



PROJET BIUM

voyager avec 200 € et 2 masques
Y a-t-il une vivante dans le globe ?

Auteur
Youva ADDAD

Client
Laure SOULIER

Table des matières

1	Résumé	2
2	Introduction	2
3	Difficultés rencontrées	2
4	Données	3
5	Extract Transform Load	5
6	Modélisation	8
6.1	Problématique	8
6.2	Schéma Conceptuelle	8
6.3	Schéma Normalisé	9
7	Outil	9
8	Analyse	9
9	Prédiction	10

1 Résumé

Dans ce travail de business intelligence nous avons réalisé ce projet centré sur la décision du choix de la destination pour les vacances en fonction de divers critères. Le coût de vie de la ville où bien la situation sanitaire du pays entre autres nos critères de sélection.

Dans ce rapport nous allons aborder les différents outils utilisés ainsi que les différentes analyses de données, la façon dont nous avons extrait les données mais en premier lieu nous allons parler de l'analyse des besoins.

Nous allons parler des différentes étapes que nous avons suivies afin d'aboutir à un résultat final qui soit le plus adéquat possible. Schéma suivi :

1. Définition des objectifs et des exigences
2. Le choix de la méthodologie et des outils à utiliser
3. Etablissement d'un programme de travail
4. Mise en place du programme de travail
5. Organiser le reporting et transmettre l'information

2 Introduction

Depuis plus d'une année maintenant le monde connaît un ralentissement dû au covid, les déplacements de chacun sont fortement très limités, il y a eu de plus un ralentissement économique dans le monde entier. Avec l'approche des vacances on peut se demander quelle ville/pays choisir qui minimisera le risque lié au covid et qui est le moins chère possible, tout en profitant bien sûr des loisirs et points d'intérêts divers et variés.

Nous allons analyser ce point sous plusieurs dimensions pour pouvoir prendre une décision la meilleure qui soit.

3 Difficultés rencontrées

S'agissant d'un sujet à l'échelle mondiale nous avons passé une partie non négligeable à la collecte de données, nous avons collecté une quantité de données suffisante pour pouvoir traiter le sujet, mais ceci était une partie très sérieuse nous avons des difficultés à trouver des données englobant le tout, des données traitant suffisamment de ce sujet, ce qui nous a permis tout de même d'exploiter de nouvelles techniques que nous détaillerons plus tard.

4 Données

Nous avons donc sélectionné les données qu'il nous fallait afin d'axer notre fait dessus, s'agissant d'un sujet qui traite du covid et de la destination la moins chère nous avons choisi de regarder :

- La situation sanitaire, le nombre cas ainsi que le nombre de vaccination.
- Le coût de la vie, typiquement le coût d'un restaurant, taxis, location
- Le prix des hôtels, le logement est une étape décisive dans le choix d'une ville
- Prix des vols
- POI, les points d'intérêt pour toutes les villes du monde, Point d'intérêt pour les restaurants, tourisme, choses à voir ou à faire

(nous avons aussi des données sur la criminalité des villes mais on les a pas exploitées).

1. `owid-covid-data.csv` ce fichier recense des informations sur le nombre cas, nombre de vaccination, nombre de mort pour chaque pays et depuis le début de la pandémie 24/02/2020 à 30/04/2021, ce fichier a été extrait depuis 'Johns Hopkins Coronavirus Resource Center' une université américaine.
2. `Cost_of_life.csv` ce fichier recense pour une ville divers coûts par exemple le coût d'un restaurant, le coût d'un café, coût fruit & légume, coût taxi ... etc, ce fichier a été scrapé depuis `numbreo.com`. Numbeo est une base de données mondiale accessible à tous sur les prix à la consommation déclarés, les taux de criminalité perçus, la qualité des soins de santé, entre autres statistiques.
3. `hostel.csv` ce fichier nous fournit pour un hôtel son nom, le prix la nuit, le nombre de lits ...etc. Nous avons scrapé ces données directement depuis Booking.
4. `vols.csv` ce fichier nous fournit le nom de l'aéroport, la date du vol, le pays destination, latitude & longitude et le prix du vol ...etc de même ce fichier a été scrapé depuis booking.
5. POI est un dossier comportant les points d'intérêt, ou chaque fichier est un csv pour un pays des points d'intérêt, ils sont répartis par catégorie & sous-catégorie dans le fichier le nom international. Nous avons récolté ce fichier dans plusieurs sources dans la majorité est les sources gouvernementales

la figure suivante montre quels sont les pays et ville que nous avons étudié :

Feuille 1



Carte basée sur les lng et lat. Les détails affichés sont associés au/à la City.

