



Projet BIUM

---

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

---

Laure Soulier

# Table des mati`eres

|          |                                  |           |
|----------|----------------------------------|-----------|
| <b>1</b> | <b>R´esum´e</b>                  | <b>2</b>  |
| <b>2</b> | <b>Introduction</b>              | <b>2</b>  |
| <b>3</b> | <b>Difficult´es rencontr´ees</b> | <b>2</b>  |
| <b>4</b> | <b>Donn´ees</b>                  | <b>3</b>  |
| <b>5</b> | <b>Extract Transform Load</b>    | <b>5</b>  |
| <b>6</b> | <b>Mod´elisation</b>             | <b>8</b>  |
| 6.1      | Probl´ematique . . . . .         | 8         |
| 6.2      | Sch´ema Conceptuelle.. . . .     | 8         |
| 6.3      | Sch´ema Normalis´e . . . . .     | 9         |
| <b>7</b> | <b>Outil</b>                     | <b>9</b>  |
| <b>8</b> | <b>Analyse</b>                   | <b>9</b>  |
| <b>9</b> | <b>Pr´ediction</b>               | <b>10</b> |

# 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 choix de la ville où bien la situation sanitaire du pays entre autre nos critères de sélection. Dans ce rapport nous abordons 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, en plus un ralentissement économique dans le monde entier. Avec l'approche des vacances on peut se demander quel 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 toutes les données traitant suffisamment de ce sujet, 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 qui 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 choisis 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 de 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 diversifié par exemple le coût d'un restaurant, le coût d'un café, coût fruit & légume, coût taxi ... etc, ce fichier a été scrappé depuis `numbeo.com`. Numbeo est une base de données mondiale accessible à tous sur les prix à la consommation d'éclairés, les taux de criminalité, 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 lit ...etc. Nous avons scrappé 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 scrappé depuis booking.
5. POI est un dossier comportant les points d'intérêt et chaque fichier est un csv pour un pays des points d'intérêt sont réparties 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 ´etudier :

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 cit   les sources et les extractions des donn  es mais nous avons pas parler de la transformation et le chargement des donn  es. Pour cela nous avons utilis   :

1. Talend pour g  rer nos donn  es
2. Dataiku pour pouvoir trait  , joindre, filtrer, remplacer, ameliore nos donn  es.
3. Notebook python
4. Tableau Desktop pour l'int  gration de nos donn  es
5. Power BI

Nous avons donc nettoy   nos donn  es avec ces outils et nous les avons converties aux formats de rapport qui conviennent, Nous avons donc appliqu   les r  gles suivante :

- Nous avons s  lectionner uniquement certaines colonnes typiquement les colonnes que nous avons jug   pertinente pour notre analyse. les colonnes ayant des valeurs nulles par exemple nous les avons pas s  lectionnes pour une meilleur coh  rence. Par exemple la base covid ayany 41 caract  ristique nous avons du enlev   celle qui ne server a rien.
- Nous avons aussi traduit les valeur cod  e par exemple les longitude et latitude ont   t   traduite en Point qui est une forme g  n  rale pour pouvoir faire une map a la suite.
- Nous avons D  river de nouvelle colonnes calcul  e par combinaison d  autre caract  ristique par exemple pour la base code de la vie nous avons cr  er une caract  ristique regroupant la somme des prix de transport, de restauration, Sim pr  pay  e.
- Nous avons trier les donn  es en fonction de certaine colonnes pour am  liorer les performances de recherche. Par exemple dans le cost of living nous les avons trier par pays.
- Nous avons joint des donn  es provenant de plusieurs base de donn  es, cette partie nous avons permit de fusionn  e et d  dupliquer les donn  es. Par exemple dans la base cost of living nous avons les villes mais pas les coordonn  es nous avons donc op  rer une jointure pour pouvoir trouv   les coordonn  es.
- Nous avons aggr  ger les donn  es par cumul, moyenne, somme. cette   tape nous a permis de ne pas prendre toute les donn  es mais seulement l  aggr  gation.
- Nous avons aussi utilis   la transposition les donn  es   tant 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 s  parer du nom avec langue nationale.
- Nous avons utilis   le select distinct pour pouvoir ne prendre qu  une seul fois un attribut par exemple.
- Nous avons cr  er de nouvelle base, par exemple pour la date, nous avons donc fait un split de dates pour pouvoir avoir une plage bien pr  cise.

Avec tout ces traitement nous avons d  finie nous table et data werehouse. Nous d  taillerons plus tard comment nous avons cr  er le datawarehouse et la fa  on dont nous l  avons alimenter.

