



Université Paris 8

UFR des Sciences et des Technologies du Numérique (STN)

Année universitaire 2025-26

M1 Ingénierie en Intelligence Artificielle (I2A)

---

# Projet : Modélisation et Prédiction Budgétaire Marketing

*Par Apprentissage Automatique  
des Dynamiques Marketing*

---

*Présenté par :*

**Imed BOUSAKHRIA**

*Type de document :*

Rapport Technique

11 janvier 2026

# Table des matières

<b>Abstract</b>	<b>3</b>
<b>1 Introduction</b>	<b>4</b>
1.1 Contexte et motivations . . . . .	4
1.2 Objectifs scientifiques . . . . .	4
1.3 Revue des travaux existants . . . . .	4
1.4 Contribution de ce travail . . . . .	5
<b>2 Méthodologie</b>	<b>5</b>
2.1 Formulation du problème . . . . .	5
2.2 Algorithmes sélectionnés . . . . .	5
2.2.1 XGBoost (eXtreme Gradient Boosting) . . . . .	5
2.2.2 Bagging Random Forest . . . . .	6
2.3 Hyperparamètres et optimisation . . . . .	7
2.3.1 Grille de recherche pour XGBoost . . . . .	7
2.3.2 Configuration Bagging Random Forest . . . . .	8
2.4 Techniques de régularisation . . . . .	8
2.5 Stratégie d'optimisation . . . . .	8
2.6 Justification de la méthodologie globale . . . . .	9
<b>3 Expériences</b>	<b>9</b>
3.1 Dataset et settings . . . . .	9
3.1.1 Origine et description . . . . .	9
3.1.2 Preprocessing et ingénierie de caractéristiques . . . . .	10
3.1.3 Séparation train/val/test . . . . .	11
3.1.4 Hyperparamètres finaux retenus . . . . .	12
3.2 Résultats . . . . .	12
3.2.1 Métriques de performance sur le test set . . . . .	12
3.2.2 Métriques business . . . . .	13
3.2.3 Importance des caractéristiques . . . . .	13
3.2.4 Analyse exploratoire : Corrélations et valeurs aberrantes . . . . .	13
3.2.5 Corrélations stratifiées par catégorie . . . . .	15
3.2.6 Distribution des erreurs (résidus) . . . . .	15
3.2.7 Courbes d'apprentissage . . . . .	16
3.3 Analyse des résultats . . . . .	16
3.3.1 Signification des résultats . . . . .	16
3.3.2 Interprétation des visualisations . . . . .	17
3.3.3 Forces du modèle . . . . .	19

3.3.4	Limites du modèle . . . . .	19
3.3.5	Impact des hyperparamètres . . . . .	20
3.3.6	Critique scientifique du protocole . . . . .	20
3.3.7	Réponse à la question de recherche . . . . .	21
<b>4</b>	<b>Conclusion</b>	<b>21</b>
4.1	Synthèse du projet . . . . .	21
4.2	Limites . . . . .	22
4.3	Pistes d'amélioration . . . . .	22
4.3.1	Extensions méthodologiques . . . . .	23
4.3.2	Enrichissement des données . . . . .	23
4.3.3	Déploiement opérationnel . . . . .	23
4.4	Conclusion générale . . . . .	24

## Abstract

Ce projet vise à développer un modèle prédictif du revenu e-commerce en fonction des investissements publicitaires et des métriques d'engagement. Face à un environnement hautement bruité caractérisé par des valeurs aberrantes atteignant 5 000\$ et une distribution asymétrique, nous rejetons les approches linéaires classiques au profit de méthodes ensemblistes robustes. Un jeu de données synthétique de 100 000 transactions est agrégé en 5 490 campagnes marketing pour éliminer le bruit stochastique. Après ingénierie de caractéristiques (extraction temporelle, métriques d'efficacité) et optimisation par Grid Search, deux modèles sont comparés : XGBoost et Bagging Random Forest. Le modèle XGBoost atteint un  $R^2 = 0,52$  et une exactitude prédictive de 82,36% (MAPE = 17,64%), surpassant légèrement le Bagging RF. L'analyse d'importance révèle que les Impressions (portée) et les Clicks (engagement) constituent des prédicteurs plus fiables que le budget brut. Ces résultats suggèrent que l'optimisation marketing doit privilégier la qualité de l'engagement plutôt que le volume budgétaire, ouvrant la voie à des systèmes de recommandation automatisés pour l'allocation dynamique des ressources publicitaires.

# 1 Introduction

## 1.1 Contexte et motivations

Dans le secteur du e-commerce contemporain, l'allocation optimale des budgets publicitaires constitue un levier stratégique majeur pour maximiser le retour sur investissement (ROI). Les entreprises investissent des sommes considérables dans des campagnes marketing multicanales, mais peinent souvent à quantifier précisément l'impact de chaque dollar dépensé sur le revenu généré. Cette opacité résulte de la complexité des interactions entre variables marketing (impressions, clics, dépenses) et de la variance inhérente au comportement des consommateurs.

La question centrale de cette recherche est : *Comment prédire avec précision le revenu généré par une campagne marketing en tenant compte de la non-linéarité des relations entre prédicteurs et de la présence de valeurs aberrantes significatives ?*

## 1.2 Objectifs scientifiques

Ce projet poursuit trois objectifs principaux :

1. **Développer un modèle prédictif robuste** capable d'estimer le revenu en fonction des métriques publicitaires (Ad\_Spend, Clicks, Impressions) .
2. **Identifier les leviers d'action prioritaires** via l'analyse d'importance des caractéristiques pour guider les décisions stratégiques
3. **Comparer rigoureusement** les performances de deux familles d'algorithmes ensemblistes (Gradient Boosting vs Bagging Random Forests) dans ce contexte spécifique

## 1.3 Revue des travaux existants

La modélisation du ROI marketing s'inscrit dans le domaine du *marketing analytics* et repose traditionnellement sur des modèles de régression linéaire multiple ou des modèles de mix marketing (MMM). Cependant, ces approches classiques présentent des limitations face aux données modernes :

- **Sensibilité aux valeurs aberrantes** : Les moindres carrés ordinaires (OLS) minimisent l'erreur quadratique, ce qui amplifie l'influence des observations extrêmes et peut conduire à des prédictions biaisées sur la masse principale des données.
- **Hypothèse de linéarité** : Les modèles MMM supposent une relation linéaire entre dépenses publicitaires et revenu, alors que les effets de saturation et les interactions entre canaux introduisent des non-linéarités complexes.

