

1. Projet commun réalisé par Jiang Chang, Christopher Dumas, Michael Thomas et l'Auteur

Résumé

Ce rapport présente la conception et mise en œuvre d'un pipeline de données complet, depuis l'extraction jusqu'à la visualisation, dans le cadre d'un projet de fin de formation. Nous avons opté pour une architecture ELT (Extract-Load- Transform) avec une *stack* technique incluant Fivetran pour l'extraction, DBT pour la transformation, BigQuery comme entrepôt de données, et Tableau pour la visualisation. L'orchestration a été assurée par Airflow, permettant d'automatiser l'ensemble des *workflows*, y compris l'exécution de modèles de Machine Learning (Prophet pour les prévisions de Chiffre d'affaires, Apriori pour l'analyse des associations de produits et XGBoost pour l'analyse des causes de pannes de machines).

Ce travail a permis de valider des compétences techniques variées (SQL, Python, orchestration, *dataviz*) tout en soulignant l'importance cruciale de la qualité des données et de la collaboration d'équipe dans un projet data. Le pipeline conçu répond aux exigences initiales et pourra facilement s'adapter à de nouveaux besoins analytiques.

Table des matières

Introduction et organisation								
1	Conception et Sélection de la Stack Technique 1.1 Modèle - Extract-Load-Transform 1.2 Data Warehouse - BigQuery 1.3 Extraction et Chargement - Outils Managés ou Scripts Python 1.4 Transformation - DBT Cloud 1.5 Orchestration - Airflow 1.6 Visualisation - Tableau	4 5 5 5 6 6 7						
2	Implémentation du Pipeline 2.1 Configuration du data warehouse - BigQuery 2.2 Configuration du Chargement et de l'Extraction - Fivetran 2.3 Transformation de données - DBT Cloud 2.3.1 Configuration de Base 2.3.2 Organisation du Projet DBT 2.3.3 Analyse de code 2.3.4 Configuration avancée 2.4 Scripts Python - Machine Learning et Géolocalisation Clients 2.4.1 Prévision du Chiffre d'Affaires avec Prophet 2.4.2 Analyse des associations de catégories de produits avec l'algorithme	8 9 10 11 11 13 18 20 21						
	Apriori 2.4.3 Géolocalisation des Clients avec leur adresse IP 2.5 Orchestration avec Airflow 2.5.1 Analyse de code : DAG 2.5.2 Analyse de code : start_dbt_build_job.py	28 32 32 35						
3	Tableaux de bords, analyse et recommandations 3.1 Tableaux de Bord	39 40 40 42 43 43						
Co	onclusion	45						
A	Code DBT A.1 Configuration	46 51 54						
В	Scripts Machine Learning et Geolocalisation	59						
C	Orchestration Airflow	65						
D	Dashboards	70						

Introduction et organisation

Le projet en fin de formation s'inscrivait dans le cadre d'une première mission professionnelle chez Carttrend, entreprise spécialisée dans l'e-commerce. Face à un marché concurrentiel, l'équipe *Data & Insights* avait été chargée de construire un pipeline complet pour exploiter des données jusqu'alors sous-utilisées et dispersées dans des fichiers Excel et Google Sheets. L'objectif était alors de transformer ces sources brutes en analyses pertinentes et en tableaux de bord clairs et conviviaux, facilitant ainsi la prise de décisions stratégiques dans quatre domaines critiques : l'analyse des ventes, l'efficacité marketing, l'optimisation logistique et la satisfaction client.

Nous avons rapidement identifié trois impératifs sous-jacents. Premièrement, concevoir des tableaux de bord intuitifs, concis et régulièrement mis à jour. Nous avons donc choisi très tôt de limiter volontairement le nombre de visualisations pour éviter la surcharge cognitive des utilisateurs, privilégiant l'usage quotidien à l'exhaustivité. De plus, il nous a semblé essentiel de proposer une solution clé en main nécessitant une maintenance minimale par des équipes non nécessairement techniques. Enfin, il fallait construire une architecture évolutive, capable de gérer un gros volume de données, mais aussi facilement adaptable à de nouveaux axes d'analyse sans refonte majeure.

Notre équipe, composée de quatre collaborateurs aux compétences diverses, a commencé par analyser les forces et faiblesses de chacun et la manière dont celles-ci pouvaient se compléter au sein du groupe. Certaines forces majoritaires rapidement identifiées, comme une préférence pour SQL dans les opérations de transformation des données, une aisance avec l'écosystème Google Cloud Platform (BigQuery) et une maîtrise de l'outil de visualisation Tableau, ont également orienté le choix de la *stack* technique. Cette diversité de compétences a permis une répartition efficace des tâches, avec un travail parallèle sur la majeure partie du pipeline (principalement pour des raisons pédagogiques), complété par des spécialisations ponctuelles sur des sujets spécifiques comme le Machine Learning ou l'orchestration Airflow.

Un calendrier initial (Figure 1) a été rapidement établi pour structurer les dix jours du projet, dans un effort conscient d'éviter l'embourbement dans la partie technique du pipeline et ainsi réserver suffisamment de temps à la visualisation et à l'analyse. L'objectif était également de prévenir le feature creep du pipeline, en n'ajoutant des fonctionnalités supplémentaires qu'une fois le cœur du travail réalisé, permettant ainsi d'en évaluer plus précisément la pertinence. Des réunions stand-up biquotidiennes ont été mises en place – le matin pour attribuer les tâches et résoudre les blocages immédiats et le soir pour évaluer les résultats de la journée, documenter

collectivement les avancées et ajuster les objectifs du lendemain. Cette organisation, combinée à des revues d'étape bihebdomadaires avec l'enseignant jouant le rôle du chef de projet, a assuré une progression fluide sans recours au *crunch*.

Mardi 20/05	Mercredi 21/05	Jeudi 22/05	Vendredi 23/05	S a	D i	Lundi 26/05	Mardi 27/05	Mercredi 28/05	Jeudi 29/05	Vendredi 30/05
				m e d i	m a n c h e					
Organisation Générale Debut Nettoyage Données (DBT) Configuration Airbyte	Nettoyage Données (DBT Staging)	Conception Marts Implementation Marts (DBT Marts)	Debut Dataviz (Tableau) Machine Learning (Python)	2 4 / 0 5	2 5 / 0 5	Dataviz (Tableau) Fin Machine Learning (Python)	Dataviz (Tableau)	Fin Dataviz (Tableau) Analyses et Recommanda- tions	Orchestration (Airflow) Finitions et Debugage	Finitions et Debugage

FIGURE 1 – Plan d'organisation initial

Chapitre 1

Conception et Sélection de la Stack Technique

Ce chapitre présente les choix fondamentaux qui ont guidé la conception de notre pipeline de données pour Carttrend. Dans le but de transformer des sources disparates en informations exploitables, nous avons dû prendre plusieurs décisions techniques stratégiques concernant l'architecture globale. La conception de notre pipeline de données a ainsi débuté par une réflexion fondamentale sur son architecture.

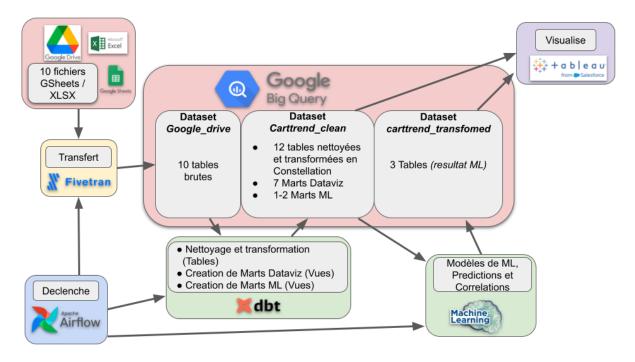


FIGURE 1.1 – Diagrame de l'architecture finale séléctionnée

1.1 Modèle - Extract-Load-Transform

Après analyse, nous avons opté pour un modèle ELT (Extract-Load-Transform) plutôt qu'une approche ETL traditionnelle. En premier lieu, le modèle ELT offre une meilleure scalabilité pour faire face à une potentielle croissance du volume de données chez Carttrend. Contrairement à un modèle ETL où les transformations sont appliquées avant le chargement, travailler en ELT nous permet de charger d'abord les données brutes dans le *data warehouse*, puis d'y appliquer les transformations nécessaires. Cette approche réduit les goulots d'étranglement potentiels lors des phases d'extraction et de chargement. De plus, cette architecture préserve une flexibilité essentielle pour l'analyse des données. En conservant les données originales dans le *data warehouse* avant transformation, nous maintenons la possibilité de retraiter ces informations selon de nouveaux axes d'analyse pas toujours identifiés immédiatement. Enfin, ce choix technique correspondait aux compétences globales de notre équipe, plus à l'aise avec SQL qu'avec Python pour les opérations de transformation complexes. L'approche ELT nous a permis d'implémenter l'ensemble de nos logiques métier en SQL via DBT, tout en bénéficiant de la puissance de calcul distribuée du *data warehouse*. Cette orientation ELT a constitué le fil directeur dans le choix de notre stack technique, qui est présentée dans les sections suivantes.

1.2 Data Warehouse - BigQuery

BigQuery de Google s'est naturellement imposé comme socle de notre infrastructure. En tant que *data warehouse* SaaS, BigQuery réduit en effet la nécessité d'interventions manuelles pour résoudre les problèmes de provisionnement de capacité et abstrait la gestion des ressources physiques (puissance de calcul, espace de stockage, accès réseau...). D'autre part, notre équipe manifestait une aisance particulière avec l'écosystème Google Cloud Platform. Cette familiarité technique, couplée à la large adoption de BigQuery dans l'industrie, a permis de réduire significativement les difficultés de configuration des différents outils du reste de la *stack* (DBT, Fivetran, scripts Python...). Enfin, le choix de BigQuery répondait à une exigence de sécurité des données. Cette solution délègue les problématiques sensibles de chiffrement et de contrôle d'accès à Google, nous permettant de se concentrer sur les données plutôt que sur l'infrastructure sous-jacente.

1.3 Extraction et Chargement - Outils Managés ou Scripts Python

Pour l'extraction et le chargement des données, nous avons délibérément choisi d'utiliser une solution managée (Airbyte Cloud ou Fivetran, le choix exact de l'outil n'ayant pas été fait à ce stade) plutôt que de développer des scripts Python spécialisés. Ces outils managés "cloud-native" offrent en effet une fiabilité supérieure et une maintenance simplifiée, avec une scalabilité immédiate, alors que des scripts Python maison nécessiteraient des ajustements et

une optimisation régulière pour s'adapter aux besoins évolutifs. Sur le plan de la sécurité, ces solutions permettent – comme pour BigQuery – de déléguer les problématiques complexes de chiffrement et de gestion des accès aux experts, plutôt que de les faire reposer sur notre propre implémentation. Enfin, leur configuration via une interface simple et intuitive réduit considérablement le temps de développement tout en assurant une maintenance simplifiée à long terme.

1.4 Transformation - DBT Cloud

Nous avons choisi DBT Cloud (Data Build Tool) comme solution de transformation des données. Cet outil s'intègre parfaitement avec BigQuery et permet d'exécuter les transformations directement dans le *data warehouse*, exploitant ainsi les performances de l'infrastructure Google tout en réduisant les mouvements de données. La flexibilité de matérialisation (*view*/table) offre un contrôle précis sur le compromis entre puissance de calcul et espace de stockage. L'architecture modulaire de DBT (*staging*, *marts*) facilite l'adaptation du pipeline à la fois aux nouvelles sources de données et aux besoins évolutifs des outils d'analyse (Machine Learning, visualisation...). La documentation intégrée et les tests faciles à rédiger garantissent la qualité des données après transformation. Enfin, l'utilisation du SQL – particulièrement adapté aux compétences de notre équipe – combinée aux fonctionnalités collaboratives de l'outil, ont définitivement confirmé notre choix.

1.5 Orchestration - Airflow

Pour l'orchestration de notre pipeline, deux solutions ont été envisagées. La première, plus basique, aurait utilisé les fonctionnalités natives d'orchestration de Fivetran, permettant un déclenchement régulier des synchronisations (par exemple quotidiennes) et l'exécution automatique d'un job DBT après chaque chargement. Cette approche, plus simple à mettre en œuvre et à maintenir par des équipes moins techniques, aurait pu suffire à répondre au cahier des charges. Les scripts Python de Machine Learning auraient alors été exécutés séparément via des *jobs 'cron'* sur un serveur dédié.

Nous avons finalement opté pour Airflow, solution plus complexe initialement mais offrant une gestion bien plus fine des dépendances entre tâches, comme par exemple l'enchaînement conditionnel de l'entraînement des modèles de machine learning seulement après un transfert et une transformation réussis des données. D'autre part, Airflow offre une flexibilité accrue pour l'ajout de nouvelles fonctionnalités, comme l'intégration ultérieure d'un script de géolocalisation des clients via leurs adresses IP. Ces atouts justifient a nos yeux la complexité initiale de mise en place, et nous avons tenté de remédier à la difficulté de maintenance en utilisant une architecture de code simple.

1.6 Visualisation - Tableau

Notre choix s'est porté sur Tableau pour la visualisation des données, principalement en raison de ses avantages en contexte gratuit. La version de démonstration de Tableau Cloud autorisait une collaboration efficace entre membres de l'équipe, avec des fonctionnalités de partage et d'édition conjointe supérieures à ce qu'offre Power BI sans licence payante. L'intégration native avec BigQuery a constitué un autre atout décisif, permettant un accès simple aux données et une mise à jour des *dashboards* automatique et régulière. Enfin, la familiarité préexistante de l'équipe avec Tableau a accéléré notre productivité, nous permettant de se concentrer sur la création de visualisations claires et efficaces.

Conclusion

Cette première partie a présenté les choix architecturaux fondamentaux qui structurent notre pipeline de données. L'architecture ELT retenue, reposant sur une *stack* technique cohérente (BigQuery, Fivetran, DBT Cloud, Airflow et Tableau), forme un pipeline robuste et évolutif. Un diagramme global de cette architecture illustrant les interactions entre ces différents outils (*Figure 2.5*), est proposé en page 4.

Chapitre 2

Implémentation du Pipeline

Cette section présente la concrétisation technique de nos choix architecturaux en un pipeline opérationnel. Nous détaillons les étapes clés de construction : configuration de BigQuery, intégration des solutions d'extraction (Fivetran), transformation via DBT, enrichissement par scripts Python, et orchestration globale avec Airflow.

2.1 Configuration du data warehouse - BigQuery

La mise en place de notre environnement BigQuery a débuté par la création d'un compte Google unique doté des droits d'administrateur sur Google Cloud Platform. Après avoir initialisé un projet BigQuery, nous avons structuré l'espace de stockage en trois *datasets* distincts (*Figure 2.1*), chacun dédié à une étape spécifique du pipeline : données brutes, données transformées, et résultats des modèles de Machine Learning.

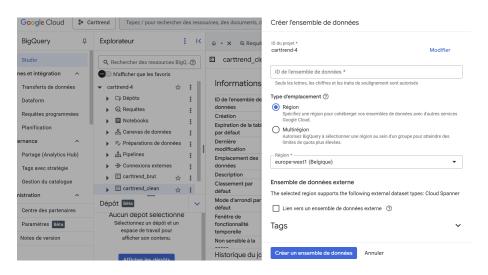


FIGURE 2.1 – Creation des datasets sur BigQuery

L'automatisation des processus a nécessité l'établissement de comptes de service individuels pour chaque outil (Figure 2.2). DBT et les scripts Python de Machine Learning exécutés via Airflow bénéficient ainsi d'identifiants spécifiques avec des permissions limitées à leurs fonctionnalités respectives. Bien que des comptes aient été initialement configurés pour Airbyte durant la phase de tests, ceux-ci ont été soigneusement désactivés après l'adoption définitive de Fivetran comme solution d'extraction.

= F	Filter Enter property name or value									
	Email	Status	Name ↑							
	<u>airbyte-gdrive-read-access@carttrend-diam.gserviceaccount.com</u>	Disabled	Airbyte GDrive Read Access							
	garttrend-airbyte@carttrend- 4.iam.gserviceaccount.com	Disabled	carttrend airbyte							
	dbt-read-write@carttrend- 4.iam.gserviceaccount.com		DBT-Read-Write							
	 <u>™achine-learning@carttrend-liam.gserviceaccount.com</u> 		Machine Learning							

FIGURE 2.2 – Comptes de services

Pour accélérer l'implémentation, nous avons opté pour une authentification OAuth via le compte administrateur Google afin d'autoriser Fivetran et Tableau à accéder aux données. Bien que pratique, cette approche présente des risques de sécurité en accordant des privilèges excessifs à ces outils – qui n'ont pas besoin d'un accès administrateur complet – tout en créant une dépendance à un compte administrateur unique.

Cette configuration permet d'ores et déjà d'utiliser BigQuery dans un pipeline fonctionnel et performant, mais nécessiterait des améliorations avant un déploiement à plus grande échelle.

2.2 Configuration du Chargement et de l'Extraction - Fivetran

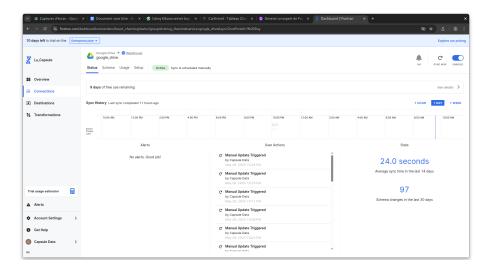


FIGURE 2.3 – Interface de Gestion de Fivetran

Initialement, Airbyte Cloud avait retenu notre attention pour son aspect open source et son interface simplifiée. Cependant, des difficultés techniques sont rapidement apparues, notamment avec le connecteur Google Drive développé par la communauté, qui ne permettait pas d'accéder efficacement aux fichiers partagés – c'est-à-dire ceux dont le compte Google utilisé

par Airbyte n'était pas propriétaire. Bien qu'un contournement ait été envisagé (chargement manuel de chaque fichier dans une source Airbyte distincte), cette solution à été jugée peu flexible. Face à ces difficultés, nous avons initialement procédé à un chargement manuel des données dans BigQuery, permettant ainsi à une partie de l'équipe de commencer les étapes de nettoyage et de transformation, tandis que l'autre pouvait opérer la migration vers Fivetran sereinement.

La configuration de Fivetran (Figure 2.4) a débutée par la création d'un compte d'essai, et la connexion à BigQuery et Google Drive via OAuth. Bien que cette méthode d'authentification présente des limites (comme mentionné précédemment), elle a permis une intégration rapide. Fivetran s'est montré remarquablement efficace pour lire sans difficulté les fichiers Google Drive partagés, traitant de manière transparente et uniforme les formats Google Sheets et Excel. L'outil a même fait preuve d'intelligence en supprimant automatiquement les diacritiques des noms de colonnes incompatibles avec BigQuery. En prévision de l'automatisation, nous avons généré une clé API et étudié la documentation de l'API REST¹ afin de pouvoir déclencher les synchronisations via Airflow, une fonctionnalité essentielle pour l'orchestration future du pipeline.

Au final, Fivetran s'est imposé comme une solution optimale pour l'extraction et le chargement. Ses nombreux connecteurs pré-intégrés, sa maintenance simplifiée et son potentiel d'automatisation en font un composant clé garantissant l'évolutivité et l'adaptabilité future de notre infrastructure de données.

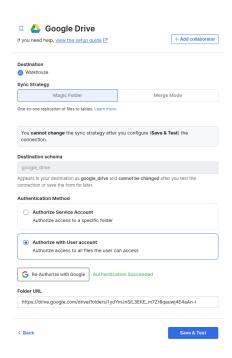


FIGURE 2.4 – Configuration de Fivetran

2.3 Transformation de données - DBT Cloud

La transformation des données est une étape capitale durant laquelle les informations brutes sont nettoyées, enrichies et converties en une structure exploitable pour l'analyse. Cette section explore en détail notre mise en œuvre de DBT (Data Build Tool), outil que nous avons retenu pour effectuer cette transformation. Nous commencerons par la configuration de base, puis nous décrirons l'organisation du projet. Cette dernière sera illustrée par l'explication détaillée de quelques modèles représentatifs. Enfin, nous aborderons la configuration avancée de DBT, permettant notamment son interfaçage avec Airflow.

^{1.} Disponible sur le site de Fivetran, https://fivetran.com/docs/rest-api/api-reference/connections

2.3.1 Configuration de Base

La configuration intiale de notre environnement DBT a suivi plusieurs étapes. Nous avons débuté par la création d'un compte d'essai sur DBT Cloud, suivie de l'ouverture d'un compte GitHub hébergeant un dépôt dédié au *versioning* de notre code de transformation. Le projet DBT a ensuite été configuré pour lui permettre d'accéder aux données brutes stockées dans BigQuery ainsi qu'à un dataset de sortie, via un compte de service spécifique. Enfin, nous avons établi la connexion entre le projet DBT Cloud et notre dépôt GitHub.

