



# Rapport de projet **Data Mining**

Yasmine SAIDI

10 mars 2021

## Table des matières

<b>Introduction</b>	<b>4</b>
<b>1 Visualisation de données avec R</b>	<b>5</b>
1.1 Compréhension des données . . . . .	5
1.2 Préparation des données . . . . .	5
1.3 Exploration de données . . . . .	6
1.4 Analyse des résultats . . . . .	7
<b>2 Classification avec Python</b>	<b>8</b>
2.1 Compréhension du métier . . . . .	8
2.1.1 Détermination des objectifs stratégiques et opérationnelles . . . . .	8
2.1.2 Détermination des objectifs du Data Mining . . . . .	8
2.2 La compréhension des données . . . . .	8
2.2.1 Description des données . . . . .	8
2.2.2 Exploration des données . . . . .	10
2.3 Préparation des données . . . . .	13
2.4 Modélisation . . . . .	15
2.4.1 Sélection des techniques de modélisation . . . . .	15
2.4.2 Explication des techniques de modélisation . . . . .	16
2.4.3 Construire le modèle . . . . .	16
2.5 Evaluation . . . . .	17
2.5.1 Comparaison des modèles . . . . .	17
2.5.2 Amélioration . . . . .	18
2.5.3 Importance des variables . . . . .	18
2.6 Déploiement . . . . .	19
<b>3 Topic Modeling</b>	<b>20</b>
3.1 Prétraitement des données . . . . .	20
3.1.1 Tokenisation . . . . .	20
3.1.2 Fréquence des mots . . . . .	20
3.1.3 Stemming : Normalisation des mots . . . . .	21
3.1.4 Lemmatisation : Base des mots . . . . .	21
3.1.5 Stop words . . . . .	22
3.1.6 Part of speech tagging (POS) . . . . .	22
3.1.7 Reconnaissance des différentes entités . . . . .	22
3.1.8 Chunking . . . . .	23
3.2 Visualisation . . . . .	23
3.3 Topic modeling avec LDA . . . . .	24
3.3.1 La fréquence de mots . . . . .	24
3.3.2 Gensim filters . . . . .	24



3.3.3	Gensim doc2bow . . . . .	24
3.3.4	Running LDA using Bag of Words . . . . .	25
3.3.5	TF-IDF . . . . .	25
3.3.6	Running LDA using TF-IDF . . . . .	26
3.3.7	Evaluation du modèle . . . . .	27

<b>Conclusion</b>	<b>28</b>
-------------------	-----------

## Table des figures

1	Visualisation avec R - Ajout de deux colonnes . . . . .	5
2	Visualisation avec R - Calcul des moyennes . . . . .	6
3	Visualisation avec R - Diagramme à bâton . . . . .	6
4	Visualisation avec R - Diagramme linéaire - Margarine A . . . . .	7
5	Visualisation avec R - Diagramme linéaire - Margarine B . . . . .	7
6	Description des attributs . . . . .	10
7	La matrice de corrélation . . . . .	10
8	le sexe des malades cardiaques . . . . .	11
9	le taux de cholestérol et glucose des malades cardiaques . . . . .	11
10	l'âge des malades cardiaques . . . . .	12
11	Equilibrage du jeu de données . . . . .	12
12	Valeurs manquantes . . . . .	13
13	Valeur anormale - poids . . . . .	13
14	Valeurs anormales - ap_hi . . . . .	14
15	Valeurs anormales - ap_lo . . . . .	14
16	Nettoyage - remplacement de poids anormale par la moyenne . . . . .	14
17	Nettoyage - suppression des lignes avec valeurs anormales . . . . .	15
18	Sélection des données pour la modélisation . . . . .	15
19	la construction des modèles . . . . .	17
20	Comparaison des modèles . . . . .	18
21	Importance des variables . . . . .	18
22	Tokenisation . . . . .	20
23	Fréquence des mots . . . . .	21
24	Stemming . . . . .	21
25	Stop words . . . . .	22
26	Part of speech tagging (POS) . . . . .	22
27	Named entity recognition . . . . .	23
28	Chunking . . . . .	23
29	Nuage de mots . . . . .	23
30	La fréquence des mots . . . . .	24
31	Gensim doc2bow . . . . .	24
32	Running LDA using Bag of Words . . . . .	25
33	TF-IDF . . . . .	26
34	LDA en utilisant TF-IDF . . . . .	26

# Introduction

Le Data Mining est une technique parmi les composants de la Big Data. Il permet d'analyser les données volumineuses. Cette année on étudie cette matière, nous avons vu les notions basiques ainsi que les principaux algorithmes supervisé et non-supervisé. Pendant le TP, nous avons découvert le langage R et la construction des modèles et leurs évaluations avec le langage python et l'application web Jupyter.

Dans ce projet on va mettre en place les notions que nous avons acquies pendant les séances de cours et de TP. Dans un premier temps je vais analyser des données avec R et dans le deuxième chapitre je vais créer un modèle de classification. Et enfin, je vais aborder le sujet d'analyse de topic modeling.

# 1 Visualisation de données avec R

## 1.1 Compréhension des données

Nous avons disposé d'un fichier csv 'Cholesterol\_R'. Le jeu de données est composé de 18 enregistrements. Pour chaque enregistrement, nous disposons de 5 informations : ID, Before, After4weeks, After8weeks, Margarine.

Champs	Signification
ID	Identifiant de la personne.
Before	Le taux de cholestérol avant le commencement de régime.
After4weeks	Le taux de cholestérol après 4 semaines du commencement de régime.
After8weeks	Le taux de cholestérol après 8 semaines du commencement de régime.
Margarine	La marque de margarine utilisée pendant de régime.

TABLE 1 – Visualisation avec R - Dictionnaire de données

## 1.2 Préparation des données

J'ai choisi de visualiser la moyenne de la variation de niveau de cholestérol pour chaque marque de margarine et la variation de taux de cholestérol chez chaque personne. Donc il fallait faire un peu de calcul.

1. AJOUT DE COLONNES : J'ai ajouté deux colonnes qui représente la variation de niveau de cholestérol pour les deux cas après 4 semaines et après 8 semaines.

i..ID	Before	After4weeks	After8weeks	Margarine	VariationAfter4weeks	
1	1	6.42	5.83	5.75	B	0.59
2	2	6.76	6.20	6.13	A	0.56
3	3	6.56	5.83	5.71	B	0.73
4	4	4.80	4.27	4.15	A	0.53
5	5	8.43	7.71	7.67	B	0.72
6	6	7.49	7.12	7.05	A	0.37
VariationAfter8weeks						
1			0.67			
2			0.63			
3			0.85			
4			0.65			
5			0.76			
6			0.44			

FIGURE 1 – Visualisation avec R - Ajout de deux colonnes

2. CALCUL DES MOYENNES : J'ai calculer la moyenne de la variation de tous les gens après avoir filtrer les données selon la marque de la margarine.

```
> meanA  
[1] 0.4855556 0.5466667  
> meanB  
[1] 0.6466667 0.7111111
```

FIGURE 2 – Visualisation avec R - Calcul des moyennes

Du coup pour la margarine A, le niveau de cholestérol est diminué d'un taux 0.48 après 4 semaines et de 0.54 après 8 semaines. Pour la margarine B, il y'a une diminution de 0.64 après 4 semaines et 0.71 après 8 semaines.