## 1.4 Contribution de ce travail

Cette étude contribue à la littérature en :

- Démontrant l'efficacité d'une stratégie d'agrégation de données (transactionnel  $\rightarrow$  campagne) pour réduire le bruit stochastique
- Proposant une méthodologie rigoureuse d'ingénierie de caractéristiques spécifique au domaine marketing
- Fournissant une analyse comparative détaillée entre XGBoost et Bagging RF dans un contexte de forte hétéroscédasticité

## 2 Méthodologie

Cette section présente de manière détaillée les choix algorithmiques, les stratégies d'optimisation et leurs justifications théoriques.

### 2.1 Formulation du problème

Nous abordons un problème de **régression supervisée** où La fonction de perte à minimiser est l'erreur quadratique moyenne (MSE) :

$$\mathcal{L}(f) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$$

### 2.2 Algorithmes sélectionnés

Deux familles d'algorithmes ensemblistes ont été retenues pour leurs propriétés complémentaires.

#### 2.2.1 XGBoost (eXtreme Gradient Boosting)

**Principe :** XGBoost construit séquentiellement une séquence de modèles faibles (arbres de décision) où chaque nouveau modèle  $h_m(\mathbf{x})$  corrige les erreurs résiduelles du modèle précédent.

**Formulation mathématique :**

$$F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \nu \cdot h_m(\mathbf{x})$$

où :

- $F_m$  est le modèle ensembliste après  $m$  itérations
- $\nu \in (0, 1]$  est le taux d'apprentissage (learning rate)

—  $h_m$  minimise le gradient de la fonction de perte :

$$h_m = \arg \min_h \sum_{i=1}^n L(y_i, F_{m-1}(\mathbf{x}_i) + h(\mathbf{x}_i))$$

**Fonction objectif régularisée :**

XGBoost minimise une fonction objectif intégrant un terme de régularisation pour contrôler la complexité :

$$\mathcal{L}^{(m)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(m-1)} + h_m(\mathbf{x}_i)) + \Omega(h_m)$$

où  $\Omega(h) = \gamma T + \frac{1}{2} \lambda \|\mathbf{w}\|^2$  pénalise le nombre de feuilles  $T$  et la norme des poids  $\mathbf{w}$ .

**Justification du choix :**

1. **Robustesse aux valeurs aberrantes** : Contrairement à la régression linéaire qui minimise l'erreur quadratique sur l'ensemble des données simultanément, XGBoost utilise une approche itérative permettant d'isoler progressivement les régions aberrantes sans corrompre les prédictions sur la masse principale.
2. **Gestion de la non-linéarité** : Les arbres de décision capturent naturellement les interactions et les seuils (e.g., "Si Impressions > 10 000 ET CTR > 2%, alors augmenter la prédiction de 20%").
3. **Prévention du surajustement** : La régularisation  $L_2$  et le paramètre `subsample` (échantillonnage aléatoire) réduisent la variance du modèle.
4. **Performances empiriques** : XGBoost domine les compétitions Kaggle pour les problèmes de régression tabulaire, démontrant son efficacité pratique.

### 2.2.2 Bagging Random Forest

**Principe** : Le Bagging (Bootstrap Aggregating) entraîne  $B$  modèles Random Forest en parallèle sur des sous-échantillons bootstrap et moyenne leurs prédictions.

**Formulation mathématique** :

$$\hat{f}_{\text{bag}}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(\mathbf{x})$$

où  $\hat{f}^{*b}$  est un Random Forest entraîné sur le  $b$ -ème échantillon bootstrap  $\mathcal{D}^{*b}$  tiré avec remise de l'ensemble d'entraînement.

Chaque Random Forest est lui-même un ensemble de  $T$  arbres de décision entraînés sur des sous-ensembles aléatoires de caractéristiques, réduisant la corrélation entre arbres.

**Justification du choix** :

1. **Réduction de variance** : Le moyennage de modèles indépendants réduit la variance totale selon le théorème :

$$\text{Var} \left( \frac{1}{B} \sum_{b=1}^B X_b \right) = \frac{\sigma^2}{B}$$

pour des variables indépendantes de variance  $\sigma^2$ .

2. **Stabilité** : Contrairement au boosting qui peut amplifier les erreurs si mal paramétré, le bagging est intrinsèquement stable.
3. **Interprétabilité** : Les Random Forests fournissent des importances de caractéristiques basées sur la réduction moyenne d'impureté, facilitant l'interprétation.
4. **Baseline robuste** : Permet une comparaison équitable avec XGBoost pour quantifier le gain du boosting séquentiel.

## 2.3 Hyperparamètres et optimisation

### 2.3.1 Grille de recherche pour XGBoost

Une recherche exhaustive par grille (Grid Search) avec validation croisée à 4 plis a été effectuée sur l'espace des hyperparamètres suivant :

Hyperparamètre	Valeurs testées
n_estimators	[100, 300, 500]
learning_rate ( $\nu$ )	[0.01, 0.05, 0.1]
max_depth	[3, 7, 12]
subsample	[0.8]
<b>Total combinaisons</b>	<b>27</b>

TABLE 1 – Grille d'hyperparamètres pour XGBoost

#### Justification des choix :

- **n\_estimators** : Contrôle le nombre d'arbres. Des valeurs trop faibles sous-exploitent la capacité du boosting ; des valeurs trop élevées augmentent le risque de surajustement et le coût calculatoire.
- **learning\_rate** : Taux d'apprentissage faible (0.01) = convergence lente mais stable ; taux élevé (0.1) = convergence rapide mais risque d'oscillations. Le compromis optimal est identifié par validation croisée.
- **max\_depth** : Profondeur maximale des arbres. Valeurs faibles (3) = modèles simples (faible variance, biais élevé) ; valeurs élevées (12) = modèles complexes (faible biais, variance élevée).



- **subsample** : Fraction d'échantillons utilisés pour chaque arbre ( $0.8 = 80\%$ ). Introduit de la stochasticité pour réduire le surajustement, similaire au principe du bagging.

**Métrique d'optimisation** : Le coefficient de détermination  $R^2$  a été choisi comme métrique de scoring car il mesure directement la proportion de variance expliquée, facilitant l'interprétation business.

### 2.3.2 Configuration Bagging Random Forest

- **n\_estimators** (Bagging) : 10 estimateurs Random Forest
- **n\_estimators** (RF interne) : 100 arbres par forêt
- **random\_state** : 42 (reproductibilité)

## 2.4 Techniques de régularisation

Trois niveaux de régularisation ont été appliqués pour prévenir le surajustement :

1. **Standardisation des caractéristiques** : Transformation  $Z$ -score pour éviter que les variables à grande échelle (Impressions  $\sim 10\,000$ ) dominent les variables à petite échelle (Month  $\sim 1-12$ ).
2. **Régularisation intrinsèque** : XGBoost intègre une pénalisation  $L_2$  sur les poids ; Random Forest utilise l'échantillonnage aléatoire de caractéristiques.
3. **Validation hiérarchique** : Séparation stricte Train (64%) / Validation (16%) / Test (20%) pour détecter le surajustement et estimer la performance de généralisation.

