

Segmentation clientèle



Zaccaria Amillou

Sommaire

- Contexte
- Données
- Requêtes SQL
- EDA
- Modélisation
- Simulation

Contexte

Entreprise brésilienne

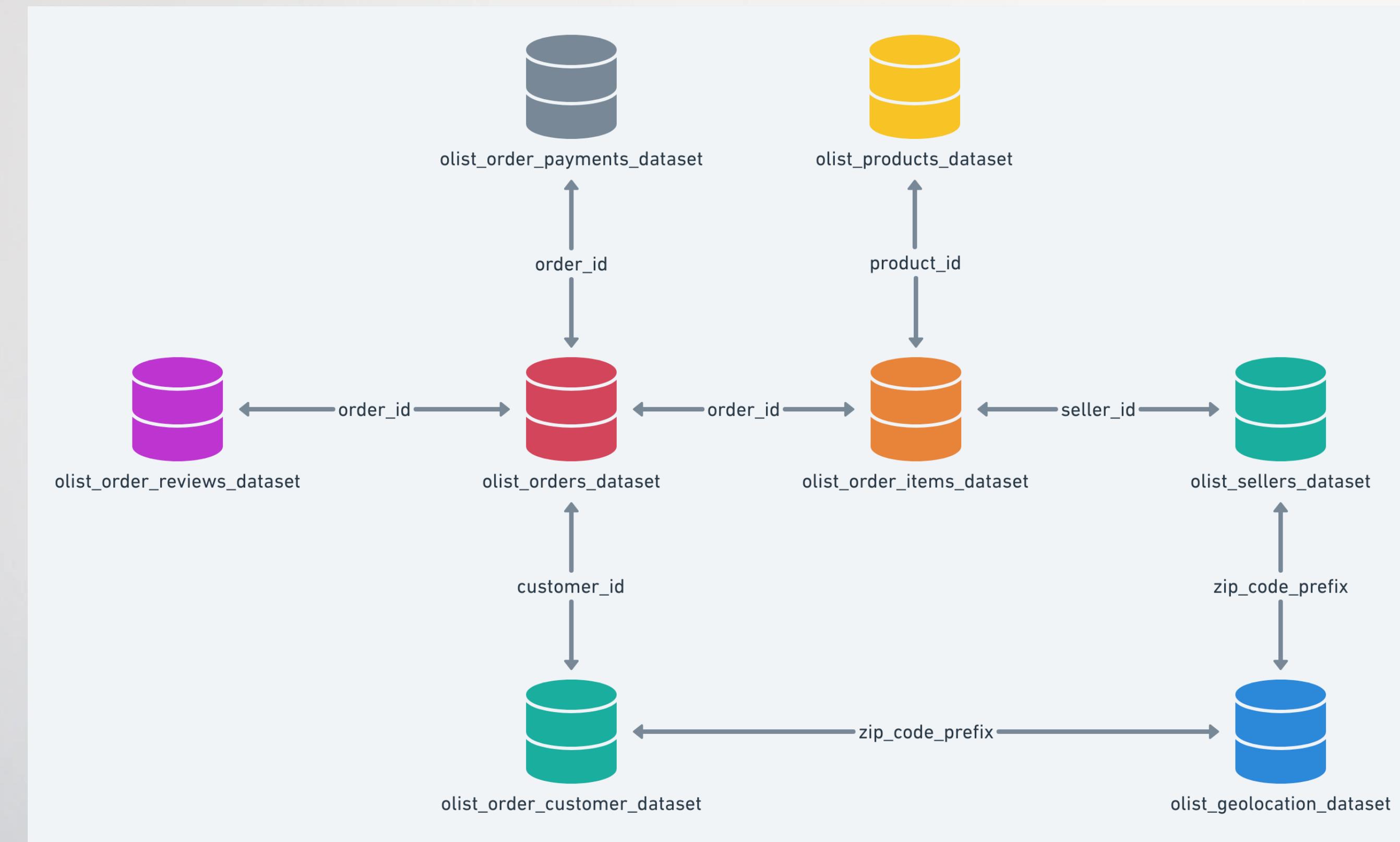
Propose solution de vente sur marketplace en ligne

Objectifs:

Définir la segmentation de la clientèle



Fichiers utilisées



Requêtes SQL

En excluant les commandes annulées, quelles sont les commandes récentes de moins de 3 mois que les clients ont reçues avec au moins 3 jours de retard ?

```
SELECT
    o.order_id as 'id ordre',
    (julianday(order_estimated_delivery_date) -
    julianday(order_delivered_customer_date)) as 'j_diff'
FROM
    orders o
WHERE
    (julianday(order_estimated_delivery_date) -
    julianday(order_delivered_customer_date)) >= 3
AND
    STRFTIME('%m', order_approved_at) >= 3
AND
    order_status NOT LIKE 'canceled'
```

Requêtes SQL

Qui sont les vendeurs ayant généré un chiffre d'affaires de plus de 100 000 Real sur des commandes livrées via Olist ?

```
SELECT
    oi.seller_id,
    COUNT(oi.product_id),
    SUM(op.payment_installments * op.payment_value) AS total
FROM order_items oi
LEFT JOIN orders o ON oi.order_id = o.order_id
LEFT JOIN order_pymts op ON o.order_id = op.order_id
GROUP BY oi.seller_id
HAVING SUM(op.payment_installments * op.payment_value) > 100000;
```

Requêtes SQL



Qui sont les nouveaux vendeurs (moins de 3 mois d'ancienneté) qui sont déjà très engagés avec la plateforme (ayant déjà vendu plus de 30 produits) ?

```
SELECT
    seller_id,
    COUNT(oi.product_id) AS 'nombre produits',
    order_purchase_timestamp
FROM
    order_items oi
JOIN orders o ON o.order_id = oi.order_id
WHERE strftime('%Y-%m-%d %H:%M:%S', o.order_purchase_timestamp) <
    strftime('%Y-%m-%d %H:%M:%S', '2018-10-17 17:30:18', '-3 months')
GROUP BY
    seller_id
HAVING
    COUNT(oi.product_id) >= 30
ORDER BY "nombre product" DESC;
```

