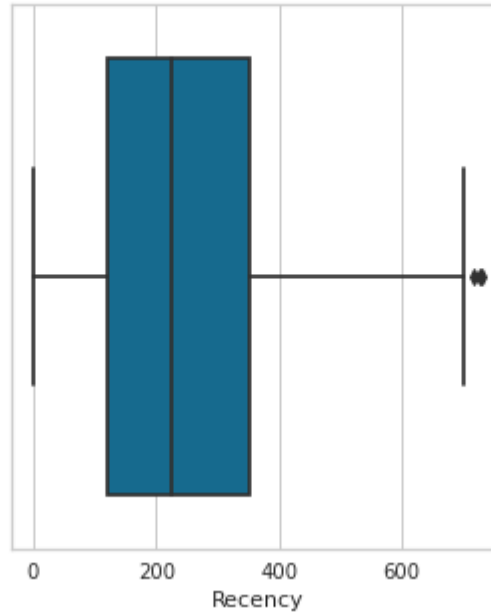
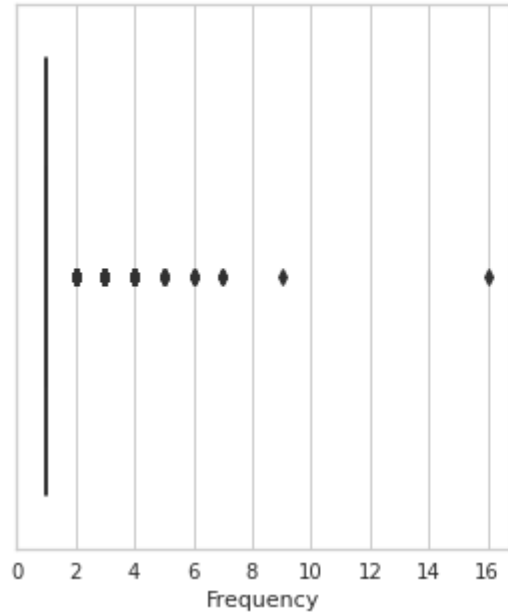
Variables de classification

Fenêtre temporelle : Septembre 2016 à Septembre 2018 (100 % du data-set)

RFM score

95 419 clients

RFM variables descriptions



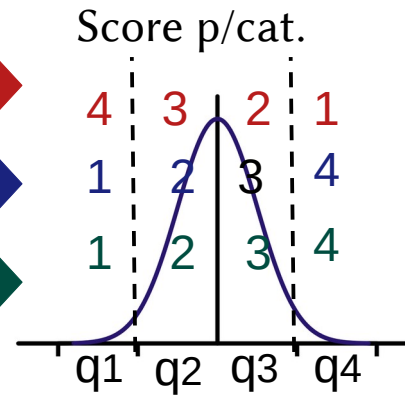
RFM score

Recency

Frequency

Monetary

95 419 clients



RFM score

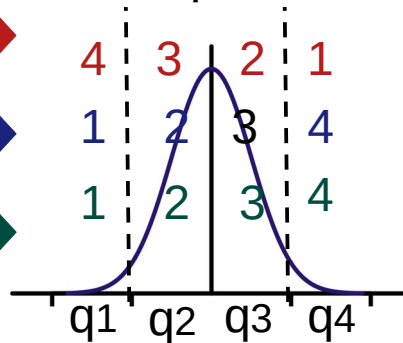
Recency

Frequency

Monetary

95 419 clients

Score p/cat.



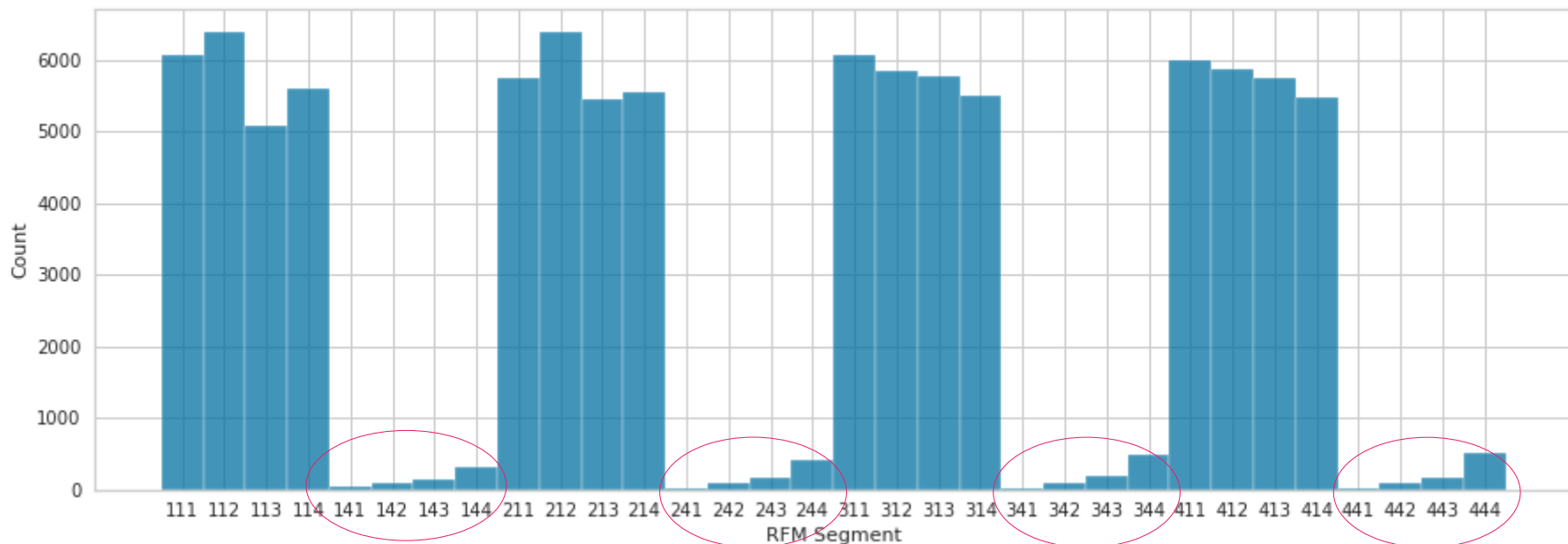
Note :

Frequency > 1 → score = 4

32 types de clients différents

RFM segment score

Clients avec frequency > 1



RFM score

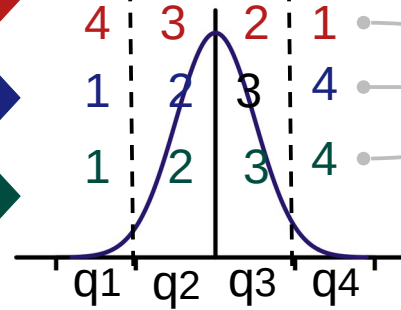
95 419 clients

Recency

Frequency

Monetary

Score p/cat.



Σ

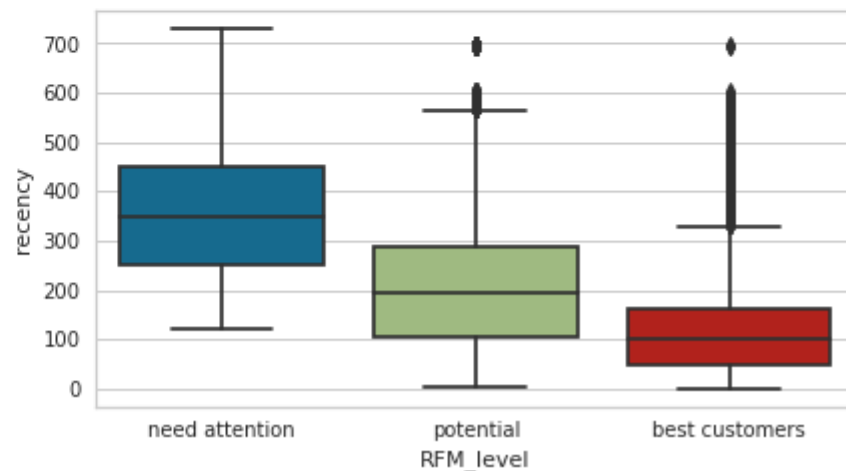
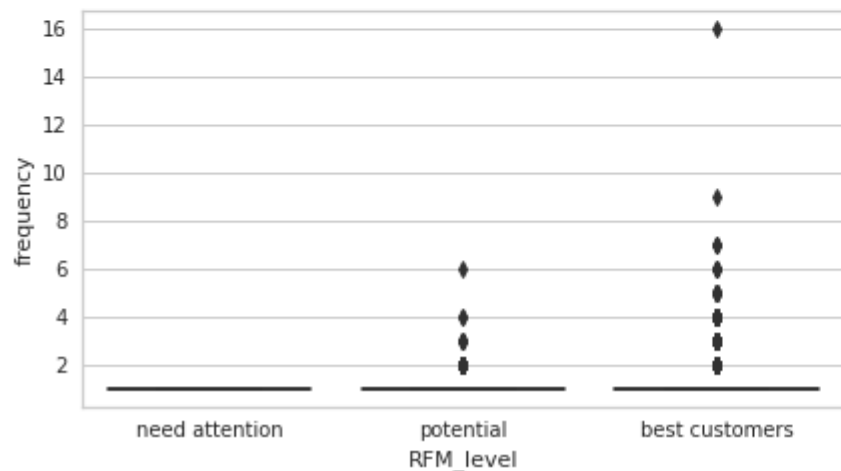
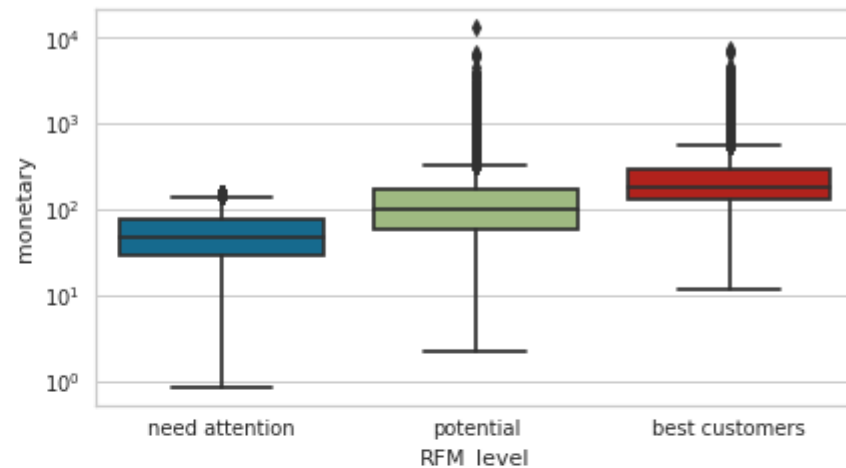
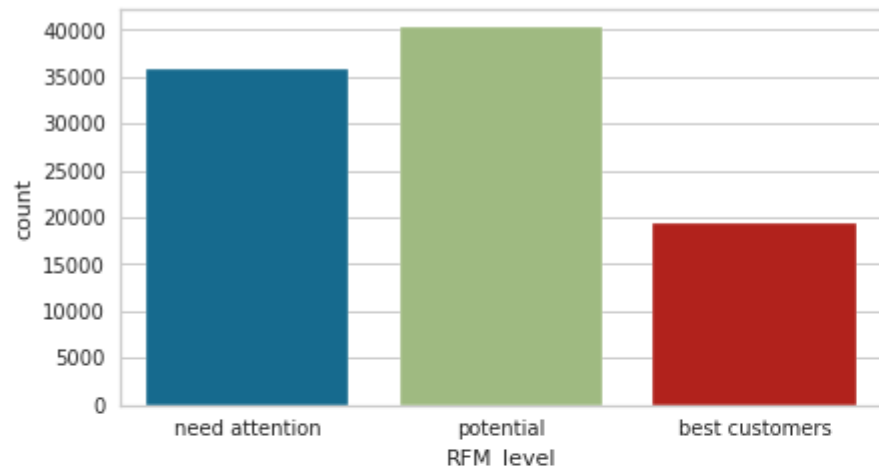
RFM SCORE