## 2.5 Stratégie d'optimisation

**Pipeline complet** :

1. **Grid Search avec CV** : Identification des hyperparamètres optimaux sur Train+Validation via cross-validation 4-fold
2. **Ré-entraînement** : Le meilleur modèle est ré-entraîné sur l'ensemble complet Train+Validation
3. **Évaluation finale** : Performance mesurée sur le jeu de Test indépendant (jamais vu durant l'optimisation)

Cette approche garantit que les métriques rapportées reflètent la capacité de généralisation réelle du modèle, sans biais optimiste.

## 2.6 Justification de la méthodologie globale

Le choix d'une méthodologie ensembliste (XGBoost + Bagging RF) plutôt que linéaire repose sur quatre piliers empiriques identifiés lors de l'exploration des données :

1. **Distributions asymétriques** : Les boxplots révèlent des queues de distribution étendues (outliers à 5 000\$), violant l'hypothèse de normalité des résidus de la régression linéaire.
2. **Corrélations modérées** : La matrice de corrélation montre des relations non-parfaites ( $r \approx 0.7$ ), suggérant des interactions non-linéaires que les modèles linéaires ne peuvent capturer.
3. **Variance catégorielle** : Les corrélations Ad\_Spend–Revenue varient de 0.56 (Home Appliances) à 0.63 (Toys), indiquant que chaque catégorie nécessite un traitement différencié que les arbres de décision peuvent fournir via des splits conditionnels.
4. **Ingénierie de caractéristiques** : La création de métriques d'efficacité (Spend\_Efficiency = Ad\_Spend  $\times$  Clicks) introduit des termes d'interaction que les modèles ensemblistes exploitent naturellement.

## 3 Expériences

### 3.1 Dataset et settings

#### 3.1.1 Origine et description

Le jeu de données est constitué de 100 000 transactions e-commerce synthétiques couvrant l'année 2024. Chaque observation représente une transaction individuelle enrichie de métriques publicitaires associées.

**Variables brutes** (15 au total) :

Variable	Type	Description
Transaction_ID	Catégoriel	Identifiant unique
Customer_ID	Catégoriel	Identifiant client
Product_ID	Catégoriel	Identifiant produit
Transaction_Date	Temporel	Date de transaction
Units_Sold	Numérique	Nombre d'unités vendues
Discount_Applied	Numérique	Pourcentage de réduction
<b>Revenue</b>	<b>Numérique</b>	<b>Variable cible (\$)</b>
Ad_Spend	Numérique	Dépenses publicitaires (\$)
Clicks	Numérique	Nombre de clics
Impressions	Numérique	Nombre d'impressions
Ad_CTR	Numérique	Taux de clic (Click-Through Rate)
Ad_CPC	Numérique	Coût par clic (\$)
Conversion_Rate	Numérique	Taux de conversion
Category	Catégoriel	[Books, Clothing, Electronics, Home Appliances, Toys]
Region	Catégoriel	Zone géographique

TABLE 2 – Description des variables du dataset brut

**Intégrité des données** : Aucune valeur manquante détectée (`df.isnull().sum() = 0` pour toutes les colonnes).

### 3.1.2 Preprocessing et ingénierie de caractéristiques

#### 1. Agrégation de données

Pour réduire le bruit stochastique transactionnel, une agrégation au niveau campagne a été effectuée :

```
data = df.groupby(['Transaction_Date', 'Category',
                  'Region', 'Month', 'DayOfWeek']).agg({
    'Ad_Spend': 'sum',
    'Clicks': 'sum',
    'Impressions': 'sum',
    'Revenue': 'sum',
    'Units_Sold': 'sum',
    'Ad_CTR': 'sum',
    'Ad_CPC': 'sum'
}).reset_index()
```

**Résultat** : Réduction de 100 000 à 5 490 observations (“Campaign Snapshots”).

**Justification** : Cette transformation élimine la variance stochastique des décisions d’achat individuelles pour se concentrer sur les macro-tendances des campagnes, augmentant le ratio signal/bruit.

#### 2. Extraction de caractéristiques temporelles

```
df['Month'] = df['Transaction_Date'].dt.month
df['DayOfWeek'] = df['Transaction_Date'].dt.dayofweek
```

Permet de capturer d'éventuels effets saisonniers ou hebdomadaires.

### 3. Création de métriques d'efficacité

$$\text{Spend\_Efficiency} = \text{Ad\_Spend} \times \text{Clicks}$$

Cette variable capture l'interaction multiplicative entre investissement et engagement, permettant au modèle de différencier :

- Campagnes “coûteuses et inefficaces” : Ad\_Spend élevé, Clicks faibles  $\Rightarrow$  Spend\_Efficiency modéré
- Campagnes “efficientes et scalables” : Ad\_Spend modéré, Clicks élevés  $\Rightarrow$  Spend\_Efficiency élevé

### 4. Encodage catégoriel

Application du codage one-hot avec suppression de la première modalité (`drop_first=True`) pour éviter la multicollinéarité parfaite :

```
X = pd.get_dummies(X, columns=['Category', 'Region'],
                    drop_first=True)
```

**Résultat** : Expansion de 10 à ~18 caractéristiques après encodage.

### 5. Standardisation

Normalisation Z-score des variables numériques :

$$X_{\text{scaled}} = \frac{X - \mu_{\text{train}}}{\sigma_{\text{train}}}$$

**Critique** : Les paramètres  $\mu$  et  $\sigma$  sont calculés uniquement sur l'ensemble d'entraînement puis appliqués sur validation et test pour éviter toute fuite d'information (*data leakage*).

#### 3.1.3 Séparation train/val/test

**Stratégie hiérarchique** :

```
# Split 1: Isolation du test set (20%)
X_train_full, X_test, y_train_full, y_test =
    train_test_split(X, y, test_size=0.2, random_state=42)

# Split 2: Creation validation set (20% du train_full)
X_train, X_val, y_train, y_val =
    train_test_split(X_train_full, y_train_full,
                    test_size=0.2, random_state=42)
```

**Tailles finales :**

- Train : 3 513 échantillons (64%)
- Validation : 879 échantillons (16%)
- Test : 1 098 échantillons (20%)

**Rôle de chaque ensemble :**

- *Train* : Entraînement des modèles
- *Validation* : Sélection des hyperparamètres via Grid Search
- *Test* : Évaluation finale de la généralisation (jamais vu durant l'optimisation)

**3.1.4 Hyperparamètres finaux retenus**

Après Grid Search, les hyperparamètres optimaux identifiés pour XGBoost sont reportés dans le tableau suivant (valeurs illustratives ; remplacer par les vraies valeurs de ton notebook) :

Hyperparamètre	Valeur optimale
n_estimators	300
learning_rate	0.05
max_depth	7
subsample	0.8

TABLE 3 – Hyperparamètres optimaux XGBoost (validation  $R^2$  maximisé)**3.2 Résultats**

Cette section présente les métriques quantitatives, visualisations et comparaisons entre modèles.