FIGURE 1 – Les villes en point noir

5 Extract Transform Load

Nous avons précédemment cité les sources et les extractions des données mais nous avons pas parler de la transformation et le chargement des données. Pour cela nous avons utilisé :

1. Talend pour gérer nos données
2. Dataiku pour pouvoir traité, joindre, filtrer, remplacer, améliorer nos données.
3. Notebook python
4. Tableau Desktop pour l'intégration de nos données
5. Power BI

Nous avons donc nettoyé nos données avec ces outils et nous les avons converties aux formats de rapport qui conviennent, Nous avons donc appliqué les règles suivante :

- Nous avons sélectionner uniquement certaines colonnes typiquement les colonnes que nous avons jugé pertinente pour notre analyse. Les colonnes ayant des valeurs nulles par exemple nous les avons pas sélectionnées pour une meilleur cohérence. Par exemple la base covid ayant 41 caractéristique nous avons enlevé celle qui ne servir à rien.
- Nous avons aussi traduit les valeur codée par exemple les longitude et latitude ont été traduites en Point qui est une forme générale pour pouvoir faire une map à la suite.
- Nous avons Dérivé de nouvelles colonnes calculées par combinaison d'autres caractéristiques. par exemple pour la base coût de la vie nous avons créé une caractéristique regroupant la somme des prix de transport, de restauration, Sim prépayé.
- Nous avons trier les données en fonction de certaines colonnes pour améliorer les performances de recherche. Par exemple dans le cost of living nous les avons trier par pays.
- Nous avons joint des données provenant de plusieurs bases de données, cette partie nous a permis de fusionner et déduplicer les données. Par exemple dans la base cost of living nous avons les villes mais pas les coordonnées nous avons donc opéré une jointure pour pouvoir trouver les coordonnées.
- Nous avons agréger les données par cumul, moyenne, somme. cette étape nous a permis de ne pas prendre toutes les données mais seulement l'aggrégation.
- Nous avons aussi utilisé la transposition des données étant dans le mauvais sens.
- Nous avons fractionner des colonnes en plusieurs colonnes. par exemples dans les POI nous avons le nom du monument en internationale séparer du nom avec langue nationale.
- Nous avons utilisé le select distinct pour pouvoir ne prendre qu'une seule fois un attribut par exemple.
- Nous avons créé de nouvelles bases, par exemple pour la date, nous avons donc fait un split de dates pour pouvoir avoir une plage bien précise.

Avec tout ces traitements nous avons défini nos tables et data warehouse. Nous détaillerons plus tard comment nous avons créé le datawarehouse et la façon dont nous l'avons alimenté.