L'utilisation partagée d'un seul compte DBT Cloud partagé par toute l'équipe a cependant révélé des limites organisationnelles. Si le travail sur les modèles SQL restait praticable en parallèle grâce à leur découpage en fichiers distincts, la gestion simultanée des fichiers schema.yaml – essentiels pour la documentation et les tests – posait davantage problème. Pour prévenir les risques d'écrasement des modifications, nous avons instauré un système de rotation informel lors de l'édition de ces fichiers critiques, ce qui a légèrement affecté notre productivité durant les phases de documentation et d'écriture de tests.

2.3.2 Organisation du Projet DBT

Les données sources étaient déjà relativement bien organisées et ne nécessitaient que peu de normalisation supplémentaire. Ce constat nous a conduits à opter pour une architecture simplifiée en seulement deux niveaux : *staging* et *marts*.

Le dossier *staging* regroupe l'ensemble des modèles décrivant des transformations initiales appliquées en une seule étape. Ces modèles sont à la fois responsables du nettoyage des données brutes (par exemple la gestion des valeurs nulles, la standardisation des formats ou le renommage des colonnes) mais aussi d'une normalisation légère (par exemple, l'unification des canaux publicitaires présents dans plusieurs tables sources). A l'issue de cette étape, chaque modèle génère une vue dans BigQuery, optimisant ainsi l'utilisation du stockage. Nous avons délibérément évité une normalisation excessive à ce stade pour éviter de complexifier inutilement les requêtes ultérieures en multipliant les jointures, l'organisation initiale des données source ne le nécessitant pas. Les données transformées adoptent donc une structure proche du modèle en constellation (*Figure 2.5*), avec :

- Trois tables de fait principales (Campagnes, Posts et Commandes),
- Des tables de dimensions clairement identifiées (ex. Produits, Entrepôts, Canaux...),
- Quelques tables hybrides conservant une nature mixte (ex. Satisfaction, Entrepots_Machines...).

Le dossier *marts* organise les données transformées selon leurs finalités. D'un côté, les modèles dédiés au machine learning, comme *mart_Previsions_Series_Temporelles_CA.sql*, mettent en œuvre un *feature engineering* spécifiquement conçu pour les algorithmes respectifs. De l'autre, les modèles orientés visualisation, à l'image de *mart_Commentaires.sql*, concentrent les agrégations et calculs d'indicateurs métiers essentiels pour alimenter les *dashboards* Tableau.

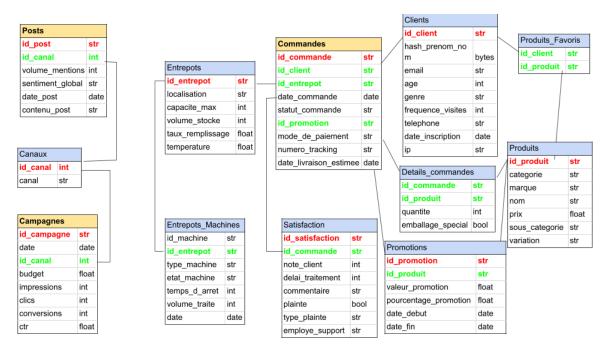


FIGURE 2.5 – Diagramme des vues générées par les modèles staging

Contrairement aux vues utilisées dans le dossier *staging*, tous les modèles de type *mart* sont matérialisés en tables dans BigQuery pour accélérer l'accès aux données, améliorant ainsi notablement la réactivité des *dashboards* tout en réduisant la puissance de calcul lors des requêtes fréquentes.

Cette architecture présente cependant des limites structurelles qui sont apparues au fil du développement. Le modèle *mart_Commandes.sql* illustre particulièrement ces contraintes, ayant progressivement accumulé trop de logique de transformation (phénomène de *scope creep*) et aurait mérité un traitement intermédiaire distinct. D'autre part, bien que parfaitement fonctionnelle pour ce projet de taille modérée, l'approche combinant nettoyage et transformation en une seule étape réduit la modularité du code. Une évolution future gagnerait à introduire un dossier *"intermediate"* (conforme aux bonnes pratiques recommandées par la documentation officielle de DBT ¹) qui permettrait d'isoler la construction des modèles complexes, de découpler clairement nettoyage des sources, normalisation et enrichissement métier, et de faciliter la réutilisation des modèles intermédiaires. Enfin, une normalisation plus poussée – qui bénéficierait aussi d'un niveau intermédiaire – pourrait s'avérer nécessaire si la structure des données sources venait à se complexifier ou si de nouvelles sources venaient s'ajouter aux existantes.

Il convient donc de considérer l'architecture actuelle comme une base solide mais perfectible, appelée à évoluer en fonction des besoins futurs.

 $^{1. \ \} we \ \ structure \ \ our \ \ dbt \ \ projects, \ \ https://docs.getdbt.com/best-practices/how-we-structure/1-guide-overview$

2.3.3 Analyse de code

Cette section présente quelques exemples détaillés de modèles et de fichiers de configuration. Ces exemples complets, ainsi que d'autres sont disponibles, largement commentés ¹.

Clients.sql

Ce modèle, situé dans le dossier *staging*, a pour objectif de nettoyer et d'enrichir la table brute des clients. Il effectue notamment un important travail de standardisation des numéros de téléphone et de validation des données.

La première partie est une *Common Table Expression* (CTE) relativement complexe nommée "tel" utilisée dans la standardisation les numéros de téléphone. Elle extrait et le numéro principal au format 000-000-0000, ainsi que l'indicatif pays (comme +033 pour la France) avec un ajout de zéros si nécessaire, afin de préparer la standardisation au format international 2).

```
1 WITH tel AS (
 2
    SELECT
 3
      id client,
      -- Formatage standard du numéro (123-456-7890)
 4
 5
      REGEXP_REPLACE (numero_telephone, r'^*.*?(\d{3})\D?(\d{4}).*$',
           r' \1-\2-\3') AS numero_telephone,
 6
       -- Left-Padding du code pays (ex: +033 pour la France)
 7
      CASE
 8
        WHEN LENGTH(ac.area_code) = 3 THEN CONCAT('+', ac.area_code)
 9
        WHEN LENGTH(ac.area_code) = 2 THEN CONCAT('+0', ac.area_code)
10
        WHEN LENGTH (ac.area_code) = 1 THEN CONCAT ('+00', ac.area_code)
11
        ELSE '+033'
                      -- Valeur par défaut
12
      END as area code
13
    FROM {{source("google_drive",'carttrend_clients_clients')}}
14
    JOIN (
15
       -- Sous-Requete pour l'extraction du code pays depuis le numéro
          original
16
      SELECT
17
      id_client,
18
      REGEXP EXTRACT (numero telephone, r'^+?([1-9]{0,3}).*$') AS area code
      FROM {{source("google_drive", 'carttrend_clients_clients')}}
19
20
    ) ac USING(id_client)
21)
```

La requête principale se concentre sur la sélection et la validation des différents champs clients. Elle inclut notamment l'anonymisation des noms et prénoms via un hachage SHA256 ainsi que la validation des emails à l'aide d'une expression régulière volontairement permissive, pour minimiser les faux négatifs qui pourraient éliminer des adresses valides. Elle vérifie en outre que les adresses IP respectent le standard IPv4³, même s'il est important de noter que

^{1.} cf. Annexe A, page 46

^{2.} ITU-T E.164 $\S6.1$, https://www.itu.int/rec/dologin_pub.asp?lang=e&id=T-REC-E.164-201011-I!!PDF-F

^{3.} RFC 791, §2.3 https://www.rfc-editor.org/rfc/rfc791#page-7

l'expression régulière ne vérifie pas que chaque segment numérique est inférieur à 256.

```
-- Identifiant client conservé
2.3
24
    id_client,
2.5
26
    -- Anonymisation des données personnelles avec hashing SHA256
27
    sha256(CONCAT(prenom_client, nom_client)) as hash_prenom_nom,
28
29
    -- Validation de l'email
30
31
      WHEN REGEXP_CONTAINS (email, r'^[\w\.-]+@[\w-]+.[A-Za-z]+$') THEN email
32
      ELSE NULL
33
    END as email,
34
35
    -- Suppression des âges négatifs
    CASE WHEN age > 0 THEN age END as age,
36
37
38
    -- Calcul de la tranche d'âge
39
    CASE
      WHEN age > 0 AND age <= 25 THEN ' <= 25'
40
41
      WHEN age > 25 AND age <= 40 THEN '25-40'
      WHEN age > 40 AND age <= 55 THEN '40-55'
42
43
      WHEN age > 55 AND age <= 70 THEN '55-70'
      WHEN age > 70 THEN '>70'
44
45
      ELSE NULL
    END as tranche_age,
46
47
48
    -- Normalisation du genre
49
    CASE
50
      WHEN lower (genre) = 'homme' THEN 'Homme'
      WHEN lower(genre) = 'femme' THEN 'Femme'
51
52
      ELSE NULL
53
    END as genre,
54
55
    -- Validation et correction de la fréquence de visites
56
    CASE
57
     WHEN frequence_visites > 0 THEN frequence_visites
58
      ELSE 0
59
    END as frequence_visites,
60
61
     -- Construction du numéro de téléphone complet
62
     -- Indicatif + Numero, format final +000-123-456-7890
63
    CONCAT(tel.area_code, '-', tel.numero_telephone) as telephone,
64
    -- Validation de la date d'inscription
65
66
    CASE
      WHEN DATE_DIFF(date_inscription, CURRENT_DATE(), DAY) <= 0</pre>
67
68
      THEN date_inscription
69
      ELSE NULL
70
    END as date_inscription,
71
72
    -- Calcul de l'ancienneté en jours
73
    CASE
74
      WHEN DATE_DIFF(date_inscription, CURRENT_DATE(), DAY) <= 0</pre>
75
      THEN DATE_DIFF(CURRENT_DATE(), date_inscription, DAY)
```

```
76 ELSE NULL
77
                                             END as anciennete_jours,
 78
79
                                               -- Validation des adresses IP (format IPv4)
80
                                             CASE
81
                                                                WHEN REGEXP_CONTAINS (adresse_ip, r' \wedge d\{1,3\} \land d\{1,3\} 
                                                                                                 ′)
                                                                THEN adresse_ip
82
                                                                ELSE NULL
83
84
                                             END as ip
8.5
86 FROM {{source("google_drive",'carttrend_clients_clients')}}
87 JOIN tel USING(id_client)
88 WHERE id_client IS NOT NULL -- Exclusion des clients sans identifiant
```

Ce modèle fournit une base de données clients propre et standardisée, prête à être utilisée pour toutes les analyses ultérieures.

Canaux.sql

Ce second exemple illustre la normalisation des données communes à plusieurs tables sources. Il s'agit ici de créer une référence unifiée des canaux de communication utilisés dans deux différentes sources : la table "Posts" et la table "Campagnes". Le modèle est d'abord constitué d'une CTE nommée "canal", utilisant opération UNION DISTINCT qui agrège les noms des canaux de communication entre les deux sources tout en garantissant l'élimination des doublons. La requête principale ajoute ensuite à cette CTE une clé primaire auto-incrémentée ("id_canal") grâce a la fonction ROW_NUMBER(). Le modèle produit une table de dimension avec deux colonnes :

- "id_canal": Identifiant unique numérique (clé primaire),
- "canal" : Nom original du canal de communication.

```
1 WITH canal AS (
    -- Extraction des canaux depuis les posts réseaux sociaux
 2
 3
    SELECT canal_social
 4
    FROM {{source("google_drive",'carttrend_posts_posts')}}
 5
 6
    UNION DISTINCT -- Combine les deux en éliminant les doublons
 7
 8
     -- Extraction des canaux depuis les campagnes marketing
 9
    SELECT canal
10
    FROM {{source('google_drive', 'carttrend_campaigns_campagnes')}}
11)
12 SELECT
13
    -- Génération d'un ID unique pour chaque canal,
14
    -- un chiffre séquentiel auto-incrémenté
15
    ROW_NUMBER() OVER() as id_canal,
16
17
    -- Conservation du nom original
18
    canal_social as canal
19
```

```
20 FROM canal
21 -- Tri alphabétique pour une meilleure lisibilité
22 ORDER BY canal ASC
```

En centralisant les valeurs communes, le processus de normalisation supprime les redondances et on garantit la cohérence des données. Pour en bénéficier, une mise à jour des tables sources "Posts" et "Campagnes" s'impose : elles devront désormais référencer la table dimension via des clés étrangères ("id_canal") plutôt que de stocker directement les noms de canaux.

mart_Commentaire.sql

Il s'agit cette fois-ci d'un *mart*, transformant les commentaires bruts des clients en données structurées pour alimenter un nuage de mots interactif dans le *dashboard* Satisfaction Client. Il extrait les termes significatifs et calcule leurs indicateurs clés :

- Fréquence : le nombre de fois que ce mot apparait. Détermine la taille du mot dans le nuage,
- Note moyenne : la note moyenne du commentaire dans lequel le mot apparait. Détermine la couleur du mot dans le nuage et permet de filtrer celui-ci afin de n'afficher que les mots apparaissant dans une certaine plage de notes.

La requête s'articule en deux étapes principales, à nouveau avec l'aide de CTEs. La première, nommée "mots", décompose chaque commentaire en mots tout en préservant la note associée. Les fonctions imbriquées UNNEST (SPLIT (' ')) permettent cette transformation d'une chaine de caractères à une colonne de mots individuels.

```
1 WITH
2  mots AS (
3   SELECT
4   note_client, -- Note associée (1-5)
5   mot
6  FROM {{ref("Satisfaction")}}, -- Table staging nettoyée
7  UNNEST(SPLIT(commentaire, ' ')) AS mot -- Désagrégation en mots
8  ),
```

Le résultat de cette première CTE est ensuite filtré pour retirer les mots vides puis agrégée par mots pour supprimer les doublons. Les mots sont ensuite nettoyés d'éventuels caractères spéciaux et normalisés en majuscule, puis fréquence et note moyenne sont calculées.

```
9
    mots_aggreges AS (
10
       SELECT
         UPPER (REGEXP_REPLACE (mot, r'[^\w]', '')) AS mot,
11
12
         COUNT(*) AS frequence,
13
         ROUND (AVG (note_client), 1) AS note_moyenne
14
       FROM mots
       WHERE mot != ''
15
16
         AND lower (mot) NOT IN (
           'a', 'about', 'above', 'after', 'again', 'against', 'all', 'am', '
17
              an', 'and'
         ) -- liste de mots vides ecourtee pour plus de lisibilite
18
```

```
19 GROUP BY mot
20 )
```

Le résultat final est enfin trié par fréquence décroissante, principalement à des fins de *debu-* gage

```
21 SELECT *
22 FROM mots_aggreges
23 ORDER BY frequence DESC
```

Cette préparation des données est essentielle pour la réalisation d'une visualisation relativement complexe dans Tableau comme celle de cet exemple. Grâce à cette dernière, les commentaires des clients pourront être analysés d'un coup œil, tout en préservant le contexte dans lequel ils ont été écrits.

mart_Prevision_Series_Temporelles_CA.sql

Ce modèle de type *mart* prépare les données historiques de chiffre d'affaires pour alimenter le modèle de prévision de séries temporelles (Prophet). Il garantit un format adapté ses besoins spécifiques (deux colonnes, l'une contenant la date sans interruption et l'autre la série à prédire, ici le chiffre d'affaires). Le modèle utilise à nouveau deux CTEs pour assembler la table finale.

La première CTE, nommée "cal", génère l'ensemble des dates entre la première et la dernière commande enregistrée, créant ainsi un calendrier continu sans interruption.

```
1 WITH
 2
    cal AS (
 3
       SELECT date
 4
       FROM UNNEST (GENERATE_DATE_ARRAY (
 5
         -- Date minimale
 6
         (SELECT MIN(date commande) FROM {{ref('mart Commandes')}}),
 7
         -- Date maximale
 8
         (SELECT MAX(date_commande) FROM {{ref('mart_Commandes')}})
 9
       ) AS date
10
    ),
```

La seconde CTE, "ca", calcule le chiffre d'affaires journalier à partir des commandes validées, en excluant explicitement les celles qui ont été annulées et ne représentent donc pas un revenu réel. L'agrégation se fait par date de commande, dont la granularité est le jour.

```
ca AS (
11
12
       SELECT
13
         date_commande AS date,
14
         SUM(prix_apres_remise) AS ca
                                       -- Somme des montants apres remise
       FROM {{ref('mart_Commandes')}}
15
       WHERE statut commande != 'Annulée'
16
                                          -- Filtre des commandes annulees
17
       GROUP BY date_commande
18
```

La requête finale assemble enfin ces deux CTE via une jointure gauche, garantissant que toutes les dates du calendrier apparaîtront dans le résultat final. Aux jours sans commandes sons assignés le chiffre 0 en valeur de CA plutôt que NULL, format préféré par Prophet.

```
19 SELECT
20 cal.date AS ds, -- Colonne 'ds' requise par Prophet
21 ROUND(IFNULL(ca.ca, 0), 2) as y -- Colonne 'y' requise par Prophet
22 FROM cal
23 -- Jointure gauche sur les dates pour garder toutes les dates
24 LEFT JOIN ca USING(date)
25 ORDER BY date -- Tri chronologique obligatoire
```

Grâce à cette préparation minutieuse, Prophet peut ingérer les données de chiffre d'affaires et ajuster un modèle de prédiction avec un minimum de transformation de données préalable.

2.3.4 Configuration avancée

La configuration du projet DBT a par la suite été modifiée pour garantir la qualité des données et de la documentation, autoriser l'automatisation du processus et gagner en performance de requêtage.

Le fichier *dbt_project.yml* a tout d'abord été configuré pour implémenter les stratégies de matérialisation mentionnées dans le paragraphe 2.3.2 en page 11.

Une attention particulière a été portée aux fichiers *schema.yaml* qui jouent un double rôle essentiel. Ils documentent précisément chaque modèle (description des tables et colonnes) tout en définissant des tests rigoureux pour assurer la qualité des données. Ces tests vérifient divers aspects critiques : intégrité des relations via les clés étrangères, présence obligatoire de données, unicité des valeurs, ou encore leur conformité à des plages acceptables. Nous avons étendu ces capacités de validation grâce à des extensions (notamment dbt_utils) installées via *packages.yml*, permettant l'implémentation de contrôles plus avancés.

Le bloc de code suivant est un extrait commenté de l'un de ces fichiers, situé dans le dossier "staging", et qui illustre le sérieux avec lequel les modèles ont été documentés et validés. Chaque colonne qui le nécessite bénéficie d'une explication claire de son rôle et format (ex. valeur de promotion exprimée en euros ou dates au format ISO¹). Il est aussi possible d'observer les vérifications plus complexes évoquées ci-dessus. Par exemple, la validation de l'identifiant d'une promotion via une expression régulière garantit le respect du format 'P' suivi de trois chiffres. Le test de validité du pourcentage de promotion, configuré comme un avertissement plutôt qu'une erreur, montre une attention particulière aux besoins métier en permettant des promotions à zéro pour-cent sans bloquer le traitement. La contrainte de clé étrangère assure quant à elle la cohérence relationnelle avec la table des produits.

^{1.} ISO 8601, https://www.iso.org/iso-8601-date-and-time-format.html

```
8
      tests:
 9
       - dbt_utils.expression_is_true:
10
           # Format clé primaire P000
11
         expression: 'REGEXP_CONTAINS(id_promotion, r"^P\d{3}$")'
12
       - dbt_utils.expression_is_true:
13
           # Alerte en cas de promotion négative
14
         expression: 'pourcentage_promotion > 0'
1.5
         config:
16
           severity: warn
17
       columns:
18
         - name: id_promotion
19
           tests:
20
           - not_null # Garantie d'existence de la clé primaire
21
           - unique # Garantie d'unicite de la clé primaire
22
         - name: id_produit
23
           tests:
2.4
           - relationships:
25
             to: ref('Produits') # Controle de cle etrangere
             field: id_produit
26
27
         - name: valeur_promotion
28
           description: en euro
29
         - name: pourcentage_promotion
30
           description: promotion en pourcentage (float entre 0 et 1)
31
         - name: date_debut
32
           description: Date de debut de la promotion 'yyyy-mm-dd'
33
           tests:
34
           - not_null # Garantie d'existence d'une date de début
35
         - name: date_fin
36
           description: Date de fin de la promotion 'yyyy-mm-dd'
     # [...]
```

Cet effort visant à garantir une documentation des modèles de qualité et l'intégrité des données a porté ses fruits, détectant au plus tôt des erreurs dans les données brutes qui auraient pu faire boule de neige et entrainer des problèmes catastrophiques par la suite.