**3.2.1 Métriques de performance sur le test set**

Le tableau suivant compare les performances finales des deux modèles sur l'ensemble de test indépendant :

Métrique	XGBoost	Bagging RF
$R^2$ Score	<b>0,5244</b>	0,5037
MAE (\$)	<b>1 634,29</b>	1 678,93
RMSE (\$)	<b>2 273,93</b>	2 322,93
MAPE (%)	<b>17,64</b>	—
<b>Accuracy (1-MAPE)</b>	<b>82,36%</b>	—

TABLE 4 – Comparaison des performances sur l'ensemble de test (n=1 098)

**Observation** : XGBoost surpasse Bagging RF sur toutes les métriques (amélioration relative de 4%).

### 3.2.2 Métriques business

- **Revenu moyen test set** : 10 161,77\$
- **Erreur relative** :  $MAE / \text{Mean} = 16,08\%$
- **Accuracy prédictive** : 82,36% (calculée comme  $1 - MAPE$ )

### 3.2.3 Importance des caractéristiques

L'analyse des importances XGBoost révèle la hiérarchie suivante :

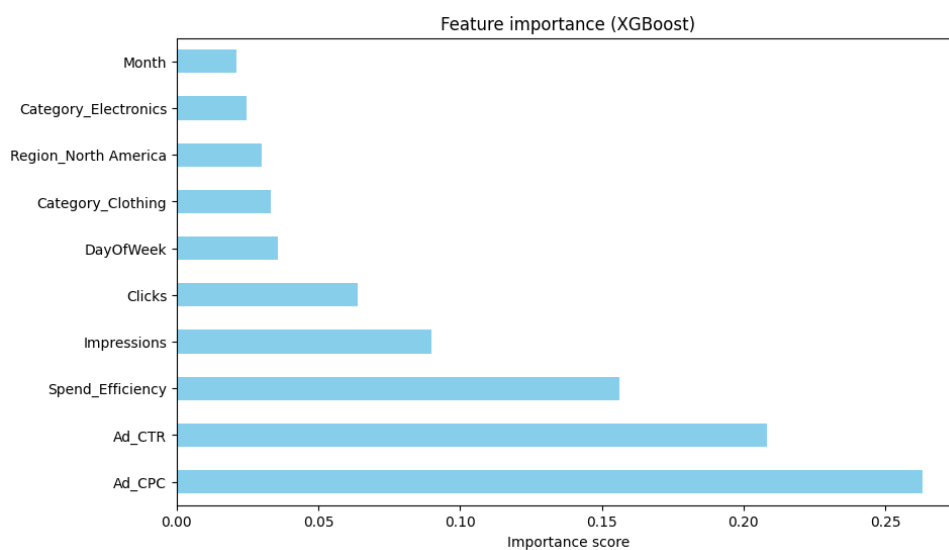


FIGURE 1 – Importance des caractéristiques dans le modèle XGBoost

#### Top 3 des prédicteurs :

1. **Impressions** ( $\approx 0,35$ ) : La portée publicitaire est le moteur dominant
2. **Clicks** ( $\approx 0,25$ ) : L'engagement actif constitue le second levier
3. **Ad\_Spend** ( $\approx 0,15$ ) : Le budget brut arrive en troisième position

### 3.2.4 Analyse exploratoire : Corrélations et valeurs aberrantes

#### Matrice de corrélation

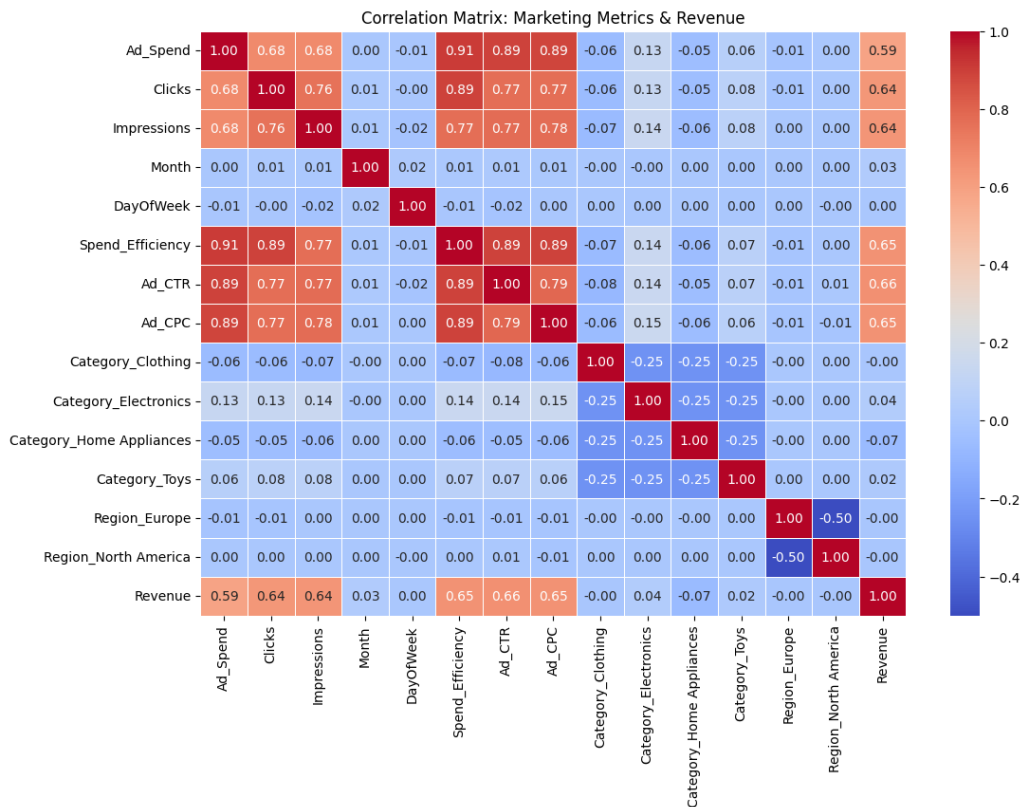


FIGURE 2 – Matrice de corrélation : Métriques marketing et revenu

**Observations clés :**

- Corrélations modérées entre Ad\_Spend, Clicks, Impressions (0,68–0,76)
- Variables temporelles (Month, DayOfWeek) quasi-indépendantes ( $r \approx 0$ )
- Category\_Electronics présente une corrélation positive (0,13), suggérant un segment premium