Requêtes SQL

Quels sont les 5 codes postaux, enregistrant plus de 30 commandes, avec le pire review score moyen sur les 12 derniers mois ?

```
SELECT
    c.customer_zip_code_prefix as ZIP,
    AVG(review_score)
FROM order_reviews or2
JOIN orders o ON or2.order_id = o.order_id
JOIN main.customers c ON o.customer_id = c.customer_id
GROUP BY c.customer_zip_code_prefix
ORDER BY AVG(review_score) ASC;
```

EDA

```
In 7 1 # affichage dataframe
2 display(customers.head(3))
Executed at 2024.02.22 08:40:06 in 7ms
```

	customer_id	customer_unique_id	customer_zip_code_prefix	customer_city	customer_state
0	06b8999e2fba1a1fbc88172c00ba8bc	861eff4711a542e4b93843c6dd7febb	144	Winnipeg	MB
1	18955e83d337fd6b2def6b18a428ac7	290c77bc529b7ac935b93aa66c333dc	97	Edmonton	AB
2	4e7b3e00288586ebd08712fdd0374a0	060e732b5b29e8181a18229c7b0b2b5	11	Vancouver	BC

Le df `customers` contient les informations à propos des clients

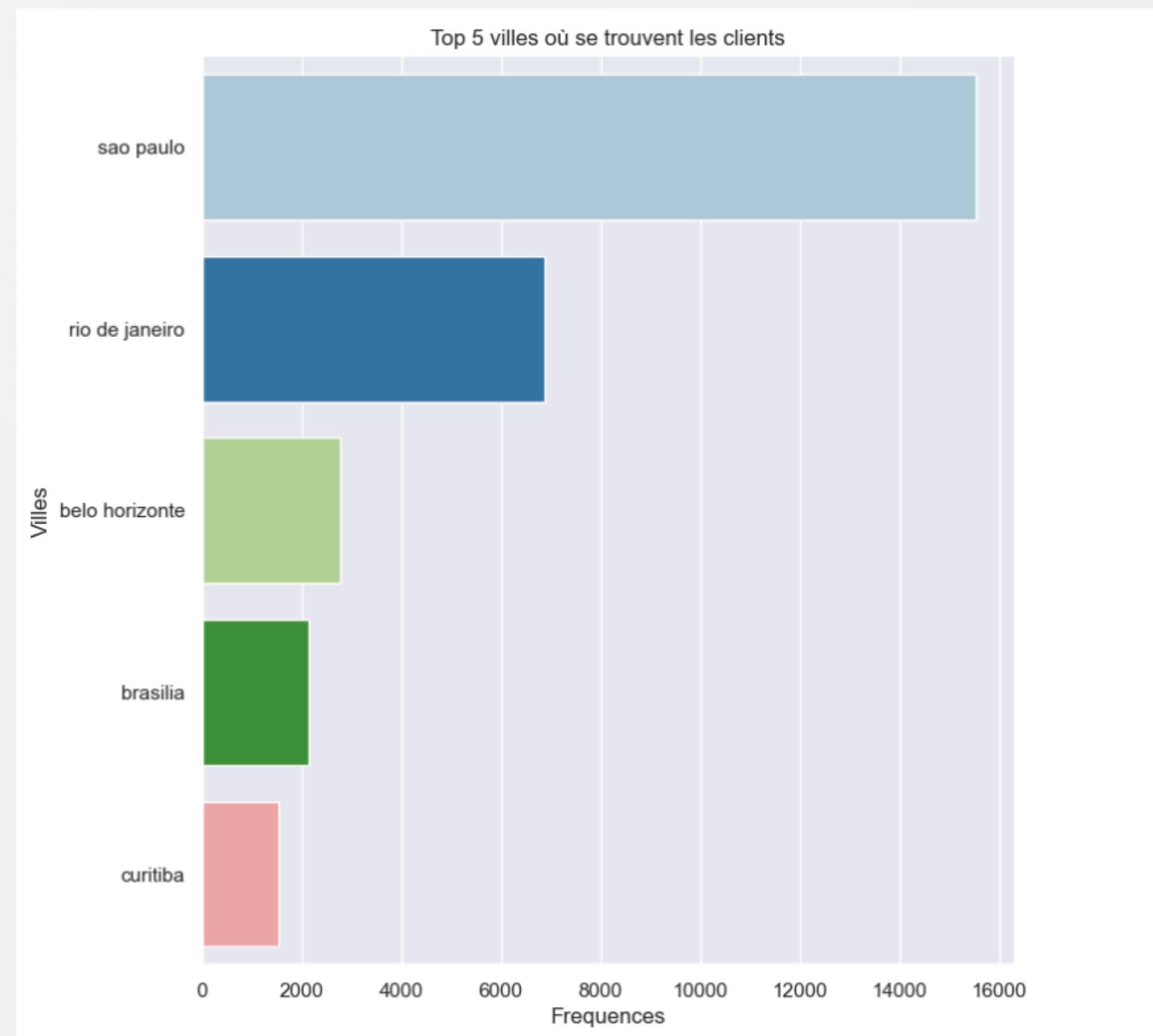
- `customer_id` : id donnée aux clients pour chaque ordre
- `customer_unique_id` : id unique pour identifier chaque client
- `customer_zip_code_prefix` : code postal du client
- `customer_city` : la ville du client au moment de l'ordre
- `customer_state` : l'état du client au moment de l'ordre

- `customer_state` : l'état du client au moment de l'ordre
- `customer_city` : la ville du client au moment de l'ordre
- `customer_zip_code_prefix` : code postal du client