- Need attention: [min, Q1);
- Potential: [Q1, Q3);
- Best customers: [Q3, max]

Rocio Isorna

RFM score



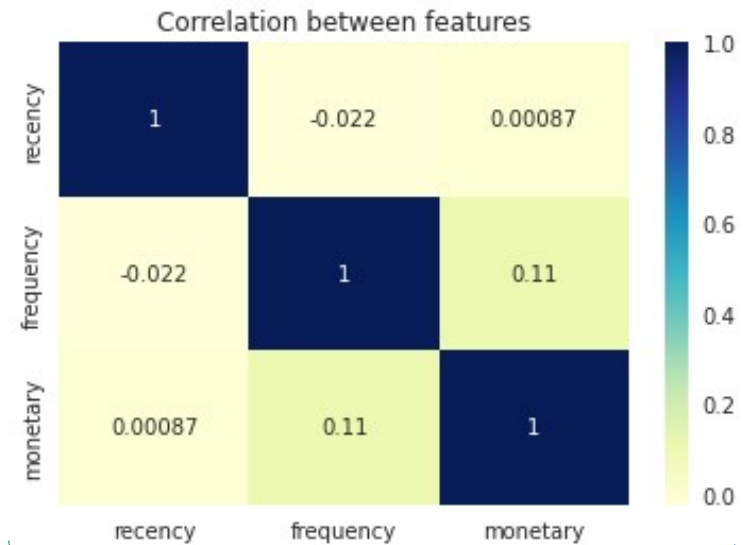


Apprentissage non-supervisé

K-Mean

95 419 clients

Features : RFM



StandardScaler()

DBSCAN et Agglomerative clustering

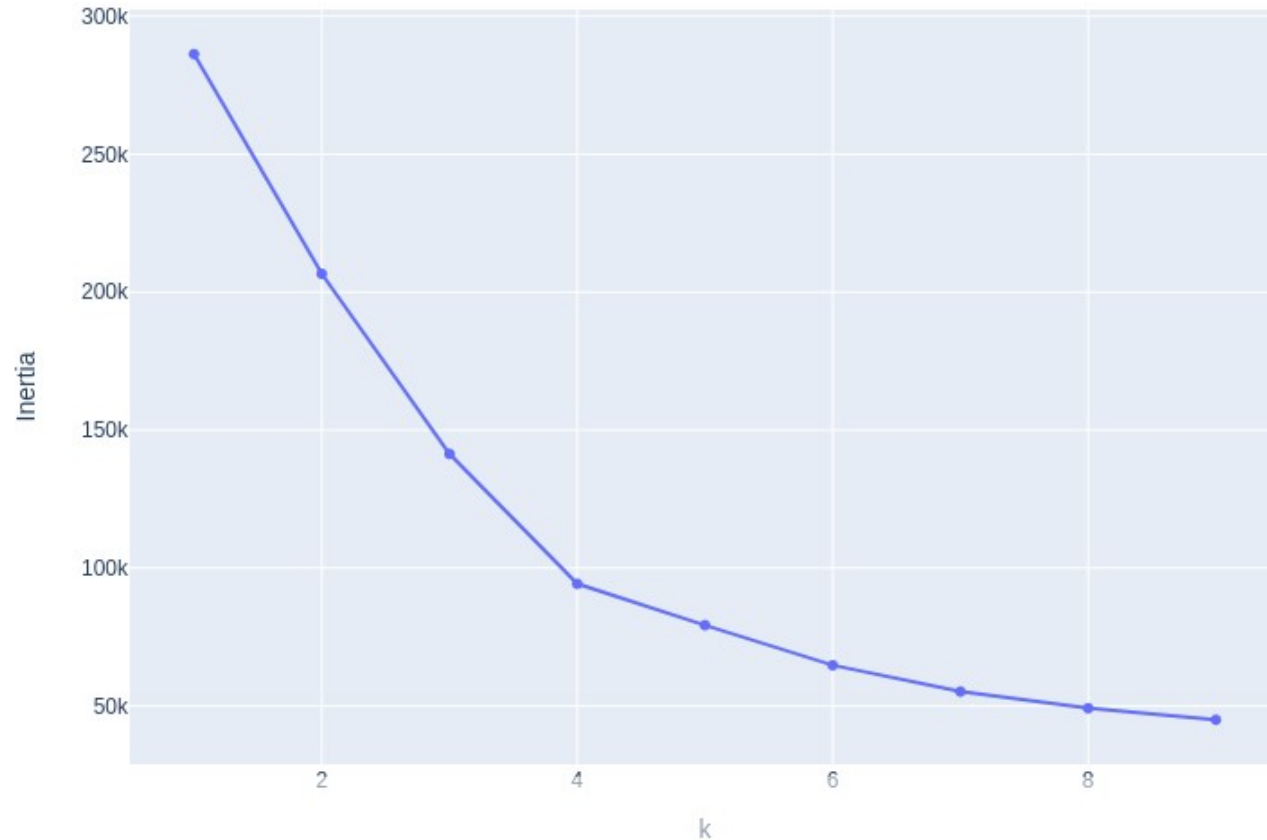
10 000 clients
random



K-Mean

Paramètre K : Méthode du coude

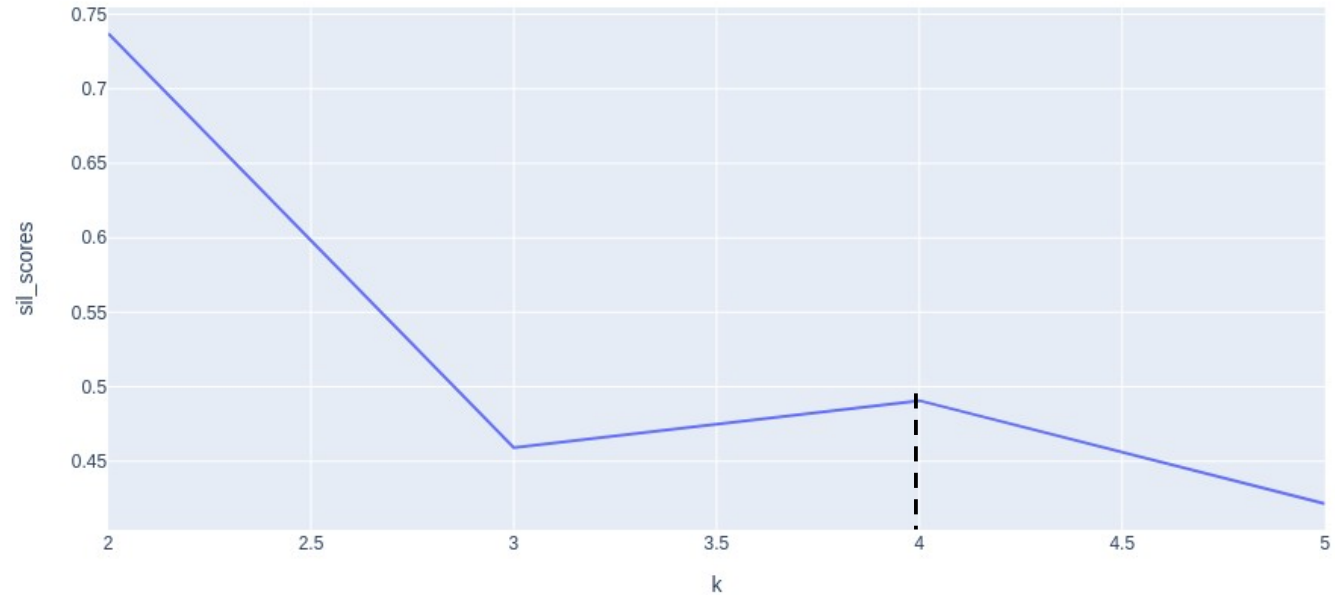
Elbow Method showing the optimal number of clusters