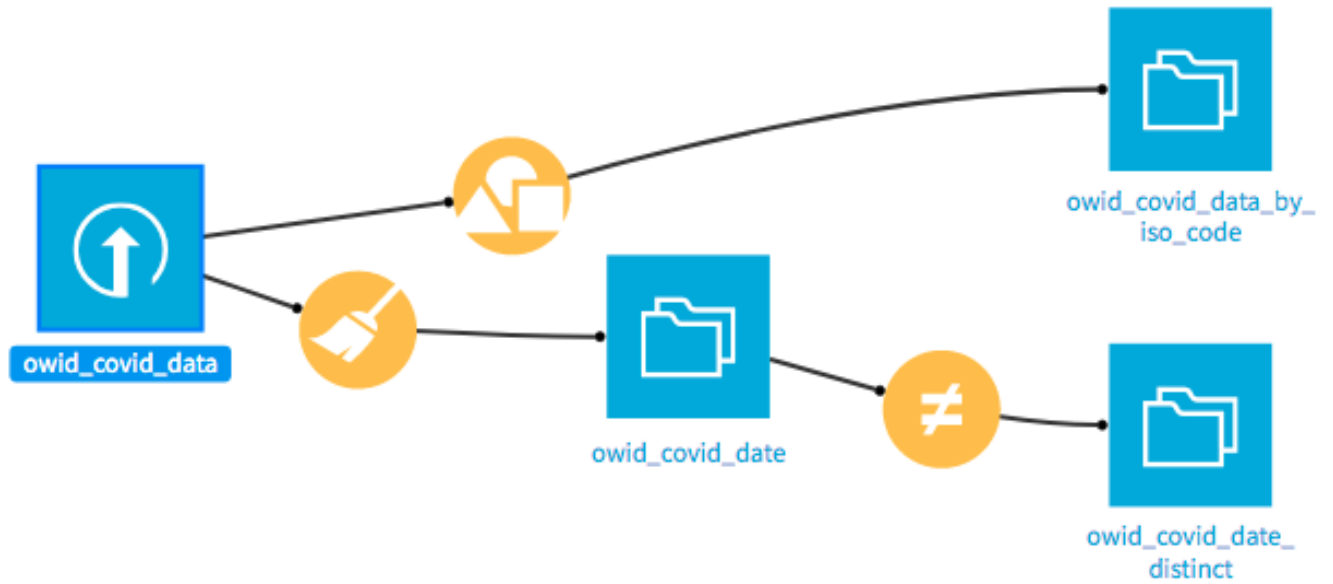


FIGURE 2 – Traitement des données covid

Ici donc nous avons deux paths, un pour la préparation des données (enlever les colonnes pas utiles ...Etc) et le select les valeurs distinct. De l'autre côté du chemin nous avons sélectionné les iso_code de pays.