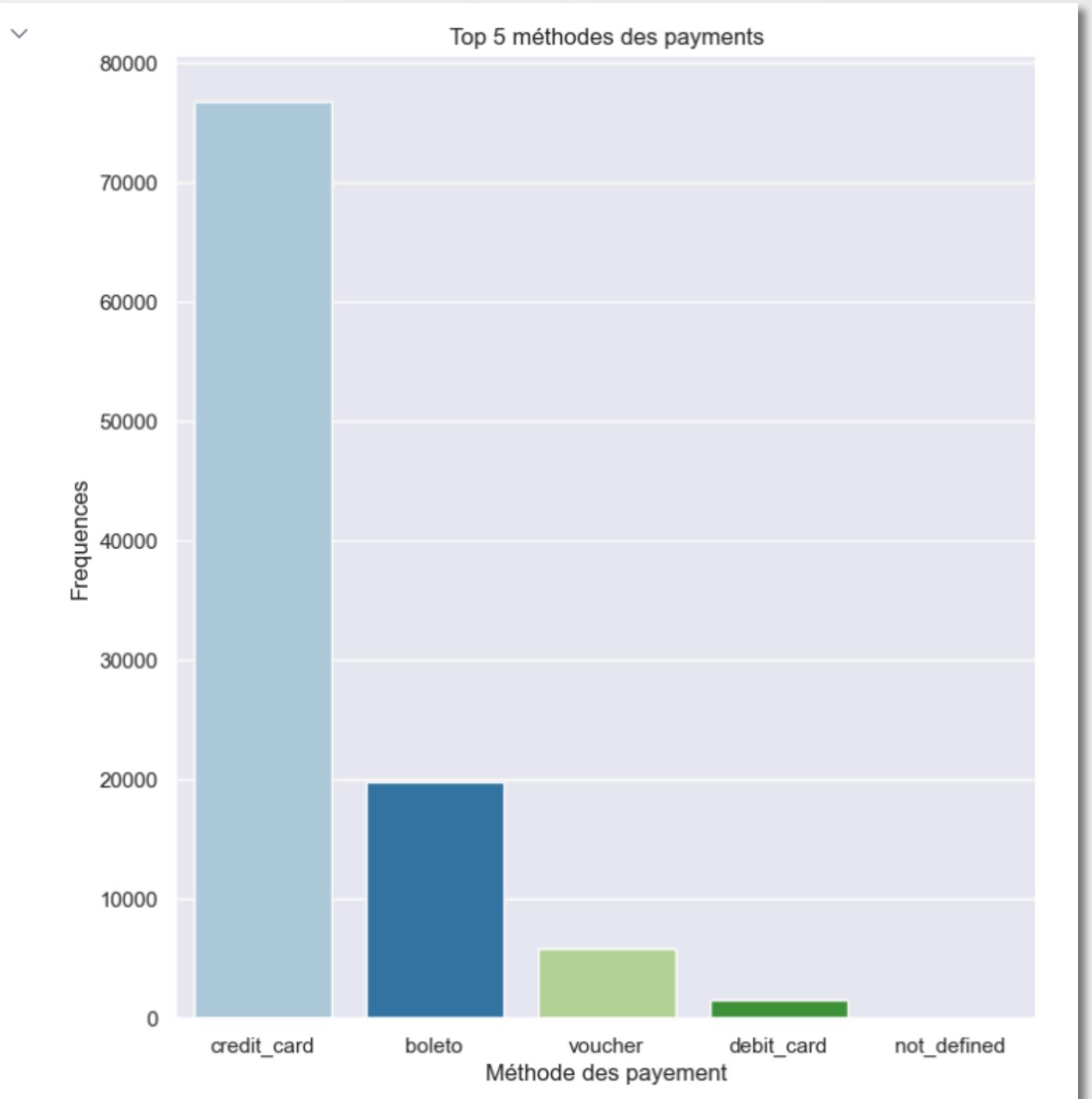
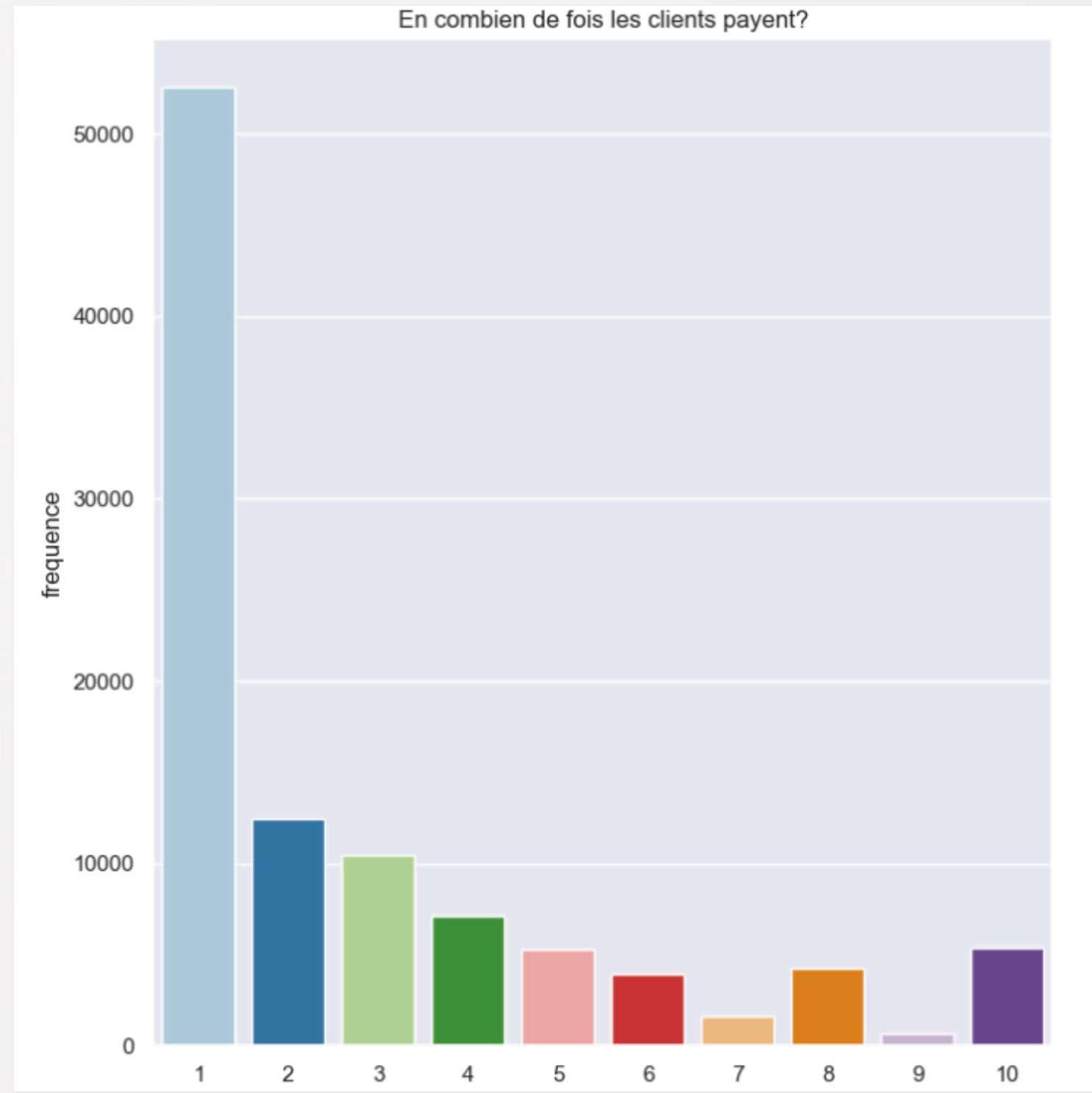