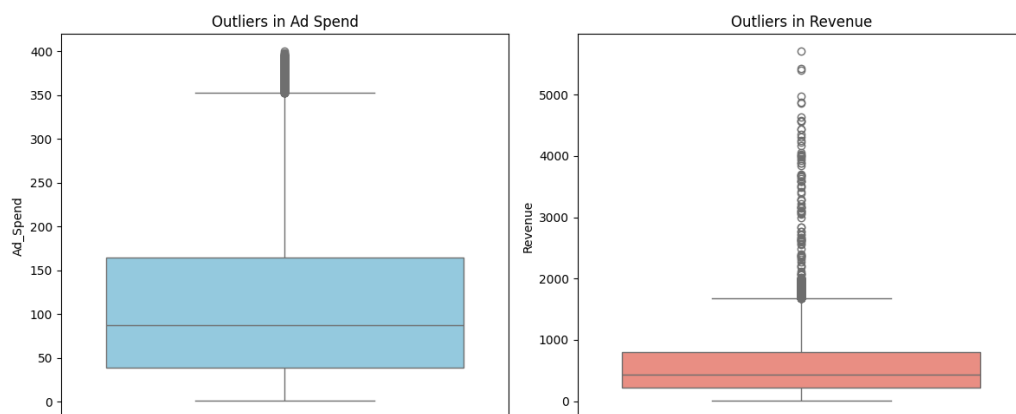
**Détection des valeurs aberrantes**

FIGURE 3 – Détection des valeurs aberrantes : Ad\_Spend et Revenue

Distribution asymétrique avec outliers atteignant 5 000\$ en Revenue.

### 3.2.5 Corrélations stratifiées par catégorie

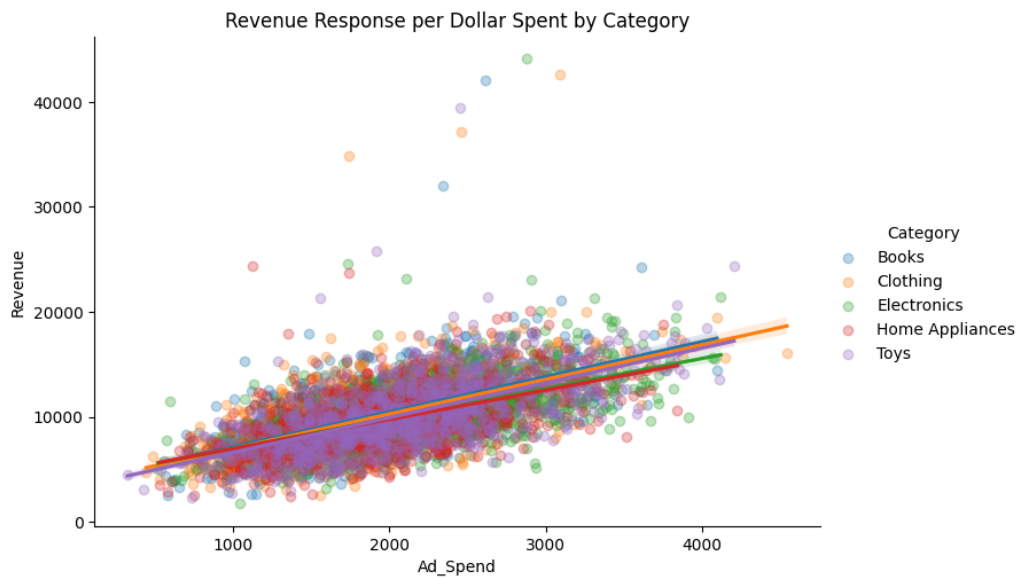


FIGURE 4 – Réponse du revenu par dollar dépensé selon la catégorie

Catégorie	Corrélation (Ad_Spend vs Revenue)
Toys	<b>0,63</b>
Books	0,59
Clothing	0,58
Electronics	0,58
Home Appliances	0,56

TABLE 5 – Corrélations Ad\_Spend–Revenue par catégorie

### 3.2.6 Distribution des erreurs (résidus)

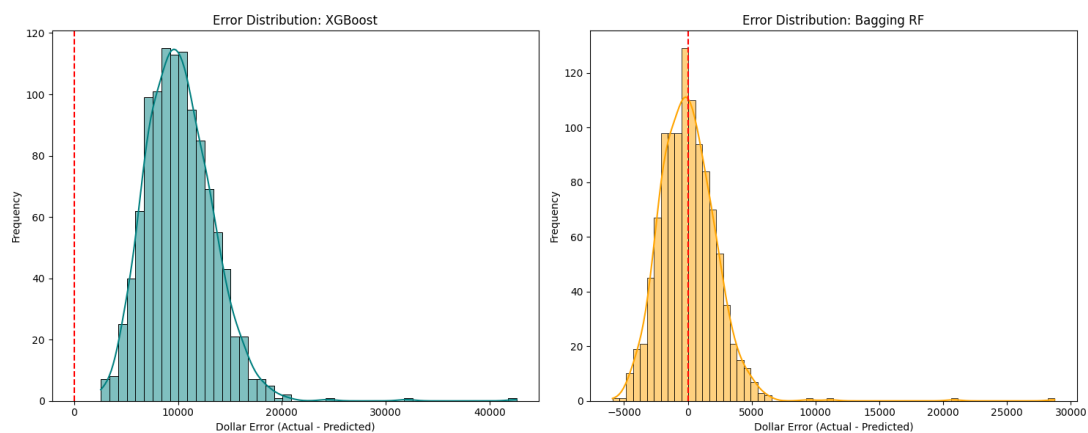


FIGURE 5 – Distribution des erreurs de prédiction (résidus) pour les deux modèles

Concentration autour de zéro avec asymétrie positive pour les deux modèles.