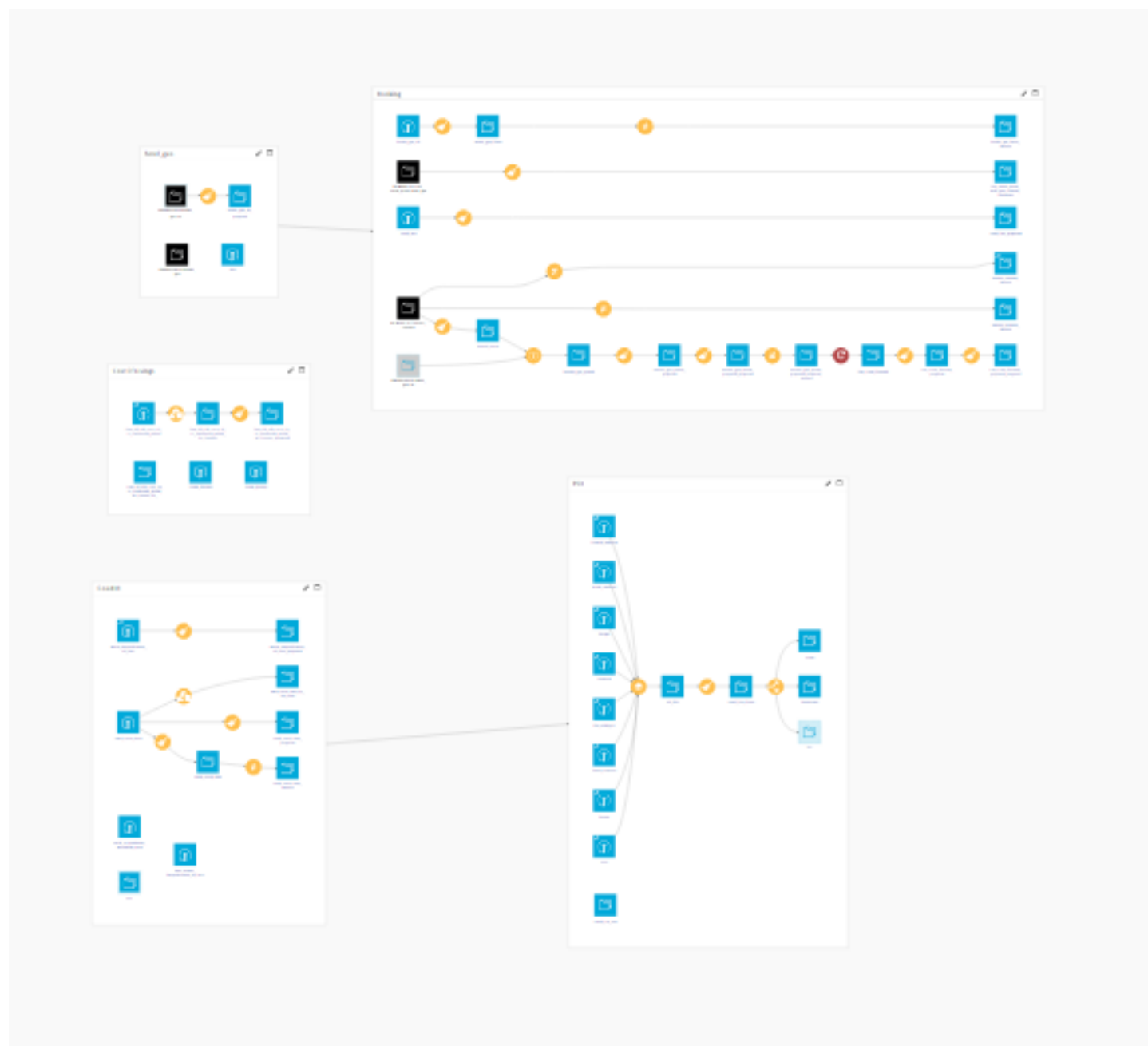


FIGURE 3 – Général des traitements effectués

6 Modélisation

6.1 Problématique

Nous souhaitant pouvoir prendre une décision de voyage pour cette été. Pour pouvoir répondre efficacement a cette problématique notre schéma doit donc répondre :

- le coût d'une ville, nourriture, transport et déplacement...etc
- la situation sanitaire en fonction du temps du nombre de vaccin, du nombre de test, et le nombre de cas.
- le prix d'un vol a une destination en fonction du temps.
- le prix d'un hotel d'une destination en fonction du temps.