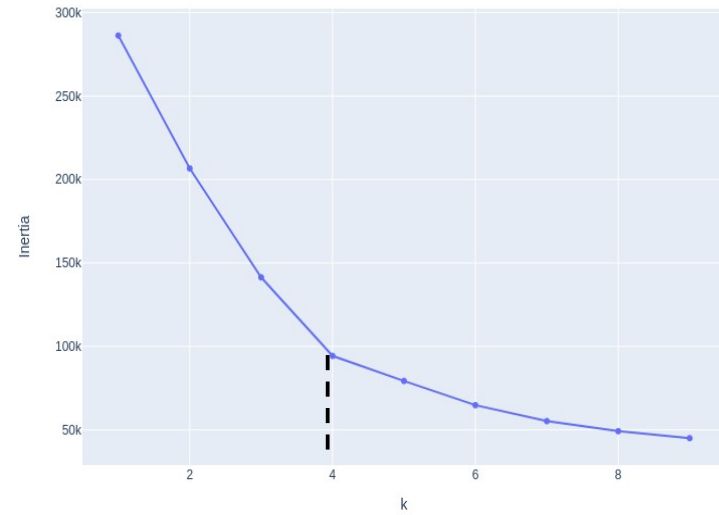
Paramètre K : Méthode du score

Score method for each K



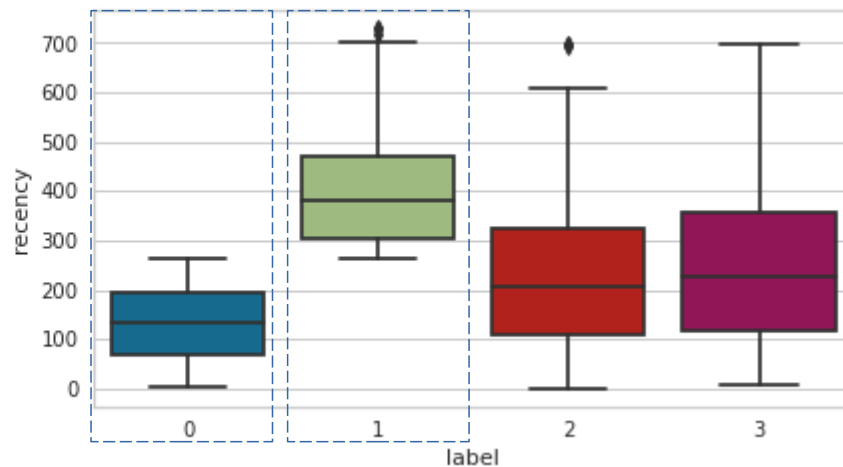
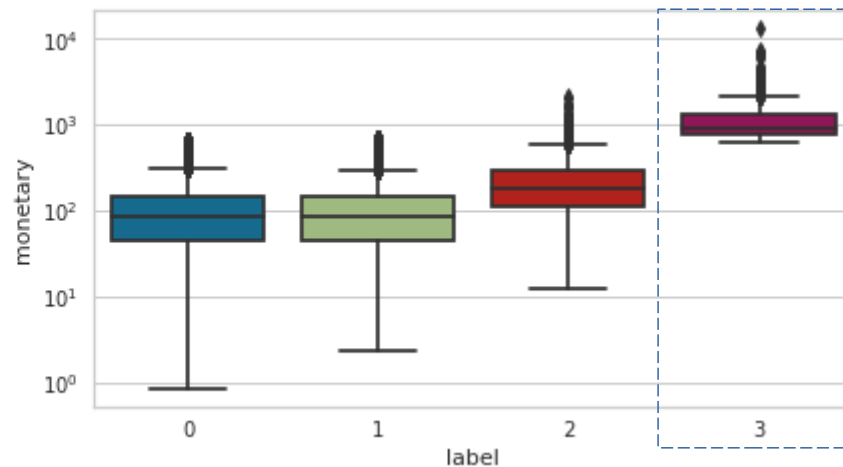
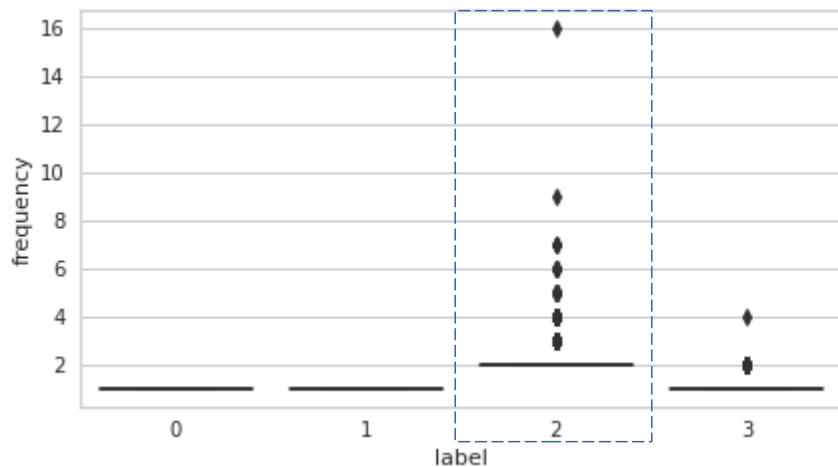
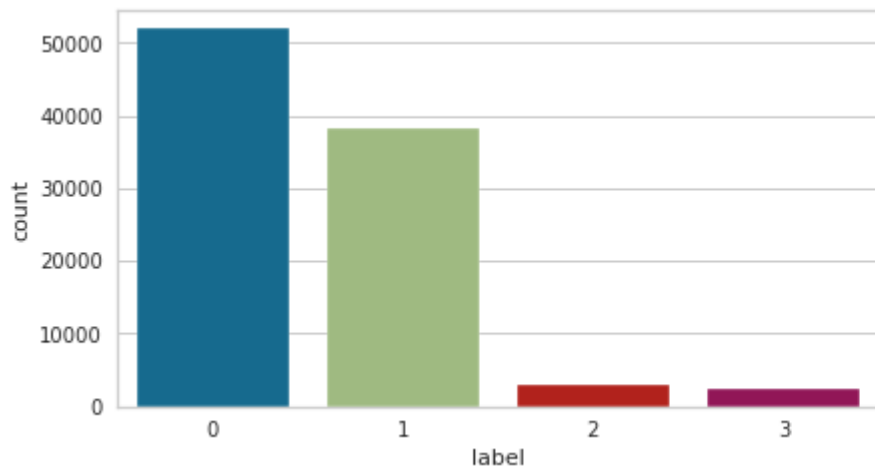
Paramètre K : Méthode du coude