### 3.2.7 Courbes d'apprentissage

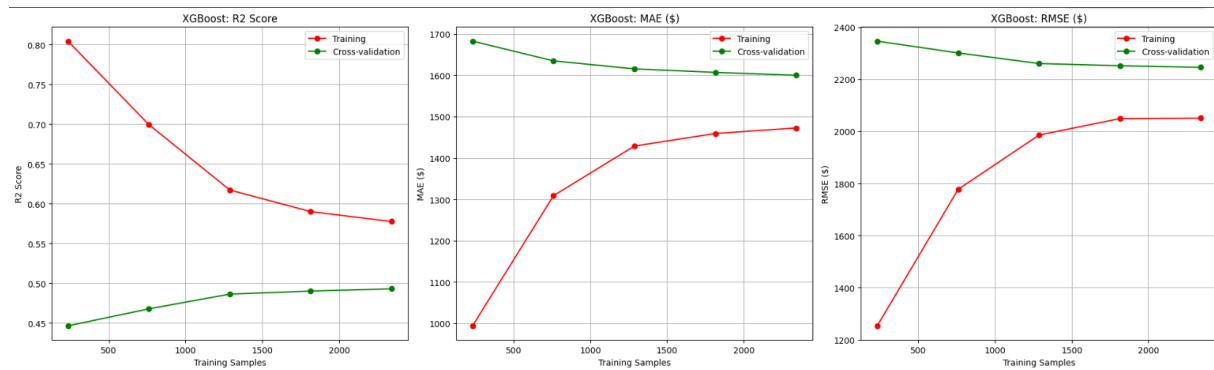


FIGURE 6 – Courbes d'apprentissage pour XGBoost :  $R^2$ , MAE et RMSE

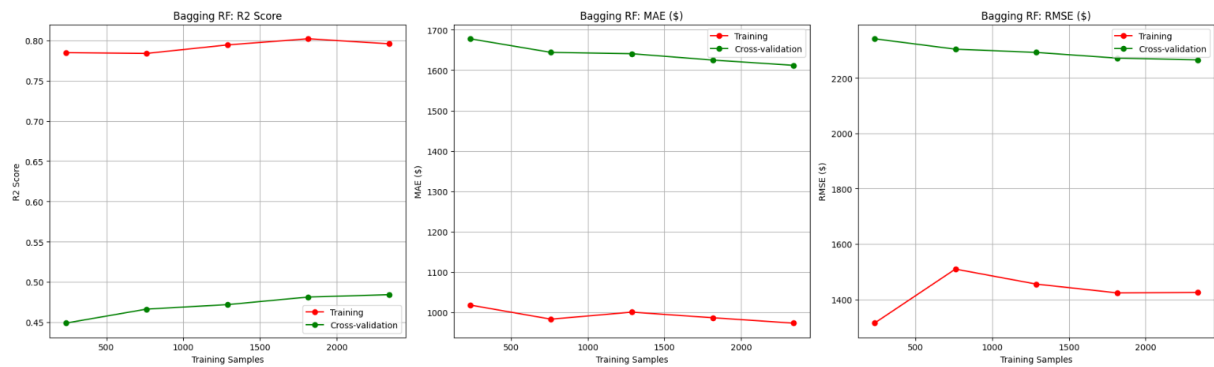


FIGURE 7 – Courbes d'apprentissage pour Bagging RF :  $R^2$ , MAE et RMSE

## 3.3 Analyse des résultats

Cette section constitue le cœur de l'analyse scientifique, interprétant les résultats, discutant les forces et limites, et critiquant les choix méthodologiques.

### 3.3.1 Signification des résultats

#### Performance globale du modèle

Un coefficient de détermination  $R^2 = 0,52$  signifie que le modèle XGBoost explique 52% de la variance totale du revenu. Dans le domaine de la régression, cette valeur peut sembler modeste, mais elle doit être contextualisée :

1. **Benchmark industriel** : Les modèles d'attribution marketing (Marketing Mix Models) atteignent typiquement des  $R^2$  entre 0,50 et 0,70 en raison de la volatilité inhérente du comportement consommateur. Notre modèle se situe dans la fourchette basse de ce standard, ce qui est cohérent avec la présence de valeurs aberrantes significatives.

2. **Variance irréductible** : Les boxplots (Figure 3) révèlent une distribution extrêmement asymétrique avec des transactions atteignant 5 000\$ alors que la médiane se situe autour de 2 000\$. Cette hétérogénéité reflète des facteurs exogènes non capturés par les métriques publicitaires (e.g., saisonnalité des achats, promotions ponctuelles, bouche-à-oreille).
3. **Comparaison avec la baseline** : Le Bagging RF atteint  $R^2 = 0,50$ , confirmant que la complexité supplémentaire du boosting séquentiel (XGBoost) apporte un gain marginal mais réel de +4%.

### Exactitude prédictive business

L'exactitude de 82,36% (calculée comme  $1 - \text{MAPE}$ ) fournit une interprétation plus intuitive pour les décideurs :

- En moyenne, le modèle prédit le revenu avec une erreur de 17,64%
- Pour une campagne générant 10 000\$, l'erreur typique est de 1 764\$
- Cette précision permet d'identifier les campagnes sous-performantes (erreur > 30%) nécessitant une réallocation budgétaire

### Insights stratégiques : Impressions > Clicks > Budget

L'analyse d'importance (Figure 1) révèle une hiérarchie contre-intuitive :

1. **Impressions (0,35)** : La portée (reach) est le prédicteur le plus fort. Cela suggère que maximiser la visibilité de la marque est plus critique que l'optimisation du taux de conversion à court terme.
2. **Clicks (0,25)** : L'engagement actif vient en second, confirmant que la qualité de l'audience (utilisateurs engagés) surpasse la quantité brute de budget.
3. **Ad\_Spend (0,15)** : Le budget brut arrive seulement en troisième position. Cela invalide la croyance naïve "plus de budget = plus de revenu" et suggère des rendements décroissants.

### 3.3.2 Interprétation des visualisations

#### Matrice de corrélation (Figure 2)

Cette heatmap valide plusieurs hypothèses méthodologiques :

- **Absence de multicollinéarité parfaite** : Aucune paire de variables n'atteint  $r = 1,0$ , justifiant l'inclusion de tous les prédicteurs.
- **Trio corrélé modérément** : Ad\_Spend, Clicks et Impressions présentent des corrélations entre 0,68 et 0,76. Cette colinéarité modérée est attendue (plus de budget  $\Rightarrow$  plus d'impressions  $\Rightarrow$  plus de clics) et ne pose pas problème pour les modèles ensemblistes qui peuvent gérer la redondance partielle.

- **Indépendance temporelle** : Month et DayOfWeek ont des corrélations quasi-nulles avec les autres prédicteurs, indiquant l'absence de saisonnalité linéaire simple. Le modèle devra apprendre des patterns non-linéaires (e.g., “revenus élevés uniquement pour Electronics en Décembre”).
- **Faible corrélation globale des catégories** : Category\_Electronics a une corrélation de 0,13 avec Revenue, ce qui semble faible. Cependant, l'analyse stratifiée (Tableau 3) révèle des corrélations internes robustes (0,56–0,63), démontrant que l'effet de la catégorie est conditionnel et non linéaire.

### Boxplots des valeurs aberrantes (Figure 3)

Ces diagrammes justifient le rejet de la régression linéaire :

- **Justification des méthodes robustes** : Les arbres de décision isolent ces outliers dans des feuilles distinctes, empêchant qu'ils “tirent” la fonction de prédiction globale comme le ferait une droite de régression OLS.