### 1.3 Exploration de données

J'ai choisi le diagramme à bâton pour visualiser les données. J'ai constitué le diagramme comme suit :

- sur l'axe horizontal, la durée de traitement ;
- Sur l'axe vertical, le taux moyenne de variation de niveau de cholestérol chez les patients.

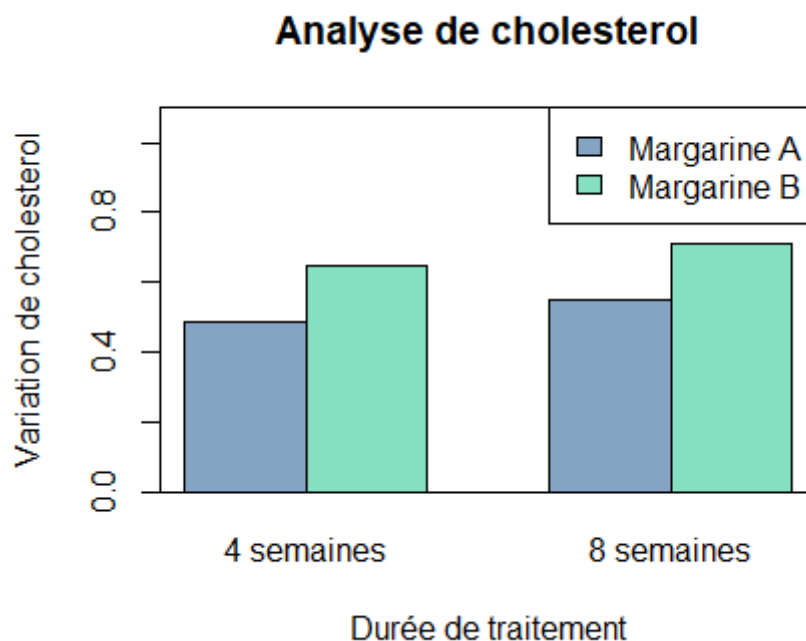


FIGURE 3 – Visualisation avec R - Diagramme à bâton

Pour visualiser le taux de cholestérol chez chaque personne dans les trois périodes j'ai choisi un diagramme linéaire.

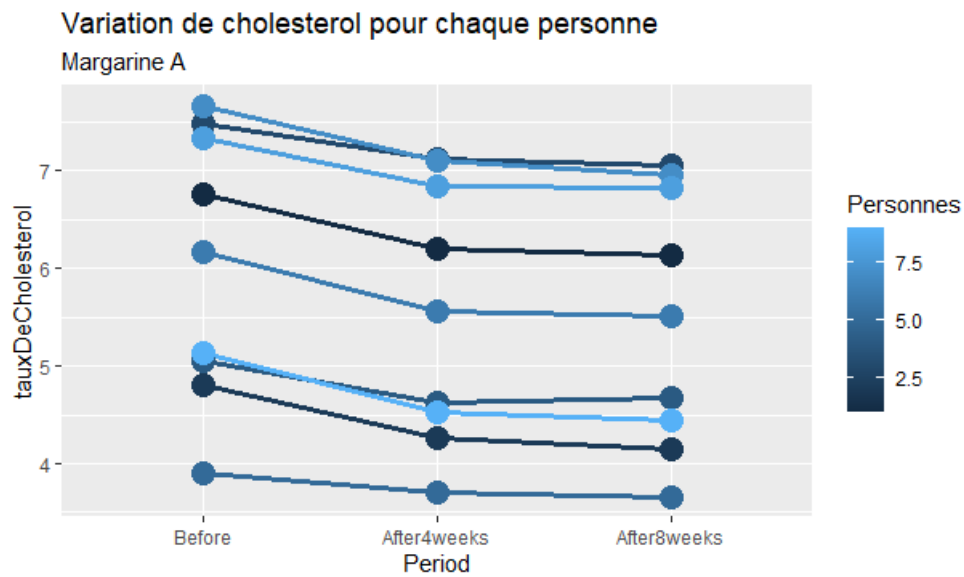


FIGURE 4 – Visualisation avec R - Diagramme linéaire - Margarine A

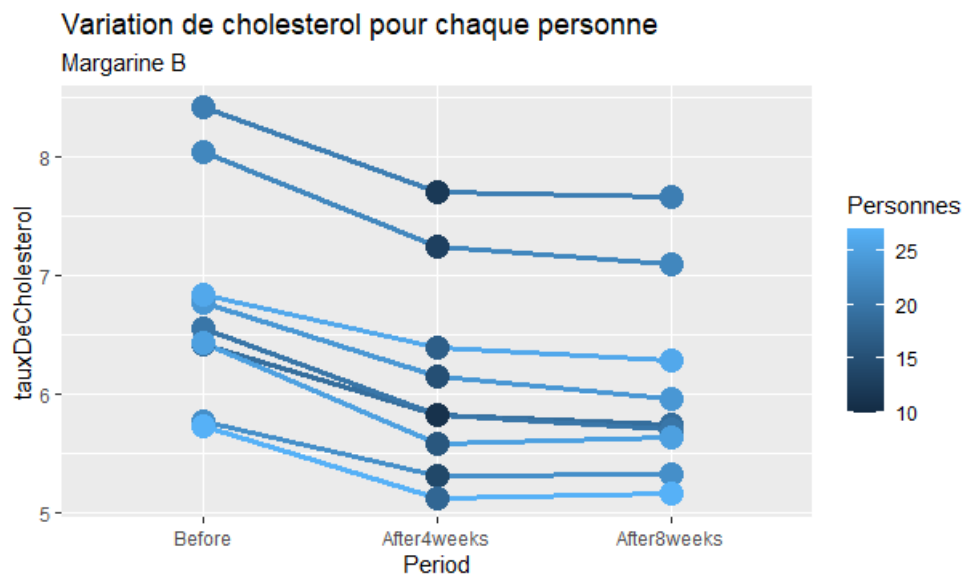


FIGURE 5 – Visualisation avec R - Diagramme linéaire - Margarine B

## 1.4 Analyse des résultats

D'après les figure précédentes, il apparaît que les deux margarine permettent de baisser le taux de cholestérol mais la margarine B permet de diminuer davantage le niveau de cholestérol chez les gens.



## 2 Classification avec Python

Le but de cette partie est de faire une classification supervisée sur les données de patients. On veut prédire si une personne a une maladie cardiaque selon plusieurs données : l'âge, la hauteur, le poids, ...

Dans ce projet je vais suivre la méthode CRISP-DM qui signifie « Cross-industry standard process for data mining ». Chaque sous-section va être une étape de cette méthode. Donc nous aurons quatre sections dans cette partie qui sont : la compréhension métier, la compréhension des données, la préparation des données, la modélisation et enfin l'évaluation.