Enfin. l'automatisation des builds via Airflow a requis une configuration spécifique. Après avoir créé un environnement de production dédié dans DBT Cloud (Figure 2.6) – isolant ainsi les exécutions planifiées des développements en cours nous avons configuré un job minimaliste exécutant la commande dbt build. Il assure la reconstruction complète des modèles tout en exécutant l'ensemble des tests décrits dans la section précédente. Pour permettre son déclenchement par Airflow, nous avons généré un token d'accès personnel (PAT) héritant des permissions de notre unique compte utilisateur DBT, puis étudié la documentation de l'API REST¹ de DBT Cloud.

FIGURE 2.6 – Environements du projet

ricing and sign up for a plan at any time during your trial 田 ⟨⟩ Environments Environments Environment variables ③ Deployment PROD X Latest (Formerly known as Versionless) Development DEV

^{1.} Disponible sur le site de DBT, nous avons utilisé la version 2, https://docs.getdbt.com/ dbt-cloud/api-v2#/

Conclusion

de données et en Machine Learning.

L'intégration de DBT dans notre pipeline de données s'est globalement déroulée avec aisance, malgré les défis collaboratifs mentionnés précédemment. La structure modulaire du projet a permis à l'équipe de travailler simultanément sur différents modèles, augmentant ainsi notre productivité. Les modèles spécialisés du dossier *marts*, bien que complexes pour certains, ont permis une grande adaptabilité – la plupart ayant été fréquemment ajustés et optimisés pour s'adapter aux évolutions des besoins en visualisation

DBT s'est donc imposé comme un composant central de notre pipeline. Son architecture modulaire garantit la maintenabilité du code, et que son approche basée sur SQL a facilité son adoption par l'équipe. La solution mise en place, bien qu'ayant révélé certaines limites, constitue une base opérationnelle, solide et prête à être adaptée aux besoins futurs.

2.4 Scripts Python - Machine Learning et Géolocalisation Clients

Cette section présente les modèles de Machine Learning développés en Python et intégrés au pipeline de données. Ces composants démontrent comment enrichir un pipeline ELT classique avec des fonctionnalités analytiques avancées. Chaque script suit une organisation similaire : extraction des données mises à jour depuis BigQuery, traitement et génération de nouvelles données, puis réinjection des résultats dans BigQuery.

Pour répondre au cahier des charges, et au regard des données disponibles, nous avons déployé trois algorithmes :

- Prévisions de chiffre d'affaires utilisant Prophet,
- Analyse des associations de catégories de produits via l'algorithme Apriori,
- Détection des causes de pannes machines avec XGBoost.

En complément, nous avons implémenté un système de géolocalisation des clients par adresse IP, en exploitant une API gratuite aux capacités limitées. Bien que secondaire, ce système soulève des défis techniques intéressants (mécanisme de *caching* et intégration d'une longue tâche dans un DAG Airflow, entre autres).

Les sections suivantes détailleront chaque composant, à l'exception du modèle de *Gradient Boosting* développé par un collaborateur. Les implémentations finales commentées sont disponibles ¹.

^{1.} cf. Annexe B, page 59

2.4.1 Prévision du Chiffre d'Affaires avec Prophet

L'utilisation de Prophet (développé par Facebook) pour la prévision du chiffre d'affaires s'est imposée principalement par à sa facilité de mise en œuvre et d'intégration dans notre pipeline. La bibliothèque propose une API claire et bien documentée, fonctionnant naturellement avec des *DataFrames* Pandas, ce qui simplifie l'interfaçage avec nos processus d'extraction depuis BigQuery et de réinjection des résultats.

Son adoption large a permis de facilement trouver des exemples pour nous inspirer. Par ailleurs, l'interprétation aisée des résultats – présentés sous forme de décomposition en série de Fourier montrant la tendance globale, la saisonnalité annuelle et hebdomadaire – facilite l'analyse des résultats. La possibilité d'ajouter simplement des régresseurs externes a également été considérée dans notre évaluation, même si cette fonctionnalité n'a finalement pas été utilisée dans notre cas particulier. La flexibilité offerte par ses ajustements manuels ciblés permet une adaptation fine aux spécificités de nos données, comme nous le verrons dans l'analyse de code suivante. Enfin, sa rapidité d'exécution préserve la performance de notre pipeline Airflow. L'ensemble de ces caractéristiques ont fait pencher la balance en faveur Prophet plutôt que d'autres solutions plus modernes et peut-être plus performantes (Neural Prophet, Réseaux Neuronaux de type LSTM...) aussi considérées.

Analyse de Code

Les données de chiffre d'affaires sont préparées par un *mart* dédié ¹. Le chargement des données dans un *DataFrame* Pandas s'effectue donc par l'envoi d'une simple requête à BigQuery.

```
1 # imports standards retires pour plus de lisibilite
 2 credentials = service_account.Credentials.from_service_account_file(
 3
    os.path.dirname(__file__) + '/Compte Service Machine Learning.json',
 4)
 5 sql = '''
 6
    SELECT *
 7
    FROM 'carttrend_clean.mart_Prevision_Series_Temporelles_CA'
 9 df = pandas_gbq.read_gbq(
10
11
    project id="carttrend-4",
12
    credentials=credentials
13 )
14 # Conversion de 'dbdate' en 'datetime'
15 df['date'] = pd.to_datetime(df['date'])
```

L'entraînement du modèle et la génération des prédictions suivent une implémentation standard de Prophet.

```
16 from prophet import Prophet
17 # Initialisation du modele Prophet
18 model = Prophet()
```

^{1.} cf. §A.10, page 57

```
# Entrainement du modele (il attend des noms de colonnes specifiques)
model.fit(df.rename(columns={'date': 'ds', 'chiffre_affaire': 'y'}))
# Creation d'un dataframe avec 30 jours supplementaires
future = model.make_future_dataframe(periods=30, freq='D')
# Prediction du chiffre d'affaire sur ces 30 jours
forecast = model.predict(future)
```

Les résultats peuvent facilement être visualisés (Figure 2.7).

25 model.plot(forecast)

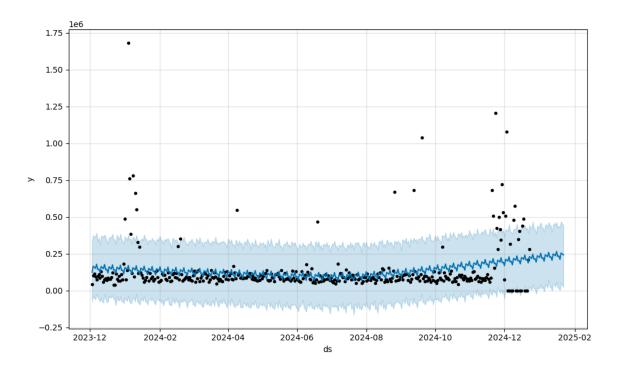


FIGURE 2.7 – Visualisation de l'analyse Prophet brute

L'analyse de ces résultats révèle des performances décevantes du modèle : intervalle de confiance très large et détection erronée d'une tendance haussière. L'examen des données montre en réalité une stabilité globale du chiffre d'affaires, perturbée par deux périodes atypiques :

- La première quinzaine de janvier,
- Entre fin novembre et mi-décembre.

Ces événements ne présentent pas de caractère saisonnier et se distinguent par leur irrégularité (pics de chiffre d'affaires, alternés avec des jours sans vente pour la deuxième période).

Pour améliorer les prédictions, nous avons appliqué l'une des stratégies recommandées dans la documentation de Prophet¹, en traitant ainsi ces périodes comme des vacances ponctuelles. Les résultats sont visibles en Figures 2.8 et 2.9.

16 vacances = pd.DataFrame([

^{1.} Spécifiquement la partie *Treating COVID-19 lockdowns as a one-off holidays* de la section *Handling Shocks*, https://facebook.github.io/prophet/docs/handling_shocks.html

```
17 {
     # Premiere periode de pics, du 30/12/23 pour 20 jours
18
19
     'holiday': 'choc 1',
     'ds': pd.to_datetime('2023-12-30'),
20
21
     'lower_window': 0,
     'upper_window': 20,
22
23 },
24 {
25
     # Deuxieme periode de pics, du 20/11/24 pour 28 jours
26
     'holiday': 'choc 2',
27
     'ds': pd.to_datetime('2024-11-20'),
     'lower_window': 0,
2.8
     'upper_window': 28,
29
30 }
31 ])
32
33 model = Prophet (holidays=vacances)
34 # Le reste du code est inchange
```

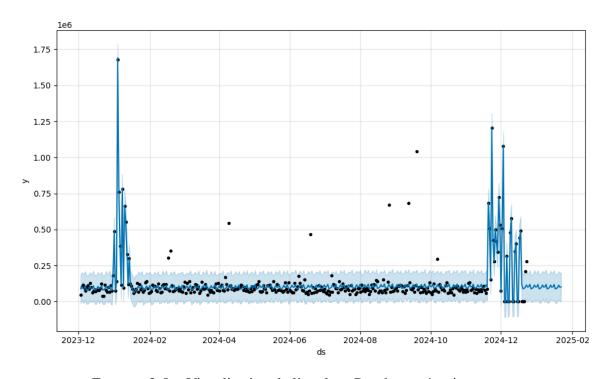


FIGURE 2.8 – Visualisation de l'analyse Prophet après ajustements

Les résultats après ajustement montrent des améliorations limitées. Bien que le modèle parvienne désormais à neutraliser correctement les deux périodes atypiques et à identifier une tendance globale plate plus conforme aux résultats attendus, ses performances restent globalement médiocres. L'intervalle de confiance reste large, reflétant l'incapacité du modèle à réaliser une prédiction fiable. La détection d'une légère saisonnalité hebdomadaire est plausible, mais trop faible pour être réellement utilisable.

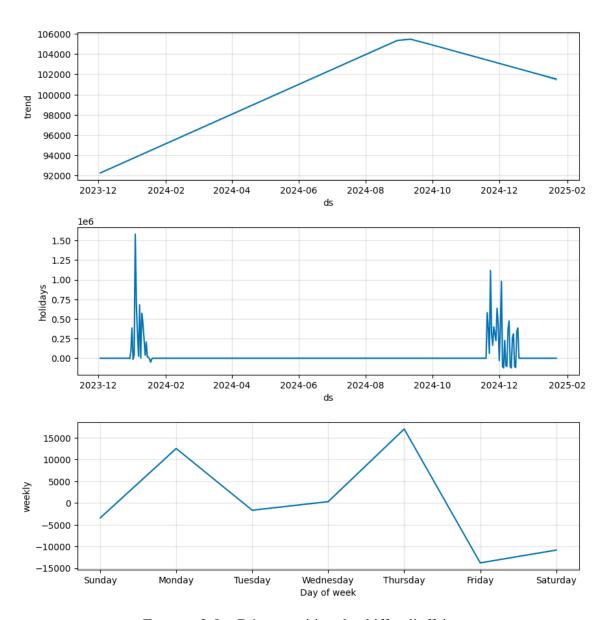


FIGURE 2.9 – Décomposition du chiffre d'affaires

Interprétation

Bien que techniquement correcte dans son implémentation, le modèle se révèle très décevant en terme qualité de prédiction. Cette contre-performance s'explique entièrement par les caractéristiques de nos données d'entrée. Généré aléatoirement pour les besoins de l'exercice, notre jeu de données présente des pics imprévisibles et une absence de tendances saisonnière. Pour la même raison, il est impossible de trouver des régresseurs externes pertinents pour améliorer la performance du modèle. Enfin, le volume historique disponible est limité (à peine plus d'une année de données).

Avec des données réelles, plusieurs pistes d'amélioration pourraient être explorées. L'intégration de régresseurs externes (campagnes marketing, indicateurs économiques généraux, vacances scolaires...) pourrait améliorer le modèle Prophet existant. Il serait aussi possible de

recourir à des algorithmes plus récents comme NeuralProphet ou, si des moyens techniques significatifs étaient disponibles, plus sophistiqués (réseaux neuronaux à mémoire par exemple LSTMs ou Transformers). Enfin, simplement attendre l'accumulation d'un volume suffisant de données historiques avant de déployer un système de prédiction pourrait aider à réaliser des prédictions plus fiables.

Cette expérience souligne avant tout le caractère purement académique de cet exercice. Un modèle aux performances aussi limitées ne devrait en aucun cas dépasser le stade exploratoire, et encore moins servir de base à des décisions stratégiques. Ce travail aura cependant permis de démontrer l'intégration technique complète de Prophet dans notre pipeline de données. Plus fondamentalement, il rappelle l'importance cruciale d'une analyse rigoureuse des données, tant en amont qu'en aval de toute modélisation.

2.4.2 Analyse des associations de catégories de produits avec l'algorithme Apriori

L'objectif initial de cette analyse était d'identifier des groupes de produits fréquemment achetés ensemble afin de pouvoir générer des recommandations pertinentes. Cependant, la structure particulière de nos données – chaque commande ne contient qu'un seul produit – a immédiatement posé problème pour une analyse classique des paniers d'achat. Face à cette contrainte, nous avons adapté notre approche en considérant l'ensemble des achats historiques par client comme un unique "panier étendu dans le temps". L'algorithme Apriori a alors été appliqué pour trouver des combinaisons de produits fréquemment associées au sein de ces historiques clients, même lorsque les achats étaient espacés dans le temps.

Malheureusement, cette analyse a été peu concluante, les métriques *support*, *lift* et *confidence* restant toutes décevantes. Cette impasse nous a conduits à modifier la granularité de notre analyse en appliquant Apriori au niveau des catégories de produits plutôt qu'aux produits individuels.

Analyse de Code

Comme pour les autres modèles, un *mart* dédié a été implémenté. Chaque ligne représente un *Array* SQL de catégories de produits achetés par client, dans l'ordre chronologique. L'obtention des données se fait toujours de la même manière.

```
# imports standards retires pour plus de lisibilite
credentials = service_account.Credentials.from_service_account_file(
    'Compte Service Machine Learning.json',
)
sql = '''
SELECT *
FROM 'carttrend_clean.mart_Apriori_Combinaisons_Produits'
'''
```

```
9 df = pandas_gbq.read_gbq(
10 sql,
11 project_id="carttrend-4",
12 credentials=credentials
13 )
14 df
```

```
categorie
0 [Jouets, Mode, Sports, Beaute]
1 [Mode, Livres]
2 [Mode, Maison, Livres, Electronique, Beaute]
3 [Mode, Electronique, Livres, Sports]
4 [Mode, Maison, Sports]
```

Cette *dataframe* est ensuite transformée au format attendu par l'objet *apriori* d'*mlxtend*, soit une liste de liste de paniers. Celle-ci est ensuite encodée, puis le modèle est entrainé.

```
from mlxtend.preprocessing import TransactionEncoder
from mlxtend.frequent_patterns import apriori

# Transformation des donnees en format liste de listes
paniers = [list(row[0]) for row in df.itertuples(index=False)]

# Encodage des "paniers" au format one-hot
te = TransactionEncoder()
te_ary = te.fit(transactions).transform(transactions)
dfml = pd.DataFrame(te_ary, columns=te.columns_)

# Execution de l'algorithme Apriori avec seuil minimal de support
frequent_itemsets = apriori(dfml, min_support=0.1, use_colnames=True)
frequent_itemsets
```

```
support itemsets
   0.122124 (Alimentation)
1
   0.116323 (Beaute)
  0.100320 (Jouets)
2
  0.897580 (Livres)
3
  0.234247 (Maison)
5
  0.962392 (Mode)
  0.638328 (Sports)
6
   0.564513 (Electronique)
7
   0.108822 (Livres, Alimentation)
8
   0.117023 (Mode, Alimentation)
```

À titre exploratoire, il est intéressant d'essayer d'extraire des règles d'association (ex : X est souvent acheté avant Y...)

La sortie ci-dessous a été altérée pour augmenter la lisibilité.

```
antecedents consequents ac supp cs supp supp conf lift
```

```
40
                                               0.225
                                                        0.130
                                                                 0.228
    (Sports, ...)
                     (Maison, ...)
                                     0.571
                                                                         1.016
43
    (Maison, ...)
                     (Sports, ...)
                                     0.225
                                               0.571
                                                        0.130
                                                                 0.580
                                                                         1.016
15
    (Sports, ...)
                     (Maison)
                                     0.571
                                               0.234
                                                        0.136
                                                                 0.238
                                                                         1.016
                     (Sports, ...)
18
                                     0.234
                                               0.571
                                                         0.136
                                                                 0.580
                                                                         1.016
   (Maison)
                                                                 0.213
17 (Sports)
                     (Livres, ...)
                                     0.638
                                               0.209
                                                         0.136
                                                                         1.016
```

Les résultats n'étant que peu concluant (voir section suivante), et la question posée par l'énoncé ne portant que sur les catégories de produits achetées ensemble, il a finalement été décidé de limiter l'analyse aux associations courantes. Ces dernières ont été sélectionnées si elles étaient composées d'au moins deux éléments, et si elles jouissaient d'un support important. La visualisation sur le tableau de bord est enfin simplifiée en transformant la sortie du script en une table de chaines de caractères.

```
29 # Filtrage des resultats :
30 # - Combinaisons d'au moins 2 categories
31 # - Support minimum de 40%
32 frequent_itemsets['length'] = frequent_itemsets['itemsets'].apply(
    lambda x: len(x)
33
34)
35
36 frequent_itemsets = frequent_itemsets[
     (frequent_itemsets['length'] >= 2) &
37
38
     (frequent_itemsets['support'] >= 0.4)
39 ][['support', 'itemsets']]
40
41 # Formatage des resultats (frozenset -> string), facilite la visualisation
42 frequent_itemsets['itemsets'] = frequent_itemsets['itemsets'].apply(
    lambda x: ' - '.join(list(x))
44)
45
46 frequent_itemsets
```

```
support itemsets
13  0.864573  Mode - Livres
14  0.571214  Sports - Livres
15  0.507401  Livres - Electronique
19  0.613523  Mode - Sports
... ...
```

Interprétation

Cette analyse visait initialement à identifier des combinaisons de catégories de produits fréquemment achetées ensemble pour mieux comprendre les comportements d'achat. L'application de l'algorithme Apriori a effectivement permis d'extraire certaines associations significatives en termes de support, comme les paires "Mode + Livres" (86,5%) ou "Sports + Livres" (57,1%), révélant ainsi des co-occurrences régulières dans les historiques d'achats. Cependant, une analyse plus approfondie des métriques d'association tempère considérablement ces résultats. Si certaines règles présentent une confiance élevée (comme "Électronique + Mode" à 96%), les valeurs de lift systématiquement proches de 1 indiquent en réalité une absence de corrélation significative entre les catégories.

Ici encore, l'analyse est limitée par nos données générées aléatoirement. Mais au-delà de ces résultats, cette étude révèle aussi les limites de performance de l'algorithme Apriori, auxquelles nous avons été confrontés. Pour exploiter un plus gros volume de données, l'investissement dans une puissance de calcul plus importante ou l'exploration d'un algorithme plus efficace, comme FP-Growth, pourrait s'avérer nécessaire. Sur le plan technique, ce travail a néanmoins validé l'intégration fonctionnelle d'Apriori dans notre pipeline. Il met aussi en lumière l'importance d'une analyse rigoureuse des métriques pour éviter les interprétations erronées.

2.4.3 Géolocalisation des Clients avec leur adresse IP

La géolocalisation des clients à partir de leur adresse IP n'était pas une exigence initiale du cahier des charges, mais elle représentait un défi technique intéressant ainsi qu'une illustration de l'intégration d'une nouvelle source de donnée au pipeline. Dans l'absence d'information directe sur la position géographique des clients dans les données, une API de géolocalisation gratuite utilisant leurs adresses IP a été utilisée. Toutefois, comme souvent avec les services gratuits, des limitations importantes étaient à prendre en compte, notamment un temps minimal entre deux requêtes fixé à 0,2 seconde (*rate limit*). Afin de respecter cette contrainte et pour éviter des appels redondants, un système de (*caching*) a été mis en place. Il permet de stocker localement les résultats préalablement obtenus pour une durée définie.