Elbow Method showing the optimal number of clusters





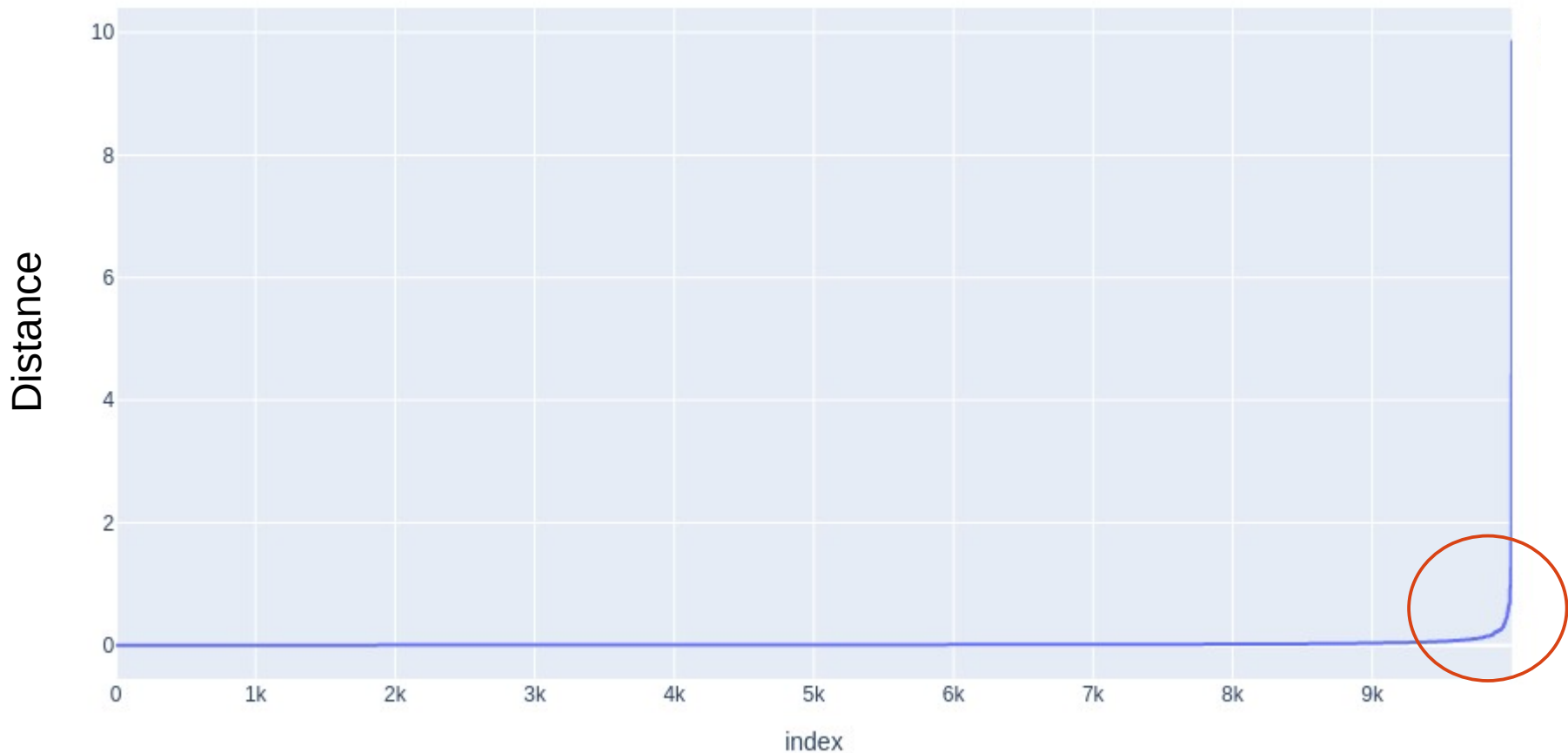
K-Mean clusters





DBSCAN : min sample et epsilon

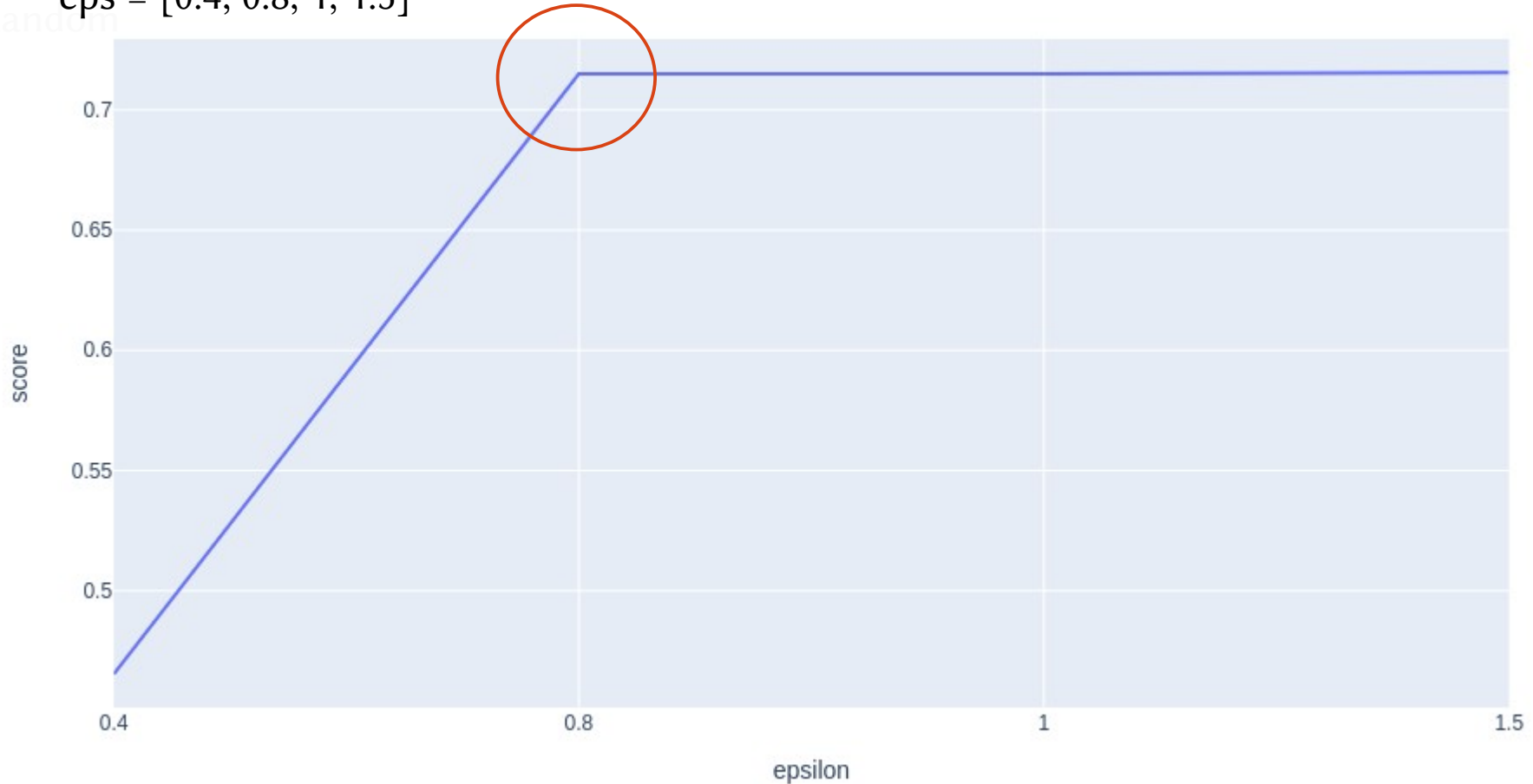
- $\text{min_samples} = 2 * \text{n_features}$
- Calcul de distance euclidien entre un point et les min_samples les plus proches





DBSCAN : min sample et epsilon

- $\text{min_samples} = 2 * \text{n_features}$
- $\text{eps} = [0.4, 0.8, 1, 1.5]$

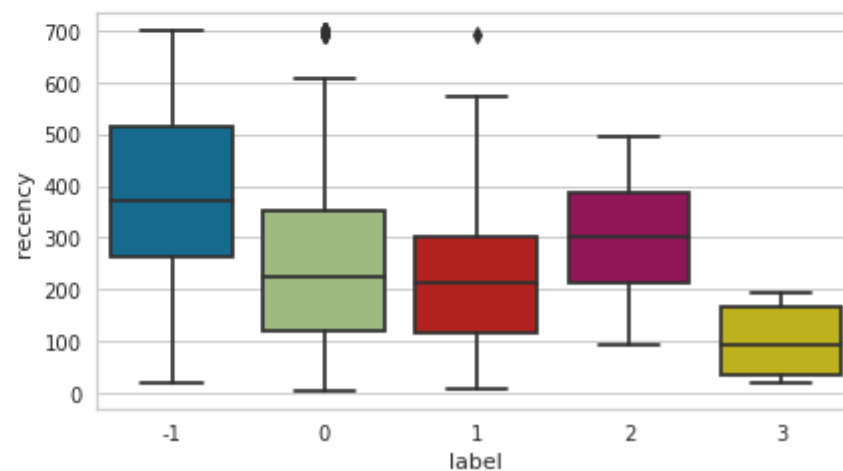
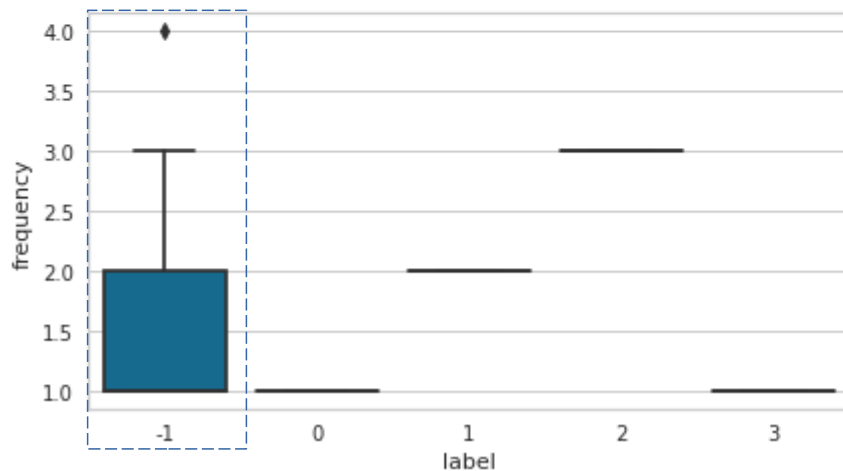
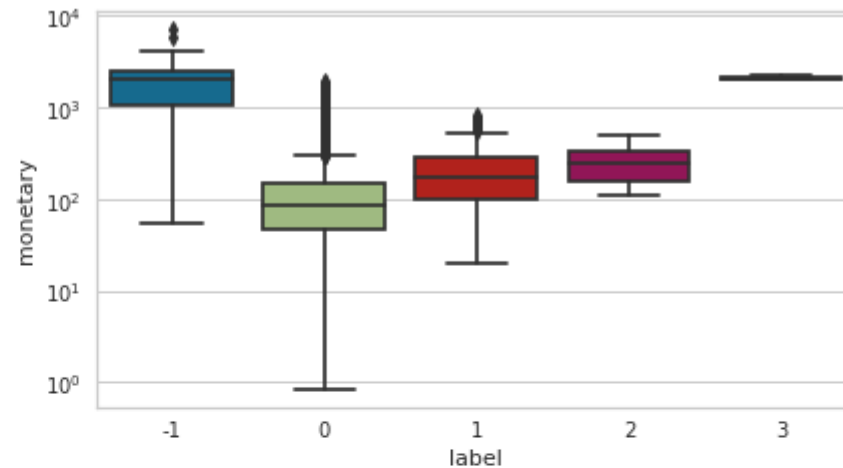
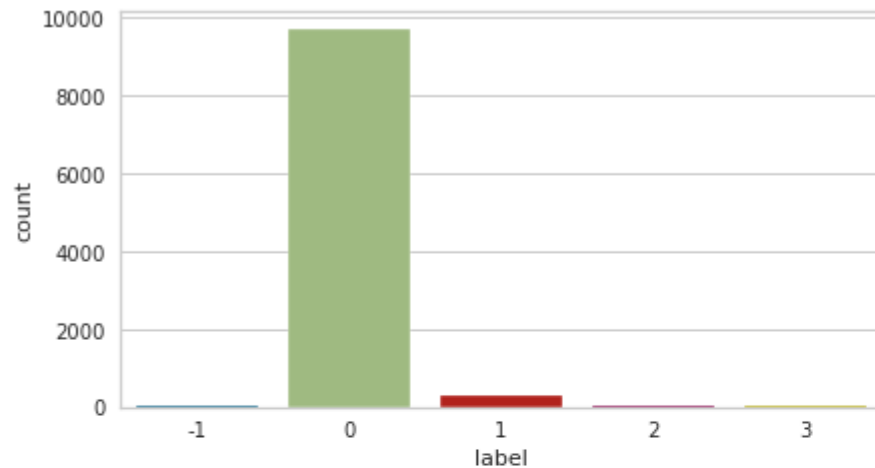