### 2.1 Compréhension du métier

Cette première phase consiste à bien comprendre les éléments métiers et la problématique qu'on vise à résoudre. Je vais commencer par déterminer les objectifs stratégiques et opérationnels. Puis, je vais traduire l'objectif stratégique en concepts de Data mining.

#### 2.1.1 Détermination des objectifs stratégiques et opérationnelles

Selon un article publié sur le site de la fondation pour la recherche médicale <https://www.frm.org> : « On recense chaque année en France près de 70 000 décès liés à l'insuffisance cardiaque, et plus de 150 000 hospitalisations. ». Ces chiffres nous donnent une idée sur la propagation de ces maladies. Donc une étude de ce sujet sera très importante.

#### 2.1.2 Détermination des objectifs du Data Mining

Puisque l'objectif stratégique est clairement défini, il convient maintenant de le traduire en concepts de Data Mining. On va opter à une technique prédictive puisque le but est de prédire si une personne peut avoir une maladie cardiaque ou non à partir des informations présentes. Parmi les techniques prédictives, on trouve les arbres de décision et les réseaux de neurones.

### 2.2 La compréhension des données

La phase de compréhension des données de CRISP-DM implique l'étude des données disponibles pour le Data mining. On doit tout d'abord décrire les données. La deuxième étape est l'exploration des données et comme étape finale vérifier la qualité des données.

#### 2.2.1 Description des données

Pour chaque enregistrement, on a :

Code	Signification	Type
id	L'identifiant de la personne	nombre entier
age	L'âge de la personne en jours	nombre entier
gender	Le sex de la personne : 1 : Si c'est une femme 2 : Si c'est un homme	variable qualitative : nombre entier
height	La longueur de la personne en cm	nombre entier
Weight	Le poids de la personne en kg	nombre réel
ap_hi	la pression artérielle systolique	nombre entier
ap_lo	la pression sanguine diastolique	nombre entier
cholesterol	le taux de cholestérol 1 : si le taux est normal 2 : si le taux est supérieur au taux normal 3 : si le taux est très supérieur au taux normal	variable qualitative nombre entier
gluc	le taux de glucose 1 : si le taux est normal 2 : si le taux est supérieur au taux normal 3 : si le taux est très supérieur au taux normal	variable qualitative nombre entier
smoke	Est que la personne fume ? 1 : si oui 0 : sinon	variable booléenne
alco	Est que c'est une personne alcoolique ? 1 : si oui 0 : sinon	variable booléenne
active	Est que c'est la personne pratique une activité sportive ? 1 : si oui 0 : sinon	variable booléenne
cardio	Est que c'est la personne a une maladie cardiaque ? 1 : si oui 0 : sinon	variable booléenne

TABLE 2 – Dictionnaire de données

La commande describe nous permet d'afficher des informations sur les attributs numérique tels que la valeur moyenne, la valeur maximale, la valeur minimale, ...

	id	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	smoke
count	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000	70000.000000
mean	49972.419900	19468.865814	1.349571	164.359229	74.205690	128.817286	96.630414	1.366871	1.226457	0.088129
std	28851.302323	2467.251667	0.476838	8.210126	14.395757	154.011419	188.472530	0.680250	0.572270	0.283484
min	0.000000	10798.000000	1.000000	55.000000	10.000000	-150.000000	-70.000000	1.000000	1.000000	0.000000
25%	25006.750000	17664.000000	1.000000	159.000000	65.000000	120.000000	80.000000	1.000000	1.000000	0.000000
50%	50001.500000	19703.000000	1.000000	165.000000	72.000000	120.000000	80.000000	1.000000	1.000000	0.000000
75%	74889.250000	21327.000000	2.000000	170.000000	82.000000	140.000000	90.000000	2.000000	1.000000	0.000000
max	99999.000000	23713.000000	2.000000	250.000000	200.000000	16020.000000	11000.000000	3.000000	3.000000	1.000000

FIGURE 6 – Description des attributs

Cette matrice nous donne plusieurs informations utiles. Le plus jeune a 29 ans et la plus maigre personne a 10kg seulement!! Ce qui apparaît illogique.

## 2.2.2 Exploration des données

- **La matrice de corrélation** : La matrice de corrélation permet de trouver les features fortement corrélées pour les supprimer en fin d'augmenter la robustesse et la stabilité du modèle et aussi simplifier les résultats obtenus par le modèle.

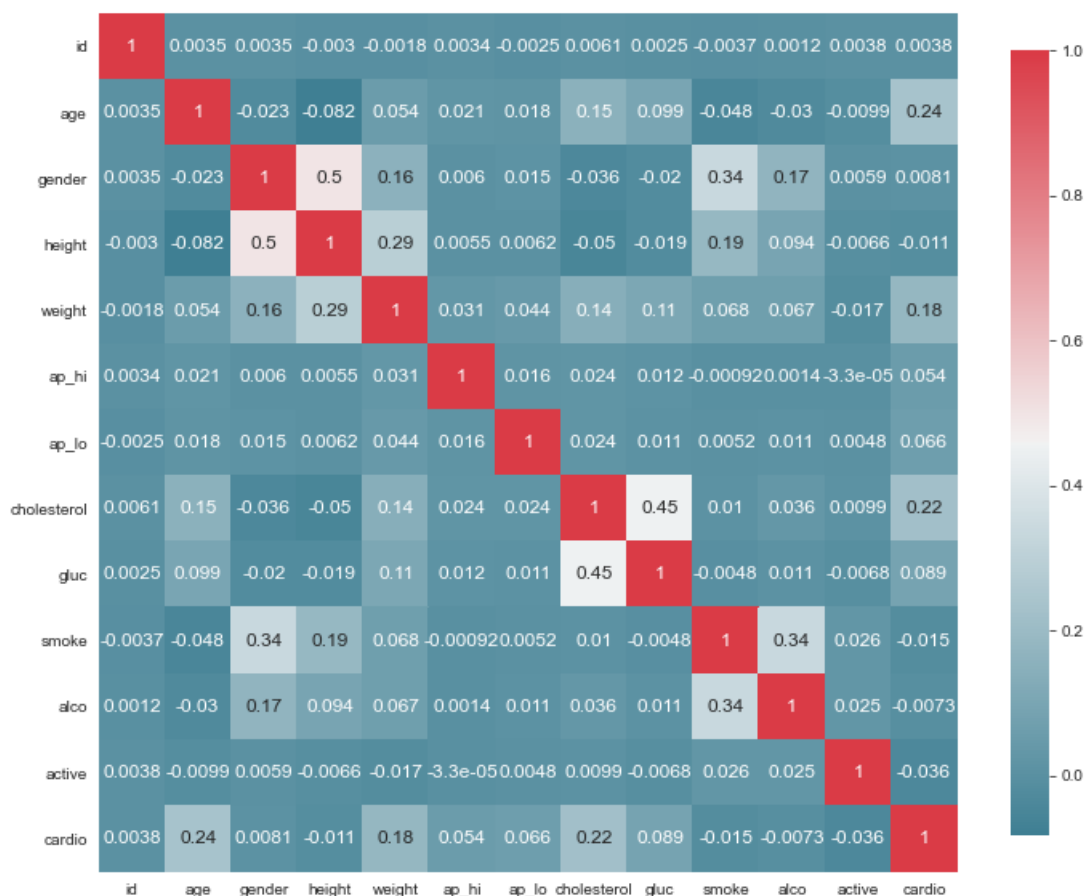


FIGURE 7 – La matrice de corrélation

D'après cette matrice, les deux variables height and gender sont très corrélées. Du coup on peut supprimer l'un des deux. Mais la variable gender est plus significative que la variable height. Donc on peut dispenser de cette variable pour créer notre modèle.

