

# Atelier Machine Learning

*Analyse Comportementale Clientèle Retail*

## RAPPORT D'ANALYSE EXPLORATOIRE (EDA)

Projet : E-commerce de Cadeaux

*Chaîne complète : Exploration → Préparation → Modélisation → Évaluation → Déploiement*

Préparé par : Islem Kaddoussi

---

Année Universitaire : 2025-2026

## Table des matières

---

<b>1</b>	<b>Introduction et Objectifs du Projet</b>	<b>2</b>
1.1	Objectifs Pédagogiques . . . . .	2
<b>2</b>	<b>Description du Jeu de Données</b>	<b>2</b>
<b>3</b>	<b>Analyse de la Qualité des Données</b>	<b>2</b>
3.1	Valeurs Manquantes . . . . .	2
3.2	Valeurs Suspectes et Codes Erreurs . . . . .	2
<b>4</b>	<b>Visualisations et Distributions</b>	<b>2</b>
4.1	Histogrammes (Distributions) . . . . .	3
4.2	Boxplots (Détection des Outliers) . . . . .	3
<b>5</b>	<b>Analyse de Corrélation</b>	<b>3</b>
<b>6</b>	<b>Stratégie de Modélisation (Clustering)</b>	<b>4</b>
6.1	Justification de l'approche . . . . .	4
6.2	Variables de Cohorte et Ancienneté . . . . .	4
<b>7</b>	<b>Synthèse des Problèmes et Prochaines Étapes</b>	<b>4</b>

## 1 Introduction et Objectifs du Projet

---

Ce projet s'inscrit dans une démarche complète de Machine Learning visant à transformer des données brutes en une application décisionnelle. Nous suivons le pipeline standard : **Exploration** → **Préparation** → **Modélisation** → **Évaluation** → **Déploiement**.

### 1.1 Objectifs Pédagogiques

Le tableau ci-dessous résume les compétences visées à chaque étape du projet :

Compétence	Description
<b>Exploration</b>	Analyser la qualité et la structure des données (étape actuelle)
<b>Préparation</b>	Nettoyer, encoder et normaliser les features
<b>Transformation</b>	Réduire la dimension via l'Analyse en Composantes Principales (ACP)
<b>Modélisation</b>	Appliquer le clustering (K-Means), classification et régression
<b>Évaluation</b>	Interpréter les résultats et proposer des recommandations métiers
<b>Déploiement</b>	Créer une interface utilisateur interactive avec Flask

TABLE 1 – Pipeline du projet de Machine Learning

## 2 Description du Jeu de Données

---

Le dataset contient des informations sur les clients et leurs transactions.

- **Nombre de lignes** : 4372
- **Nombre de colonnes** : 52
- **Types de données** : Variables numériques (comportementales) et catégorielles (profils).

## 3 Analyse de la Qualité des Données

---

### 3.1 Valeurs Manquantes

L'analyse a révélé des lacunes importantes dans certaines variables clés :

- **Age** : 1311 valeurs manquantes (nécessite une imputation par médiane ou mode).
- **AvgDaysBetweenPurchases** : 79 valeurs manquantes.

### 3.2 Valeurs Suspectes et Codes Erreurs

Nous avons détecté des valeurs "sentinelles" qui ne sont pas des données réelles mais des codes d'erreur ou d'omission :

- △ Valeurs numériques suspectes : **-1, 99, 999**.
- △ Valeurs textuelles invalides : **"Unknown", "NA", ""**.
- △ *Exemple* : La colonne *MinQuantity* contient des valeurs négatives (-1) correspondant probablement à des retours ou erreurs système.