DBSCAN : clusters

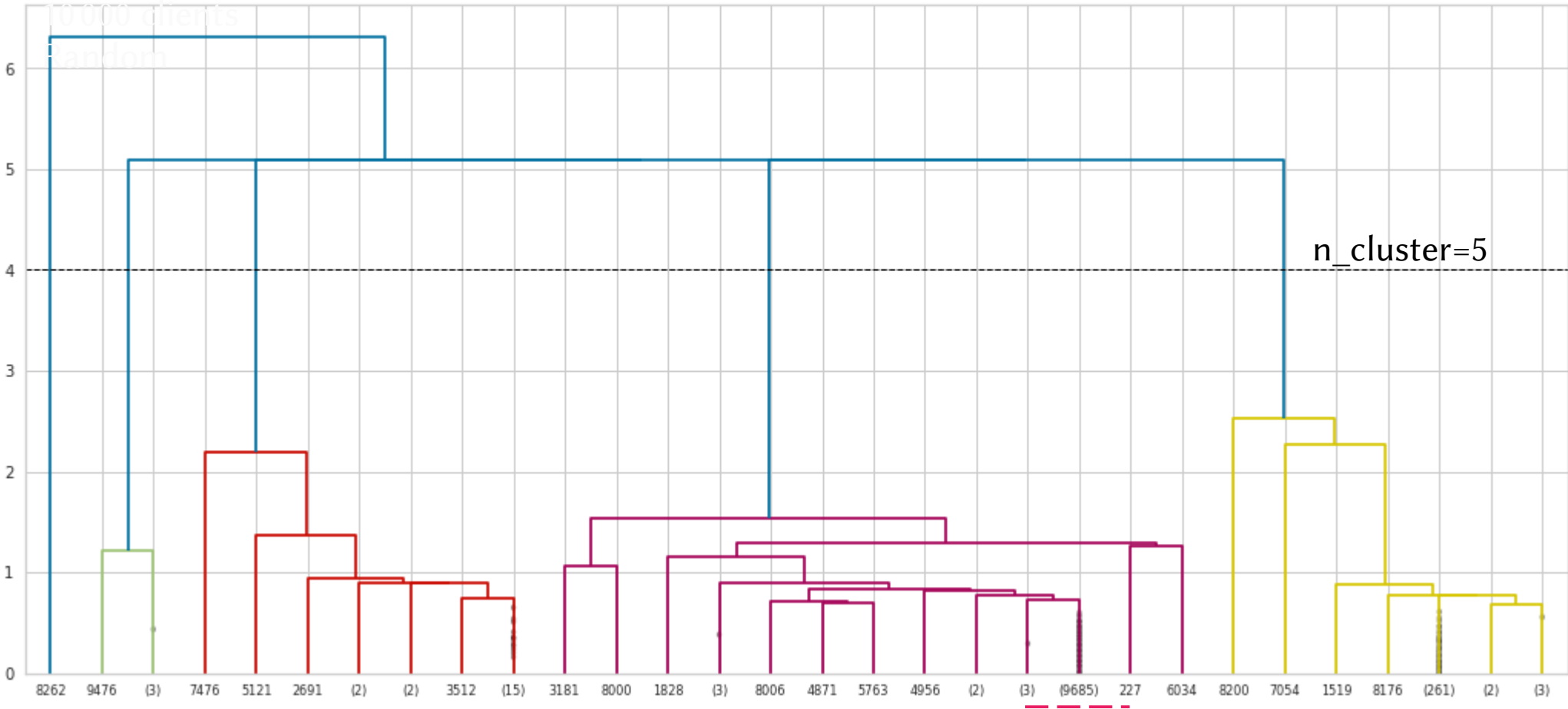
Biblio

10
Ra



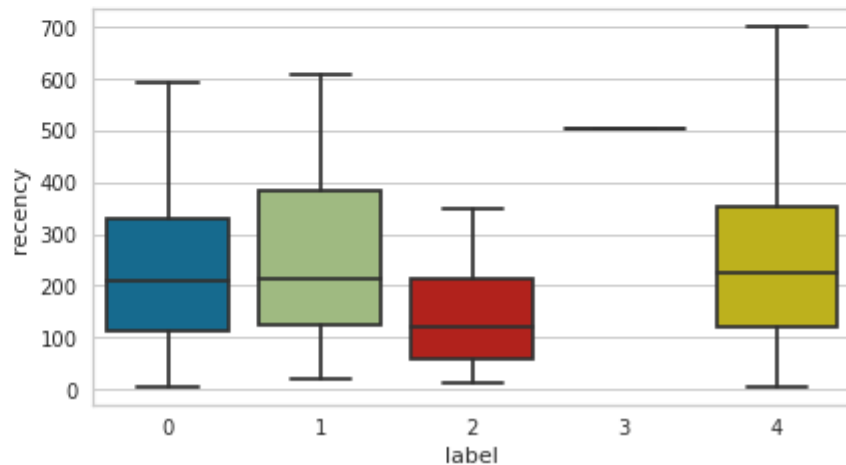
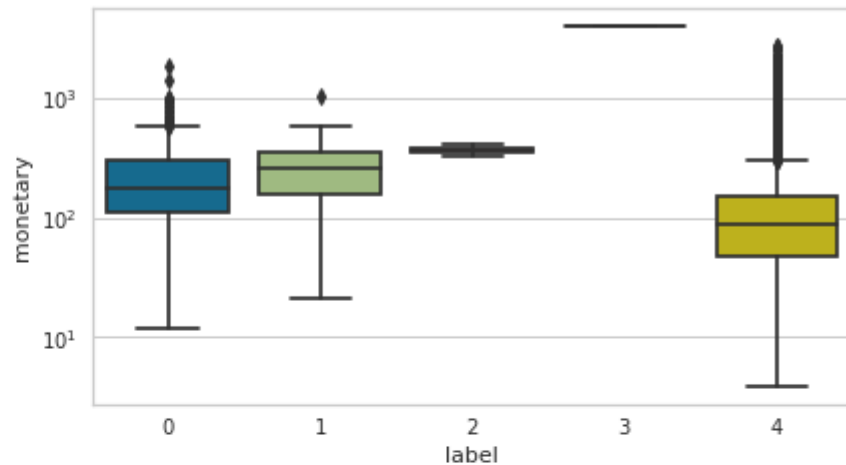
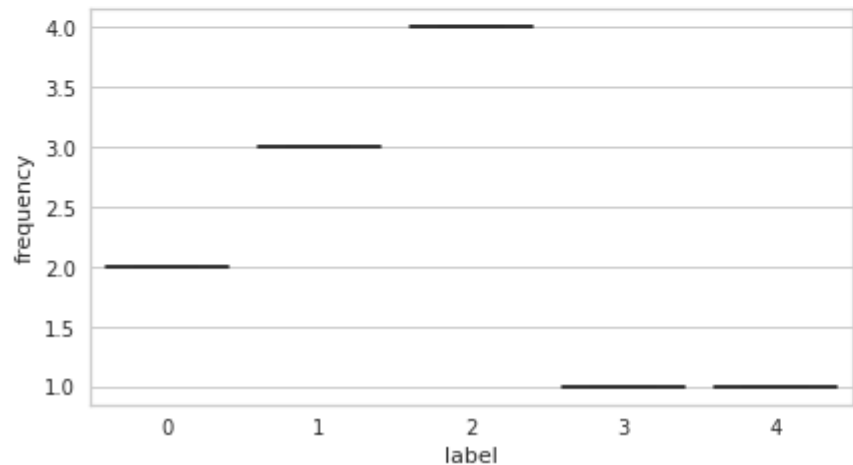
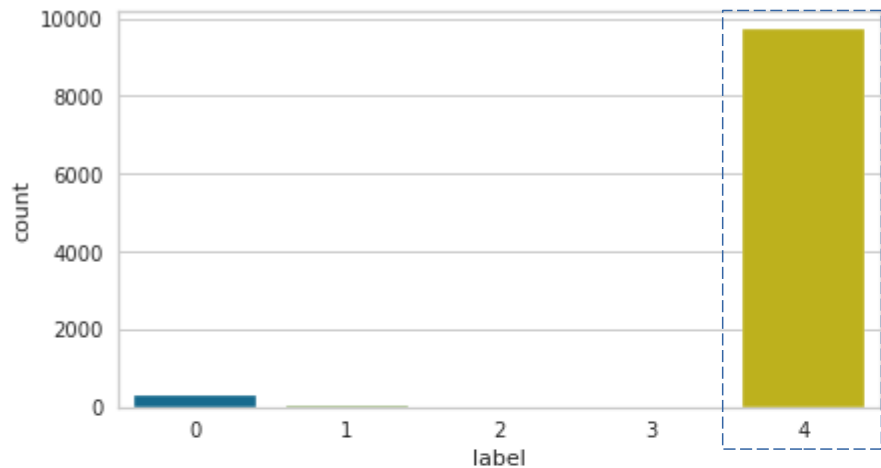