- **Diagramme circulaire représentant le sexe des malades cardiaques :**

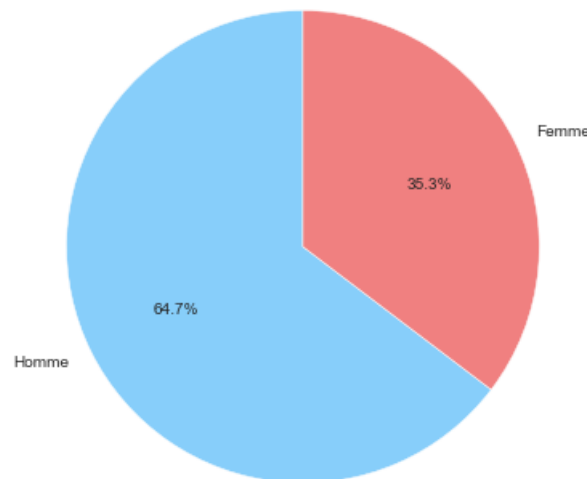


FIGURE 8 – le sexe des malades cardiaques

Ce diagramme montre que les femmes sont plus touchées aux maladies cardiaques que les hommes.

- **Diagramme en barres de taux de cholestérol et glucose des malades cardiaques :**

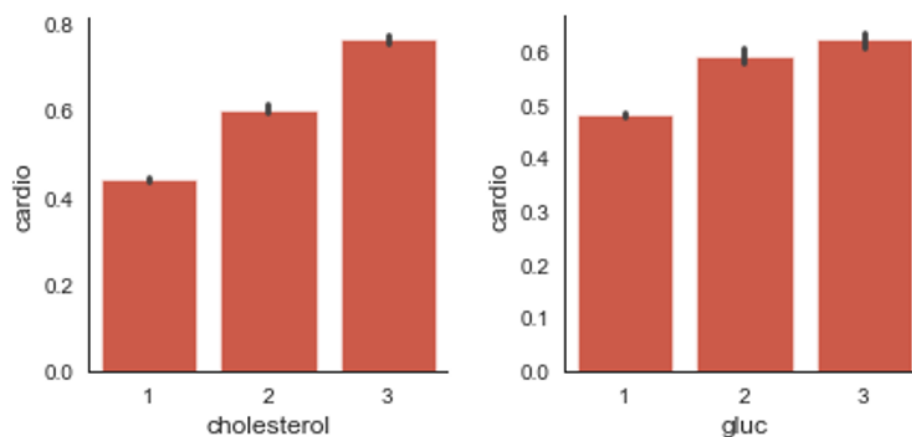


FIGURE 9 – le taux de cholestérol et glucose des malades cardiaques

Cette figure montre que le cholestérol est une bonne variable de prédiction puisque la probabilité d'avoir des maladies cardiaques varie largement selon le niveau de cholestérol.

— Histogramme d'âge des gens qui ont des maladies cardiaques :

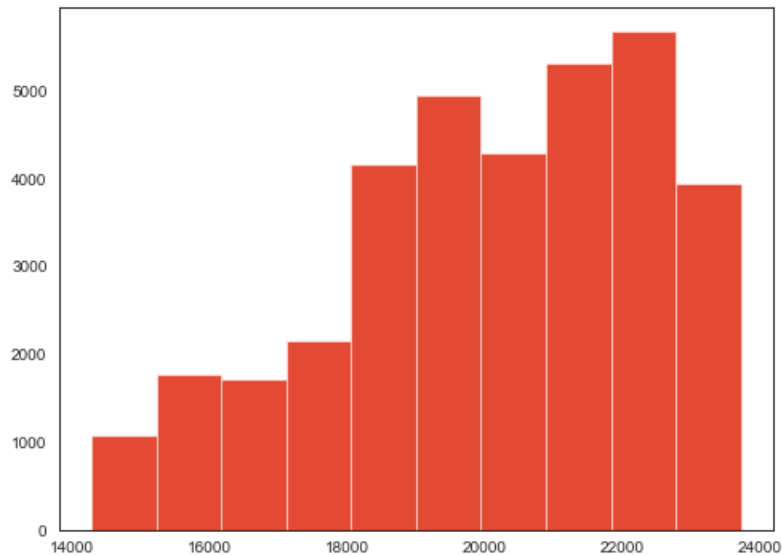


FIGURE 10 – l'âge des malades cardiaques

les gens âgées sont plus touchées par les maladies cardiaques

- **Équilibrage de jeu de données** : Pour construire un modèle robuste, les classes de la variable cible doit être équilibrées. Du coup j'ai exploré cette information pour s'assurer que les données ne sont pas biaisées.

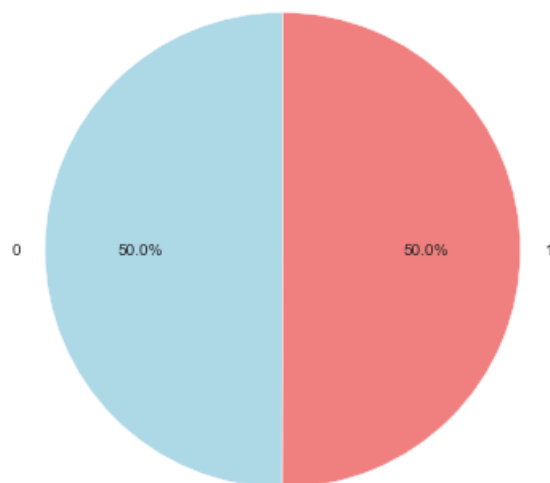


FIGURE 11 – Équilibrage du jeu de données

Heureusement les données sont parfaitement équilibrées.

## 2.3 Préparation des données

La préparation des données est l'un des aspects les plus importante et les plus coûteux en temps du Data Mining.

### Nettoyage des données

#### – Valeurs manquantes :

En utilisant la commande `isnull`, on peut savoir le nombre des valeurs manquantes dans chaque caractéristiques.

D'après la figure suivante, on peut déduire qu'il y a aucune valeur manquante dans le jeu de données.

```
id          0
age         0
gender      0
height      0
weight      0
ap_hi       0
ap_lo       0
cholesterol 0
gluc        0
smoke       0
alco        0
active      0
cardio      0
dtype: int64
```

FIGURE 12 – Valeurs manquantes

#### – Valeurs erronées :

Comme déjà vu dans la section 2.2.1, j'ai remarqué une valeur anormale de poids

	id	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	smoke	alco	active	cardio
57858	82567	18804	2	165	10.0	180	1100	2	2	0	0	1	1

FIGURE 13 – Valeur anormale - poids

une personne de 51ans avec seulement 10kg est une chose normale. Puisque l'âge moyen de la population est 53ans donc on peut remplacer la valeur erronée par la valeur moyenne des poids 74kg.