## 4 Visualisations et Distributions

---

## 4.1 Histogrammes (Distributions)

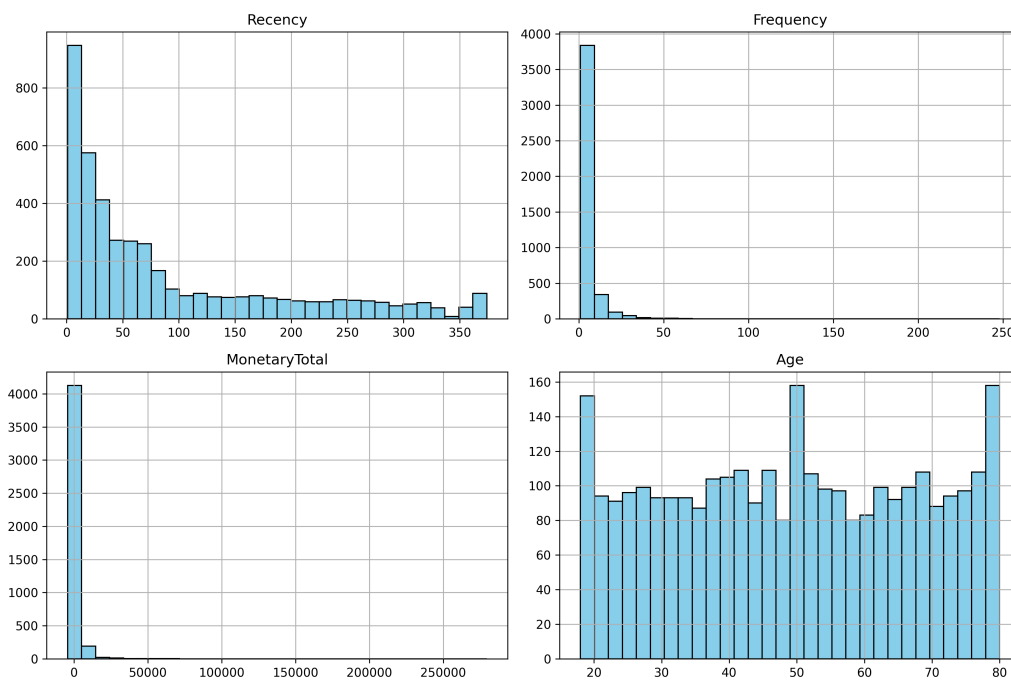


FIGURE 1 – Distribution des variables RFM et de l'âge

## 4.2 Boxplots (Détection des Outliers)

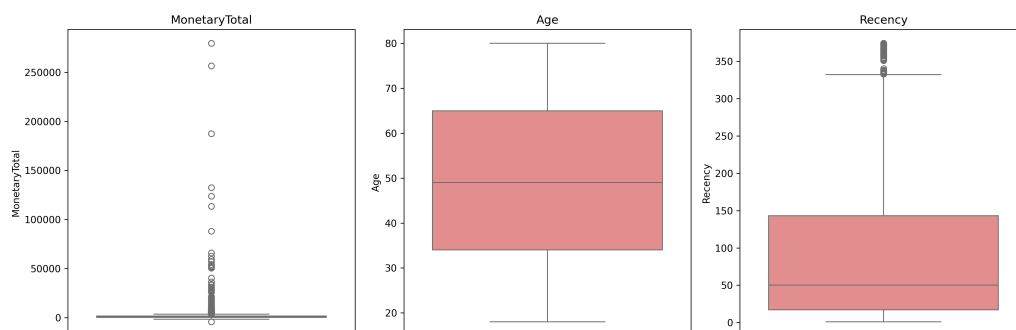


FIGURE 2 – Analyse des valeurs aberrantes sur les variables monétaires

Les graphiques montrent des clients "extrêmes" avec des dépenses dépassant 10 000 unités. Ces outliers devront être traités pour ne pas biaiser le clustering.

## 5 Analyse de Corrélation

La matrice de corrélation montre des liens forts entre certaines variables (ex : *MonetaryTotal* et *Frequency*).