Agglomerative clustering : dendrogramme





Agglomerative clustering : clusters



Conclusions

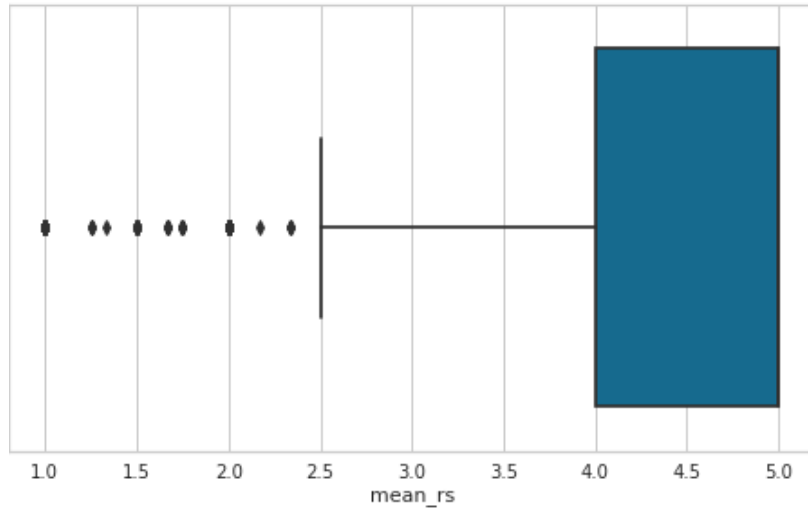
- La méthode K-Mean donnée des clusters plus équilibrés avec des définition métier plus marquées.
- En terme de mémoire et temps de calcul, la méthode K-Mean est plus performante
- La méthode Agglomerative n'as pas de fonction predict