En analysant la pression artérielle systolique(`ap_hi`), j'ai remarqué qu'il y'a des valeurs anormales (1000, 11,, 14, 12, -150,160020). Donc j'ai décidé de visualiser le nuage des points cette caractéristique.

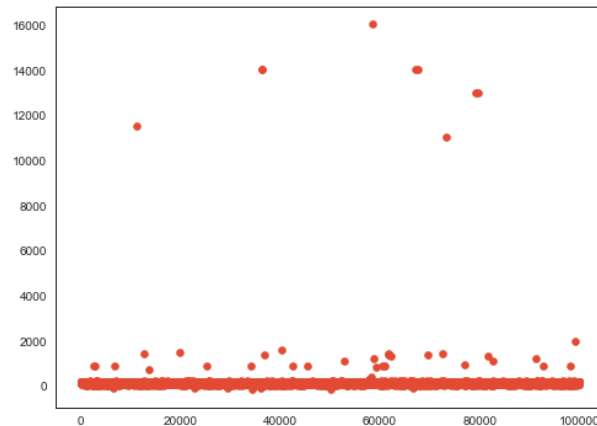


FIGURE 14 – Valeurs anormales - ap\_hi

D'après la figure, il y'a 6 valeurs anormales. Puisque sont peu nombreux j'ai décidé de les supprimer.

De même pour la pression sanguine diastolique.

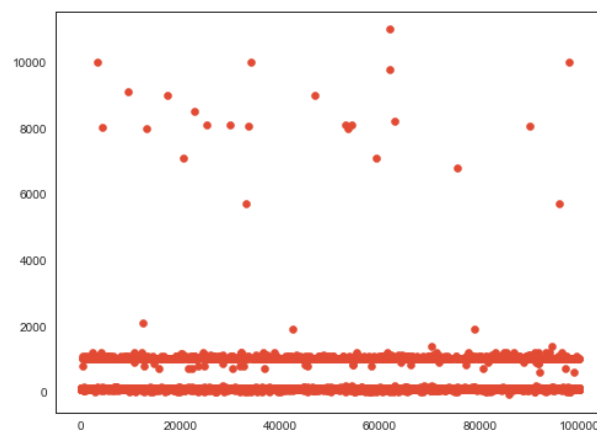


FIGURE 15 – Valeurs anormales - ap\_lo

#### — Nettoyage des données :

J'ai remplacé la valeur erronée du poids par la moyenne des poids.

	id	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	smoke	alco	active	cardio
57858	82567	18804	2	165	74.206607	180	1100	2	2	0	0	1	1

FIGURE 16 – Nettoyage - remplacement de poids anormale par la moyenne

J'ai supprimé les enregistrements qui contiennent des valeurs anormales de ap\_hi et ap\_lo.

	id	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	smoke	
count	69966.000000	69966.000000	69966.000000	69966.000000	69966.000000	69966.000000	69966.000000	69966.000000	69966.000000	69966.000000	6
mean	49974.100163	19469.001143	1.349527	164.357974	74.207008	127.108867	93.783738	1.366907	1.226424	0.088143	
std	28852.056580	2467.426418	0.476824	8.209866	14.394833	28.278689	108.322281	0.680298	0.572244	0.283505	
min	0.000000	10798.000000	1.000000	55.000000	11.000000	-150.000000	-70.000000	1.000000	1.000000	0.000000	
25%	25007.250000	17664.000000	1.000000	159.000000	65.000000	120.000000	80.000000	1.000000	1.000000	0.000000	
50%	50001.500000	19703.000000	1.000000	165.000000	72.000000	120.000000	80.000000	1.000000	1.000000	0.000000	
75%	74891.750000	21327.000000	2.000000	170.000000	82.000000	140.000000	90.000000	2.000000	1.000000	0.000000	
max	99999.000000	23713.000000	2.000000	250.000000	200.000000	2000.000000	1900.000000	3.000000	3.000000	1.000000	

FIGURE 17 – Nettoyage - suppression des lignes avec valeurs anormales

### Sélection des données

Après avoir effectué le nettoyage des données on doit maintenant choisir les données pertinentes pour nos objectifs de Data Mining.

J'ai éliminé une seule variable qui est la longueur puisqu'elle a une grande dépendance avec la variable genre.

	age	gender	weight	ap_hi	ap_lo	cholesterol	gluc	smoke	alco	active	cardio
0	18393	2	62.0	110	80	1	1	0	0	1	0
1	20228	1	85.0	140	90	3	1	0	0	1	1
2	18857	1	64.0	130	70	3	1	0	0	0	1
3	17623	2	82.0	150	100	1	1	0	0	1	1
4	17474	1	56.0	100	60	1	1	0	0	0	0

FIGURE 18 – Sélection des données pour la modélisation

J'ai pris 70% des enregistrements pour l'apprentissage et 30% pour la validation.

## 2.4 Modélisation

La quatrième étape est la modélisation qui est au cœur de tout projet d'apprentissage automatique.

Cette étape est responsable des résultats qui devraient satisfaire ou aider à atteindre les objectifs du projet. Bien que ce soit la partie glamour du projet, c'est aussi la plus courte dans le temps.

### 2.4.1 Sélection des techniques de modélisation

Il y a plusieurs types de modélisation. Le choix du modèle le plus adéquat sera généralement basé sur les critères suivants :

- Les types de données disponibles pour l'exploration : Les champs sont numériques.



- Les objectifs de data mining : Notre problème est un problème de classification supervisé. Nous avons plusieurs méthodes de classification supervisé : k plus proches voisins, Les forêt aléatoires, support vector machine, Régression logistique, ... Nous avons vu quelques méthodes dans le cours et d'autres dans le TP.

On va entraîner tous ces modèles et on va choisir celle qui donne la meilleure accuracy.

#### 2.4.2 Explication des techniques de modélisation

- Random Forests : C'est un algorithme qui effectue l'apprentissage sur plusieurs arbres de décision entraînés sur des sous ensembles du jeu de données différents.
- Support Vector Machines : Ce sont une famille d'algorithmes permettant de résoudre des problèmes de classification, de régression et de détection d'anomalie. Ils ont pour but de séparer les données en classes à l'aide d'une frontière de telle façon que la distance entre les différents groupes de données et la frontière qui les sépare soit maximale.
- Gradient Boosting : est une technique qui permet de résoudre des problèmes de classification et de régression. Elle donne un poids à chaque individu et au fur à mesure les poids sont corrigés.
- KNN : L'algorithme de plus proche voisins permet de construire des ensemble d'individus à partir des données d'apprentissage et pour chaque nouveau individu on cherche la plus proche classe pour lui associe cette classe.
- Gaussian naive bayes : C'est un classifieur linéaire basé sur le théorème de Bayes avec une forte indépendance (dite naïve) des hypothèses.
- Régression logistique : C'est un modèle de régression binomiale permet d'associer à un vecteur de variables aléatoires une variable aléatoire binomiale.