EDA - Analyses

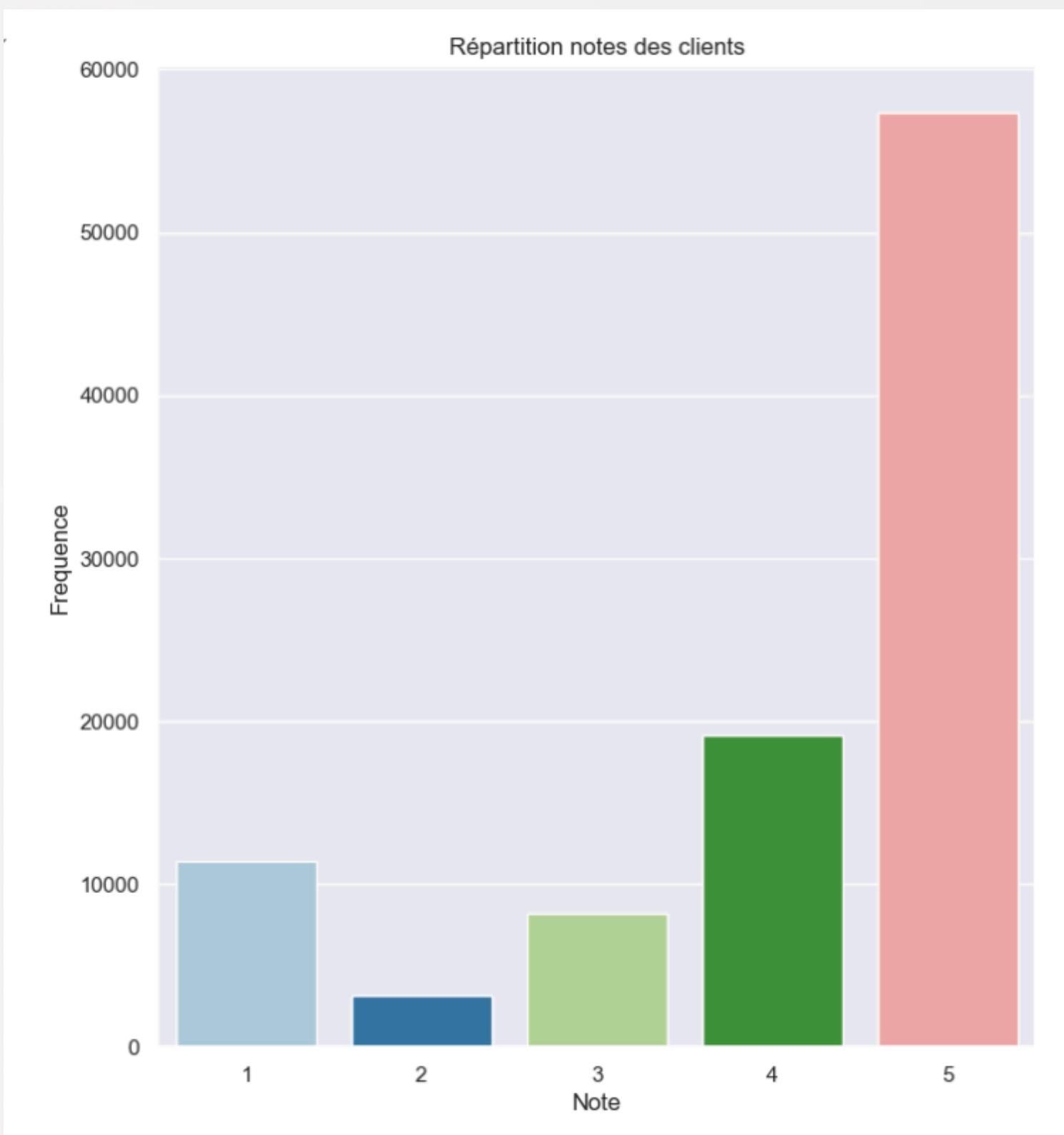
```
In 36 1 def graphique_top_count(data, colonne):
2     """ Creation bar plot des top 10 des iteration d'une colonne données
3     ARGS:
4         data: dataframe
5         colonne: colonne à itérer"""
6     # count d'itérations des typologies de payments
7     data_count = data[colonne].value_counts().reset_index()
8     data_count.columns = [colonne, 'count']
9     # top 10
10    top_10_count = data_count.sort_values(by='count', ascending=False).head(10)
11    # creation graphique
12    sns.barplot(
13        top_10_count,
14        x=colonne,
15        y='count',
16        palette='Paired')
Executed at 2024.02.22 09:26:50 in 2ms
```



