

Projet L3

Ingénierie des langues

Fouille de données

GOEHRYS Martial
16711476

19 mars 2024

Table des matières

1	Phase 1 : Récupération des données	5
1.1	Récolte des données	5
1.2	Pré-traitement	6
1.2.1	Importation	6
1.2.2	Extraction des corps des mails	6
1.2.3	Nettoyage	7
1.2.4	Mise en base	12
1.3	Données de la phase 1	17
2	Phase 2 : Traitement des données	22
2.1	Génération des données statistiques	22
2.2	Sortie de la partie 2 : Stats d'exploitation	24
2.3	Analyses statistiques	25
2.4	Traitement du langage	32
2.4.1	Lemmatisation	35
2.4.2	Vectorisation TF-IDF	38
3	Phase 3	43
A	Développement visualisation distribution de Zipf	44
B	Déploiement des bases de données	50
B.1	ElasticSearch	50
B.1.1	Conteneurisation	50
B.1.2	Initialisation de l'index	54
B.2	PostgreSQL	55
B.2.1	Conteneurisation	55
B.2.2	Initialisation de la base de données	56
C	Sorties PSQL et la liste de mots	59
D	Tableau des choix technologiques	73

E	Modèles	74
E.1	Naïves Bayes	74
F	Bibliographie	79
G	Sitotec	79
G.1	Corpus	79
G.2	Modules	79
G.3	Modèles	80
H	Codes sources	80
H.1	Github	80
H.2	Analyse statistiques de la phase 1	80

Introduction

Ce projet a pour but de développer un modèle permettant de catégoriser des emails en spam ou ham. La définition d'un spam dans le dictionnaire *Larousse* est :

"Courrier électronique non sollicité envoyé en grand nombre à des boîtes aux lettres électroniques ou à des forums, dans un but publicitaire ou commercial."

Il est possible d'ajouter à cette catégorie tous les mails indésirables comme les tentatives d'hameçonnage permettant de soutirer des informations personnelles à une cible.

L'objectif est de travailler uniquement sur les données textuelles issues du corps du mail. Nous avons donc comme point de départ les éléments suivants :

- langue : anglais
- corpus : monolingue écrit
- type : e-mail

Déroulé Le développement de ce projet s'articule autour de 3 phases majeures

- Phase 1 : Récupération des données
- Phase 2 : Analyse des caractéristiques
- Phase 3 : Construction du modèle

Phase 1 La phase 1 concerne la récolte des informations et les traitements minimums nécessaires pour la mise en base. Les objectifs de traitement de cette phase sont :

- Extraire les corps des mails et éliminer les méta-données superflues
- Éliminer les mails non anglais
- Éliminer les mails en doublons
- Éliminer les parties de textes non pertinentes (liens, réponses, certaines ponctuations)

Cette phase se termine avec la mise en base des documents dans un index Elasticsearch et le stockage du nombre de liens dans une base PSQL.

Phase 2 La phase 2 vise à générer les données statistiques et numériques à partir des corps de mails. A l'issue de cette phase, il est alors possible d'analyser manuellement les statistiques générés.

Phase 3 La phase 3 regroupe tous les opérations d'exploitation des données et vise à développer et à créer un modèle de classement des mails et d'en évaluer les performances.

Afin de conserver une certaine cohérence dans le déroulé entre les phases et au vu du temps que j'ai pris pour réaliser ce projet, l'ensemble des étapes est automatisé avec Python. Seule la récolte initiale des mails a été réalisée à la main.

Le schéma ci-dessous donne une vue synthétique des étapes du projet



FIGURE 1 – Schéma des grandes étapes

1 Phase 1 : Récupération des données

La suite d'opération de cette phase vas permettre d'extraire un maximum d'information d'un email en essayant de ne dénaturer ni le fond ni la forme. Il pourra ensuite être stocké avec sa catégorie d'appartenance. Durant cette phase nous allons également initialiser les bases de données en créant les index (ES) et les tables (PSQL et SQLITE).

Ci-dessous le schéma général de cette phase.