#### 2.4.3 Construire le modèle

En utilisant la fonction fit, j'ai pu entraîner tous les modèles.

### 3.2.1 Random Forests Model

```
Entrée [31]: modelRFC.fit( train_data , train_data_ )  
Out[31]: RandomForestClassifier()
```

### 3.2.2 Support Vector Machines

```
Entrée [32]: modelSVC.fit( train_data , train_data_ )  
Out[32]: SVC()
```

### 3.2.3 Gradient Boosting Classifier

```
Entrée [33]: modelGBC.fit( train_data , train_data_ )  
Out[33]: GradientBoostingClassifier()
```

### 3.2.4 K-nearest neighbors

```
Entrée [34]: modelKNC.fit( train_data , train_data_ )  
Out[34]: KNeighborsClassifier(n_neighbors=3)
```

### 3.2.5 Gaussian Naive Bayes

```
Entrée [35]: modelGNB.fit( train_data , train_data_ )  
Out[35]: GaussianNB()
```

### 3.2.6 Logistic Regression

```
Entrée [36]: modelLR.fit( train_data , train_data_ )  
Out[36]: LogisticRegression()
```

FIGURE 19 – la construction des modèles

## 2.5 Evaluation

### 2.5.1 Comparaison des modèles

J'ai utilisé 70% pour faire l'apprentissage du modèle et le 30% des données on va l'utiliser pour valider le modèle. On a visualiser l'accuracy de tous les modèle pour qu'on puisse faire une comparaison et choisir le meilleur modèle.

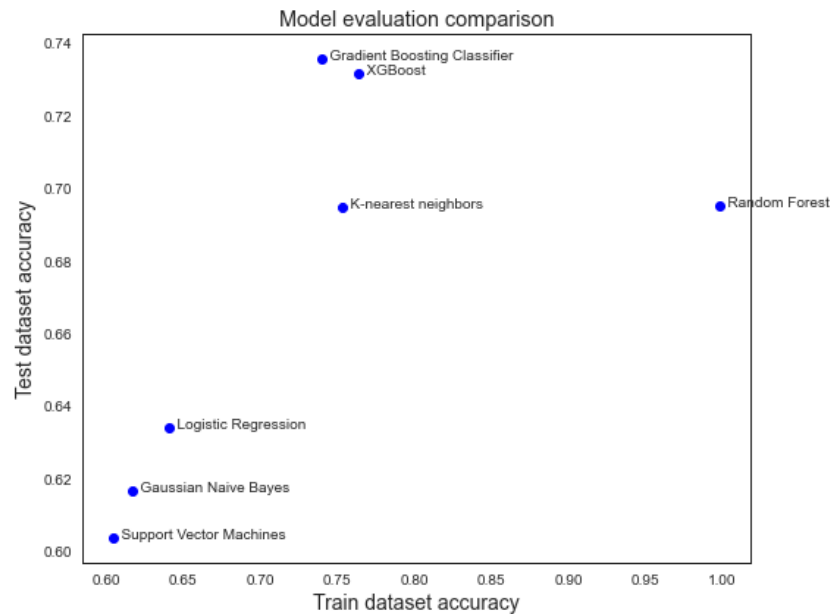


FIGURE 20 – Comparaison des modèles

Le modèle qui donne une valeur d'accuracy la plus haute est Gradient boosting Classifier.

### 2.5.2 Amélioration

J'ai essayé plusieurs paramètres pour les deux algorithmes KNN et random forest et aussi un autre algorithme qui est XGBoost. Et dans tous les cas c'était l'algorithme Gradient boosting Classifier qu'avait toujours le meilleurs score.

### 2.5.3 Importance des variables

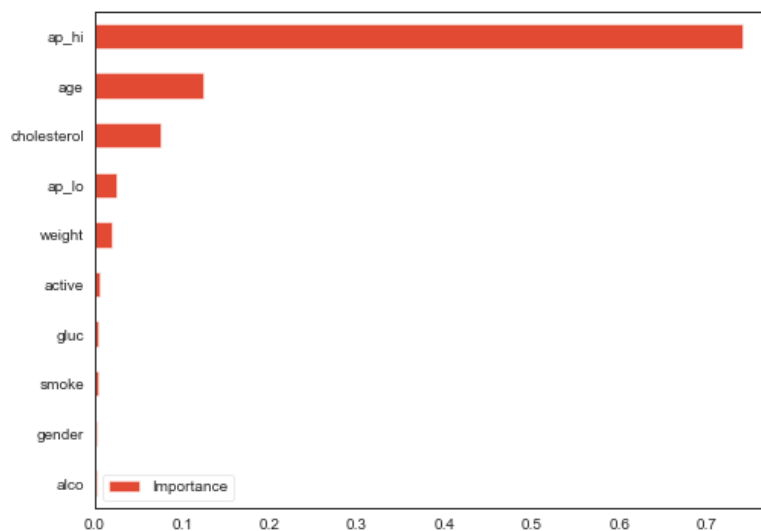


FIGURE 21 – Importance des variables



C'est la pression artérielle systolique qui joue un grand rôle pour prédire la variable cible selon le modèle de Gradient boosting Classifier.

## **2.6 Déploiement**

C'est la dernière étape de CRISP-DM. J'ai enregistré les résultats de la validation (identifiant de la personne et la valeur de la variable cible).

### 3 Topic Modeling

Dans cette partie, on va faire le topic modeling qui consiste à savoir le sujet ou le thème abstrait dans un document.

#### 3.1 Prétraitement des données

On va suivre les huit étapes vu dans les séances de TP depuis la Tokenisation jusqu'à Chunking. Je vais utiliser seulement la colonne texte puisque c'est elle qui donne plus d'informations.

##### 3.1.1 Tokenisation

La tokenisation est le process de découper un texte en tokens(mots). C'est la fonction `word_tokenize()` qui permet de faire cette tâche.

```
4250624
Out[30]: ['If',
          'I',
          'smelled',
          'the',
          'scent',
          'of',
          'hand',
          'sanitizers',
          'today',
          'on',
          'someone',
          'in',
          'the',
          'past',
          ',',
          'I',
          'would',
          'think',
          'they']
```

FIGURE 22 – Tokenisation

Le document contient 4250624 mots.

##### 3.1.2 Fréquence des mots

`FreqDist` nous permet d'avoir la fréquence de chaque mots et la fonction `most_common` donne les plus fréquents mots dans le document.

```
Out[33]: [(' ', 266985),  
          (':', 208423),  
          ('https', 177118),  
          ('the', 87362),  
          ('@', 85902),  
          (',', 80194),  
          ('COVID19', 77745),  
          ('.', 76884),  
          ('to', 71257),  
          ('of', 57334)]
```

FIGURE 23 – Fréquence des mots

Il paraît que parmi les tokens les fréquents on trouve : « : », « the », « to » et « of ». D'où vient l'importance de l'étape de stop words.