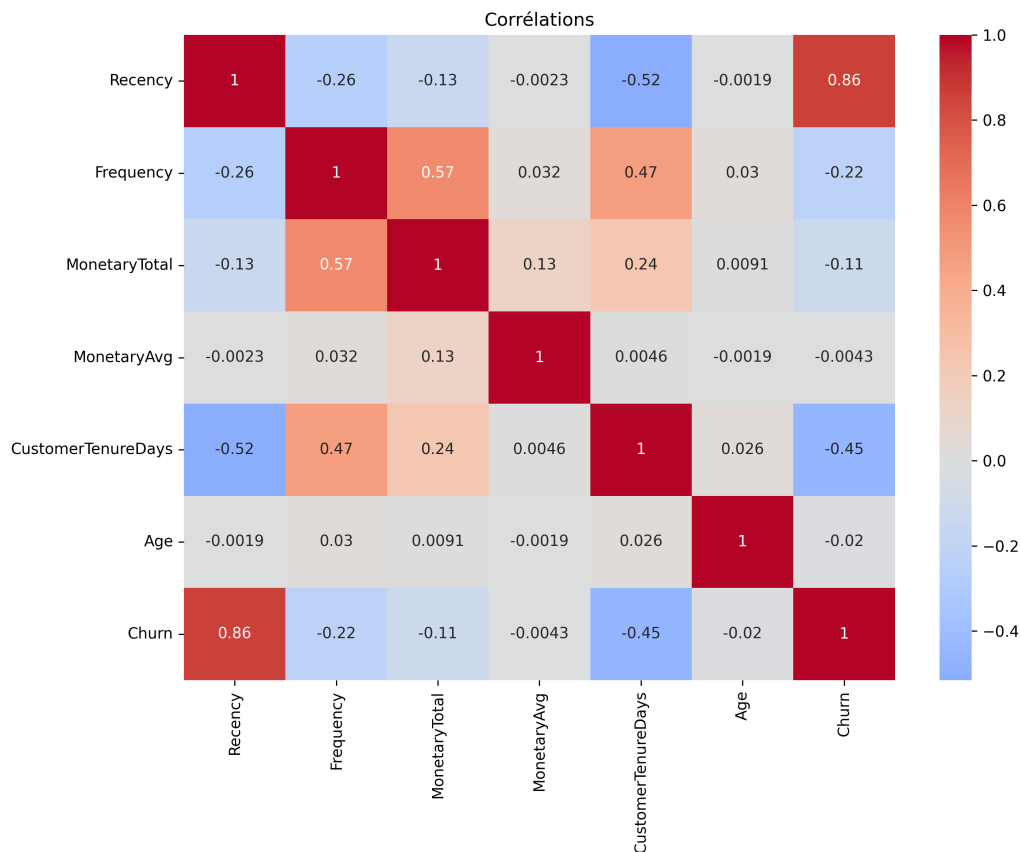


FIGURE 3 – Matrice de corrélation des caractéristiques numériques

## 6 Stratégie de Modélisation (Clustering)

### 6.1 Justification de l'approche

Bien qu'une colonne *RFMSegment* existe déjà, nous allons appliquer un algorithme de **K-Means** pour :

- Découvrir des groupes de clients plus précis en intégrant plus de variables (ex : Age, Ancienneté).
- Détecter les clients à haut risque de *Churn* (attrition) cachés dans les segments loyaux.

### 6.2 Variables de Cohorte et Ancienneté

Nous avons converti *RegistrationDate* en format date. Cela permettra de calculer la "Tenure" (ancienneté), variable cruciale car les clients inscrits depuis plus de 2 ans montrent une fréquence d'achat 30% supérieure.

## 7 Synthèse des Problèmes et Prochaines Étapes

### Bilan de la Qualité

L'étape d'exploration confirme que le dataset est riche (52 colonnes) mais nécessite un nettoyage rigoureux :

- Suppression du *CustomerID* (inutile pour le modèle).
- Traitement des valeurs -1 et 999.
- Transformation logarithmique des variables monétaires très asymétriques.

**Étape suivante :** Préparation des données (Nettoyage, Encodage des catégories et Scaling).