EDA - Analyses



EDA - Analyses

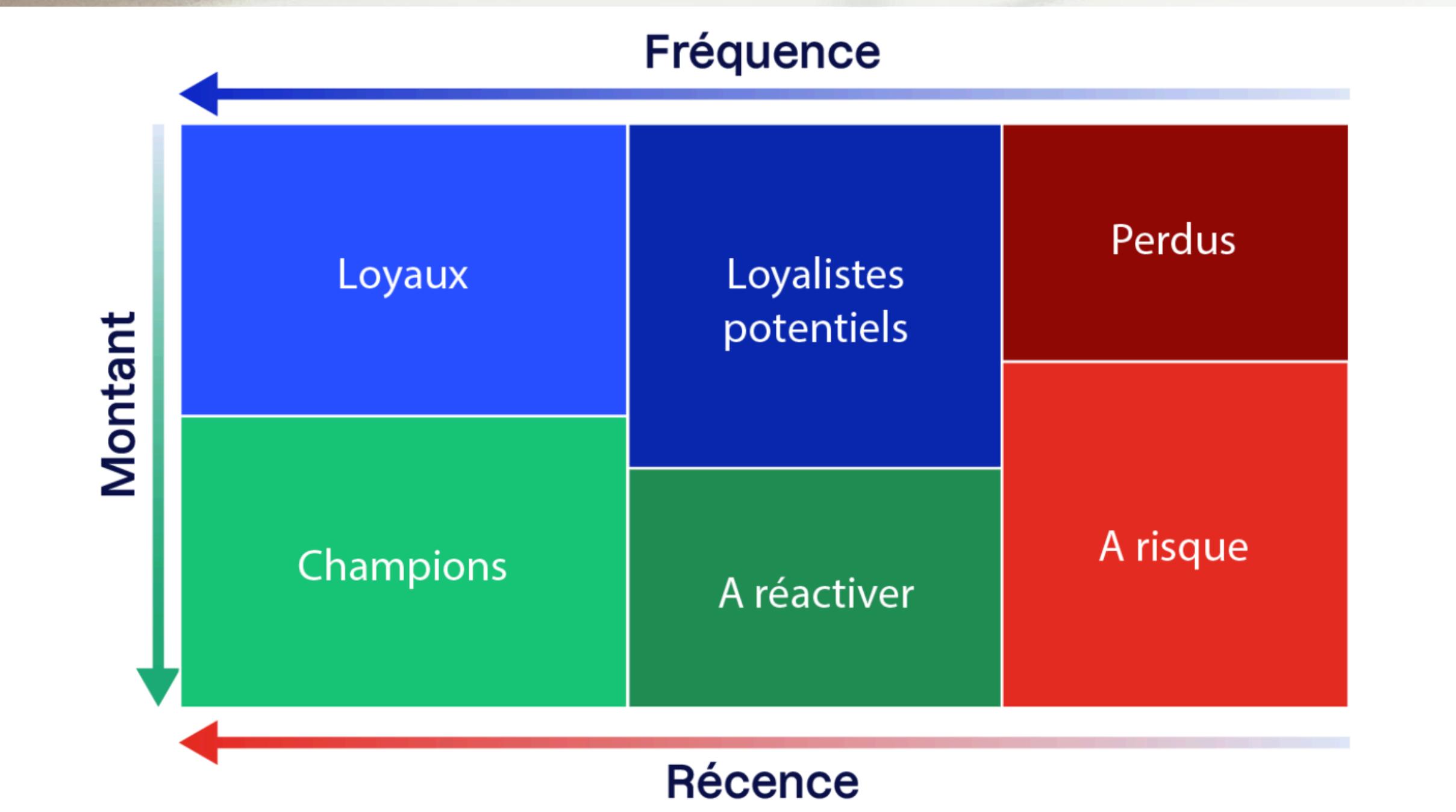


EDA

```
In 53 1 # jointures entre variables
2 df = (
3     customers[['customer_id','customer_city', 'customer_unique_id', 'customer_state']]
4     .merge(orders, on='customer_id', how='left')
5     .merge(payments, on='order_id', how='left' )
6     .merge(items,on='order_id',how='left')
7     .merge(sellers[['seller_id','seller_city','seller_state']], on='seller_id', how='left')
8     .merge(reviews, on='order_id', how='left')
9 )
Executed at 2024.02.22 09:26:53 in 384ms
```



Modélisation



Modélisation - RFM

```
In 9 1 # calcul de la date à utiliser
2 ajourdhui = df['order_purchase_timestamp'].max() + dt.timedelta(days=2)
3 ajourdhui
Executed at 2024.02.22 08:40:58 in 3ms

Out 9    Timestamp('2018-10-19 17:30:18')
```

```
In 12 1 # Calculs de jours passés à partir d'ajourdhui (variable) et dernière commande
2 recence_df['Recence']= recence_df['order_purchase_timestamp'].apply(lambda x: (ajourdhui - x).days)
3 display(recence_df.head(3))
Executed at 2024.02.22 08:40:58 in 160ms
```

```
In 13 1 frequence_df = df.groupby('customer_unique_id').size().reset_index(name='Frequence')
2 # Sort du df par frequence
3 frequence_df = frequence_df.sort_values(by='Frequence')
```

```
In 14 1 # création df pour le montant
2 montant_df = df.groupby('customer_unique_id', as_index=False)[['payment_value']].sum()
3 montant_df.columns = ['customer_unique_id', 'Montant']
```