### 3.1.3 Stemming : Normalisation des mots

Dans cette étape on associe chaque mot à sa racine et on supprime les terminologies. Il existe plusieurs méthodes pour le stemming. J'ai choisi la fonction stem.

```
Out[38]: ['If',  
          'I',  
          'smell',  
          'the',  
          'scent',  
          'of',  
          'hand',  
          'sanit',  
          'today',  
          'on',  
          'someone',  
          'in',  
          'the',  
          'past',  
          ',',  
          'I',  
          'would',  
          'think',  
          'they',  
          'were',  
          '.']
```

FIGURE 24 – Stemming

Par exemple, on peut voir que le mot « smelled » est remplacé par « smell ».

### 3.1.4 Lemmatisation : Base des mots

La lemmatisation est par définition une action consistant à l'analyse lexicale d'un texte avec pour but de regrouper les mots d'une même famille. La lemmatisation prend en compte : le genre (masculin ou féminin), le nombre (singulier ou pluriel), la possession (moi, toi, eux...) et le mode (indicatif, impératif...). La fonction lemmatize permet de faire ce travail.

### 3.1.5 Stop words

Dans cette étape on supprime les mots inutiles comme par exemple : « the », « a », « : ». Mais avant j'ai mis tous les mots en minuscule.

```
Out[75]: ['smell',  
          'scent',  
          'hand',  
          'sanit',  
          'today',  
          'someone',  
          'past',  
          'would',  
          'think',  
          'intox',  
          'that_',  
          '//t.co/qzvybrogb0',  
          'hey',  
          'yanke',  
          'yankeespr',  
          'mlb',  
          'would',  
          "n't",  
          'made',  
          'sen',
```

FIGURE 25 – Stop words

### 3.1.6 Part of speech tagging (POS)

Dans cette partie on attribue à chaque mot sa nature (un nom, un verbe, un adjectif,...).

```
[('smell', 'NN')]  
[('scent', 'NN')]  
[('hand', 'NN')]  
[('sanit', 'NN')]  
[('today', 'NN')]  
[('someone', 'NN')]  
[('past', 'NN')]  
[('would', 'MD')]  
[('think', 'NN')]  
[('intox', 'NN')]  
[('that_', 'NN')]  
[('//t.co/qzvybrogb0', 'NN')]  
[('hey', 'NN')]  
[('yanke', 'NN')]  
[('yankeespr', 'NN')]  
[('mlb', 'NN')]  
[('would', 'MD')]  
[('n't', 'RB')]  
[('made', 'VBN')]  
[('sen', 'NN')]
```

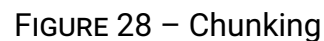
FIGURE 26 – Part of speech tagging (POS)

### 3.1.7 Reconnaissance des différentes entités

Cette étape consiste à organiser les mots dans des classes telles que noms de personnes, noms d'organisations ou d'entreprises, noms de lieux, quantités, distances, valeurs, dates, etc.



Le Chunking est un processus d'extraction de phrases à partir d'un texte non structuré, ce qui signifie analyser une phrase pour identifier les constituants (groupes de noms, verbes, groupes de verbes, etc.).



Le nuage de mots est une méthode de visualisation qui permet d'avoir les mots les plus fréquents dans un document. Plus le mot est fréquent plus il est de grand taille.





Cette représentation nous donne une idée sur les mots les plus utilisés dans les tweets sur covid. On trouve : covid19, covid, pandem, mask, ...

### 3.3 Topic modeling avec LDA

LDA est une méthode non supervisée de Natural Language Processing. Il permet de présenter un documents sous forme de thèmes pour le synthétiser.

#### 3.3.1 La fréquence de mots

La première étape consiste à construire un dictionnaire où la clé est la fréquence du mot et la valeur est le mot.

```
smell 155  
scent 13  
hand 1189  
sanit 479  
today 5227  
someon 842  
past 672  
would 2491  
think 2422  
intox 2  
that... 279
```

FIGURE 30 – La fréquence des mots

#### 3.3.2 Gensim filters

J'ai gardé seulement les 100000 mots plus fréquents qui apparaît dans moins de 15 documents et plus que un seul document.

#### 3.3.3 Gensim doc2bow

Dans cette phase, on attribue à chaque mot dans un document sa fréquence.

```
Word 8 ("would") appears 1 time.  
Word 9 ("made") appears 1 time.  
Word 10 ("player") appears 1 time.  
Word 11 ("respect") appears 1 time.  
Word 12 ("yanke") appears 1 time.
```

FIGURE 31 – Gensim doc2bow

Ce modèle a généré 9 ensembles de document. Pour chaque ensemble j'ai affiché la probabilité de TF-IDF.

### 3.3.4 Running LDA using Bag of Words

J'ai choisi d'extraire 6 topics et pour chaque topic d'afficher les 10 plus fréquents mots-clés. Un nuage de mots a donné le résultat suivant :