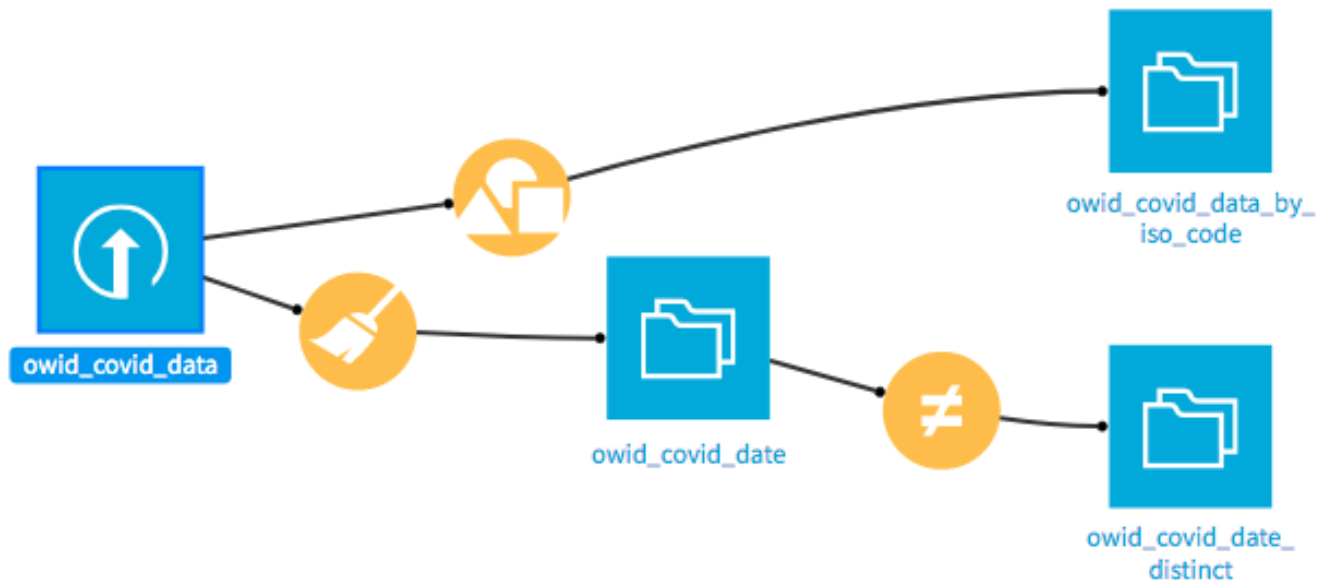


Figure 2 – Traitement des données covid

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.