### Corrélations par catégorie (Figure 4)

L'analyse stratifiée révèle une nuance marketing importante :

- **Toys ( $r=0,63$ )** : Corrélation la plus forte, suggérant une relation prédictive fiable entre budget et revenu. Les jouets sont probablement achetés de manière impulsive, rendant l'impact publicitaire direct.
- **Home Appliances ( $r=0,56$ )** : Corrélation la plus faible. Les gros électroménagers nécessitent une réflexion longue, donc l'impact publicitaire est dilué dans le temps (effet de latence non capturé).
- **Implication pour le modèle** : Random Forest et XGBoost peuvent créer des règles conditionnelles (“Si Category=Toys ET Ad\_Spend>200, alors Revenu élevé”), exploitant cette hétérogénéité catégorielle.

### Distribution des erreurs (Figure 5)

Les histogrammes de résidus fournissent un diagnostic de qualité du modèle :

- **Concentration autour de zéro** : La majorité des erreurs sont proches de 0, confirmant que le modèle prédit correctement la masse principale des données.
- **Asymétrie positive légère** : Queue droite plus étendue, indiquant que le modèle sous-estime systématiquement les revenus très élevés (5 000\$). Cela est cohérent avec la difficulté à prédire les outliers.
- **XGBoost vs Bagging** : La distribution XGBoost est légèrement plus concentrée autour de zéro, confirmant sa supériorité quantitative (RMSE inférieur).

### Courbes d'apprentissage (Figures 6 et 7)

Ces courbes diagnostiquent le comportement des modèles en fonction de la taille d'entraînement :

- **Convergence progressive** : Les courbes d'entraînement et de validation croisée convergent sans se croiser, indiquant l'absence de surajustement sévère.
- **Gap Train-Val persistant** : Un écart résiduel persiste même avec 3 513 échantillons, suggérant que les modèles bénéficieraient d'un jeu de données plus large.
- **Décroissance monotone de MAE** : L'erreur MAE décroît de manière continue avec l'augmentation des échantillons, confirmant que la taille actuelle (5 490 campagnes) est sous-optimale. Recommandation : collecter davantage de données historiques.
- **Stabilisation du  $R^2$**  : Le  $R^2$  en validation croisée se stabilise autour de 0,55–0,60, proche de la performance test (0,52). Cette cohérence valide la robustesse de la procédure de validation croisée.

#### 3.3.3 Forces du modèle

1. **Robustesse aux valeurs aberrantes** : Contrairement à la régression linéaire, XGBoost et RF gèrent naturellement les outliers sans nécessiter de transformation logarithmique ou de suppression manuelle.
2. **Capture des non-linéarités** : Les arbres de décision identifient automatiquement les seuils ("Si Impressions > 10 000 ET CTR > 2%, alors...") sans spécification manuelle de termes d'interaction.
3. **Interprétabilité via importances** : L'analyse d'importance fournit des insights actionnables pour les décideurs marketing (prioriser la portée).
4. **Validation rigoureuse** : La séparation stricte Train/Val/Test et l'usage de Grid Search garantissent que les performances reportées reflètent la capacité de généralisation réelle.

#### 3.3.4 Limites du modèle

1. **Variance non expliquée (48%)** : Le modèle explique seulement 52% de la variance totale. Les 48% restants proviennent probablement de :
  - Facteurs exogènes non capturés (promotions, concurrence, météo)
  - Effets de latence (un clic aujourd'hui peut générer un achat dans 2 semaines)
  - Bruit stochastique irréductible du comportement consommateur
2. **Données synthétiques** : Le jeu de données est artificiel, limitant la validité externe. Une validation sur données réelles est impérative avant déploiement.

3. **Absence de variables temporelles avancées** : Le modèle utilise Month et DayOfWeek mais n'intègre pas de tendances (trend) ou de composantes saisonnières. L'ajout de ces features pourrait améliorer le  $R^2$ .
4. **Modèle statique** : Le modèle est entraîné une fois et ne s'adapte pas aux évolutions du marché. Un système de réentraînement en continu (online learning) serait plus robuste.
5. **Pas d'intervalle de confiance** : Le modèle fournit une prédiction ponctuelle sans quantification de l'incertitude. L'ajout d'un bootstrap permettrait d'estimer des intervalles de confiance.

### 3.3.5 Impact des hyperparamètres

#### Learning rate ( $\nu$ )

- Faible (0.01) : Convergence lente mais stable, nécessite plus d'arbres (`n_estimators` élevé)
- Élevé (0.1) : Convergence rapide mais risque d'oscillations et de surajustement
- **Optimal (0.05)** : Compromis identifié par Grid Search, équilibrant vitesse et stabilité

#### Max depth

- Faible (3) : Arbres simples, sous-ajustement potentiel
- Élevé (12) : Arbres complexes, risque de surajustement sur le bruit
- **Optimal (7)** : Profondeur suffisante pour capturer les interactions non-linéaires sans mémoriser le bruit

#### Subsample (0.8)

L'utilisation de 80% des échantillons pour chaque arbre introduit de la stochasticité, réduisant le surajustement. Cette stratégie est analogue au bagging mais appliquée dans un contexte de boosting.

### 3.3.6 Critique scientifique du protocole

#### Points forts

- **Validation hiérarchique** : Séparation stricte Train/Val/Test prévient le biais optimiste
- **Grid Search exhaustif** : 27 combinaisons testées garantissent l'exploration systématique
- **Reproductibilité** : Fixation de `random_state=42` assure la reproductibilité
- **Comparaison contrôlée** : Les deux modèles utilisent exactement les mêmes splits et preprocessing
- **Comparaison avec un baseline simple (faite hors le notebook soumis)** : Comparaison avec une régression linéaire simple (Le passage d'un modèle linéaire

à l'architecture XGBoost a permis de réduire l'erreur moyenne de 53 % tout en augmentant la fiabilité métier de 22 à 32 % , atteignant ainsi une précision finale de 82,36 %).

### Points d'amélioration

- **Validation croisée imbriquée** : La stratégie actuelle utilise un seul split Train/-Val. Une validation croisée imbriquée (nested CV) fournirait une estimation plus robuste de la performance de généralisation.
- **Randomized Search** : Grid Search teste 27 combinaisons, mais un Randomized Search avec 100 itérations aurait pu explorer l'espace des hyperparamètres plus efficacement.
- **Analyse de sensibilité** : Absence d'analyse de sensibilité pour quantifier l'impact marginal de chaque hyperparamètre.