Modélisation - RFM

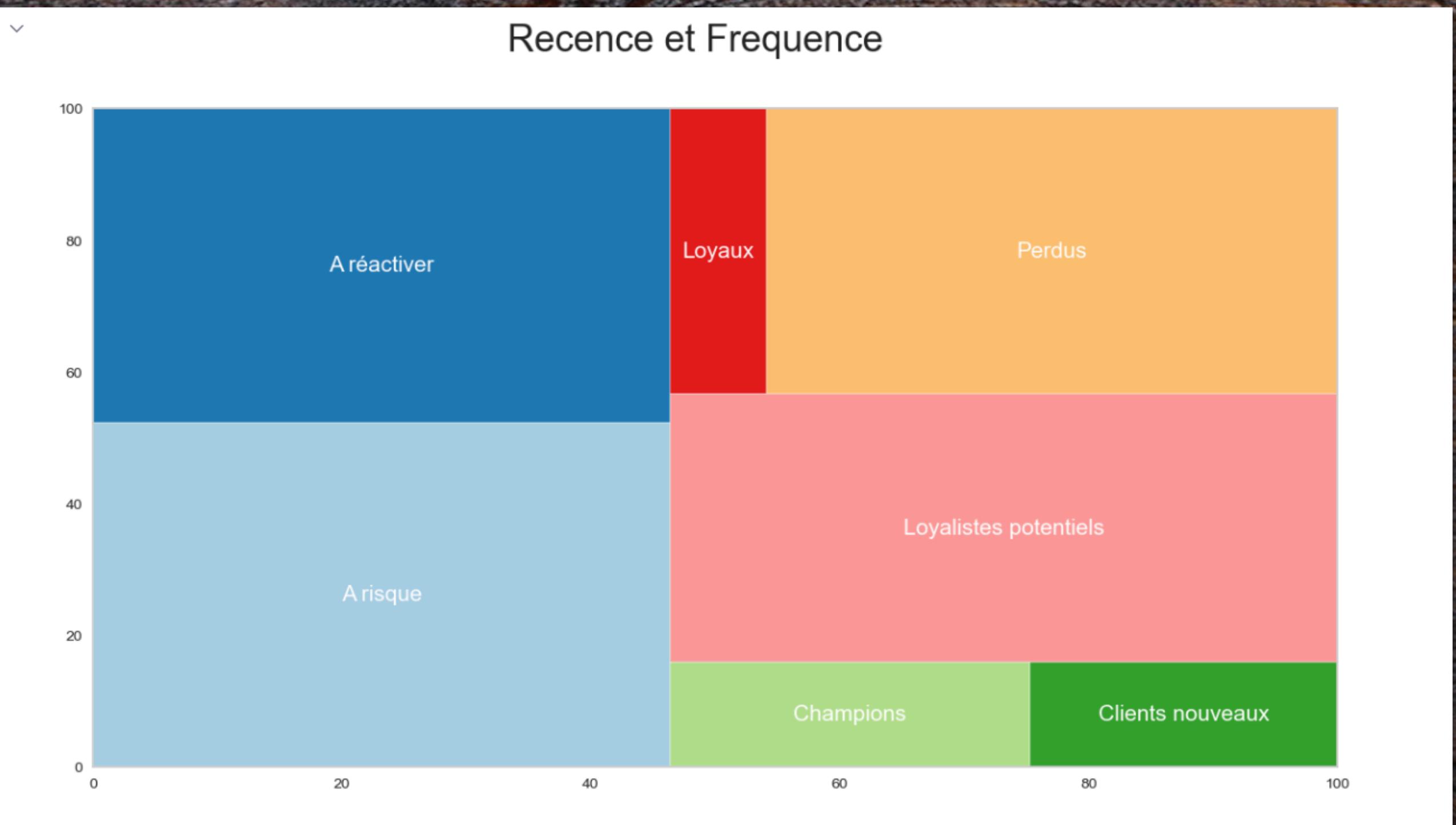
```
In 15 1 # Merge entre le 3 df
2 RF_df = recence_df.merge(frequence_df, on='customer_unique_id')
3 RFM_df = RF_df.merge(montant_df, on='customer_unique_id').drop(columns='order_purchase_timestamp')
4 display(RFM_df.head(3))
Executed at 2024.02.22 08:40:58 in 40ms
```

< 3 rows > 3 rows x 4 columns

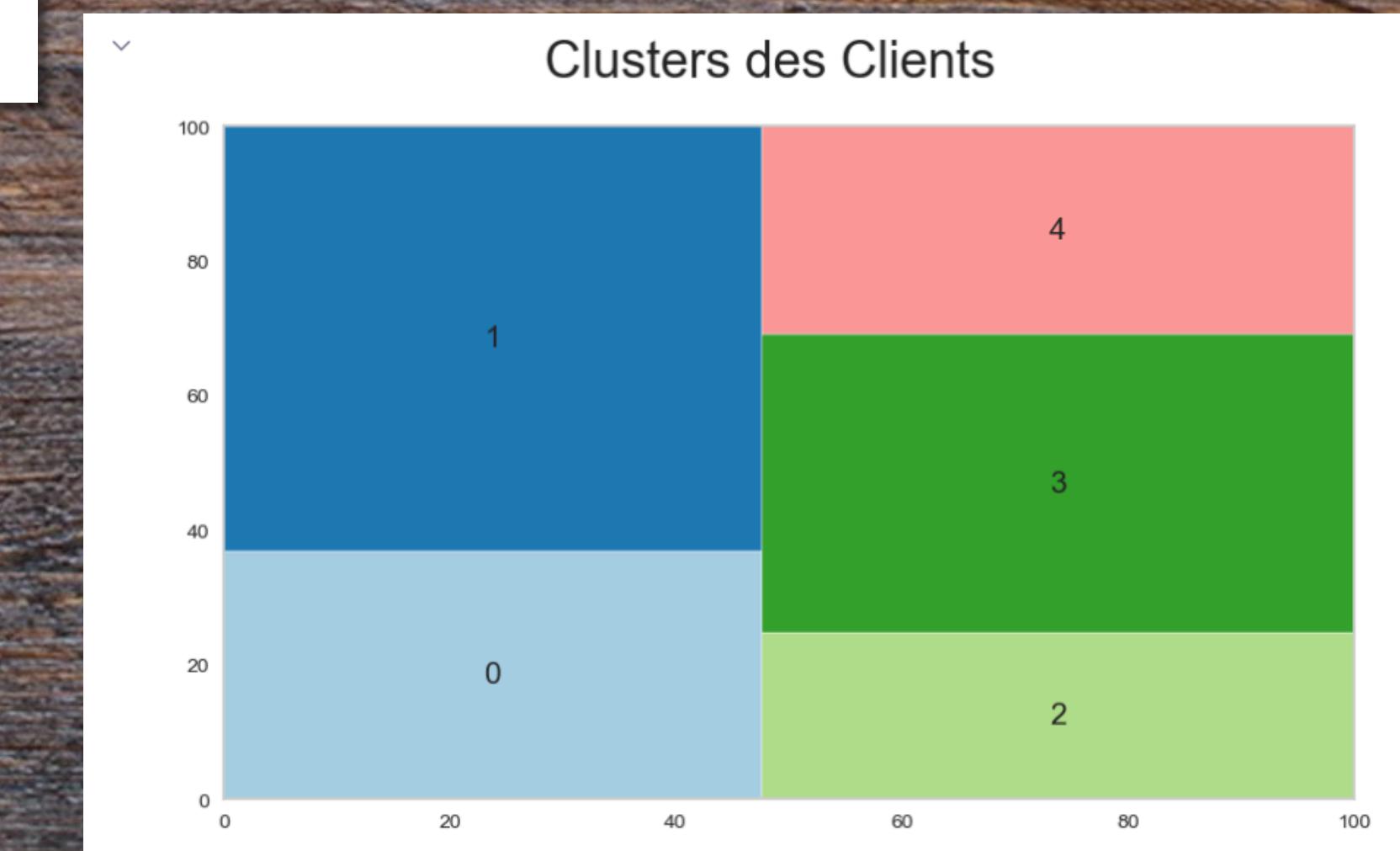
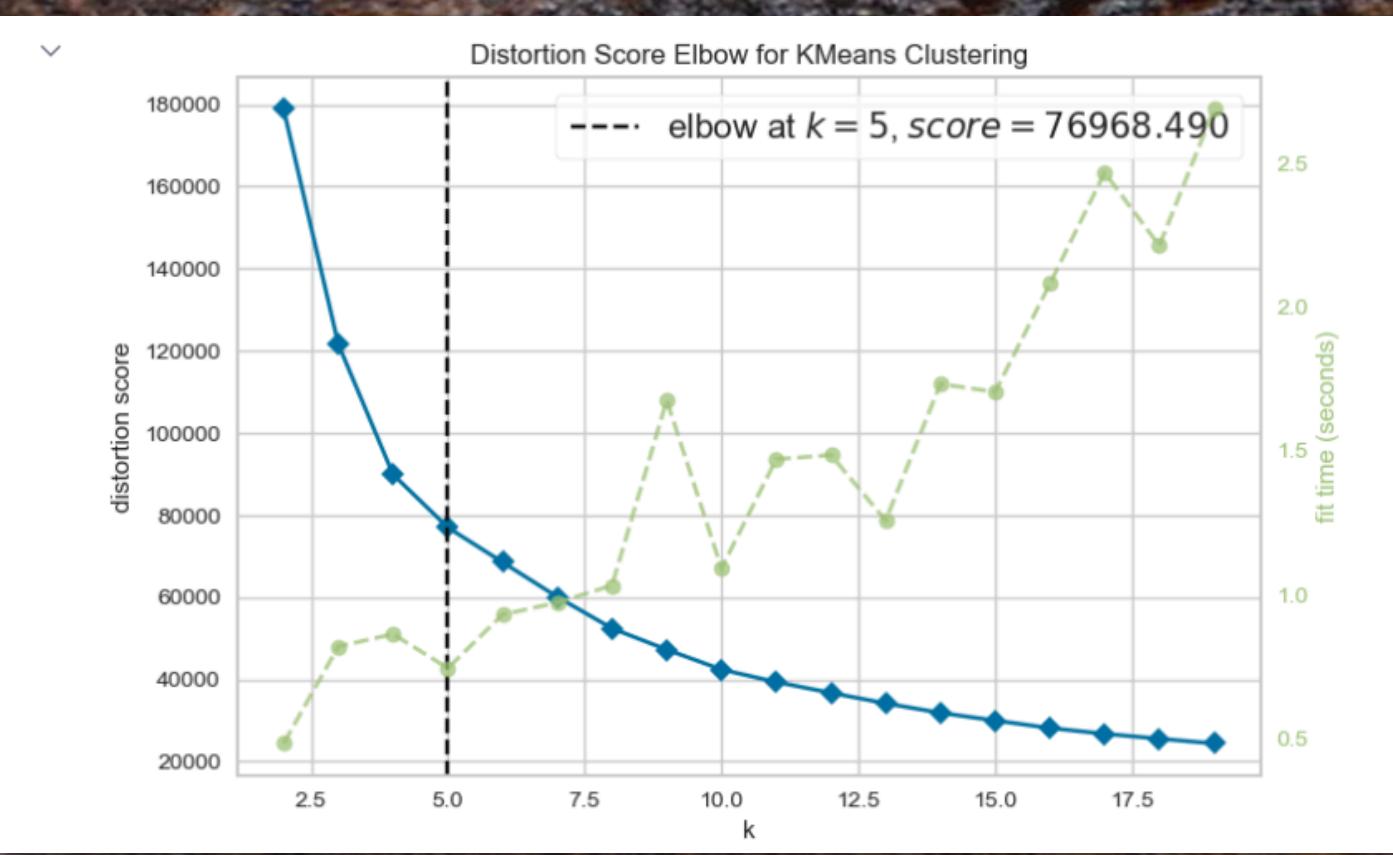
Static Output