Analyse de Code

Les données sont directement prélevées sur la table Clients nettoyée, en sélectionnant les colonnes *id_client* et *ip*.

```
1 # imports standards retires pour plus de lisibilite
 2 credentials = service account.Credentials.from service account file(
 3
     'Compte Service Machine Learning.json',
 4)
 5 sql = '''
    SELECT id client, ip
 6
 7
    FROM 'carttrend-4.carttrend_clean.Clients'
 8 ′′′
 9 df = pandas_gbq.read_gbq(
10
11
    project_id="carttrend-4",
    credentials=credentials
12
13)
14 df
```

Le cache est implémenté en utilisant le module pickle de Python ¹. Beaucoup d'objets Python (et en particulier les DataFrame Pandas) peuvent être sérialisés de cette manière puis stockée sur le disque dur, pour être chargées en mémoire par la suite. Il est ainsi possible de connaître la date de dernière mise à jour du cache en consultant, via le système d'exploitation, la date de modification du fichier. Si elle est jugée trop ancienne (ici, supérieure à une semaine), le cache est ignoré (et ensuite écrasé par les nouvelles données, comme expliqué ci-après).

```
15 # Initialisation du DataFrame de résultats
16 client_geo = pd.DataFrame(columns=['id_client', 'pays', 'region', 'ville'])
17
18 # Gestion du cache
19 file_age = None
20 max_age = datetime.timedelta(days=7) # Durée de validité du cache
21 if os.path.isfile(os.path.dirname(__file__) + '/client_geo.pkl'):
    filetime = os.path.getmtime(
22
       os.path.dirname(__file__) + '/client_geo.pkl'
23
24
    file_age = datetime.timedelta(seconds=time.time()-filetime)
25
    if file_age < max_age: # Utilisation du cache si valide</pre>
2.6
27
      client_geo = pd.read_pickle(
         os.path.dirname(__file__) + '/client_geo.pkl'
28
29
```

On construit ensuite une DataFrame contenant uniquement les clients qui ne sont pas déjà présents dans le cache. Si ce dernier a été ignoré, la DataFrame "client_geo" est vide, auquel cas l'ensemble des clients est ajouté à "df_missing".

```
30 df_missing = pd.concat([df, client_geo]).drop_duplicates(
31    subset = ['id_client'],
32    keep=False
33 )
```

On appelle ensuite l'API ip2location pour remplir la DataFrame "client_geo". L'API ne permettant pas, à ma connaissance, de traitement en lot, une requête est effectuée pour chaque adresse IP (donc pour chaque client). Si le serveur ne répond pas ou si une erreur se produit, on ignore le client en cours (dont l'IP contient peut etre une erreur) et l'on essaye avec le client suivant, dans la limite de 5 tentatives.

```
34 retries = 5 # Nombre de tentatives en cas d'échec
35 for row in df_missing.itertuples(index=False):
36
    try:
       # Appel à l'API ip2location
37
38
      response = requests.get(
       f'https://api.ip2location.io/?key=7149995C158ABB38F6BCE6505B006082&ip={
39
          row.ip}&format=json',
40
      timeout=10
41
       )
42
43
       # Si bonne reponse...
44
      if response.status_code == 200:
45
         data = response.json()
46
         # ...ajout des données de géolocalisation aux résultats
```

^{1.} pickle fait parti de la bibliothèque standard de Python, https://docs.python.org/3/library/ pickle.html

```
47
         client_geo = pd.concat([
48
           client_geo,
49
             pd.DataFrame([[row.id_client,
50
             data.get('country_name'),
51
             data.get('region_name'),
52
             data.get('city_name')]],
53
             columns=client geo.columns)
         ], ignore_index=True)
54
55
       # Sinon...
56
       else:
57
         # ...si il y a encore des tentatives disponibles...
58
         if retries:
           retries -= 1
59
60
           time.sleep(1) # Attend plus longtemps en cas d'echec
61
           continue #...on continue avec le client suivant
62
         else:
63
           raise Exception(f'API error: {response.status_code}')
64
65
    except Exception as e:
66
       print(f"Error processing IP {row.ip}: {str(e)}")
67
      break
68
    time.sleep(0.2) # Respect du rate limiting de l'API
69
```

Enfin, si le cache était obsolète, on l'écrase avec les nouvelles données.

```
70 if not file_age or file_age > max_age:
71   client_geo.to_pickle(
72   os.path.dirname(__file__) + '/client_geo.pkl'
73 )
```

Interprétation

Les résultats obtenus, affichés dans le *dashboard* "Ventes - Produits" ¹, correspondent bien aux attentes : les adresses IP des clients étant aléatoires dans les données, les pays les plus représentés ne reflètent pas une réalité commerciale, mais simplement la répartition mondiale des plages d'adresses IP. Les États-Unis et la Chine, qui disposent de larges blocs d'adresses, apparaissent donc très fréquemment dans les résultats. Par ailleurs, même avec des données plus réalistes, la précision de cette approche reste limitée. L'utilisation éventuelle par les clients de VPNs, les erreurs de l'API gratuite ou encore l'évolution constante de l'allocation des adresses IP peuvent facilement fausser les résultats.

D'un point de vue technique, ce script introduit aussi une contrainte non négligeable dans le pipeline : à chaque rafraîchissement complet du cache, il ajoute environ 0,2 seconde par client au temps d'exécution du DAG, en raison de la limitation imposée par l'API utilisée. Cette limitation aurait pu être partiellement contournée par une implémentation plus fine du cache : la table de sortie aurait pu inclure une date de dernière mise à jour pour chaque ligne dans une colonne supplémentaire. On aurait alors la possibilité de ne mettre à jour que les clients dont les données sont devenues obsolètes, grâce à cette granularité accrue. Cela n'aurait pas

^{1.} cf. Annexe D.5, page 75

entièrement résolu le problème de performance, mais aurait permis de limiter la forte surcharge du pipeline lors de la reinitialisation complète des données. Cette approche aurait en outre l'avantage d'éviter le chargement d'un pickle depuis le disque dur, opération qui peut être problématique en termes de sécurité ¹.

Au final, cet exercice reste majoritairement académique. Une entreprise comme Carttrend disposerait très probablement des adresses postales de ses clients (ne serait-ce que pour pouvoir honorer les commandes...), bien plus précises et robustes que les données dérivées d'IP. Toute-fois, cette expérimentation constitue une bonne démonstration de l'intégration d'une source de données externe dans un pipeline de traitement, avec les adaptations techniques et les précautions que cela implique.

Conclusion

Dans l'ensemble, les travaux réalisés autour des scripts Python pour l'analyse prédictive, l'exploration des comportements d'achat et la géolocalisation des clients ont permis de valider des aspects techniques clés du pipeline, tout en mettant en évidence les limites imposées par la nature synthétique des données utilisées. Le modèle de prévision basé sur Prophet, a montré des performances très faibles, l'analyse des associations de produits n'a pu produire que des résultats peu pertinents, révélant l'absence de signal réel dans les données. La géolocalisation des clients, bien qu'extérieure au périmètre du projet, a illustré l'intégration d'un service externe dans un pipeline de données. Là encore, la nature artificielle des données limite toute interprétation opérationnelle.

Un modèle de régression basé sur XGBoost a également été développé par un collaborateur pour tenter d'identifier les principales causes des pannes machines à partir de variables explicatives telles que la localisation de l'entrepôt, la température ou le volume de produits traités. Là encore, les résultats se sont révélés très peu probants, avec une erreur (RMSE) élevée, soulignant le caractère académique de l'exercice et l'absence de lien causalité entre les variables générées aléatoirement.

Dans un contexte réel, une étape de validation automatisée des modèles serait recommandée avant toute intégration des résultats dans BigQuery. Des métriques de performance (RMSE, précision, *recall*, etc.) pourraient être évaluées systématiquement, et seuls les résultats satisfaisants seraient envoyés vers le *data warehouse*. Cela éviterait de présenter à l'utilisateur final des prédictions et analyses peu fiables (comme celles générées dans ce projet). Cette précaution garantirait que seules les analyses pertinentes alimentent les *dashboards*.

Au final, ces travaux rappellent de manière transversale que la qualité des données conditionne directement la pertinence des analyses, et que les algorithmes, aussi perfectionnés soient-

^{1.} Un utilisateur mal intentionné ayant accès au serveur Airflow pourrait théoriquement modifier le fichier pickle de telle manière à lui permettre une exécution de code arbitraire avec les droits d'accès octroyés au DAG. Explications et exemples dans un article de Nicolas Lara: https://lincolnloop.com/blog/playing-pickle-security/

ils, ne peuvent produire de résultats robustes que s'ils s'appuient sur des fondations solides. Sur le plan pédagogique et technique, ce travail à cependant rempli ses objectifs : il démontre l'intégration de plusieurs briques analytiques dans un pipeline cohérent, avec une attention spéciale portée à la validation et l'interprétation des résultats.

2.5 Orchestration avec Airflow

L'ensemble des traitements développés dans ce projet repose sur une orchestration centralisée assurée par Apache Airflow. Ce composant joue un rôle crucial dans l'automatisation, la planification et le suivi de l'exécution des différentes tâches du pipeline de données. Il permet d'enchaîner les étapes de manière contrôlée, tout en assurant la gestion des dépendances et des erreurs éventuelles.

Le DAG principal mis en place regroupe six tâches, reflétant les différentes étapes de la chaîne analytique. Il commence par déclencher le transfert de données Fivetran, qui alimente le *data warehouse* en données brutes. Une fois ce transfert terminé, Airflow déclenche la commande dbt build, qui applique les transformations SQL pour construire les *marts* nécessaires aux analyses. Viennent enfin la tâche de géolocalisation des clients et les trois tâches distinctes de Machine Learning : la prévision du chiffre d'affaires (Prophet), l'analyse des causes de panne des machines (XGBoost), et l'analyse des comportements d'achat (Apriori).

Cette orchestration garantit un traitement complet des données, du chargement initial dans le *data warehouse* jusqu'à leur enrichissement final.

2.5.1 Analyse de code : DAG

Le DAG est implémenté de manière à simplifier l'intégration de nouvelles tâches. Celles-ci correspondent en effet à un simple appel à leur fonction respective. Chacune de ces fonctions appelle ensuite un script Python externe : le processus principal donne naissance à un sousprocessus chaque tâche, à l'aide de la fonction run () du module subprocess. Le DAG commence donc par la définition de ces fonctions.

```
from airflow import DAG
from airflow.operators.python import PythonOperator
from datetime import datetime, timedelta
import subprocess

def start_fivetran_sync():
    # Déclenche la synchronisation des données sources via Fivetran
    print("Starting Fivetran Sync Task")
    subprocess.run(['python3', '../../airflow/scripts/start_fivetran_sync.py'
    ])

def start_dbt_build():
```

```
12 # Lance les transformations DBT (staging + marts)
13
    print('Starting DBT Build Job')
     subprocess.run(['python3', '../../airflow/scripts/start_dbt_build_job.py'
14
        1)
15
16 def start_ml_prediction_ca():
    # Exécute le modèle Prophet de prévision de CA
17
    print('Starting ML Prediction CA')
18
     subprocess.run(['python3', '../../airflow/scripts/ml_prediction_ca.py'])
19
20
21 def start_ml_a_priori():
2.2
   # Exécute l'analyse des associations produits
23
    print('Starting ML APriori')
24
    subprocess.run(['python3', '../../airflow/scripts/ml_a_priori.py'])
25
26 def start_ml_regression_machines():
     # Exécute le modèle XGBoost d'analyse des pannes
27
28
    print('Starting ML Regression Machines')
29
    subprocess.run(['python3', '../../airflow/scripts/ml_regression_machines.
        py'])
30
31 def get_customer_locations():
32
    # Appel l'API ip2location pour géolocaliser les clients
33
    print('Starting Customer Location API Call Task')
     {\tt subprocess.run} \, (\, [\, '\, python 3'\, , \, \, '\, \ldots /\, \ldots /\, airflow/scripts/get\_customer\_location\, .
34
```

L'objet DAG est ensuite instancié avec les paramètres appropriés. Ici, le DAG est programmé pour être déclenché tous les jours à minuit. En cas d'échec, une seconde tentative est effectuée 30 minutes plus tard. Si celle-ci échoue à nouveau, un e-mail d'alerte est envoyé à l'adresse spécifiée, grâce à la configuration préalable du protocole SMTP.

```
35 default_args = {
36
    'owner': 'airflow',
37
    'start_date': datetime(2025, 5, 27),
    'schedule_interval': '@daily', # Exécution quotidienne à minuit
38
39
    'retries': 1, # Nombre de tentatives en cas d'échec
40
    'retry_delay': timedelta(minutes=30),
                                           # Délai entre les tentatives
     'email_on_failure': True,
41
                               # Notification par email (configuration SMTP
        prealable)
42
     'email': ['carttrenddata4@gmail.com'] # Destinataires
43 }
44
45 # Initialisation du DAG
46 \text{ dag} = DAG(
    'carttrend_data_pipeline',
47
                                 # Nom du DAG
48
    default_args=default_args,
49
    description='Pipeline ELT/ML complet pour CartTrend',
50
    catchup=False # Désactive les runs rétroactifs
51)
```

Les tâches sont ensuite définies. Chacune d'entre elles correspond à l'une des fonctions décrites ci-dessus. Le code ci-dessous est volontairement abrégé pour plus de lisibilité, mais la version complète est fournie ¹.

^{1.} cf. Annexe C.1, page 65

```
52 t1 = PythonOperator(
53
    task_id="start_fivetran_sync",
    python_callable=start_fivetran_sync,
54
55
    dag=dag,
56)
57
58 t2 = PythonOperator(
    task_id="start_dbt_build",
59
    python_callable=start_dbt_build,
60
61
     dag=dag,
62)
63
64 t3 = PythonOperator(
65
    task_id="start_ml_prediction_ca",
66
    python_callable=start_ml_prediction_ca,
67
    dag=dag,
68)
69 # ...
```

Enfin, les dépendances entre les tâches sont définies. Le pipeline commence systématiquement par t1 (extraction et chargement des données brutes), suivi de t2 (transformation via DBT). Ensuite, les trois modèles de Machine Learning (t3, t4 et t5) ainsi que la géolocalisation des clients (t6) peuvent s'exécuter en parallèle.

```
70 t1 >> t2
71 t2 >> [t3, t4, t5, t6]
```

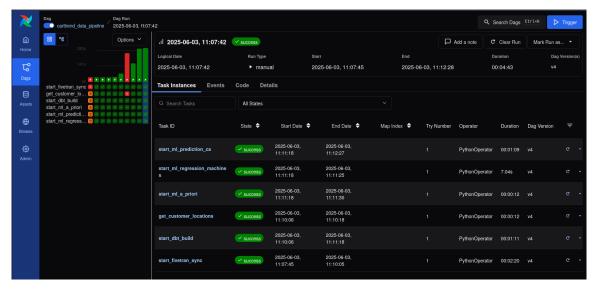


FIGURE 2.10 – Six carrés verts, le DAG s'execute avec succès!

L'architecture retenue pour ce DAG présente plusieurs avantages importants. En déléguant chaque tâche à un script Python externe via subprocess.run(), elle permet une évolution simple et modulaire du pipeline: l'ajout d'une nouvelle étape se résume à l'écriture d'un script autonome, suivi de l'ajout d'une fonction et d'une tâche dans le DAG. Cette approche favorise la séparation des responsabilités et facilite la maintenance du code sur le long terme.

Cette structure a toutefois quelques limitations. Le fait d'exécuter chaque étape dans un

sous-processus isolé rend plus difficile la transmission d'informations entre tâches. Par exemple, il n'est pas possible d'utiliser directement le mécanisme de communication natif d'Airflow (XCom). Ce manque peut cependant être partiellement contourné : l'une des possibilités est d'utiliser des print () dans les scripts appelés, et de capturer les sorties standard pour les analyser ou les *logguer* depuis le DAG. Ce type de contournement nécessite néanmoins un peu plus d'effort de mise en œuvre si des échanges de données plus complexes sont envisagés.

Par ailleurs, l'ordonnancement défini permet une exécution parallèle de plusieurs traitements indépendants, notamment les trois modèles de Machine Learning et la géolocalisation. Cela optimise le temps d'exécution global du pipeline, en tirant parti de la capacité d'Airflow à gérer plusieurs tâches en parallèle.

2.5.2 Analyse de code : start_dbt_build_job.py

start_dbt_build_job.py partage avec le script *start_fivetran_sync.py* une logique similaire : appeler un service externe, attendre la fin de son exécution, et signaler toute erreur en levant une exception. Dans ce rapport, seule l'analyse du premier est détaillée; l'autre, quasiment identique dans sa conception, est disponible ¹.

start_dbt_build_job.py est d'abord de déclencher l'exécution du job DBT configuré précédemment (cf. §Configuration avancée, page 18). Ce dernier va appliquer les modèles de transformation SQL définis dans le projet. Le script attend ensuite patiemment que le processus se termine, vérifie le code de retour, et lève une exception explicite en cas d'échec ou d'absence de réponse, permettant ainsi à Airflow d'identifier correctement une erreur et de réagir selon la configuration du DAG.

Le script commence par les imports nécessaires et la définition de quelques variables globales :

```
1 import requests
 2 import time
 3 import enum
 5 # Configuration d'accès à 1'API DBT Cloud
 6 # Tout ces identifiants ont ete crees sur l'interface DBT Cloud
 7 ACCOUNT_ID = '70471823462199' # Identifiant du compte DBT
 8 JOB_ID = '70471823467186' # Identifiant du job
 9 API_KEY = 'dbtu_57j1bQUrnaKY_Xf03UdF2kHC9q4DkqZNAmyAy2D3QAqQmxIqiM' # Token
10
11 # Enumeration des statuts possibles d'un job DBT
12 # Les codes sont documentés sur la doc DBT :
13 # https://docs.getdbt.com/dbt-cloud/api-v2#/operations/List%20Runs
14 class DbtJobRunStatus (enum.IntEnum):
15 OUEUED = 1
    STARTING = 2
16
17 RUNNING = 3
```

^{1.} cf. Annexe C.2, page 67

```
18 SUCCESS = 10
19 ERROR = 20
20 CANCELLED = 30
```

Une fonction est ensuite définie pour déclencher le *job* via l'API de DBT. Cette API fourni l'identifiant unique de l'exécution, que la fonction récupère et renvoie.

```
21 def trigger_job():
    url = f'https://jp371.usl.dbt.com/api/v2/accounts/{ACCOUNT_ID}/jobs/{
2.2
        JOB ID}/run/
23
    headers = {'Authorization' : f'Token {API_KEY}'}
    payload = {'cause' : 'Triggered from Airflow'}
24
25
    # Envoi de la requete POST a l'API DBT
26
    response = requests.post(url, json=payload, headers=headers)
27
    response.raise_for_status() # Si echec, leve une exception
28
29
30
    # Retourne l'ID d'exécution
31
    job_run_id = response.json()['data']['id']
32 return job_run_id
```

Une deuxième fonction est déclarée pour interroger l'état du *job* en cours, toujours via l'API DBT.

```
33 def get_job_status(job_run_id):
    url = f'https://jp371.usl.dbt.com/api/v2/accounts/{ACCOUNT_ID}/runs/{
        job_run_id}/'
35
    headers = {'Authorization': f'Token {API_KEY}'}
36
37
     # Envoi de la requete GET a l'API DBT
    response = requests.get(url, headers=headers)
38
39
    response.raise_for_status()
40
41
     # Retourne le statut de l'execution
42 return response.json()['data']['status']
```

Le point d'entrée principal du script se trouve dans le bloc __main__. La fonction trigger_job() est appelée pour lancer le job, et retourne l'identifiant de son exécution. Cet identifiant est ensuite passé à la fonction get_job_status(), appelée régulièrement (toutes les trois secondes) pour suivre l'avancement du job.

Il est important de noter qu'un échec temporaire à ce deuxième appel à l'API DBT (ex. *timeout* ou erreur réseau) ne signifie pas forcément que le *job* a échoué. C'est pourquoi un mécanisme de *retries* est mis en place : si la récupération du statut échoue plus cinq fois, le script considère qu'il est impossible de déterminer le résultat du job et lève une exception.

```
43 if __name__ == '__main ':
    retries = 5 # Nombre maximum de tentatives d'obtention du statut
44
45
46
    # Phase 1: Déclenchement du Job
47
    job_run_id = trigger_job()
    print(f'DBT Build Job Triggered : {job_run_id}')
48
49
50
    # Phase 2: Verification periodique du statut du Job
51
    while True:
52 # Demande l'etat du transfer a Fivetran
```

```
53
       try:
54
         status = get_job_status(job_run_id)
55
56
       # Lève une exception si la demande de statut n'aboutit pas 5 fois
57
       except Exception as e:
58
         print(f'Could not get job status : {e} - Retries : {retries}')
59
         if retries:
           retries -= 1
60
61
         else:
62
           raise Exception('Could not get job status')
63
64
       # Si on a une réponse de l'API...
65
       else:
66
         # ...et qu'il est fini : stop.
         if status == DbtJobRunStatus.SUCCESS:
67
68
         # ...et qu'il s'est produit une erreur : exception.
69
70
         elif status == DbtJobRunStatus.ERROR \
71
           or status == DbtJobRunStatus.CANCELLED:
72
           raise Exception("Job Failed !")
73
74
       # ...attendre 3 secs et recommencer.
75
      time.sleep(3)
76
77
   print('Job Finished')
```

Ce script assure le déclenchement fiable d'un *job* DBT à distance, avec un suivi en temps réel de son exécution. Il gère correctement les erreurs, ce qui le rend tout à fait adapté à son intégration dans notre environnement orchestré par Airflow.