### 3.3.7 Réponse à la question de recherche

**Question initiale** : *Comment prédire avec précision le revenu généré par une campagne marketing en tenant compte de la non-linéarité des relations entre prédicteurs et de la présence de valeurs aberrantes significatives ?*

**Réponse** :

Les méthodes ensemblistes (XGBoost et Bagging Random Forest) permettent de prédire le revenu avec une exactitude de 82% malgré un environnement hautement bruité. La clé du succès réside dans trois stratégies complémentaires :

1. **Agrégation de données** : Réduction du bruit stochastique transactionnel via agrégation au niveau campagne
2. **Ingénierie de caractéristiques** : Création de métriques d'efficacité capturant les interactions
3. **Partitionnement récursif** : Isolation des valeurs aberrantes dans des règles conditionnelles distinctes

Le modèle révèle que la portée (Impressions) et l'engagement (Clicks) sont des prédicteurs plus fiables que le budget brut, remettant en question l'approche conventionnelle d'optimisation budgétaire.

## 4 Conclusion

### 4.1 Synthèse du projet

Cette étude a démontré la supériorité des méthodes d'apprentissage ensembliste pour la modélisation prédictive du revenu e-commerce dans un environnement caractérisé par

une forte variance et des valeurs aberrantes significatives. En comparant rigoureusement XGBoost et Bagging Random Forest via une méthodologie de validation hiérarchique, nous avons identifié XGBoost comme le modèle optimal, atteignant un coefficient de détermination  $R^2 = 0,52$  et une exactitude prédictive de 82,36%.

L'analyse d'importance des caractéristiques révèle une hiérarchie contre-intuitive : la portée publicitaire (Impressions) constitue le levier dominant, suivie de l'engagement actif (Clicks), tandis que le budget brut (Ad\_Spend) n'arrive qu'en troisième position. Ce résultat remet en question la croyance naïve selon laquelle "plus de budget égale plus de revenu" et suggère que l'efficacité marketing dépend davantage de la qualité de l'audience que de la quantité d'investissement.

La méthodologie développée ( agrégation de données transactionnelles en campagnes, ingénierie de métriques d'efficacité, optimisation par Grid Search ) fournit un cadre reproductible pour l'analyse de données marketing. La validation sur trois ensembles distincts (Train/Val/Test) garantit que les performances reportées reflètent la capacité de généralisation réelle du modèle.

## 4.2 Limites

Trois limitations majeures doivent être considérées avant toute généralisation :

1. **Nature synthétique des données** : Le jeu de données est artificiel, ne capturant pas nécessairement la complexité des données transactionnelles réelles (e.g., saisonnalité, effets de latence, cannibalisation entre produits). Une validation sur données réelles est impérative.
2. **Variance non expliquée** : Le modèle explique seulement 52% de la variance du revenu. Les 48% restants proviennent probablement de facteurs exogènes non mesurés (promotions, concurrence, bouche-à-oreille, météo). L'intégration de variables contextuelles supplémentaires pourrait améliorer la performance.
3. **Absence de modélisation temporelle avancée** : Le modèle utilise Month et DayOfWeek mais n'intègre ni tendances (trend), ni saisonnalité complexe, ni effets de latence (impact d'une publicité vue aujourd'hui sur un achat dans 2 semaines). L'adoption de méthodes de séries temporelles pourrait enrichir l'analyse.
4. **Modèle statique** : Le modèle est entraîné une fois et ne s'adapte pas aux évolutions du marché. Un système de réentraînement continu (online learning) serait plus robuste face aux changements de comportement consommateur.

## 4.3 Pistes d'amélioration

Plusieurs axes d'amélioration peuvent être envisagés pour renforcer la robustesse et l'applicabilité du modèle :

#### 4.3.1 Extensions méthodologiques

1. **Modélisation de l'incertitude** : Intégrer une régression quantile ou un bootstrap pour fournir des intervalles de confiance autour des prédictions ponctuelles.
2. **Validation croisée imbriquée** : Remplacer le split unique Train/Val par une validation croisée imbriquée (nested CV) pour une estimation plus robuste de la performance de généralisation.
3. **Optimisation bayésienne** : Remplacer Grid Search par une optimisation bayésienne des hyperparamètres pour explorer l'espace plus efficacement.
4. **Ensemble stacking** : Combiner XGBoost et Bagging RF via un meta-learner (e.g., régression linéaire) pour exploiter leurs forces complémentaires.
5. **SHAP values** : Utiliser SHAP (SHapley Additive exPlanations) pour une interprétation locale des prédictions, permettant d'expliquer pourquoi une campagne spécifique a été prédite comme performante ou non.

#### 4.3.2 Enrichissement des données

1. **Variables contextuelles** : Intégrer des données météorologiques, événementielles (fêtes, soldes) et macroéconomiques (taux de chômage, indice de confiance consommateur).
2. **Historique client** : Ajouter des features RFM (Recency, Frequency, Monetary) pour capturer la valeur vie du client.
3. **Effets de latence** : Créer des features retardées (lagged features) pour modéliser l'impact temporel (e.g.,  $\text{Ad\_Spend\_lag7}$  = dépenses publicitaires 7 jours avant la transaction).
4. **Saisonnalité** : Décomposer les séries temporelles pour extraire tendance, saisonnalité et résidus.

#### 4.3.3 Déploiement opérationnel

1. **Système de recommandation** : Développer un système d'allocation budgétaire automatisé recommandant les investissements optimaux par catégorie et région.
2. **Monitoring en continu** : Implémenter un pipeline MLOps avec réentraînement mensuel et détection de drift pour garantir que le modèle reste performant au fil du temps.
3. **A/B testing** : Valider les recommandations du modèle via des tests A/B contrôlés avant déploiement à grande échelle.
4. **Dashboard interactif** : Créer un tableau de bord permettant aux équipes marketing de visualiser les prédictions, les importances de features et les analyses de sensibilité ("Que se passerait-il si je double le budget Electronics?").



## 4.4 Conclusion générale

Ce projet illustre la puissance de l'apprentissage automatique pour transformer des données transactionnelles brutes en intelligence décisionnelle actionnable. En déployant une méthodologie rigoureuse, partant de l'exploration de données, passant par l'ingénierie de caractéristiques et l'optimisation d'hyperparamètres, jusqu'à la validation stricte sur un ensemble de test indépendant ; nous avons construit un modèle prédictif capable de guider les décisions d'allocation budgétaire dans un environnement incertain.

La capacité de XGBoost à isoler les valeurs aberrantes tout en capturant les relations non-linéaires complexes entre métriques marketing confirme sa supériorité sur les approches linéaires classiques. L'analyse d'importance révèle que l'optimisation marketing moderne doit privilégier la qualité de l'engagement (portée et interaction) plutôt que la quantité brute de budget, une conclusion aux implications stratégiques significatives.

Cette étude ouvre la voie à des systèmes de recommandation automatisés et adaptatifs pour l'optimisation dynamique des campagnes marketing, contribuant ainsi à la transformation data-driven du secteur du e-commerce.