Modèle final, fonction utilisateur et maintenance

K-Mean final

95 419 clients

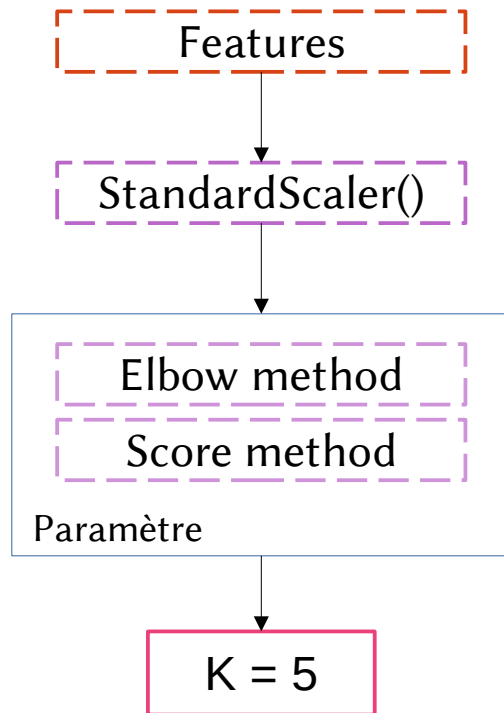
Features : RFM + Mean review score



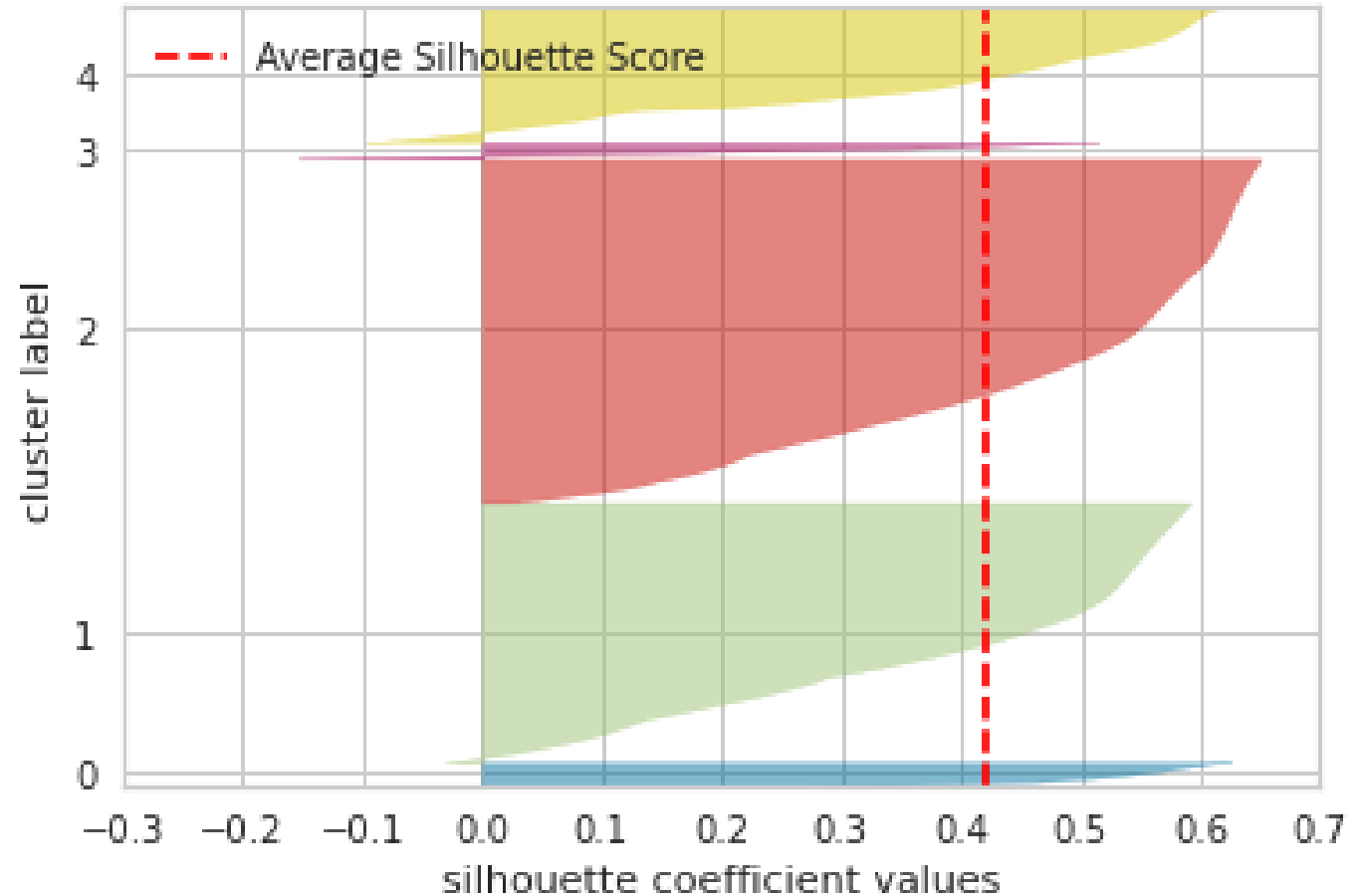
Features : correlation



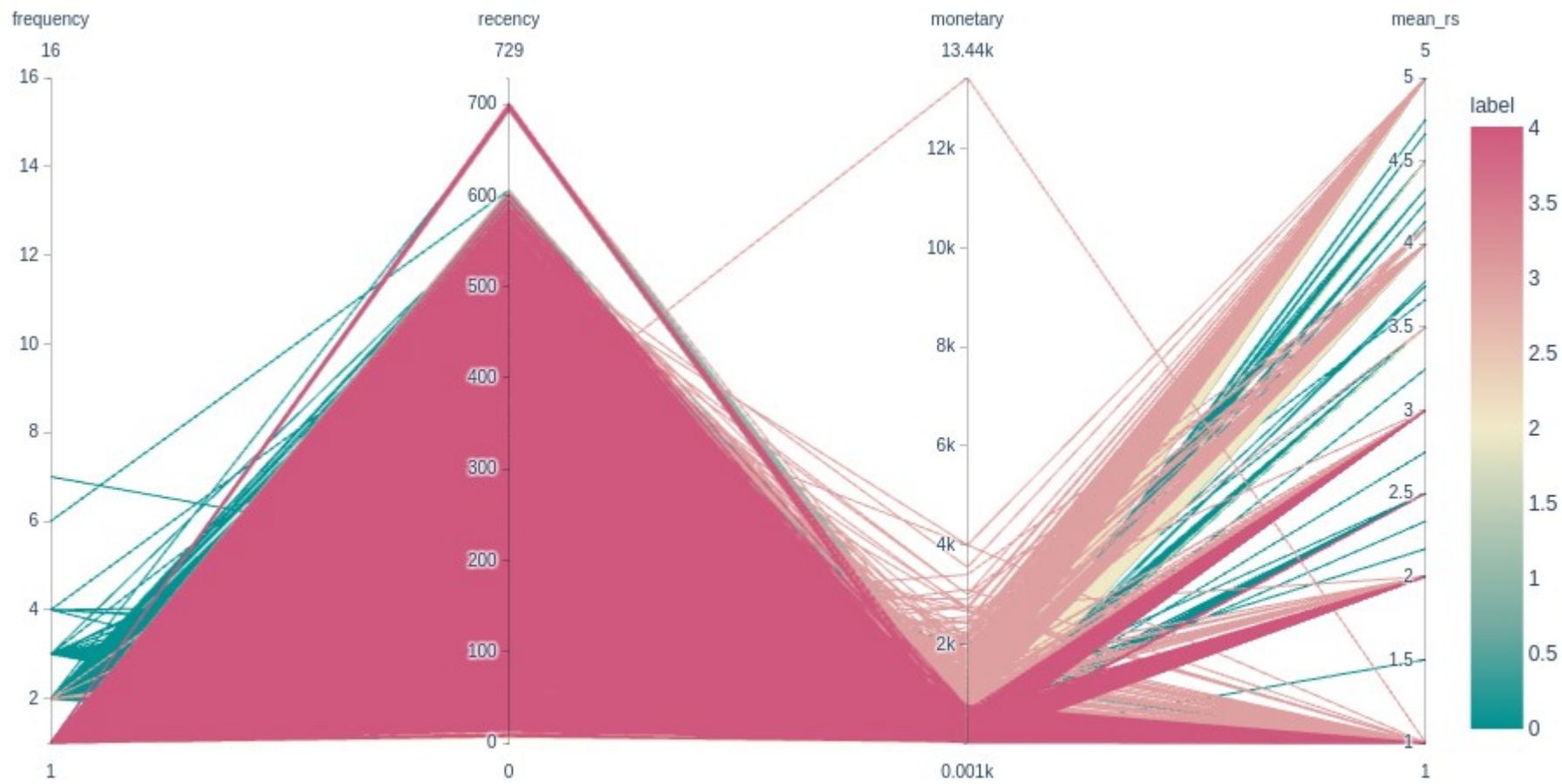
K-Mean final



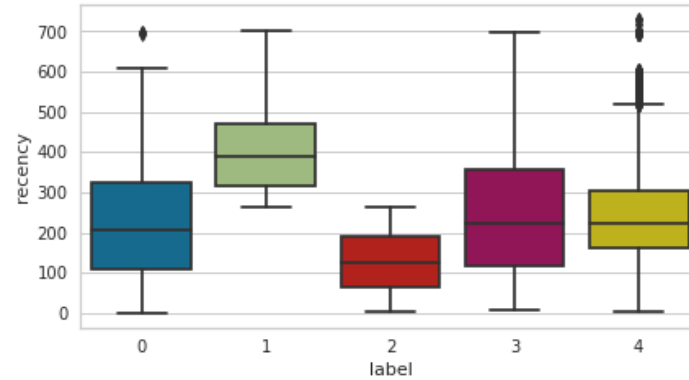
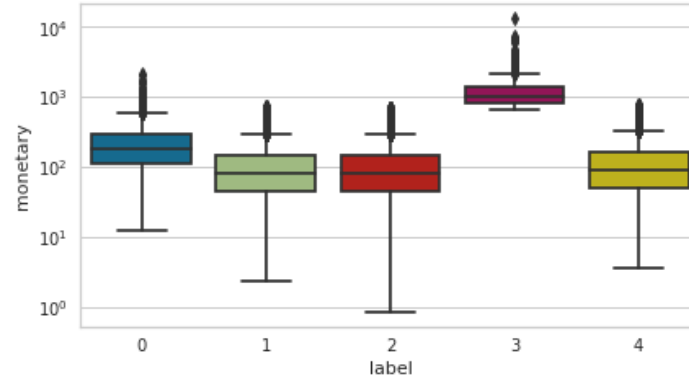
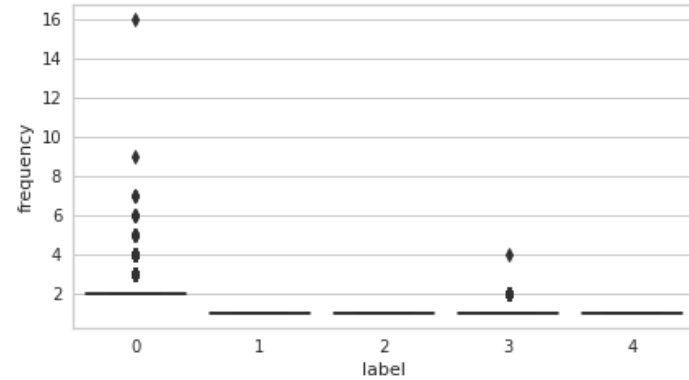
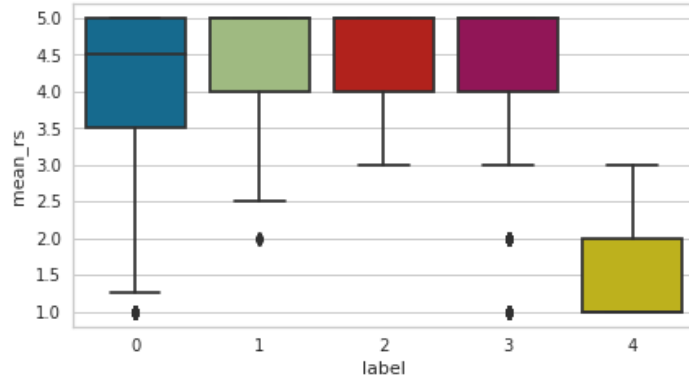
Silhouette Plot of KMeans Clustering for 95420 Samples in 5 Centers



K-Mean final



Clusters et « Persona »



Jean:
- Achète
fréquemment



Marie:
- 1 achat il y +
1 an



Pierre:
- Client récent (1
an) pas dépensier



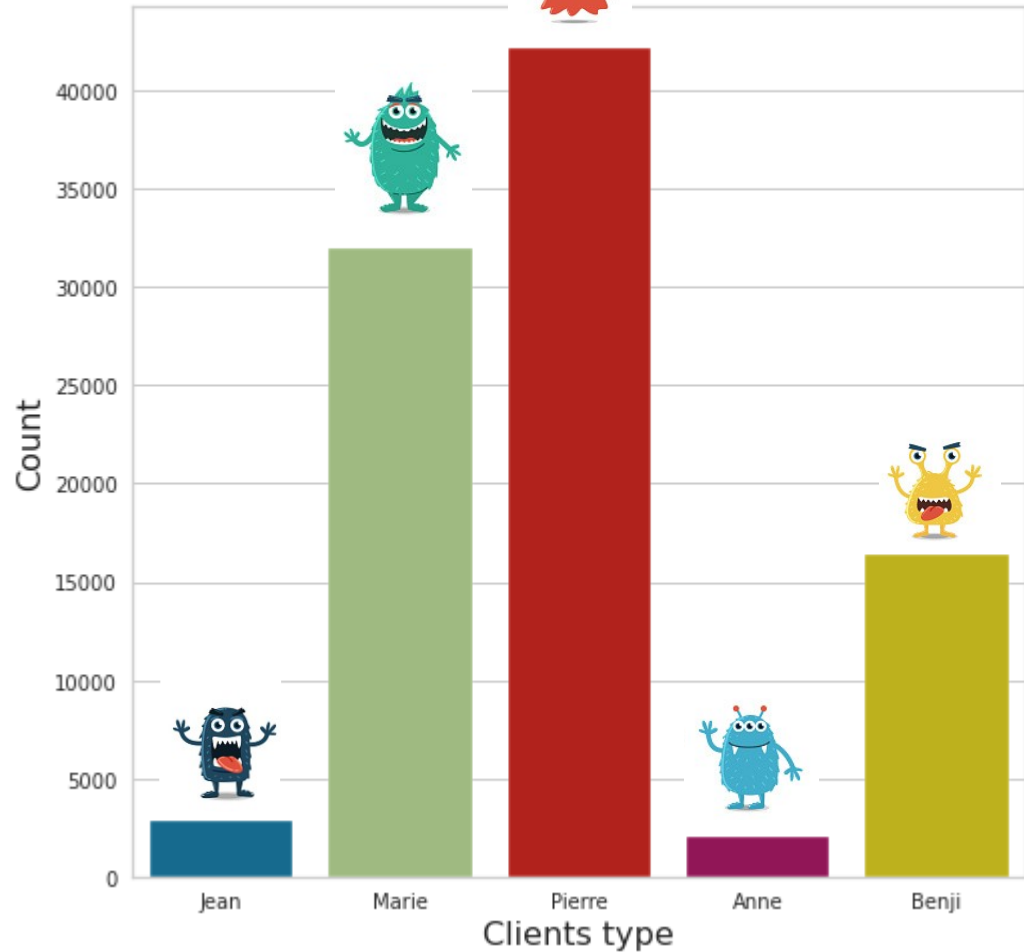
Anne:
- Ce qu'on
aime \$\$\$



Benji:
- Pas content

Clusters et « Persona »