Conclusion

L'intégration d'Airflow dans ce projet s'est révélée être un choix pertinent, notamment en termes de flexibilité et de centralisation de l'orchestration du pipeline de données. Chaque étape – extraction, transformation et enrichissement externe – est orchestrée de manière claire et séquencée, tout en permettant l'exécution parallèle des traitements indépendants. Cette orchestration aurait été plus approximative et moins facile à maintenir en cas d'utilisation des alternatives potentielle évoquées en premiere partie (§Orchestration - Airflow, page 6).

L'architecture retenue, fondée sur l'appel de scripts Python externes, est orientée vers l'évolutivité. Ajouter une tâche ou en modifier une existante ne nécessite qu'un ajustement minimal dans le DAG. En contrepartie, cette approche demande plus d'attention pour la gestion des erreurs, des logs ou encore des transferts de résultats intermédiaires.

Enfin, l'utilisation d'Airflow en mode *standalone*, avec une base SQLite et un exécuteur séquentiel, entraine certaines limitations comme l'absence d'une véritable parallélisation des tâches et une réduction des capacités de gestion des volumes importants. Ces contraintes, pas réellement limitantes dans le contexte de cet exercice, n'existeraient de toute manière pas dans un environnement de production, où une configuration soigneuse avec un *scheduler* dédié, un

exécuteur parallèle et une base de données robuste (PostgreSQL, par exemple) aurait été réalisée.

Conclusion

Cette deuxième partie a permis de détailler, étape par étape, la mise en œuvre technique du pipeline de données. Tout d'abord, l'infrastructure a été posée avec la configuration de Big-Query, utilisée comme entrepôt de données central. Fivetran assure ensuite l'extraction et le chargement des données brutes issues de sources variées. Ces données sont alors structurées, nettoyées et enrichies au moyen de modèles DBT, en facilitant ainsi l'utilisation postérieure. Des scripts Python ont ensuite été intégrés pour effectuer plusieurs traitements spécifiques : prédiction des ventes, analyse des causes de pannes, extraction d'associations fréquentes, ainsi qu'un enrichissement géographique basé sur les adresses IP. Enfin, l'ensemble du processus a été orchestré avec Airflow, garantissant un déclenchement cohérent, un suivi centralisé du pipeline et une exécution efficace des tâches.

Certaines améliorations apparaissent prioritaires si un passage en production était envisagé : remplacer les quelques utilisations d'OAuth par des comptes de service, adopter une configuration Airflow plus robuste (avec exécuteur parallèle et base de données adaptée), et revoir l'architecture des modèles DBT pour la rendre plus modulaire, lisible et maintenable. Cependant, le pipeline construit est dès à présent fonctionnel, adapté aux objectifs pédagogiques du projet, et facilement extensible. Chaque composant peut être remplacé ou amélioré indépendamment des autres, ce qui en fait une base robuste pour de futures évolutions.

La partie suivante s'appuie sur ce pipeline pour visualiser les données dans Tableau, illustrant ainsi l'importance d'un traitement amont rigoureux pour pouvoir construire des tableaux de bord fiables et fluides.

Chapitre 3

Tableaux de bords, analyse et recommandations

L'objectif principal de cette phase du projet a été de concevoir une série de *dashboards* interactifs dans Tableau, permettant une lecture claire, rapide et utile des données pour les utilisateurs finaux. L'accent a été mis sur la simplicité d'utilisation au quotidien, avec des visuels immédiatement compréhensibles et des filtres pertinents, adaptés aux besoins des différents métiers. Un effort particulier a été porté sur l'harmonisation visuelle entre les différents tableaux de bord : couleurs, respect d'une charte graphique, positionnement des filtres et légendes... Ces choix esthétiques, souvent perçus comme secondaires, contribuent pourtant fortement à l'adhésion des utilisateurs et à l'efficacité de l'analyse.

Afin d'éviter de submerger l'utilisateur d'informations, chaque *dashboard* a été volontairement limité à un nombre réduit de visualisations. Néanmoins, des filtres interactifs permettent d'explorer les données plus en profondeur, en adaptant dynamiquement les graphiques affichés selon la période, la catégorie de produits, ou encore d'autres dimensions pertinentes. Cette structure permet de répondre à deux besoins souvent contradictoires : une lecture rapide des indicateurs clés, et une capacité d'analyse plus poussée pour les utilisateurs avancés.

L'intégration des données dans Tableau a été grandement facilitée par l'utilisation des *mart* préparés dans DBT, offrant une structure prête à l'emploi pour la visualisation (un exemple est détaillé en page 16, d'autres sont disponibles 1). La connexion au *data warehouse* s'est faite facilement (mais pas forcément idéalement) par OAuth. La flexibilité du pipeline a permis des mises à jour régulières des *marts*, adaptant ainsi les données aux besoins évolutifs de la visualisation sans nécessiter de modifications profondes.

Au total, cinq *dashboards* (*Figure 3.1*) ont été conçus pour couvrir quatre fonctions métiers principales : deux pour les ventes, un pour la satisfaction client, un pour la logistique et un dernier dédié au marketing. Chacun répond à des problématiques spécifiques énumérées dans

^{1.} cf. Annexe A.3, page 54

l'énoncé du problème. Des images haute résolution sont proposées ¹, mais véritable valeur résidant dans leurs filtres interactifs et pour en mesurer tout le potentiel, le lecteur est invité à consulter les animations vidéo disponibles via les liens ci-dessous.

- Marketing: https://tinyurl.com/yc4eyf7k
- Optimisation Logistique: https://tinyurl.com/mwk9444d
- Satisfaction Client: https://tinyurl.com/zbcdjkfb
- Ventes Chiffres: https://tinyurl.com/44tb7eck
- Ventes Produits: https://tinyurl.com/2s3dv3t4

L'ensemble des tableaux de bord a fait l'objet d'une étroite collaboration entre les participants, mais nous nous concentrerons ici sur les deux dédiés aux ventes dont l'auteur était plus particulièrement responsable. Le premier se concentre sur les indicateurs financiers globaux, tandis que le second offre une granularité plus importante, proposant notamment une analyse par produits et pays.

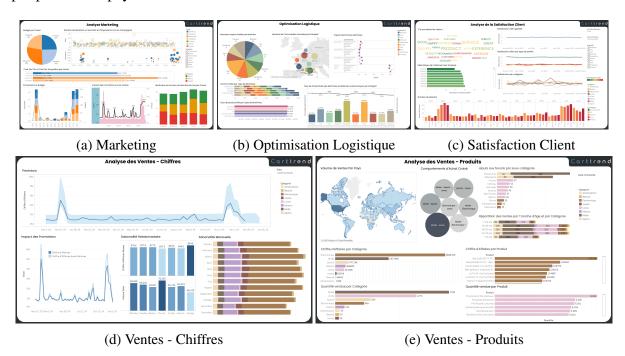


FIGURE 3.1 – Vue d'ensemble des dashboards réalisés

3.1 Tableaux de Bord

3.1.1 Ventes - Chiffres

Ce premier tableau de bord vise à fournir une vue d'ensemble claire de la performance commerciale à travers le suivi du chiffre d'affaires (CA) et de ses évolutions. L'élément central

^{1.} cf. Annexe D, page 70

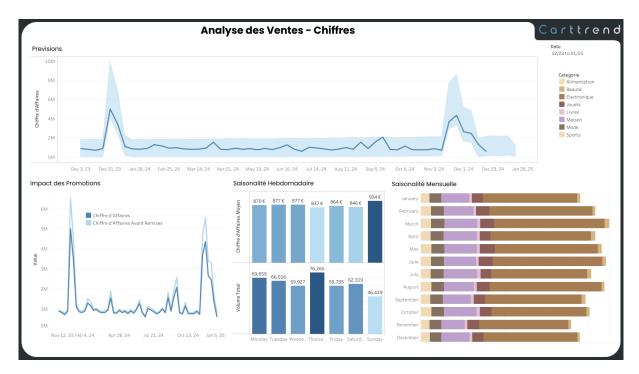


FIGURE 3.2 – Ventes - Chiffres

du *dashboard* est un graphique combiné présentant le CA historique ainsi qu'une projection de son évolution future. La prévision, quant à elle, n'est volontairement pas affichée sous forme chiffrée. Ce choix découle de la qualité jugée insuffisante du modèle prédictif¹, qui aurait pu induire les utilisateurs en erreur quant à la fiabilité des projections. Elle est donc représentée uniquement par un intervalle de confiance au sein duquel la prévision se situe. Cet intervalle est également affiché pour les données historiques, permettant ainsi de se faire une idée visuelle des performances du modèle.

Tous les graphiques présents dans ce *dashboard* sont filtrables selon la période temporelle souhaitée. L'utilisateur peut ainsi facilement isoler une plage de dates spécifique, par exemple pour observer l'impact d'actions commerciales ciblées.

Des filtres supplémentaires viennent enrichir l'analyse, permettant de croiser les données selon différents axes : catégorie de produit, mois de l'année, ou jour de la semaine. Ces filtres ne s'appliquent toutefois pas au graphique de prévision, qui ne bénéficie pas des mêmes possibilités d'interaction. Cette limitation est liée à l'architecture technique du pipeline, qui ne permet pas l'exécution en temps réel du modèle de prévision sur des données filtrées dynamiquement par l'utilisateur. Une solution plus flexible – consistant à intégrer le modèle dans un champ calculé exécuté à la volée – a bien été présentée par un développeur de Tableau lors d'une conférence technique², mais elle est inapplicable dans notre environnement Tableau Cloud. Cela limite donc l'interactivité de ce graphique par rapport au reste du *dashboard*.

^{1.} cf. §2.4.1, page 24

 $^{2. \ \} Conférence \ donnée \ par \ Nathan \ Mannheimer \ le \ 26 \ oct \ 2018. \ Enregistrement \ disponibe : \ \verb|https://youtu.be/nRtOMTnBz_Y|$

3.1.2 Ventes - Chiffres

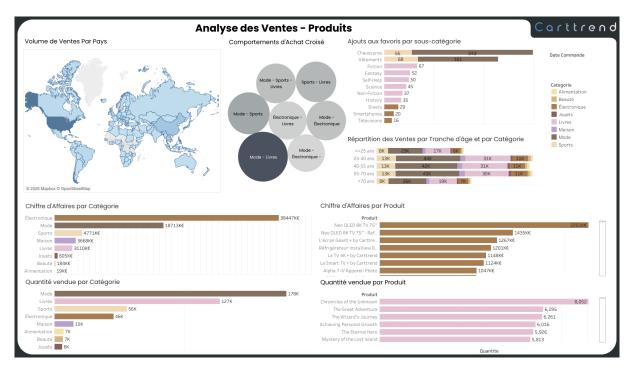


FIGURE 3.3 – Ventes - Produits

Ce deuxième tableau de bord a été conçu pour offrir une analyse plus détaillée des performances commerciales, permettant aux utilisateurs d'explorer les données sous des nouveaux angles. L'ensemble des visualisations réagit là encore dynamiquement aux filtres appliqués, avec des options étendues par rapport au premier dashboard. Outre les filtres temporels et par catégorie de produits, une carte (en haut a gauche) permet un filtre géographique par pays. Cette fonctionnalité facilite l'identification rapide de tendances régionales, bien que la fiabilité des données de géolocalisation (issues d'adresses IP) soit à considérer avec prudence ¹.

Le graphique dédié aux comportements d'achat croisés a été délibérément laissé assez vague, en raison des performances médiocres du modèle². La encore, plutôt que d'afficher des résultats peu fiables, la visualisation se concentre sur des représentations qualitatives, évitant ainsi de fournir des *insights* potentiellement trompeurs. La moitié inférieure du *dashboard* propose quant à elle quatre graphiques qui peuvent chacun être réorganisés dynamiquement du plus bas au plus haut. Cette approche permet d'identifier instantanément les meilleures et moins bonnes performances selon différents critères (CA par produit ou catégorie de produits, volume des ventes, etc.), offrant ainsi une vue synthétique et intuitive des "top" et "flop". Le graphique en haut à droite est enfin notable par l'utilisation de barres colorées par sous-catégories de produits (par exemple, la catégorie chaussures est séparée en "sport" et "mode"), permettent de visualiser la répartition des favoris au sein des grandes catégories.

^{1.} cf. §2.4.3, page 30

^{2.} cf. §2.4.2, page 27

3.2 Analyse et Recommandations

Les analyses menées dans ce projet se heurtent à une limite fondamentale : la nature aléatoire des données synthétiques utilisées. Cette caractéristique a rendu impossible l'obtention *d'insights* réellement pertinents. Les visualisations ont donc été conçues avec une approche prudente, présentant les données de manière impartiale, et évitant au maximum de forcer des conclusions non soutenues par les informations disponibles. Par exemple, certains graphiques sont délibérément laissés vagues et sans chiffres précis, et les axes ont rarement été tronqués. Cette retenue dans l'interprétation constitue un choix méthodologique conscient.

Il serait en outre très inapproprié de formuler des recommandations opérationnelles concrètes dans ce contexte. Avec un accès limité à des données peu fiables et sans possibilité de consulter l'expertise métier des équipes concernées, toute suggestion concernant par exemple la maintenance des machines ou les stratégies commerciales serait présomptueuse. Cependant, et dans le but exclusivement pédagogique de répondre à l'exercice, certaines observations été réalisées et des recommandations ont été formulées en conséquence. Elles illustrent comment les *dash-boards* pourraient orienter des décisions stratégiques.

3.3 Constats et Recommandations

Premièrement, les données semblent indiquer une possible corrélation entre les performances commerciales et les périodes promotionnelles, avec des volumes plus importants en fin d'année. Pour atténuer cette dépendance, des programmes de fidélité pourraient être envisagés, combinés à des événements commerciaux ciblés pendant les périodes creuses et une adaptation des promotions selon les jours de la semaine.

De plus, on observe des disparités entre les catégories de produits. L'électronique semble présenter un ratio CA/volume différent des autres catégories. Cette catégorie pourrait donc présenter une opportunité de *cross-selling*, tandis qu'un rééquilibrage des investissements marketing vers les catégories moins performantes pourrait d'être envisagé.

Certaines combinaisons de catégories de produits fréquemment achetés ensemble semblent se dégager. Une mise en avant stratégique de ces combinaisons, synchronisée avec les moments forts de l'année, comme "Mode + Livres" à la rentrée scolaire, ou "Mode + Sports" en janvier et avant l'été pour les nouvelles résolutions, ainsi que des recommandations intelligentes au panier pourraient potentiellement stimuler les ventes croisées.

Enfin, la répartition par âge montre une distribution gaussienne du chiffre d'affaires, centrée sur les 40-55 ans. Cibler davantage les moins de 25 ans via des canaux comme TikTok (largement sous-exploité, comme le montre le *dashboard* marketing) ¹ et développer des offres adaptées aux seniors pourrait élargir la base de clients.

^{1.} cf. Annexe D.1, page 71

Conclusion

Tableau s'est révélé être un excellent choix pour ce projet, permettant une collaboration efficace entre les membres de l'équipe tout en offrant des résultats visuels professionnels, malgré les limitations imposées par la licence d'essai.

Il est important de souligner que ces *dashboards* ne constituent pas une solution figée. Au contraire, leur force réside dans leur capacité à évoluer en fonction des besoins changeants des utilisateurs. La flexibilité de l'écosystème Tableau, couplée à la modularité du pipeline ELT (notamment grâce à DBT), permet des adaptations rapides à tous les niveaux : ajout de nouvelles visualisations, modification des indicateurs clés, ou intégration de sources de données supplémentaires. Cette agilité garantit la pérennité de ces outils face aux évolutions métiers et techniques.

Conclusion

Ce projet a été l'occasion de construire un pipeline de données complet, depuis la collecte jusqu'à l'analyse, en passant par la visualisation. Nous avons opté pour une architecture ELT moderne avec Fivetran pour l'extraction, DBT pour la transformation et BigQuery comme *data warehouse* – des choix qui se sont révélés pertinents pour leur flexibilité et leur facilité d'utilisation. L'automatisation du *workflow* avec Airflow a permis de gérer efficacement l'ensemble du processus, y compris l'exécution de nos modèles de Machine Learning (prévisions avec Prophet, analyse des associations avec Apriori et régression avec XGBoost) et la géolocalisation des clients. Grâce à Tableau, nous avons pu créer des *dashboards* clairs et interactifs qui répondent aux besoins des différents métiers. Le travail sur l'harmonisation visuelle et l'ergonomie a porté ses fruits, présentant les données synthétiquement, esthétiquement et avec intégrité.

En planifiant soigneusement les tâches et en maintenant des "stand-ups" quotidiens, nous avons pu avancer de manière constante sans avoir recours au *crunch*. La répartition des rôles selon les compétences de chacun (certains plus orientées data engineering, d'autres sur le Machine Learning ou la *dataviz*) a permis de tirer le meilleur de notre équipe.

Sur le plan pédagogique enfin, ce projet a été extrêmement formateur. Il nous a donné l'opportunité de mettre en pratique des concepts théoriques sur un cas concret, avec les contraintes et les défis d'un vrai projet data. La gestion des dépendances entre les différentes briques techniques, la nécessité de documenter proprement le code et les modèles, l'importance de tests rigoureux – autant d'aspects qui prennent tout leur sens quand on les expérimente en situation réelle.