6.2 Schéma Conceptuelle

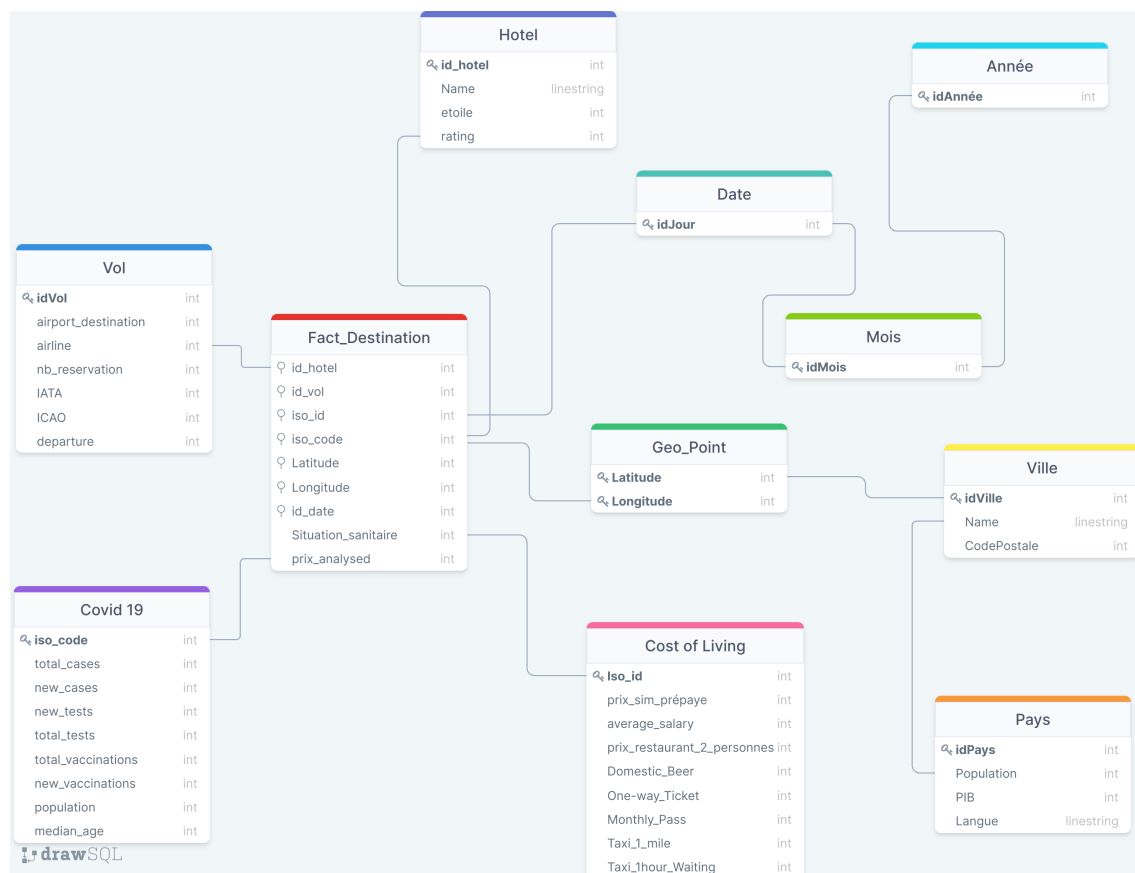


FIGURE 4 – Schéma en flocon

6.3 Schéma Normalisé

- Cost_of_living(iso_id, prix_sim_prepay, average_salary, prix_restaurants_2_personnes, Domestic.beer, ticket, monthly_pass, taxi_1_mile, taxi_1_hour).
- Covid_19(iso_code, total_cases, new_cases, new_tests, total_tests, total_vaccinations, new_vaccinations, population, median_age).
- Hotel(id_hote, Name, etoile, rating).
- Vol(id_vol, airport_destination, airline, nb_reservation, IATA, ICAO, departure).
- Date(id_jour, #id_mois).
- Mois(id_mois, #Année).
- Année(id_année).
- Geo_Point(Latitude, Longitude, #id_ville).
- Ville(id_ville, #id_pays, CodePostal).
- Pays(id_pays, Population, PIB, Langue).
- Fact_Destination(#iso_id, #iso_code, #id_hote, #id_vol, #id_jour, #Latitude, #Longitude, Situation_sanitaire, prix_analysed)

7 Outil

1. **Talend** : Nous a permis de nettoyer et d'intégrer nos données afin de les exploiter.
2. **Dataiku** : Nous a permis d'appliquer efficacement plein d'aggrégation, de jointure entre les données, unifier les données en splittant le tout, appliquer des filtres, créer de nouvelles colonnes pré-calculées.
3. **Notebook python** : Nous a permis de collecter les données via le scrapping, afin de les visualiser pour déterminer la qualité des données collectées.
4. **Tableau Desktop** : Nous a permis de faire des visualisations pour pouvoir les intégrer dans le dashboard.

8 Analyse

Nous avons effectué notre analyse avec Power BI et Tableau Desktop qui nous offre beaucoup de fonctionnalités. Nous avons effectué les analyses suivantes :

1. **Analyse Coût de la vie** : pour pouvoir voir les pays ayant un coût de vie moins cher, nous avons donc commencé à intégrer nos données dans Tableau, nous avons choisi donc de ne voir que le coût d'un restaurant, d'un taxi, SIM prépayé, salaire moyen du pays ou de la ville.
2. **Analyse situation sanitaire** : nous avons effectué un comparatif du nombre de cas, test, vaccin réalisé pour un pays en suivant la granularité Année/Mois.
3. **Analyse vols** : Ici nous analysons les prix des vols pour voir la destination la plus propice pour voyager.
4. **Analyse hôtel** : Il s'agit ici d'avoir une idée générale sur les tarifs appliqués dans ces villes en fonction du temps.

9 Prédiction

Dans cette section nous avons décidé de prédire le nombre moyen ou total de morts liées au covid dans un futur proche.

1. Lasso
2. Ridge
3. Linear Regression
4. SVM Regression

Nous avons décidé de travailler avec les SVM Regression tout simplement parce qu'ils sont très puissants, ils séparent en suivant et maximisent une marge. De plus ils permettent d'apprendre des solutions non linéaires avec utilisation d'un noyau.

Donc les SVM's est un choix naturellement bon pour pouvoir faire des prédictions sur un futur proche.