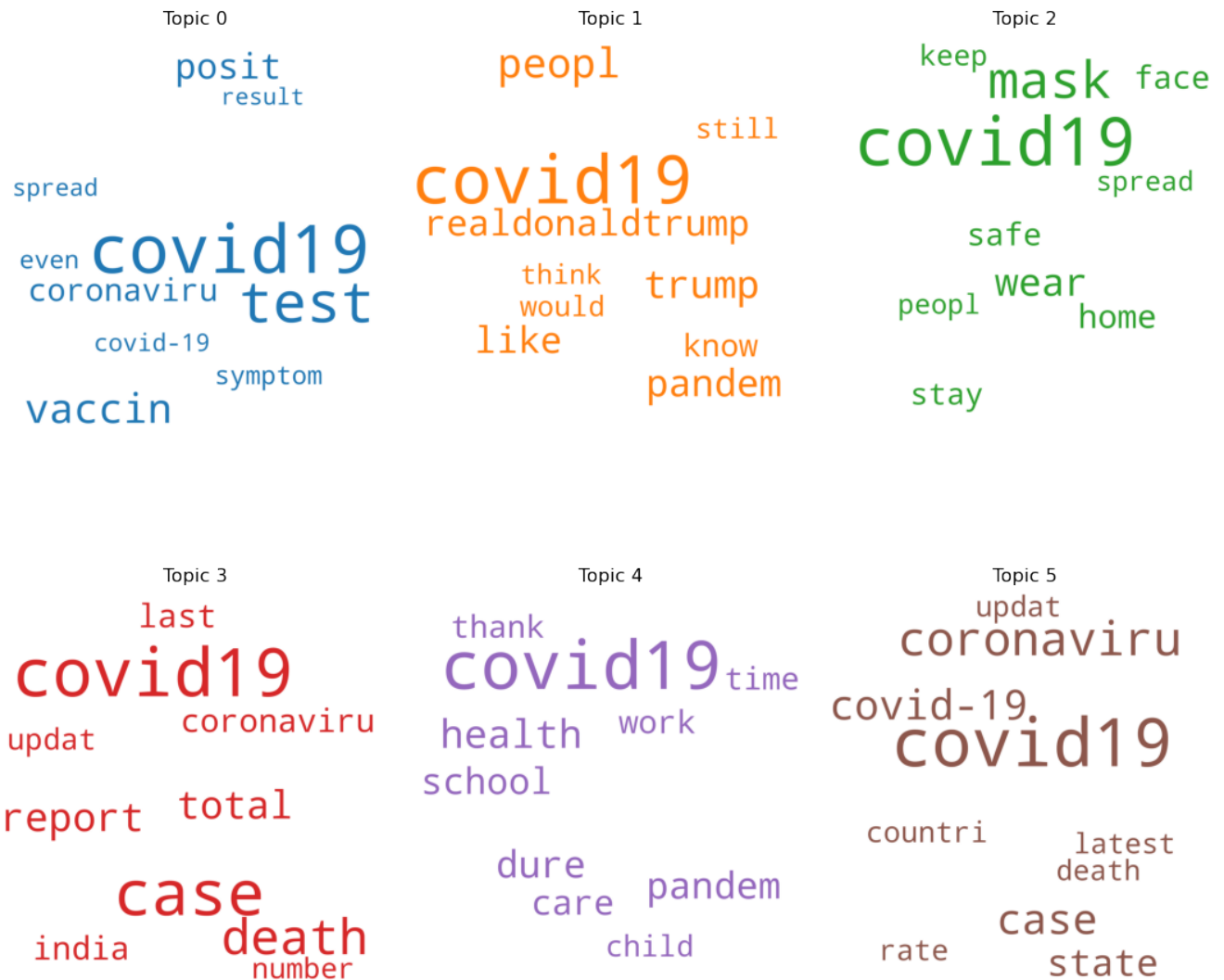


FIGURE 32 – Running LDA using Bag of Words

### 3.3.5 TF-IDF

Terme Frequency - Inverse Document Frequency (TF-IDF). Cette méthode consiste à rassembler un ensemble de documents partageant les mêmes mots clés. Deux documents ont plus de chance d'être dans le même ensemble s'ils possèdent des mots-clés en une grandes fréquence, et que ces mots sont rares dans d'autres documents.

```
[ (0, 0.3120066919227304),
  (1, 0.34386088345342114),
  (2, 0.36804123747488343),
  (3, 0.43730975990531323),
  (4, 0.33144355605282877),
  (5, 0.3978718001045878),
  (6, 0.2666917498141616),
  (7, 0.21853705022170467),
  (8, 0.26585967117591836) ]
```

FIGURE 33 – TF-IDF

Ce modèle a généré 9 ensembles de document. Pour chaque ensemble j'ai affiché la probabilité de TF-IDF.

### 3.3.6 Running LDA using TF-IDF

De même pour TF-IDF, j'ai sélectionné 6 topics et affiché 10 mots pour chaque topic.

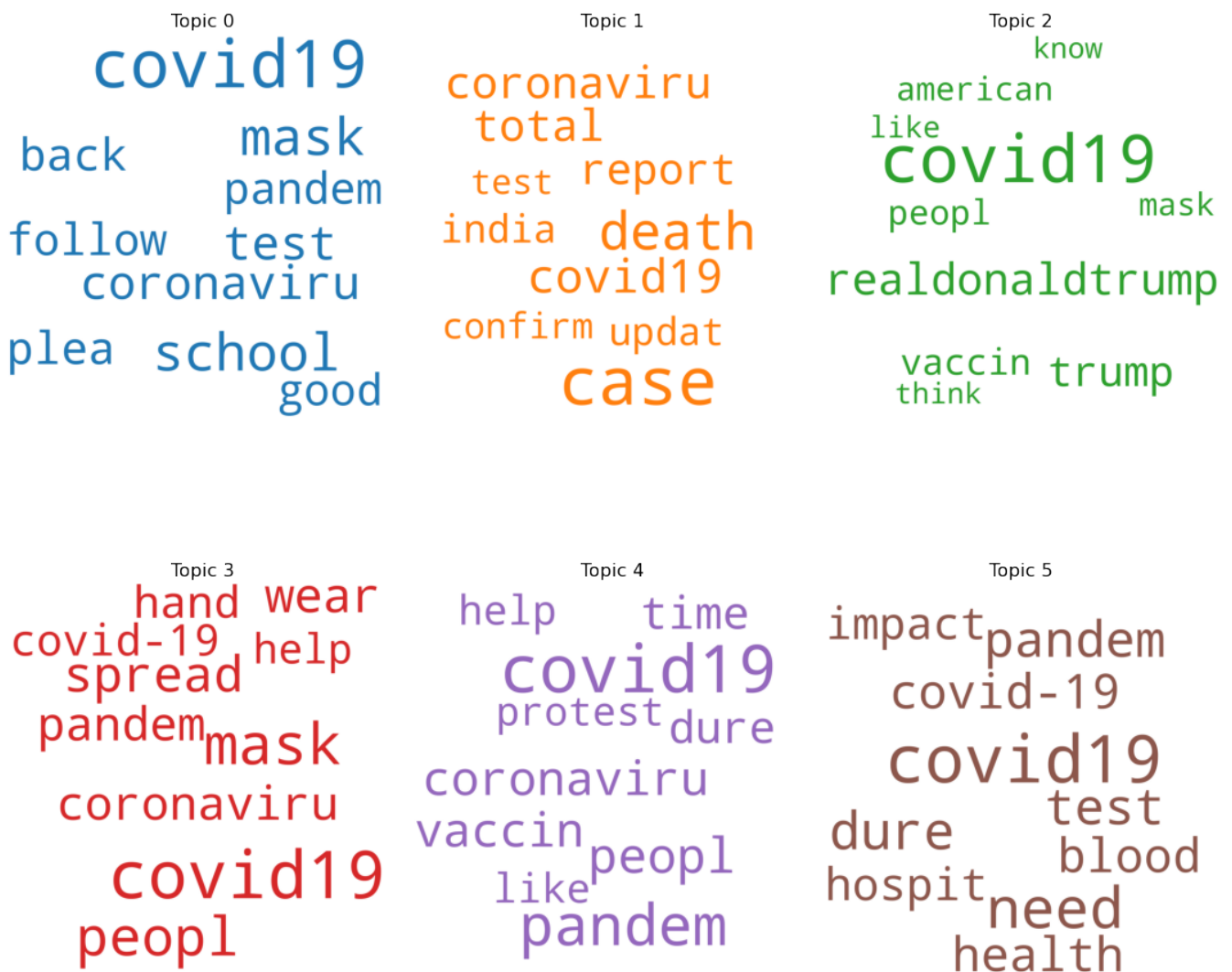


FIGURE 34 – LDA en utilisant TF-IDF



### 3.3.7 Evaluation du modèle

J'ai évalué le même document avec les deux types de modèle. Et j'ai remarqué que les deux algorithmes ont associé le document à deux topics différents.

## Conclusion

A travers ce projet, j'ai pu en effet consolider mes connaissances acquises et d'enrichir mon expérience en matière de data mining. A travers les trois sujets traités, on constate que le data mining peut être appliqué dans tous les domaines et s'avère utile pour tous les décideurs.

Pour atteindre notre objectif j'ai mis en place la démarche CRISP-DM ( Cross Industry Standard Process for Data Mining ). Il s'agit d'un modèle de processus de data mining qui décrit une approche communément utilisée par les experts en data mining pour résoudre les problèmes qui se posent à eux et les algorithmes de NLP.