Annexe A

Code DBT

A.1 Configuration

Listing A.1 – Definition des sources de donnees brutes

```
1 # [Fichier : sources.yaml]
 2 # Role : Definition des sources de donnees brutes a partir de Google Drive.
 3 # Ce fichier liste les tables sources qui seront utilisees comme input pour les
      transformations DBT.
 4
 5 version: 2
 6
7 sources:
8
    - name: google_drive # Dataset BigQuery dans lequel se trouve les tables brutes
9
      tables:
10
        - name: carttrend campaigns campagnes
11
        - name: carttrend_clients_clients
12
        - name: carttrend_commandes_commandes
13
        - name: carttrend details commandes details commandes
        - name: carttrend_entrepots_entrepots
14
        - name: carttrend_entrepots_machines_machines_entrepots
15
16
        - name: carttrend_posts_posts
        - name: carttrend_produits_produits
17
        - name: carttrend_promotions_promotions
18
19
       - name: carttrend_satisfaction_satisfaction
```

Listing A.2 – Documentation des modeles DBT "staging"

```
1 # [Fichier : staging/schema.yaml]
 2 # Role : Documentation des modeles DBT apres etape de staging (nettoyage/validation).
 3 # Decrit la structure des tables transformees, leurs relations et les tests de
      qualite associes.
 4 # Utilise des packages DBT (dbt_utils, dbt_expectations) pour les tests avances.
 5
 6 version: 2
7
 8 models:
 9
    # Campagne_Marketing -----
10
11
    - name: Campagnes
12
      description: Table des Campagnes Marketing Normalisée et Nettoyée.
13
      tests:
        - dbt_utils.expression_is_true: # Format clé primaire CAMP000
14
15
             expression: 'REGEXP_CONTAINS(id_campagne, r"^CAMP\d{3}$")'
             config:
16
17
               severity: warn # Warning plutot qu'erreur
```

```
- dbt_utils.expression_is_true: # Controle de coherence des metriques
18
            expression: 'budget > 0 and impressions > 0 and clics > 0 and conversions >
19
20
            config:
21
              severity: warn
2.2
      columns:
23
        - name: id_campagne
24
          description: Clé Primaire. Pas de test 'not_null', les campagnes n'ayant pas
              d'id sont automatiquement supprimés
25
          tests:
26
            - unique
27
        - name: date
2.8
          description: Jour (DATE) de la campagne
29
        - name: id canal
30
          description: id du Canal Marketing
31
        - name: budget
32
          description: Cout de la campagne marketing en E (float). Chiffre superieur a
33
        - name: impressions
34
          description: Cout de la campagne marketing (int). Chiffre superieur a 0
35
        - name: clics
36
          description: Nombre de Clics générés par la campagne marketing (int). Chiffre
               superieur a 0
37
        - name: conversion
          description: Nombre de Conversions générés par la campagne marketing (int).
38
              Chiffre superieur a 0
39
        - name: ctr
40
          description: CTR - Nombre de clics / nombre d'impression. Float entre 0 et 1
             ou NULL
41
42
    # Canaux ------
43
     # Pas de test car table extraite des donnée d'autres tables. Tests réalisés sur ces
44
         autres tables
45
    - name: Canaux
      description: Table de dimension contenant les canaux publicitaires (résaux
46
          sociaux)
47
      columns:
48
         - name: id_canal
49
          description: Clé Primaire. Tests non necessaire car clé crée par rownumber =
              garantie d'etre unique et non nulle
50
        - name: canal
51
          description: Nom du Canal Publicitaire (Facebook, Twitter...)
52
    # Clients ------
53
54
55
    - name: Clients
      description: Table des clients nettoyée. Les produits des favoris des clients ne
56
          sont pas normalisés
57
      tests:
        - dbt_utils.expression_is_true: # Format clé primaire C00000
58
59
            expression: 'REGEXP_CONTAINS(id_client, r"^C\d{5}$")'
60
        - dbt_utils.expression_is_true:
61
            expression: 'age < 110'
62
            config:
63
              severity: warn
64
      columns:
65
        - name: id_client
          description: Clé Primaire. Pas de test 'not_null', les clients n'ayant pas d'
66
              id valide sont automatiquement supprimés
67
          test:
            - unique
68
69
        - name: hash_prenom_nom
70
          description: Hash sha264 (pour anonymisation) du prenom et du nom du client
```

```
71
         - name: email
 72
           description: Email du client valide ou NULL
 73
         - name: age
 74
           description: Age du client (int). Chiffre superieur a 0 ou NULL
 75
         - name: tranche_age
          description: Tranche d'age. String '<=25', '25-40', '40-55', '55-70' ou '>70'
 76
 77
         - name: genre
 78
          description: Genre du client. 'Homme', 'Femme' ou NULL
 79
         - name: frequence_visites
 80
           description: Nombre de visites (int) du client sur le site. Chiffre superieur
               a 0 ou NULL
 81
         - name: telephone
 82
          description: Numero de téléphone au format '+XXX-XXX-XXXX' ou NULL
 83
         - name: date inscription
           description: Date d'inscription (date) du client ou NULL. Avant la date d'
              aujourd'hui
 85
         - name: anciennete_jours
 86
           description: Nombre de jours depuis l'inscription du client
         - name: ip
 87
 88
           description: Adresse IPv4 du client ou NULL. Au format 'XXX.XXX.XXXX.XXXX'
 89
 90
     91
 92
     - name: Details_commandes
 93
       description: Information sur la commande - produits, quantitées...
 94
 95
         - dbt_utils.expression_is_true:
 96
            expression: 'quantite > 0'
 97
       columns:
 98
         - name: id_commande
 99
          tests:
100
            - relationships:
                to: ref('Commandes')
101
                field: id_commande
102
103
         - name: id_produit
104
          tests:
105
            - relationships:
106
                to: ref('Produits')
107
                field: id_produit
         - name : quantite
108
109
           description : Nombre du même article par commande. Supérieur strictement à 0
110
         - name : emballage_special
111
           description : La commande est emballée avec un emballage spécial. Type boolé
112
113
     -----
114
     - name: Entrepots
115
116
       description: Table des Entrepots Nettoyée
117
       tests:
118
         - dbt_utils.expression_is_true:
119
            expression: 'capacite_max > volume_stocke'
120
            config:
121
              severity: warn
         - dbt_utils.expression_is_true: # Format clé primaire E000
122
             expression: 'REGEXP_CONTAINS(id_entrepot, r"^E\d{3}$")'
123
124
       columns:
125
         - name: id_entrepot
           description: Clé Primaire, au format 'EXXX'
126
127
          tests:
128
            - not_null
            - unique
129
130
         - name: localisation
131
           description: Ville dans laquelle se trouve l'entrepot
```

```
132
         - name: capacite_max
           description: Nombre de produits (int) maximum que peut contenir l'entrepot.
133
               Garantie > 0
134
         - name: volume stocke
135
           description: Nombre de produits (int) actuellement stocké dans l'entrepot.
136
         - name: taux_remplissage
137
           description: Pourcentage de stockage utilisé. float > 0
138
         - name: temperature
139
           description: Temperature moyenne mesurée dans l'entrepot
140
141
     # Entrepots_Machines ------
142
143
     - name: Entrepots_Machines
144
       description: Table qui décrit l'etat de fonctionnement des machines des entrepots
145
146
         - dbt_utils.expression_is_true: # Format clé primaire ME000000
             expression: 'REGEXP_CONTAINS(id_machine, r"^ME\d{6}$")'
147
148
149
       - name: id machine
150
         description: identifiant de la machine
151
         tests:
152
           - not_null
153
       - name: id_entrepot
154
         description: id de l'entrepot contenant la machine. Cle etrangere
155
         tests:
156
           - relationships:
157
               to: ref('Entrepots')
158
               field: id_entrepot
159
       - name: type_machine
160
         description: Fonction de la machine. Garanti 'Tri', 'Maintenance', 'Transport', '
             Stockage','Conditionnement' ou 'Autre'
       - name: etat_machine
161
         description: Etat de Fonctionnement de la machine. Garanti 'Service', 'Panne', '
162
             Maintenance' ou 'Autre'
163
       - name: temps_d_arret
164
         description: Duree de non fonctionnement de la machine
165
       - name: volume_traite
         description: le nombre de produit traités par ma machine dans le mois
166
167
       - name: date
         description: date du premier jour du mois évalué
168
169
170
     171
172
     - name : Posts
173
       columns :
174
       - name : id_post
175
         tests :
176
           - not_null
177
           - unique
178
       - name : date_post
179
         description : date de publication du post
180
181
           - dbt_expectations.expect_column_values_to_be_of_type:
182
               column_type: date
       - name : id_canal
183
         description : id du canal utilisé. Clé étrangère, mais pas de test car garantie
184
              d'exister dans la table Canaux
185
       - name : volume_mentions
186
         description : Nombre d'interactions par post (int), garanti > 0
187
       - name : sentiment global
         description : Sentiment lié au post. Peut être 'Positif', 'Neutre' ou 'Negatif'
188
              ou 'NULL'
189
       - name : contenu_post
190
         description : Contenu du post publié
```

```
191
192
193
194
     - name: Produits
195
       description: Table des Produits nettoyée
196
       tests:
197
         - dbt_utils.expression_is_true: # Format clé primaire P00000
198
             expression: 'REGEXP_CONTAINS(id_produit, r"^P\d{5}$")'
199
         - dbt_utils.expression_is_true:
200
             expression: 'prix > 0'
201
             config:
202
               severity: warn
203
       columns:
204
         - name: id_produit
205
           description: Clé Primaire.
206
           tests:
207
             - not_null
208
             - unique
209
         - name: categorie
           description: Categorie du produit. N'importe lequel de 'Livres','Alimentation
210
               ','Mode','Maison','Électronique','Jouets','Beauté','Sports' ou 'Autre'
211
           description: Marque du produit ou NULL. Premiere lettre en majuscule
212
213
         - name: nom
           description: Nom du produit ou NULL. Premiere lettre en majuscule
214
215
         - name: prix
216
           description: Prix du produit. (float) > 0
217
         - name: sous_categorie
218
           description : Nom de la sous categorie du produit ou NULL. Premiere lettre en
               majuscule
219
         - name: variation
220
           description: Nom de la variation (taille, gamme...) ou NULL. Premiere lettre
               en majuscule
221
222
     # Produits Favoris -----
     #-----
223
224
     - name: Produits_Favoris
225
       description: Table contenant les identifiants des produits que les clients ont
           ajoutés dans leur favoris
       columns:
226
227
          - name: id_client
228
           description: Clé etrangere client
229
           tests:
230
             - relationships:
231
                 to: ref('Clients')
232
                 field: id_client
         - name: id_produit
233
234
           description: Clé etrangere produit
235
           tests:
236
             - relationships:
                 to: ref('Produits')
2.37
238
                 field: id_produit
239
240
     # Promotions -----
241
     - name: Promotions
242
243
       tests:
244
         - dbt_utils.expression_is_true: # Format clé primaire P000
245
             expression: 'REGEXP_CONTAINS(id_promotion, r"^P\d{3}$")'
246
         - dbt_utils.expression_is_true:
247
             expression: 'pourcentage_promotion > 0'
248
             config:
249
               severity: warn
250
       columns:
```

```
251
         - name: id_promotion
252
           tests:
253
             - not_null
254
             - unique
255
         - name: id_produit
256
           tests:
257
             - relationships:
258
                 to: ref('Produits')
259
                 field: id_produit
260
         - name: valeur_promotion
261
           description: en euro
262
         - name: pourcentage_promotion
           description: promotion en pourcentage (float entre 0 et 1)
263
264
         - name: date debut
           description: Date de début de la promotion 'yyyy-mm-dd'
265
266
           tests:
              - not_null
267
268
          - name: date fin
269
            description: Date de fin de la promotion 'yyyy-mm-dd'
270
271
      # Satisfaction -----
272
273
      - name: Satisfaction
274
       columns:
275
          - name: id_satisfaction
276
           tests:
277
             - unique
278
             - not_null
279
         - name: id_commande # Hugues :
280
           description: id de la commande a laquelle se rapporte cet avis
281
           tests:
282
             - relationships:
                  to: ref('Commandes')
283
284
                 field: id_commande
285
         - name : note_client
286
           description : Note laissee par le client de 1 à 5
287
         - name : commentaire
           description : commentaire laissé par le client
288
289
         - name : plainte
290
           description : le commentaire fait l'objet d'une plainte de la part du client.
                Type booléen
291
         - name : delai_traitement
           description : Duree de reponse effectuee par l'equipe support
292
293
         - name : type_plainte
294
           description : Type de plainte effectuée par le client
295
          - name : employe_support
296
           description : identifiant du membre du service support
```

A.2 Exemples de Moèles "staging"

Listing A.3 – Nettoyage et Transformation des donnees brutes des entrepots

```
1 -- [Fichier : Entrepots.sql]
2 -- Role : Nettoyage et Transformation des donnees brutes des entrepots.
3 -- Effectue le nettoyage de base et le calcul des indicateurs durant l'etape staging.
4
5 SELECT
6    id_entrepot, -- Conserve l'identifiant original
7
8    -- Normalisation de la localisation (mise en majuscule de la premiere lettre)
9    INITCAP(localisation) AS localisation, -- ex. 'paris' => 'Paris'
10
11    -- Capacite toujours >= 0
```

```
12
      CASE
13
           WHEN capacite_max < 0 THEN 0
14
           ELSE capacite_max
15
       END AS capacite_max,
16
17
       -- Stock toujours >= 0
18
19
           WHEN volume_stocke < 0 THEN 0
20
           ELSE volume_stocke
21
      END AS volume_stocke,
22
23
       -- Calcul du taux de remplissage (derive)
      volume_stocke / capacite_max AS taux_remplissage,
2.4
25
       -- Simplification du nom de la colone temperature
26
27
       temperature_moyenne_entrepot AS temperature
2.8
29 -- Table contenant les donnees brutes
30 FROM {{source("google_drive", "carttrend_entrepots_entrepots")}}
```

Listing A.4 – Normalisation des canaux publicitaires depuis Posts et Campagnes

```
1 -- [Fichier : Canaux.sql]
 2 -- Normalisation des canaux publicitaires depuis Posts et Campagnes
 3 -- Crée une table de référence unique pour tous les canaux marketing
 4
 5 WITH canal AS (
 6
      -- Extraction des canaux depuis les posts réseaux sociaux
 7
      SELECT canal social
8
      FROM {{source("google_drive",'carttrend_posts_posts')}}
9
10
     UNION DISTINCT -- Combine les deux en éliminant les doublons
11
12
       -- Extraction des canaux depuis les campagnes marketing
13
      SELECT canal
14
      FROM {{source('google_drive', 'carttrend_campaigns_campagnes')}}
15)
16 SELECT
17
      -- Génération d'un ID unique pour chaque canal,
18
       -- un chiffre séquentiel auto-incrémenté
19
      ROW_NUMBER() OVER() as id_canal,
20
       -- Conservation du nom original
21
22
      canal_social as canal
23
24 FROM canal
25 -- Tri alphabétique pour une meilleure lisibilité
26 ORDER BY canal ASC
```

Listing A.5 – Nettoyage des données de posts réseaux sociaux avec intégration des canaux normalisés

```
1 -- [Fichier : Posts.sql]
 2 -- Nettoyage des données de posts réseaux sociaux avec intégration des canaux
      normalisés
 3 -- Montre l'utilisation de la table de dimension Canaux créée précédemment
 4
 5 SELECT
 6
      id_post, -- Identifiant original du post
7
      id_canal, -- Clé étrangère vers la table Canaux normalisée
8
      -- Garantit volume mentions >= 0
9
10
      CASE
11
          WHEN volume_mentions < 0 THEN 0</pre>
12
          ELSE volume_mentions
13
      END AS volume_mentions,
```

```
14
       -- Validation des valeurs de sentiment
15
16
      CASE
           WHEN sentiment_global IN ('Positif', 'Négatif', 'Neutre') THEN
17
              sentiment_global
18
           ELSE NULL
19
      END AS sentiment_global,
20
21
       -- Garantit date_post <= date actuelle
22
      CASE
23
           WHEN date_post <= CURRENT_DATE() THEN date_post</pre>
24
           ELSE CURRENT_DATE()
2.5
      END AS date_post,
26
27
       -- Texte original du post conservé sans transformation
28
      contenu_post
29
30 FROM {{source("google_drive",'carttrend_posts_posts')}} p -- Table source des posts
       -- Jointure avec la table Canaux normalisée
31
32
      JOIN {{ref('Canaux')}} c ON c.canal = p.canal_social
```

Listing A.6 – Nettoyage et normalisation des données clients

```
1 -- [Fichier : Clients.sql]
2 -- Nettoyage et normalisation des données clients avec :
3 -- 1. Anonymisation des données sensibles
 4 -- 2. Standardisation des formats
 5 -- 3. Validation des données
 6 -- 4. Enrichissement léger
8 WITH tel AS (
9
      -- CTE pour le traitement des numéros de téléphone
      SELECT
10
11
          id_client,
          -- Formatage standard du numéro (123-456-7890)
12
13
          REGEXP_REPLACE (numero_telephone, r'^*.*?(\d{3})\D?(\d{4}).*$', r'
              \1-\2-\3') AS numero_telephone,
           -- Left-Padding du code pays (ex: +033 pour la France)
14
15
          CASE
16
              WHEN LENGTH (ac.area_code) = 3 THEN CONCAT ('+', ac.area_code)
               WHEN LENGTH (ac.area_code) = 2 THEN CONCAT ('+0', ac.area_code)
17
               WHEN LENGTH(ac.area_code) = 1 THEN CONCAT('+00', ac.area_code)
18
              ELSE '+033'
                            -- Valeur par défaut
19
20
          END as area_code
21
      FROM {{source("google_drive",'carttrend_clients_clients')}}
2.2
      JOIN (
23
           -- Sous-Requete pour l'extraction du code pays depuis le numéro original
          SELECT
24
25
               id_client,
2.6
              REGEXP_EXTRACT(numero_telephone, r'^\+?([1-9]{0,3}).*$') AS area_code
27
          FROM {{source("google_drive",'carttrend_clients_clients')}}
28
      ) ac USING(id_client)
29)
30 SELECT
31
      -- Identifiant client conservé
32
      id_client,
33
34
       -- Anonymisation des données personnelles avec hashing SHA256
35
      sha256(CONCAT(prenom_client, nom_client)) as hash_prenom_nom,
36
37
      -- Validation de l'email
38
      CASE
39
          WHEN REGEXP_CONTAINS (email, r'^[\w\.-]+[\w-]+[\A-Za-z]+$') THEN email
40
          ELSE NULL
41
      END as email,
```

```
42
43
                          -- Suppression des âges négatifs
44
                          CASE WHEN age > 0 THEN age END as age,
45
46
                          -- Calcul de la tranche d'âge
47
                         CASE
48
                                        WHEN age > 0 AND age <= 25 THEN ' <= 25'
49
                                        WHEN age > 25 AND age <= 40 THEN '25-40'
50
                                        WHEN age > 40 AND age <= 55 THEN '40-55'
51
                                        WHEN age > 55 AND age <= 70 THEN '55-70'
52
                                        WHEN age > 70 THEN '>70'
53
                                        ELSE NULL
                         END as tranche_age,
54
55
                          -- Normalisation du genre
56
57
                         CASE
58
                                        WHEN lower (genre) = 'homme' THEN 'Homme'
                                        WHEN lower (genre) = 'femme' THEN 'Femme'
59
60
                                        ELSE NULL
61
                         END as genre,
62
63
                          -- Validation et correction de la fréquence de visites
64
                         CASE
65
                                        WHEN frequence_visites > 0 THEN frequence_visites
                                        ELSE 0
66
67
                         END as frequence_visites,
68
69
                          -- Construction du numéro de téléphone complet
                          -- Indicatif + Numero, format final +000-123-456-7890
70
71
                         CONCAT(tel.area_code, '-', tel.numero_telephone) as telephone,
72
                          -- Validation de la date d'inscription
73
74
                         CASE
75
                                         WHEN DATE_DIFF(date_inscription, CURRENT_DATE(), DAY) <= 0 THEN
                                                      date inscription
76
                                        ELSE NULL
77
                         END as date_inscription,
78
79
                          -- Calcul de l'ancienneté en jours
80
                         CASE
                                        WHEN DATE_DIFF(date_inscription, CURRENT_DATE(), DAY) <= 0</pre>
81
82
                                        THEN DATE_DIFF(CURRENT_DATE(), date_inscription, DAY)
83
                                        ELSE NULL
                         END as anciennete_jours,
84
85
86
                          -- Validation des adresses IP (format IPv4)
87
                         CASE
                                         WHEN REGEXP_CONTAINS (adresse_ip, r' \wedge d\{1,3\} \setminus d\{1,3\} 
88
                                                     THEN adresse_ip
89
                                        ELSE NULL
90
                         END as ip
91
92 FROM {{source("google_drive",'carttrend_clients_clients')}}
                         JOIN tel USING(id_client)
94 WHERE id_client IS NOT NULL -- Exclusion des clients sans identifiant
```

A.3 Exemples de Moèles "marts"

Listing A.7 – Analyse textuelle des commentaires clients pour *dataviz*

```
1 -- [Fichier : mart_Commentaires.sql]
2 -- Analyse textuelle des commentaires clients pour dataviz :
3 -- 1. Extraction et nettoyage des mots individuels
```

```
4 -- 2. Filtrage des stopwords
5 -- 3. Calcul des indicateurs clés (fréquence, sentiment moyen)
6
7 WITH
8
    -- Découpage des commentaires en mots individuels
9
    mots AS (
10
      SELECT
        note_client, -- Note associée (1-5)
11
12
13
      FROM {{ref("Satisfaction")}}, -- Table staging nettoyée
14
      UNNEST (SPLIT (commentaire, '')) AS mot -- Désagrégation en mots
15
    ),
16
17
    -- Agrégation et analyse des mots
18
    mots_aggreges AS (
19
      SELECT
        UPPER (REGEXP_REPLACE (mot, r'[^\w]', '')) AS mot, -- Nettoyage (majuscules,
2.0
            ponctuation)
21
        COUNT(*) AS frequence, -- Nombre d'occurrences
2.2
        ROUND (AVG (note_client), 1) AS note_moyenne -- Sentiment moyen associé
2.3
      FROM mots
24
      WHERE mot != '' -- Exclusion des chaînes vides
25
        -- Filtre des stopwords (mots courants non significatifs)
        AND lower (mot) NOT IN (
2.6
          'a', 'about', 'above', ..., 'yourselves' -- Liste écourtée de stopwords
27
              anglais
28
        )
29
      GROUP BY mot
30
    )
31
32 -- Résultat final pour dataviz
33 SELECT *
34 FROM mots_aggreges
35 -- Tri par fréquence (principalement pour debugage)
36 ORDER BY frequence DESC
```