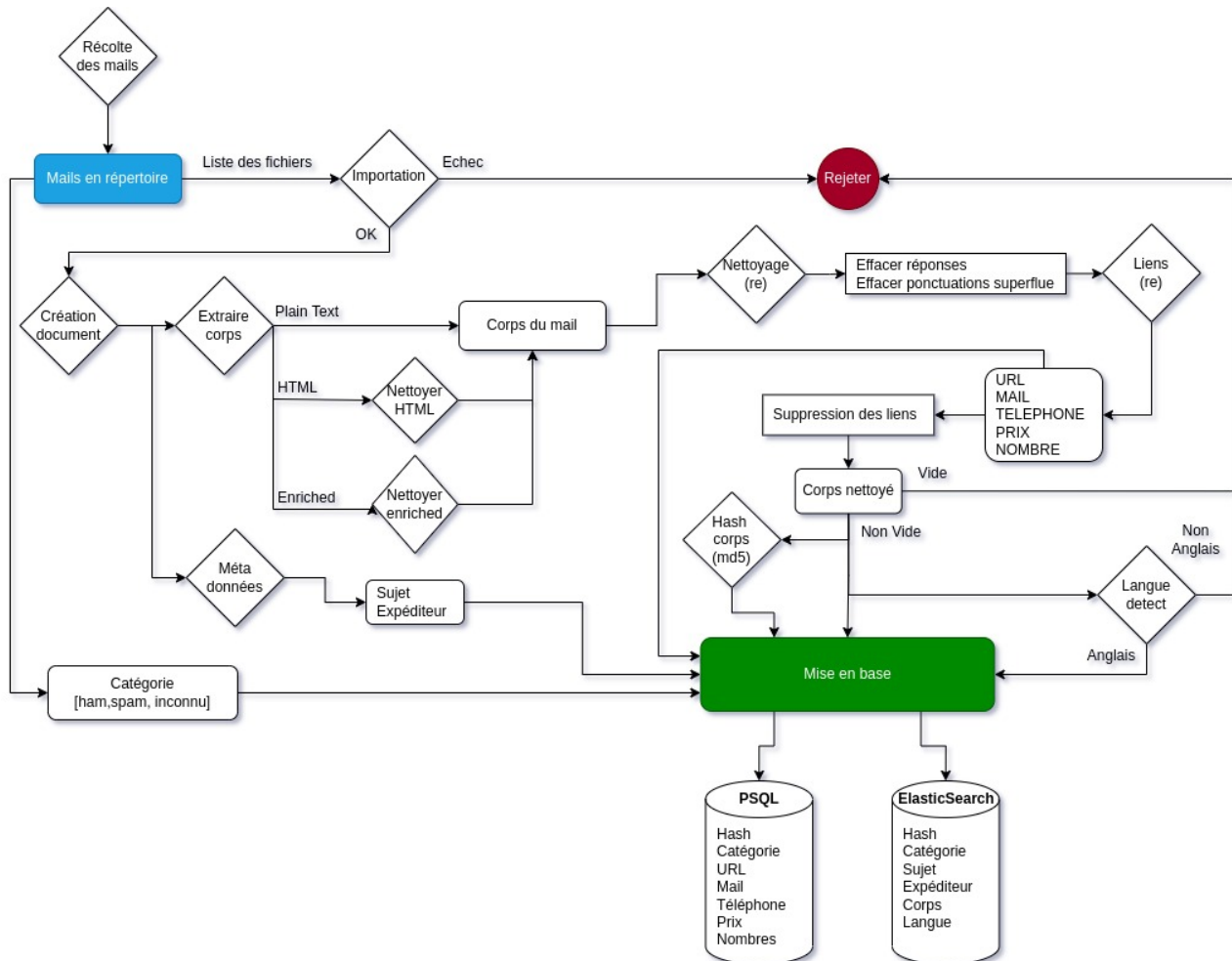


FIGURE 2 – Schéma des étapes de la phase 1

1.1 Récolte des données

Recherche de dataset Ma volonté initiale était de travailler sur des mails en français. Cependant, je n'ai pas trouvé de dataset dans cette langue. Je me suis donc retourné vers les dataset de mails en anglais.

J'ai pu alors récupérer deux dataset :

- Enron company mails (voir : G.1)
- Dataset SpamAssassin (voir : G.1)

Les mails de SpamAssassin ont l'avantage d'être pré-trié, contrairement aux mails de la compagnie Enron. Ainsi le développement du moteur se fera uniquement avec les mails du SpamAssassin afin de pouvoir vérifier les résultats de l'analyse.

Téléchargement des données Le téléchargement du dataset Enron est possible a partir du moment où l'on possède un compte sur la plateforme Kaggle. Le dataset SpamAssassin est

ouvert, il suffit de télécharger les archives de chaque catégorie.

La récolte des données a été réalisée à la main sans automatisation. Les mails sont alors stockés dans plusieurs répertoires *HAM* et *SPAM* selon leur catégorie.