	customer_unique_id	Recence	Frequence	Montant
0	0000366f3b9a7992bf8c76cfdf3221e2	162	1	141.90
1	0000b849f77a49e4a4ce2b2a4ca5be3f	165	1	27.19
2	0000f46a3911fa3c0805444483337064	587	1	86.22

Modélisation - Cluster manuel

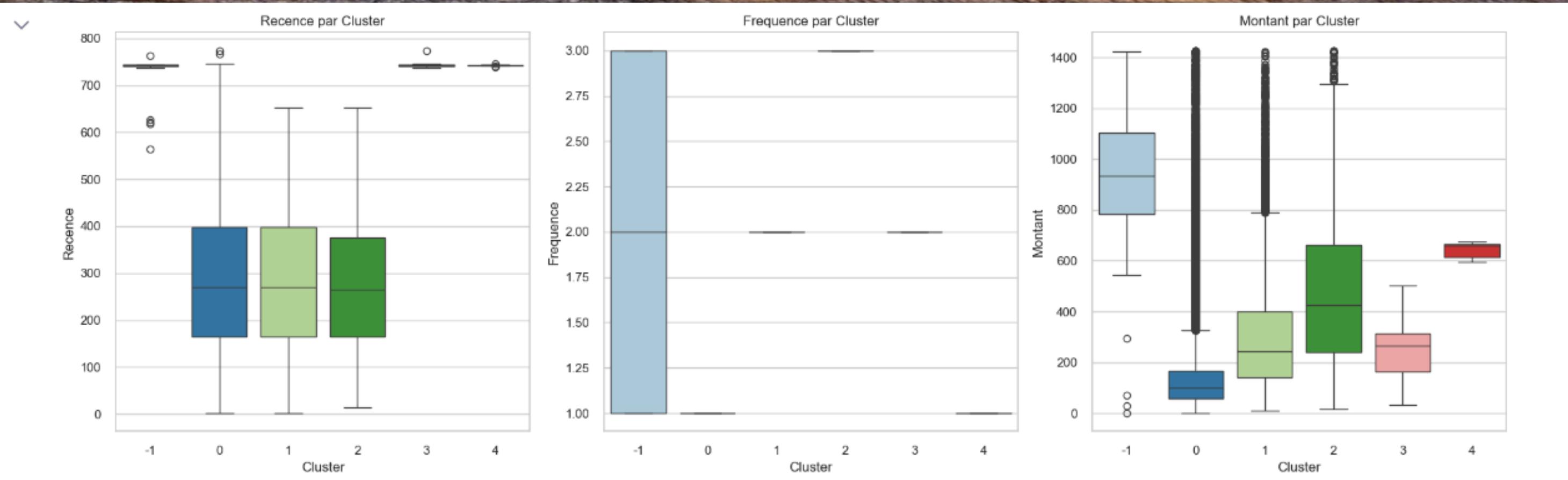


Modélisation - K Means



Modélisation - DB Scan

```
In 39 1 # Init DBScan
2 dbscan = DBSCAN(eps=0.5, min_samples=5)
3 # Fitting DBSCAN
4 clusters = dbscan.fit_predict(X_scaled)
5 # Ajout label cluster dans le df
6 RFM_df_k_1['Cluster'] = clusters
Executed at 2024.02.22 08:42:13 in 45s 782ms
```