Listing A.8 – Mart simplifiant l'analyse et la visualisation des notes des clients dans leur contexte

```
1 -- [Fichier : mart_Notes_Clients.sql]
 2 -- Mart simplifiant l'analyse et la visualisation des notes des clients dans leur
     contexte
 3 -- à l'aide des jointures multiples pour croiser les dimensions pertinentes
 4
 5 SELECT
 6
     p.id_produit, -- Identifiant produit
 7
                     -- Nom du produit
     p.nom,
     p.categorie,
                    -- Catégorie produit
 8
 9
     c.date commande, -- Date de commande
     c.date_livraison_estimee, -- Date de livraison estimée
10
      s.type_plainte, -- Type de plainte associée (le cas échéant)
11
                     -- Identifiant client
12
      c.id_client,
13
      AVG(s.note_client) AS avg_note_client, -- Note moyenne calculée
14
      e.localisation AS localisation_entrepot -- Localisation de l'entrepôt
15
16 -- Jointure des différentes tables staging
17 FROM {{ ref("Produits") }} p
18 JOIN {{ ref("Details_commandes")}} d USING(id_produit)
19 JOIN {{ ref("Commandes") }} c USING(id_commande)
20 JOIN {{ ref("Satisfaction") }} s USING(id_commande)
21 JOIN {{ ref("Entrepots")}} e USING(id_entrepot)
2.2
23 -- Groupement par dimensions d'analyse
24 GROUP BY
      p.id_produit,
25
26
    p.nom,
```

```
p.categorie,
c.date_commande,
c.date_livraison_estimee,
s.type_plainte,
c.id_client,
e.localisation

-- Tri par note moyenne (plus basse en premier pour identifier les problèmes)

ORDER BY avg_note_client ASC
```

Listing A.9 – Table complète des commandes clients

```
1 -- [Fichier : mart_Commandes.sql]
 2 -- Table complète des commandes clients avec :
 3 -- 1. Calcul des prix (initial et après remise)
 4 -- 2. Contexte client et produit
 5 -- 3. Promotions éventuelles
 7 SELECT
 8
      co.date_commande, -- Date de la commande
 9
      p.nom, -- Nom du produit
                                -- Prix unitaire de base
10
      p.prix as prix_unitaire,
11
      c.quantite, -- Quantité commandée
12
13
       -- Calcul du prix total avant remise (quantité * prix unitaire)
14
      ROUND(p.prix * dc.quantite, 2) AS prix_initial,
1.5
16
       -- Calcul du prix après application de la remise éventuelle
17
      ROUND (CASE
           WHEN pourcentage_promotion IS NOT NULL THEN
18
19
               p.prix * dc.quantite * (1 - pourcentage_promotion)
20
           ELSE p.prix * dc.quantite
2.1
      END, 2) as prix_apres_remise,
22
23
      p.categorie, -- Catégorie produit
24
      p.sous_categorie, -- Sous-catégorie
25
2.6
       -- Flag indiquant si le produit est dans les favoris du client
27
      CASE
28
           WHEN EXISTS (
29
               SELECT *
30
               FROM {{ref('Produits_Favoris')}} pf
               WHERE pf.id_client = cl.id_client AND pf.id_produit = p.id_produit)
31
32
           THEN true
33
           ELSE false
34
      END as produit_favori,
35
       -- Informations promotionnelles (avec valeurs par défaut)
36
37
       IFNULL(pr.valeur_promotion, 0) as valeur_promotion,
      IFNULL(pr.pourcentage_promotion, 0) as pourcentage_promotion,
38
39
40
       -- Informations du client
41
      cl.id_client,
42
      cl.age,
43
      cl.tranche_age,
44
      cl.genre,
45
      cl.frequence_visites,
46
      cl.date_inscription,
47
      cl.anciennete_jours,
48
49
      co.statut_commande, -- Statut de la commande
50
      p.id_produit -- ID produit pour les jointures
51
52 -- Tables sources avec jointures LEFT pour conserver toutes les commandes
53 FROM {{ref('Commandes')}} co
```

```
54 LEFT JOIN {{ref('Details_commandes')}} dc USING (id_commande)
55 LEFT JOIN {{ref('Produits')}} p USING (id_produit)
56 LEFT JOIN {{ref('Promotions')}} pr USING (id_promotion)
57 LEFT JOIN {{ref('Clients')}} cl USING (id_client)
```

Listing A.10 – Preparation des series temporelles de CA pour Prophet

```
1 -- [Fichier : mart_Prevision_Series_Temporelles_CA.sql]
 2 -- Preparation des series temporelles de CA pour Prophet :
 3 -- 1. Generation d'un calendrier continu sans trous
 4 -- 2. Aggregation du CA quotidien
 5 -- 3. Formatage specifique pour l'analyse forecasting
 6
 7 WITH
8
      -- Generation du calendrier complet entre premiere et derniere commande
9
      cal AS (
10
          SELECT date
11
          FROM UNNEST (GENERATE_DATE_ARRAY(
               (SELECT MIN(date_commande) FROM {{ref('mart_Commandes')}}), -- Date
12
               (SELECT MAX(date_commande) FROM {{ref('mart_Commandes')}}) -- Date
13
                  maximale
14
           ) AS date
15
16
17
       -- Calcul du CA quotidien (hors commandes annulees)
      ca AS (
18
          SELECT
19
2.0
               date commande AS date,
21
              SUM(prix_apres_remise) AS ca -- Somme des montants apres remise
          FROM {{ref('mart_Commandes')}}
2.2
23
          WHERE statut_commande != 'Annulée' -- Filtre des commandes annulees
2.4
           GROUP BY date_commande
25
      )
26
27 -- Jointure finale pour serie temporelle complete
28 SELECT
29
      cal.date AS ds, -- Colonne 'ds' requise par Prophet
30
      ROUND (IFNULL (ca.ca, 0), 2) as y -- Colonne 'y' requise par Prophet
31 FROM cal
32
      LEFT JOIN ca USING(date) -- Jointure sur les dates, left pour garder toutes les
          dates
33 ORDER BY date -- Tri chronologique obligatoire
```

Listing A.11 – Preparation des données pour modele XGBoost de prediction des pannes

```
1 -- [Fichier : mart_Regression_Fonctionnement_Machines.sql]
 2 -- Preparation des donnees pour modele XGBoost de prediction des pannes
 3
 4 SELECT
 5
    -- Variables Explicatives
    EXTRACT (MONTH FROM m.date) AS month, -- Saisonnalite (1-12)
 6
    e.localisation, -- Localisation de l'entrepot (categorielle)
 7
    e.temperature, -- Temperature moyenne
 8
 9
    m.volume_traite, -- Charge de travail (continue)
10
    -- L'analyse sera faite par type de machines
11
12
    m.type_machine, -- Type d'equipement (categorielle)
13
     -- Variable cible : Uptime
14
15
     -- [0-1], 1 = 100% uptime
16
    1 - m.temps_d_arret / (DATE_DIFF(
17
      DATE_ADD (m.date, INTERVAL 1 MONTH),
      m.date,
18
19
   DAY
```

```
20  ) * 24) as pourcentage_uptime,
21
22 FROM {{ref("Entrepots_Machines")}} m
23  JOIN {{ref('Entrepots')}} e USING(id_entrepot)
```

Annexe B

Scripts Machine Learning et Geolocalisation

Listing B.1 – Implementation du modele Prophet pour la prevision de CA

```
1 # [Fichier : ml_prediction_ca.py]
 2 # Implementation du modele Prophet pour la prevision de CA :
 3 # 1. Extraction des donnees depuis BigQuery
 4 # 2. Configuration du modele avec gestion des périodes particulière
 5 # 3. Entrainement et prediction
 6 # 4. Stockage des resultats dans BigQuery
8 from google.oauth2 import service_account
9 import pandas_gbq
10 from prophet import Prophet
11 import pandas as pd
12 import os
13
14 # Authentification GCP via service account
15 credentials = service_account.Credentials.from_service_account_file(
      os.path.dirname(__file__) + '/Compte Service Machine Learning.json',
16
17)
18
19 # Extraction des donnees historiques
20 sql = '''
21
      SELECT *
22
      FROM 'carttrend_clean.mart_Prevision_Series_Temporelles_CA'
23 ′′′
24 df = pandas_gbq.read_gbq(sql, project_id="carttrend-4", credentials=credentials)
25 df['date'] = pd.to_datetime(df['date'])
27 # Certaines périodes de ventes sont très particulières, et pas forcement repré
      sentatives
28 # Elles sont déclarées comme vacances et sont ignorées par le modèle
29 # cf. https://facebook.github.io/prophet/docs/handling_shocks.html
30 vacances = pd.DataFrame([
31
           'holiday': '2023',
32
33
           'ds': pd.to_datetime('2023-12-30'),
34
           'lower_window': 0,
35
           'upper_window': 20,
36
      },
37
38
           'holiday': '2024',
           'ds': pd.to_datetime('2024-11-20'),
39
40
           'lower_window': 0,
41
           'upper_window': 25,
42
       }
43 ])
44
```

```
45 # Initialisation du modele Prophet
46 model = Prophet(holidays=vacances)
47
48 # Entrainement du modele (il attend des noms de colonnes specifiques)
49 model.fit(df.rename(columns={'date': 'ds', 'chiffre_affaire': 'y'}))
50
51 # Creation d'un dataframe avec 30 jours supplementaires
52 future = model.make_future_dataframe(periods=30, freq='D')
53
54 # Prediction du chiffre d'affaire sur ces 30 jours
55 # ds : date, yhat : prediction,
56 # yhat_lower et yhat_upper : bornes de l'intervalle de confiance
57 forecast = model.predict(future)[['ds', 'yhat', 'yhat_lower', 'yhat_upper']]
58
59 # Formatage des resultats
60 forecast['ds'] = pd.to_datetime(forecast['ds'])
61 output = forecast.rename(columns={'ds': 'date'})
63 # Combinaison des données de départ et des prédictions
64 # Cela simplifie la visualisation par la suite
65 output = pd.concat([forecast.set_index('date'), output.set_index('date')], axis=1).
      reset_index()
66
67 # Envoi vers BigQuery
68 pandas_gbq.to_gbq(
69
      output,
70
      "carttrend_transform.ml_predictions_ca",
71
      project_id="carttrend-4",
72
      if_exists="replace"
73)
```

Listing B.2 – Implementation de l'algorithme Apriori pour l'analyse des associations produits

```
1 # [Fichier : ml_a_priori.py]
 2 # Implementation de l'algorithme Apriori pour l'analyse des associations produits :
 3 # 1. Extraction des transactions depuis BigQuery
 4 # 2. Preprocessing des donnees pour ML
 5 # 3. Execution de l'algorithme Apriori
 6 # 4. Filtrage et stockage des resultats
7
8 from google.oauth2 import service_account
9 import pandas_gbq
10 import pandas as pd
11 from mlxtend.preprocessing import TransactionEncoder
12 from mlxtend.frequent_patterns import apriori
13 import os
14
15 # Authentification GCP via service account
16 credentials = service_account.Credentials.from_service_account_file(
       os.path.dirname(__file__) + '/Compte Service Machine Learning.json',
17
18)
19
20 # Extraction des combinaisons produits depuis BigQuery
21 sql = '''
22
      SELECT *
2.3
      FROM 'carttrend_clean.mart_Apriori_Combinaisons_Produits'
24 '''
25 df = pandas_gbq.read_gbq(sql, project_id="carttrend-4", credentials=credentials)
27 # Transformation des donnees en format liste de listes (format attendu par mlxtend)
28 transactions = [list(row[0]) for row in df.itertuples(index=False)]
29
30 # Encodage des "paniers" au format one-hot
31 te = TransactionEncoder()
32 te_ary = te.fit(transactions).transform(transactions)
```

```
33 dfml = pd.DataFrame(te_ary, columns=te.columns_)
34
35 # Execution de l'algorithme Apriori avec seuil minimal de support
36 frequent_itemsets = apriori(dfml, min_support=0.1, use_colnames=True)
37
38 # Filtrage des resultats :
39 # - Combinaisons d'au moins 2 categories
40 # - Support minimum de 40%
41 frequent_itemsets['length'] = frequent_itemsets['itemsets'].apply(lambda x: len(x))
42 frequent_itemsets = frequent_itemsets[
       (frequent_itemsets['length'] >= 2) &
44
       (frequent_itemsets['support'] >= 0.4)
45 ][['support', 'itemsets']]
46
47 # Formatage des resultats (frozenset -> string), facilite la visualisation
48 frequent_itemsets['itemsets'] = frequent_itemsets['itemsets'].apply(
       lambda x: ' - '.join(list(x))
49
50)
51
52 # Envoi des resultats vers BigQuery
53 pandas_gbq.to_gbq(
54
      frequent_itemsets,
55
      "carttrend_transform.ml_categories_achetees_ensemble",
56
      project_id="carttrend-4",
57
      if_exists="replace"
58)
```

Listing B.3 – Modele XGBoost pour l'analyse des causes de pannes machines

```
1 # [Fichier : ml_regression_machines.py]
 2 # Modele XGBoost pour l'analyse des causes de pannes machines :
 3 # 1. Preparation des donnees (normalisation, encodage)
 4 # 2. Entrainement d'un modèle par type de machine
 5 # 3. Evaluation des performances pour debugage
 6 # 4. Export des resultats vers BigQuery
 7
 8 from google.oauth2 import service_account
9 import pandas as pd
10 import numpy as np
11 import pandas_gbq
12 from sklearn.preprocessing import StandardScaler, MinMaxScaler, LabelEncoder
13 from sklearn.model_selection import train_test_split
14 from sklearn.metrics import mean_squared_error
15 import xgboost as xgb
16 import os
17
18 # Authentification GCP via service account
19 credentials = service account.Credentials.from service account file(
20
      os.path.dirname(__file__) + '/Compte Service Machine Learning.json',
21 )
22
23 # Extraction des donnees depuis BigQuery
24 sql_machine = '''
25
      SELECT *
26
      FROM 'carttrend-4.carttrend_clean.mart_Regression_Fonctionnement_Machines'
27 ′′′
28 df_machine = pandas_gbq.read_gbq(sql_machine, project_id="carttrend-4", credentials=
      credentials)
29
30 # Normalisation de la variable cible (entre 0 et 1)
31 scaler_uptime = MinMaxScaler()
32 df_machine['pourcentage_uptime_normalise'] = scaler_uptime.fit_transform(
33
      df_machine[['pourcentage_uptime']]
34)
35
```

```
36 # Standardisation du volume traite (moyenne=0, ecart-type=1)
37 sta_volume = StandardScaler()
38 df_machine['volume_traite_standalise'] = sta_volume.fit_transform(
39
       df_machine[['volume_traite']]
40)
41
42 # Encodage des variables categorielles (localisation)
43 le = LabelEncoder()
44 df_machine['localisation'] = le.fit_transform(df_machine['localisation'])
45
46 # Extraction de la liste des types de machines
47 lst_typesMachines = df_machine['type_machine'].unique()
48
49 # Initialisation d'une DataFrame pour stocker les résultats
50 df out = pd.DataFrame(columns=[
51
       'type_machine', 'month', 'localisation', 'temperature',
       'volume_traite_standalise', 'mse', 'rmse'
52
53 ])
54
55 # Entrainement d'un modèle par type de machine
56 for typeMachine in lst_typesMachines:
57
       # Filtrage des donnees
58
       # Ne garde que celles concernant le type de machine en cours d'analyse
59
      df_filtre = df_machine[df_machine['type_machine'] == typeMachine]
60
61
       # Separation variables explicatives / variables cibles
62
      X = df_filtre.drop(columns=[
63
           'pourcentage_uptime', 'pourcentage_uptime_normalise',
64
           'type_machine', 'volume_traite'
65
66
      y = df_filtre['pourcentage_uptime_normalise']
67
68
       # Decoupage train/test
69
      X_train, X_test, y_train, y_test = train_test_split(
70
           X, y, test_size=0.2, random_state=42
71
      )
72
73
       # Configuration et entrainement du modele
74
      model = xgb.XGBRegressor(
75
           n_estimators=100, # Nombre d'arbres
           learning_rate=0.1, # Taux d'apprentissage
76
           max_depth=3 # Profondeur maximale des arbres
77
78
79
      model.fit(X_train, y_train)
80
81
       # Prediction et evaluation
82
      y_pred = model.predict(X_test)
83
      y_test_inv = scaler_uptime.inverse_transform(y_test.to_numpy().reshape(-1, 1))
84
      y_pred_inv = scaler_uptime.inverse_transform(y_pred.reshape(-1, 1))
85
      mse = mean_squared_error(y_test_inv, y_pred_inv)
86
      rmse = np.sqrt(mse)
87
      rmses.append(rmse)
88
89
       # Stockage des resultats
90
       importances = model.feature_importances_
       df_out.loc[len(df_out)] = [typeMachine] + importances.tolist() + [mse, rmse]
91
92
93 # Envoi vers BigQuery
94 pandas_gbq.to_gbq(
95
      df_out,
       "carttrend_transform.ml_cause_panne_machine",
96
      project_id="carttrend-4",
97
98
      if_exists="replace"
99)
```

Listing B.4 – Enrichissement des données clients par géolocalisation IP

```
1 # [Fichier : get_customer_location.py]
 2 # Enrichissement des données clients par géolocalisation IP :
 3 # 1. Extraction des IP clients depuis BigQuery
 4 # 2. Appels à l'API ip2location.io (limité à 1 requête/0.2s) pour obtenir la position
 5 # 3. Mise en cache des résultats pour 7 jours
 6 # 4. Sauvegarde du résultat dans BigQuery
7
8 import requests
9 from google.oauth2 import service_account
10 import pandas as pd
11 import pandas_gbq
12 import time
13 import os
14 import datetime
15
16 # Authentification GCP via service account
17 credentials = service_account.Credentials.from_service_account_file(
      os.path.dirname(__file__) + '/Compte Service Machine Learning.json',
19)
20
21 # Extraction des clients à géolocaliser et de leurs IP
22 sql = '''
2.3
      SELECT id_client, ip
2.4
      FROM 'carttrend-4.carttrend_clean.Clients'
25 ′′′
26 df = pandas_gbq.read_gbq(sql, project_id="carttrend-4", credentials=credentials)
2.7
28 # Initialisation du DataFrame de résultats
29 client_geo = pd.DataFrame(columns=['id_client', 'pays', 'region', 'ville'])
30
31 # Gestion du cache
32 file_age = None
33 max_age = datetime.timedelta(days=7) # Durée de validité du cache
34 if os.path.isfile(os.path.dirname(__file__) + '/client_geo.pkl'):
3.5
       filetime = os.path.getmtime(os.path.dirname(__file__) + '/client_geo.pkl')
       file_age = datetime.timedelta(seconds=time.time()-filetime)
36
37
       if file_age < max_age: # Utilisation du cache si valide
38
           client_geo = pd.read_pickle(os.path.dirname(__file__) + '/client_geo.pkl')
39
40 # Selection des IPs dont la position n'est pas dans le cache
41 df_missing = pd.concat([df, client_geo]).drop_duplicates(
       subset=['id client'],
42
43
      keep=False
44)
45
46 # Appels à l'API
47 retries = 5 # Nombre de tentatives en cas d'échec
48 for row in df_missing.itertuples(index=False):
49
      try:
50
           # Appel à l'API ip2location
51
           response = requests.get(
               f'https://api.ip2location.io/?key=7149995C158ABB38F6BCE6505B006082&ip={
52
                  row.ip}&format=json',
53
               timeout=10
           )
54
55
56
           # Si bonne reponse...
57
           if response.status_code == 200:
58
               data = response.json()
59
               # ...ajout des données de géolocalisation aux résultats
60
               client_geo = pd.concat([
61
                   client_geo,
62
                   pd.DataFrame([[row.id_client,
```

```
63
                                data.get('country_name'),
64
                                data.get('region_name'),
65
                                data.get('city_name')]],
66
                              columns=client_geo.columns)
67
               ], ignore_index=True)
           # Sinon...
68
69
           else:
70
               # ...si il y a encore des tentatives disponibles...
71
               if retries:
72
                   retries -= 1
73
                   time.sleep(1)
                                  # Attend plus longtemps en cas d'echec
74
                   continue #...on continue avec le client suivant
75
               else:
76
                   raise Exception(f'API error: {response.status_code}')
77
78
      except Exception as e:
79
           print(f"Error processing IP {row.ip}: {str(e)}")
80
           break
81
82
      time.sleep(0.2) # Respect du rate limiting de l'API
83
84 # Mise à jour du cache si nécessaire
85 if not file_age or file_age > max_age:
86
      client_geo.to_pickle(os.path.dirname(__file__) + '/client_geo.pkl')
87
88 # Envoi vers BigQuery
89 pandas_gbq.to_gbq(
90
     client_geo,
91
      'carttrend_transform.Clients_Geo',
92
      project_id='carttrend-4',
      if_exists='replace'
93
94)
```