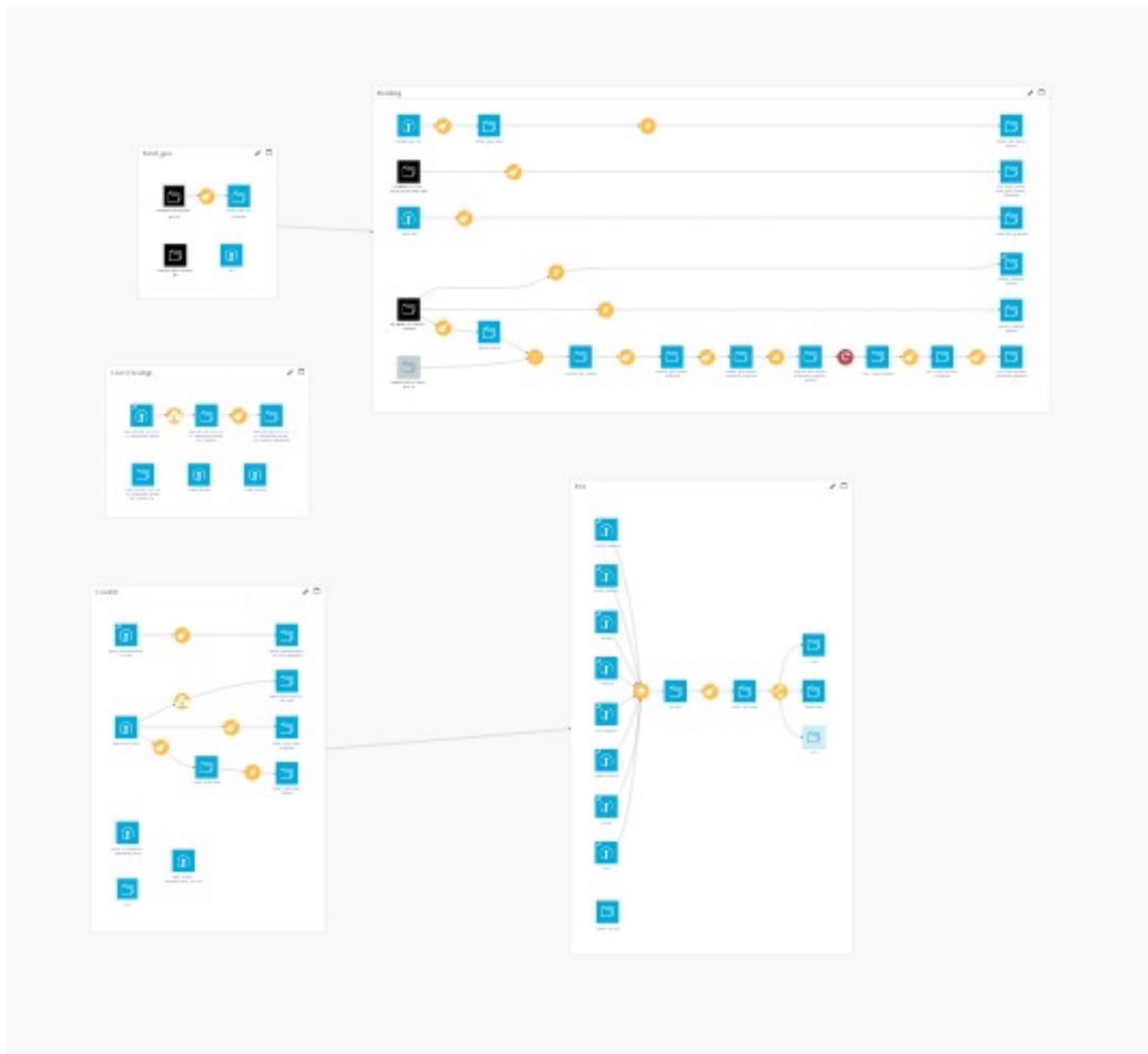


Figure 3 – G'én'eral des traitements effectu'e



## 6 Modélisation

### 6.1 Problématique

Nous souhaitant pouvoir prendre une décision de voyage pour cette période, pour répondre efficacement à cette problématique notre schéma doit donc répondre :

- le coût d'une ville, nourriture, transport et d'eplacement...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

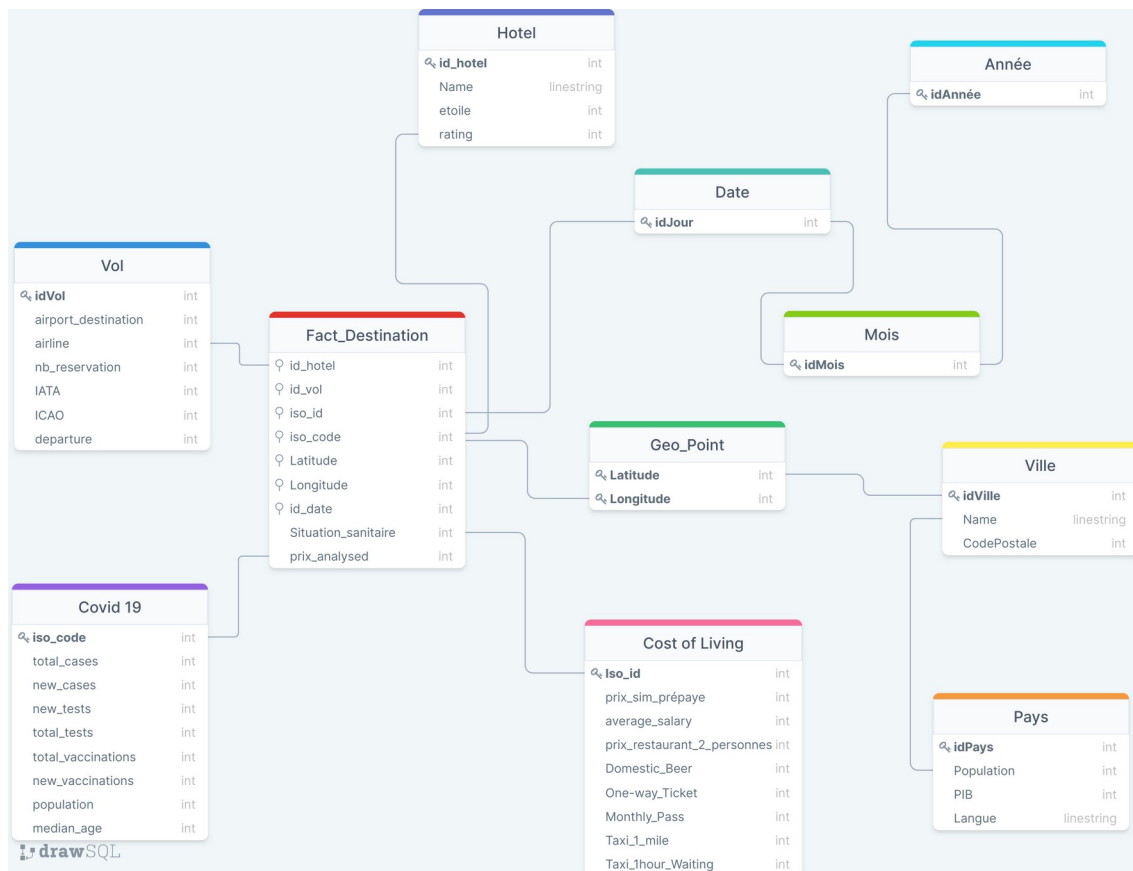


Figure 4 – Schéma en flocon

### 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 analysé)

## 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'agrégation, 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 prépayée, salaire moyen du pays ou de la ville.
2. **Analyse situation sanitaire** : nous avons effectué un comparatif du nombre de cas, test, vaccin réalisée 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 mort liée 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, sépare en suivant en maximisant une marge. 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.