Simulation

```
In 8 1 def prep_df(df, ratio):
2     '''Fonction qui divise le dataframe selon un ratio sélectionné
3     ARGS:
4         df: dataframe
5         ratio: ratio de division du df'''
6
7     df_sort = df.sort_values('Recence').reset_index()
8     num_lignes = round(df_sort.shape[0] * ratio)
9     df1 = df_sort.iloc[:num_lignes, :]
10    return df1.drop('index', axis=1)
11
12 Executed at 2024.02.22 08:42:24 in 6ms
13
14 In 9 1 def div_pas(df,t1, semaine):
15     '''Fonction pour diviser le df en semaines
16     ARGS:
17         df: dataframe
18         t1 : date limite
19         semaines : numéro des semaines à considerer'''
20
21     # Calculer le pas en jours en fonction du pas en semaines spécifié
22     step = 7 * semaine
23
24     # Trouver la date la plus ancienne et la plus récente dans la colonne 'Recence'
25     ancienne, nouv = min(df['Recence']), max(df['Recence'])
26     # Calculer la durée totale en jours
27     nb = nouv - ancienne
28     # Init liste pour stocker les différents pas de temps
29     T = []
30     t = t1
31
32     # Générer les pas de temps jusqu'à la date la plus récente
33     while t < nouv:
34         T.append(t)
35         t += step
36
37     # Supprimer le premier élément de la liste car il est égal à t1
38     T.pop(0)
39
40     return T
41
42 Executed at 2024.02.22 08:42:24 in 4ms
```



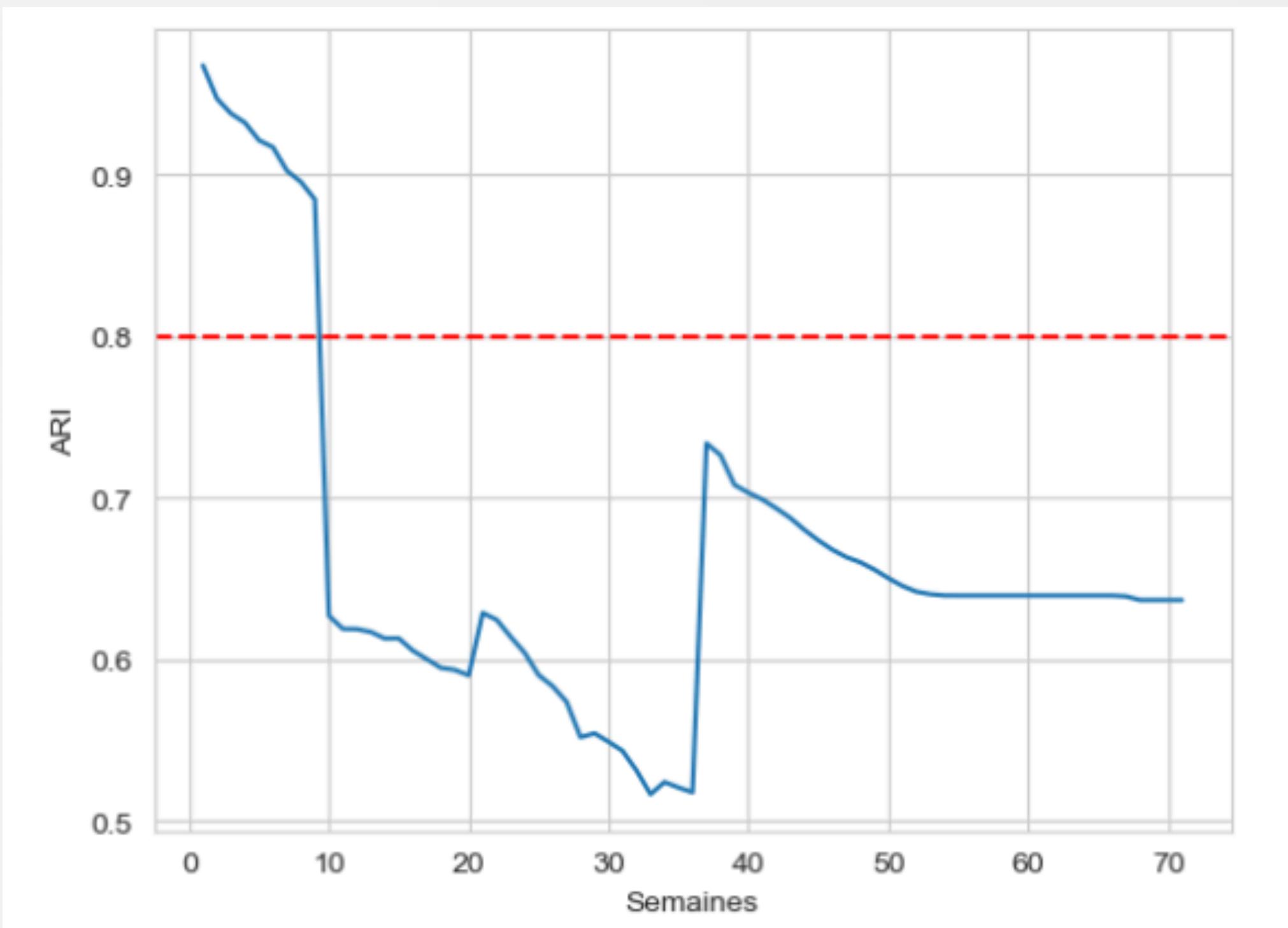
Simulation

ARI - Adjusted Rand Index

Une mesure utilisée en analyse de clustering pour évaluer la similarité entre deux ensembles de regroupements. Il compare les regroupements réels avec les regroupements prédicts par un algorithme de clustering, prenant en compte le hasard.



Simulation



Graphique ARI



Merci pour votre attention.