Annexe C

Orchestration Airflow

Listing C.1 – DAG Airflow orchestrant le pipeline ELT/ML complet

```
1 # [Fichier : carttrend_data_pipeline.py]
 2 # DAG Airflow orchestrant le pipeline ELT/ML complet :
 3 # 1. Synchronisation des données brutes via Fivetran
 4 # 2. Transformation avec DBT
 5 # 3. Exécution des modèles ML
 6 # 4. Enrichissement via API externe
 8 from airflow import DAG
 9 from airflow.operators.python import PythonOperator
10 from datetime import datetime, timedelta
11 import subprocess
12
13 # -----
14 # Définition des tâches
15 # -----
16
17 def start_fivetran_sync():
18
     # Déclenche la synchronisation des données sources via Fivetran
19
      print("Starting Fivetran Sync Task")
20
      subprocess.run(['python3', '../../airflow/scripts/start_fivetran_sync.py'])
21
22 def start_dbt_build():
23
       # Lance les transformations DBT (staging + marts)
2.4
      print('Starting DBT Build Job')
25
      subprocess.run(['python3', '../../airflow/scripts/start_dbt_build_job.py'])
26
27 def start_ml_prediction_ca():
       # Exécute le modèle Prophet de prévision de CA
2.8
29
      print('Starting ML Prediction CA')
      subprocess.run(['python3', '../../airflow/scripts/ml_prediction_ca.py'])
30
31
32 def start_ml_a_priori():
33
       # Exécute l'analyse des associations produits
34
      print('Starting ML APriori')
35
      subprocess.run(['python3', '../../airflow/scripts/ml_a_priori.py'])
36
37 def start_ml_regression_machines():
38
       # Exécute le modèle XGBoost d'analyse des pannes
39
      print('Starting ML Regression Machines')
40
      subprocess.run(['python3', '../../airflow/scripts/ml_regression_machines.py'])
41
42 def get_customer_locations():
43
       # Appel l'API ip2location pour géolocaliser les clients
44
      print('Starting Customer Location API Call Task')
45
      subprocess.run(['python3', '../../airflow/scripts/get_customer_location.py'])
```

```
46
 47 # -----
 48 # Configuration du DAG
 49 # -----
 50
 51 default_args = {
 52
       'owner': 'airflow',
 53
       'start_date': datetime(2025, 5, 27),
 54
        'schedule_interval': '@daily', # Exécution quotidienne à minuit
 55
        'retries': 1, # Nombre de tentatives en cas d'échec
 56
        'retry_delay': timedelta(minutes=30), # Délai entre les tentatives
 57
        'email_on_failure': True,
                                  # Notification par email (configuration SMTP prealable
        'email': ['carttrenddata4@gmail.com'] # Destinataires
 58
 59 }
 60
 61 # Initialisation du DAG
 62 \text{ dag} = DAG(
 63
        'carttrend_data_pipeline', # Nom du DAG
 64
       default_args=default_args,
 65
       description='Pipeline ELT/ML complet pour CartTrend',
 66
       catchup=False # Désactive les runs rétroactifs
 67 )
 68
 69 # -----
 70 # Définition des opérateurs
 71 # -----
 72
 73 t1 = PythonOperator(
 74
      task_id="start_fivetran_sync",
 75
       python_callable=start_fivetran_sync,
 76
       dag=dag,
 77)
 78
 79 t2 = PythonOperator(
       task_id="start_dbt_build",
 80
 81
       python_callable=start_dbt_build,
 82
       dag=dag,
 83)
 84
 85 t3 = PythonOperator(
 86
       task_id="start_ml_prediction_ca",
 87
       python_callable=start_ml_prediction_ca,
 88
       dag=dag,
 89)
 90
 91 t4 = PythonOperator(
       task_id="start_ml_regression_machines",
 92
 93
       python_callable=start_ml_regression_machines,
 94
       dag=dag,
 95)
 96
 97 t5 = PythonOperator(
 98
       task_id="start_ml_a_priori",
 99
       python_callable=start_ml_a_priori,
100
       dag=dag,
101)
102
103 t6 = PythonOperator(
104
       task_id="get_customer_locations",
105
       python_callable=get_customer_locations,
       dag=dag,
106
107)
108
```

```
109 # ------
110 # Orchestration des tâches
111 # ------
112
113 # Schéma d'exécution:
114 # 1. Fivetran (t1) doit terminer avant DBT (t2)
115 # 2. DBT (t2) doit terminer avant les modèles ML (t3,t4,t5) et geoloc (t6)
116 t1 >> t2
117 t2 >> [t3, t4, t5, t6]
```

Listing C.2 – Script de declenchement de synchronisation des données via l'API Fivetran

```
1 # [Fichier : start_fivetran_sync.py]
 2 # Script de declenchement de synchronisation des données via l'API Fivetran :
 3 # 1. Déclenchement de la synchronisation
 4 # 2. Vérification périodique de l'état
5 # 3. Gestion des erreurs
7 import requests
8 import time
9
10 # Déclenche une synchronisation via l'API Fivetran
11 def start_sync():
12
      url = "https://api.fivetran.com/v1/connections/haunt_charming/sync"
13
      headers = {
14
           "Accept": "application/json; version=2",
           "Authorization": "Basic
1.5
              dWY4ZDBmMUZ1WGZ5Z3M5cTpacDJFYjV1RDZzQUp0ZVUxWG1ZUGt1WGV2T1R2Z2Vyag==",
16
           "content-type": "application/json"
17
      }
18
19
       # Envoi de la requête POST pour démarrer le sync
      response = requests.post(url, json={'force' : False}, headers=headers)
2.0
21
       response.raise_for_status() # Si echec, lève une exception
22
       return True
23
24 # Récupère le statut actuel de la synchronisation
25 def get_sync_state():
2.6
      url = "https://api.fivetran.com/v1/connections/haunt_charming"
27
      headers = {
28
           "Accept": "application/json; version=2",
           "Authorization": "Basic
29
              dWY4ZDBmMUZ1WGZ5Z3M5cTpacDJFYjV1RDZzQUp0ZVUxWG1ZUGt1WGV2T1R2Z2Vyag=="
30
31
32
       # Envoi de la requete GET a l'API Fivetran
      response = requests.get(url, headers=headers)
33
       response.raise for status()
34
35
       # Retourne l'etat du transfer ('syncing', 'paused', etc.)
36
      return response.json()['data']['status']['sync_state']
37
38 if __name__ == '__main__':
39
       retries = 5 # Nombre maximum de tentatives d'obtention du statut
40
41
       # Phase 1: Démarrage de la synchronisation
42
       start_sync()
43
      print('Sync Requested')
44
45
       # Phase 2: Verification periodique du statut de la synchronisation
46
       sync_started = False
      while True:
47
48
          try:
49
               sync_state = get_sync_state()
50
           # Lève une exception si la demande de statut n'aboutit pas 5 fois
51
```

```
52
           except Exception as e:
53
               print(f'Could not get sync status : {e} - Retries : {retries}')
54
               if retries:
55
                   retries -= 1
56
               else:
57
                   raise Exception('Could not get sync status')
58
59
           # Si on a une réponse de l'API...
60
           else:
61
               # ...et qu'il est fini : stop.
62
               if sync_state != 'syncing' and sync_started:
63
64
65
               # ...et qu'il n'a pas encore commence...
               if sync_state != 'syncing' and not sync_started:
66
67
                   continue
68
69
               # ...et qu'il est toujours en cours...
               if sync_state == 'syncing':
70
                   sync_started = True
71
72
                   print('Syncing')
73
74
               # ...et que la synchronisation est un echec : exception.
75
               else:
76
                   raise Exception('Sync Failed !')
77
78
           # ...attendre 3 secs et recommencer.
79
           time.sleep(3)
80
81
       print('Syncing Done')
```

Listing C.3 – Script de declenchement du Job DBT Build via l'API DBT

```
1 # [Fichier : start_dbt_build_job.py]
 2 # Script de declenchement du Job DBT Build via l'API DBT :
 3 # 1. Déclenchement du job
 4 # 2. Vérification périodique de l'état
 5 # 3. Gestion des erreurs
 6
7 import requests
8 import time
9 import enum
10
11 # Configuration d'accès à l'API DBT Cloud
12 # Tout ces identifiants ont ete crees sur l'interface DBT Cloud
13 ACCOUNT_ID = '70471823462199' # Identifiant du compte DBT
14 JOB_ID = '70471823467186' # Identifiant du job
15 API_KEY = 'dbtu_57j1bQUrnaKY_Xf03UdF2kHC9g4DkgZNAmyAy2D3QAgQmxIqiM' # Token d'API
16
17 # Enumeration des statuts possibles d'un job DBT
18 # Les codes sont documentés sur la doc DBT :
19 # https://docs.getdbt.com/dbt-cloud/api-v2#/operations/List%20Runs
20 class DbtJobRunStatus(enum.IntEnum):
21
      QUEUED = 1
22
      STARTING = 2
23
      RUNNING = 3
      SUCCESS = 10
2.4
25
      ERROR = 20
26
      CANCELLED = 30
27
28 # Déclenche l'exécution du job DBT et retourne l'ID d'exécution
29 def trigger_job():
30
      url = f'https://jp371.usl.dbt.com/api/v2/accounts/{ACCOUNT_ID}/jobs/{JOB_ID}/run/
      headers = {'Authorization' : f'Token {API_KEY}'}
31
```

```
32
       payload = {'cause' : 'Triggered from Airflow'}
33
34
       # Envoi de la requete POST a l'API DBT
35
       response = requests.post(url, json=payload, headers=headers)
36
       response.raise_for_status() # Si echec, lève une exception
37
38
       # Retourne l'ID d'exécution
39
       job_run_id = response.json()['data']['id']
40
       return job_run_id
41
42 # Récupère le statut actuel d'une exécution DBT
43 def get_job_status(job_run_id):
44
       url = f'https://jp371.usl.dbt.com/api/v2/accounts/{ACCOUNT_ID}/runs/{job_run_id}/
45
       headers = {'Authorization': f'Token {API_KEY}'}
46
47
       # Envoi de la requete GET a l'API DBT
48
       response = requests.get(url, headers=headers)
49
       response.raise_for_status()
50
51
       # Retourne le statut de l'execution
52
       return response.json()['data']['status']
53
54
55 if __name__ == '_
                    __main__':
56
       retries = 5 # Nombre maximum de tentatives d'obtention du statut
57
58
       # Phase 1: Déclenchement du Job
59
       job_run_id = trigger_job()
60
       print(f'DBT Build Job Triggered : {job_run_id}')
61
62
       # Phase 2: Verification periodique du statut du Job
63
       while True:
64
           # Demande l'etat du transfer a Fivetran
65
           try:
66
               status = get_job_status(job_run_id)
67
           # Lève une exception si la demande de statut n'aboutit pas 5 fois
68
69
           except Exception as e:
70
               print(f'Could not get job status : {e} - Retries : {retries}')
71
               if retries:
72
                   retries -= 1
73
               else:
74
                   raise Exception('Could not get job status')
75
76
           # Si on a une réponse de l'API...
77
           else:
78
               # ...et qu'il est fini : stop.
79
               if status == DbtJobRunStatus.SUCCESS:
80
81
               # ...et qu'il s'est produit une erreur : exception.
82
               elif status == DbtJobRunStatus.ERROR or status == DbtJobRunStatus.
                  CANCELLED:
83
                   raise Exception("Job Failed !")
84
           # ...attendre 3 secs et recommencer.
85
86
           time.sleep(3)
87
88
      print('Job Finished')
```

Annexe D

Dashboards

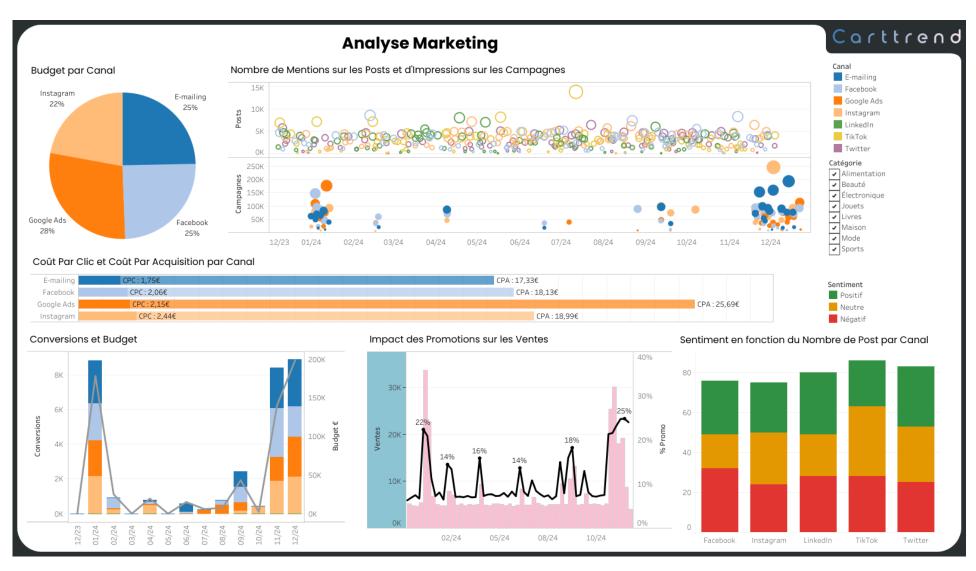


FIGURE D.1 – Analyse Marketing

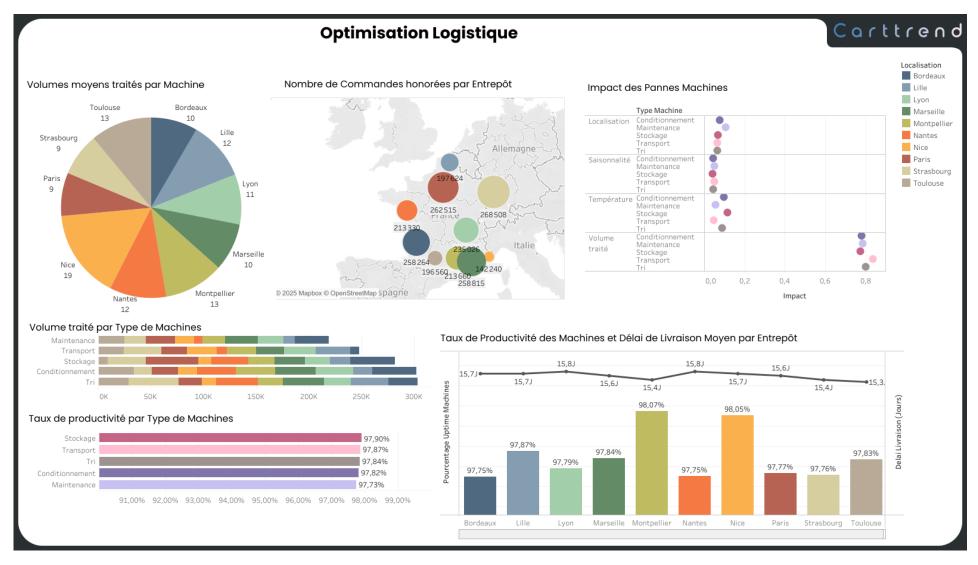


FIGURE D.2 – Optimisation Logistique

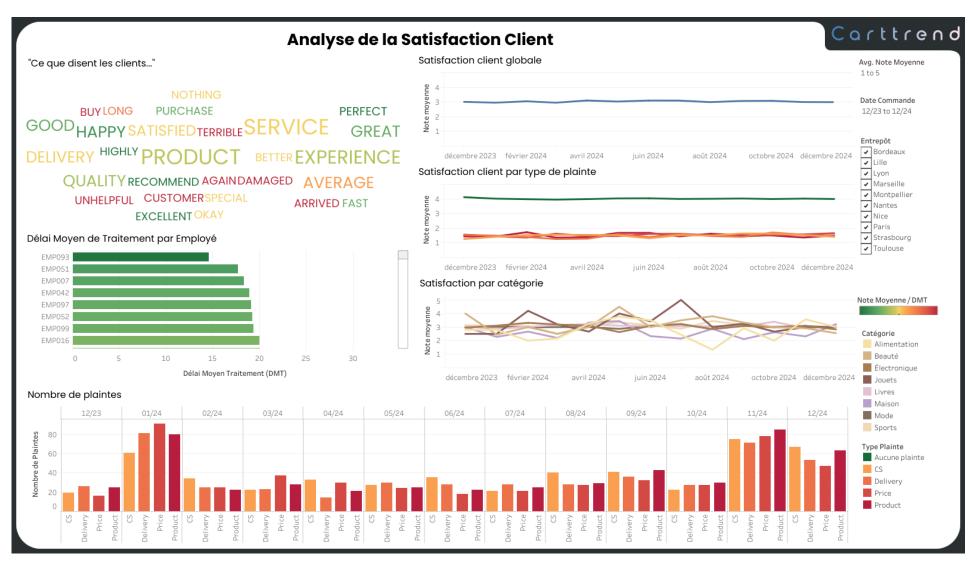


FIGURE D.3 – Analyse de la Satisfaction Client

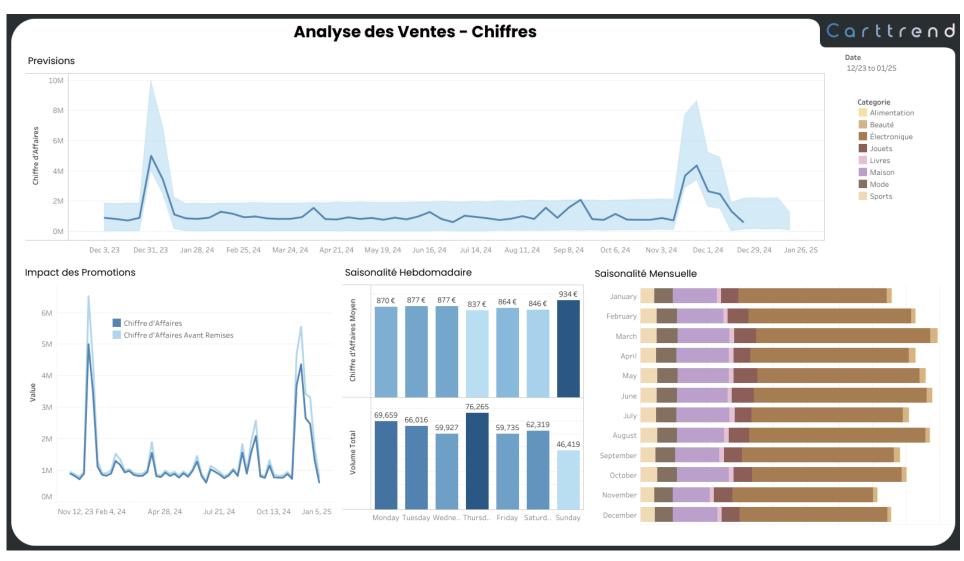


FIGURE D.4 – Ventes - Chiffres

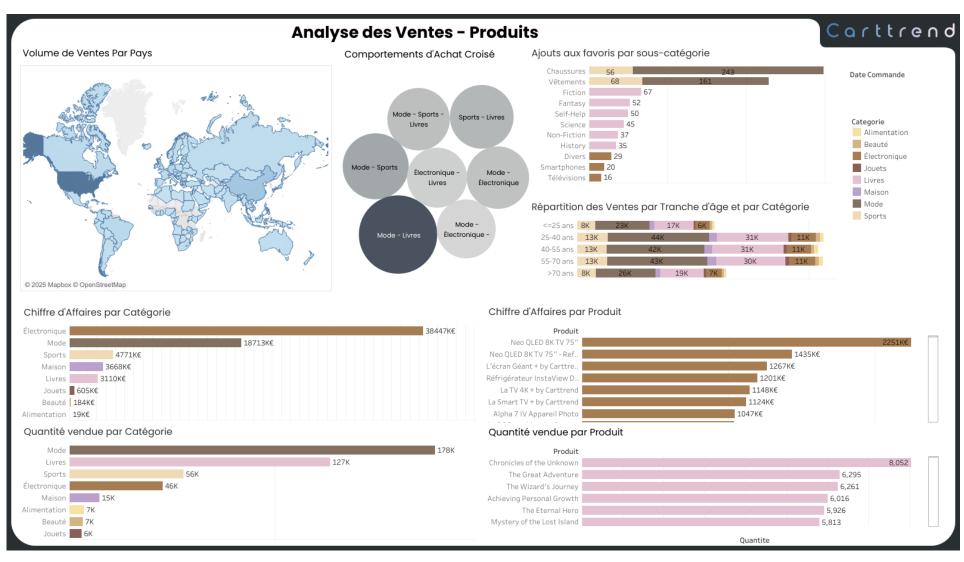


FIGURE D.5 – Ventes - Produits