Clients



Jean:
- Achète fréquemment



Marie:
- 1 achat il y + 1 an



Pierre:
- Client récent (1 an) pas dépensier



Anne:
- Ce qu'on aime \$\$\$

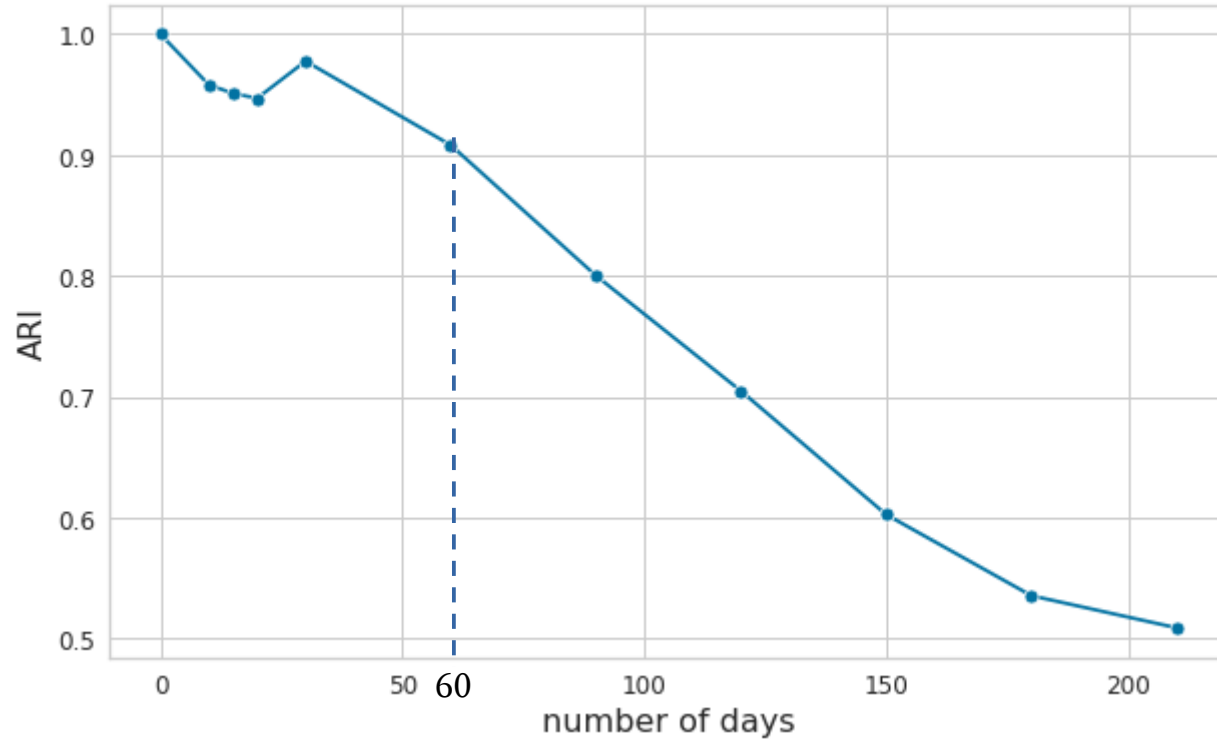


Benji:
- Pas content

orna

Maintenance :

ARI variation with number of days not considered

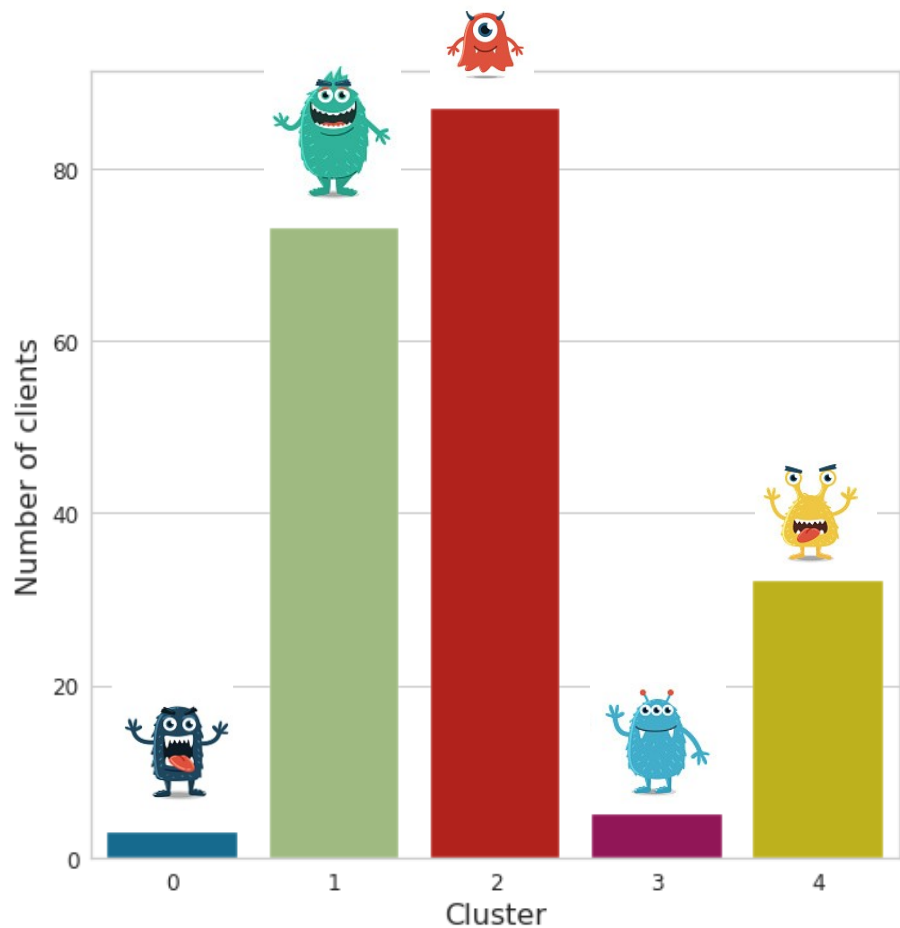


La maintenance est conseillé tout les 2 mois

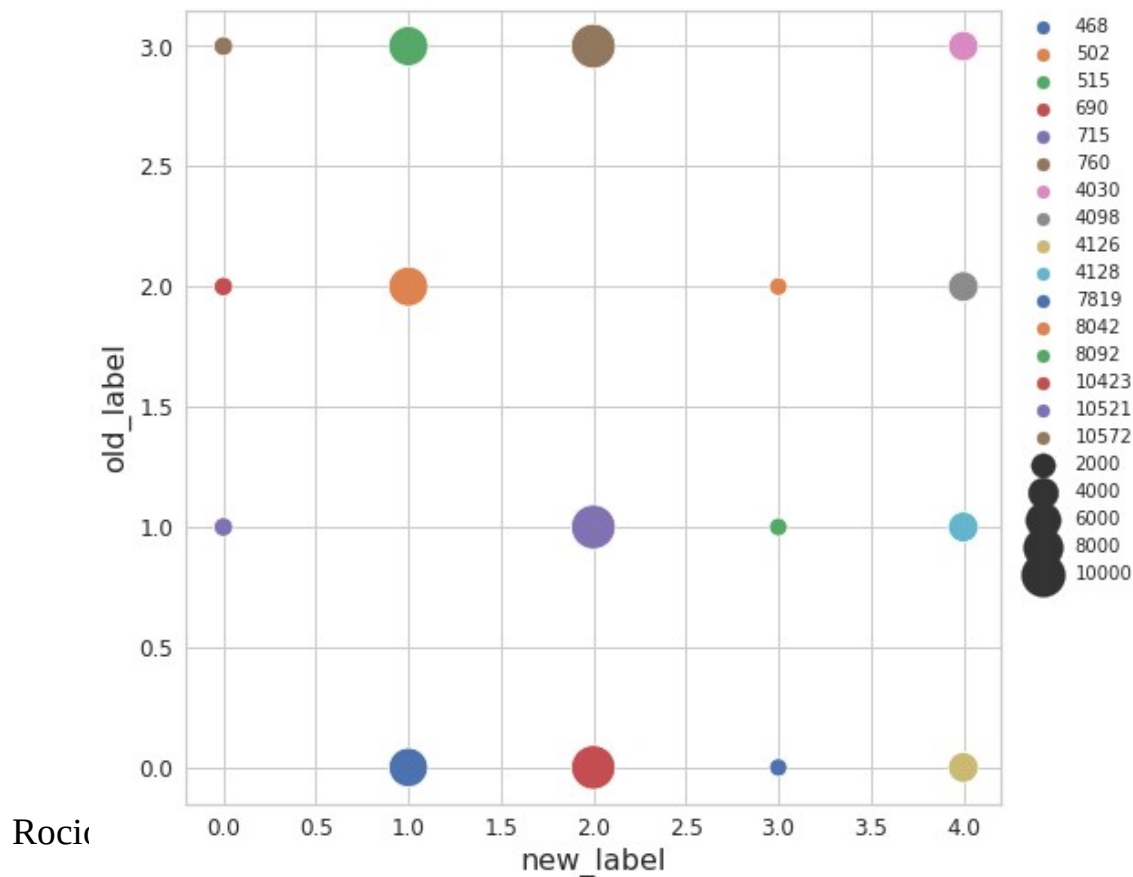
Fonction pour rajouter nouveaux clients: Example

Données d'entrée : base de données de la **dernier mise à jours** et **base des données à jour**

New clients labels/persona



Clients evolution



Conclusion:

- Cinq clusters ont été établie pour identifier le comportement des clients
- Il a été établie qu'il est pertinent de faire une maintenance de l'algorithme de clustering tout les 2 mois
- Un fonction est développe pour voir l'évolution des clients dans le temps et identifier à quel cluster appartiens les nouveaux clients

Merci de votre attention