Format :

- *Enron* - 1 fichier CSV avec tous les mails
- *SpamAssassin* - 1 fichier texte par mail

1.2 Pré-traitement

Les étapes de pré-traitement regroupent toutes les étapes et actions réalisées avant la mise en base. L'objectif de ces étapes est d'extraire le message en retirant les métadonnées du mail. Il va être possible d'effectuer certain traitement de nettoyage et de récupération d'informations primaires.

Les manipulations de messages dans Python se font principalement à l'aide du module python natif *email*.

1.2.1 Importation

La fonction *email.message_from_binary_file* permet de transformer un fichier mail en objet python manipulable :

Fonction d'importation des fichiers

```
1 def import_from_file(chemin):
2     try:
3         with open(chemin, 'rb') as data:
4             msg = message_from_binary_file(data, policy=policy.default)
5             return msg
6
7     except FileNotFoundError:
8         print("Fichier : '{}' non trouve".format(chemin), file=sys.stderr)
9         return None
```

1.2.2 Extraction des corps des mails

Une fois le fichier importé au format *EmailMessage*, il est possible d'en extraire le corps.

Le corps du mail peut être composé de plusieurs parties qui ne sont pas forcément du texte. Les parties non textuelles ne sont pas conservées.

Extraction du corps du mail

```
1 def extract_body(msg):
2     refused_charset = ['unknown-8bit', 'default', 'default_charset',
3                       'gb2312_charset', 'chinesebig5', 'big5']
4     body = ""
5
6     if msg.is_multipart():
7         for part in msg.walk():
8             if not part.is_multipart():
9                 body += extract_body(part)
10    return body
```

```

11
12     if msg.get_content_maintype() != 'text':
13         return ""
14
15     if msg.get_content_charset() in refused_charset:
16         return ""
17
18     if msg.get_content_subtype() == 'plain':
19         payload = msg.get_payload(decode=True)
20         body += payload.decode(errors='ignore')
21
22     if msg.get_content_subtype() == 'html':
23         payload = msg.get_payload(decode=True)
24         body += nettoyage.clear_html(payload.decode(errors='ignore'))
25
26     if msg.get_content_subtype() == 'enriched':
27         payload = msg.get_payload(decode=True)
28         body += nettoyage.clear_enriched(payload.decode(errors='ignore'))
29
30     return body

```

1.2.3 Nettoyage

Le nettoyage du texte utilise principalement les expressions régulières pour retirer un maximum d'éléments indésirables dans le texte. J'utilise également 2 modules externes afin de traiter le code HTML et faire la détection des mails qui ne sont pas écrits en anglais.

Par regex J'utilise le module python *re* pour réaliser les traitements suivants :

Suppression des réponses Lorsque l'on répond à un mail, le texte du message précédent est conservé dans le corps du mail. Afin de permettre la distinction avec les mails précédant le caractère '>' est ajouté en début de ligne. Je retire toutes les lignes correspondant à des réponses afin de limiter les doublons.

Nettoyage des réponses

```

1 def clear_reply(texte):
2     pattern = re.compile('^>.*$', flags=re.MULTILINE)
3     return re.sub(pattern, '', texte)

```

Suppression des ponctuations Afin de ne pas surcharger la base de données et pour se concentrer sur le texte, une grande partie des caractères de ponctuation seront retirés. L'idée est de se concentrer sur les ponctuations classiques (.,?!).

Nettoyage des ponctuations

```

1 def clear_punctuation(texte):
2     pattern_ponct = re.compile('[*#\-\_=:;<>\\[\\]"\'~)(|/$+}{@%&\\\\]',
3     flags=re.MULTILINE)
4     return re.sub(pattern_ponct, '', texte)

```

Suppression des balises pour les enriched text Certaines parties du corps de mail sont de type *enriched text*. Les balises ne sont pas pertinente dans notre analyse et sont donc retirées.

Nettoyage des balises enriched text

```
1 def clear_enriched(texte):
2     pattern = re.compile('<.*>')
3     return re.sub(pattern, '', texte)
```

Suppression des liens La présence de certaines informations comment les liens URL, les adresses mails et les numéros de téléphone ne peuvent pas être utilisés dans l'analyse textuelle. Cependant il peut être intéressant de conserver une trace de leur présence. Nous allons donc modifier ces liens qui seront comptabilisés avant d'être retirés du texte.

Nettoyage des liens

```

1 def change_lien(texte, liens):
2     pattern_mail = re.compile(' [a-zA-Z0-9_+~]+@[a-zA-Z0-9-]+\.\.[a-zA-Z0-9-]+\.' )
3
4     pattern_url1 = re.compile(' (http|ftp|https)?://([w~-]+(?:([w~-]+(?:[w~-]+)+))?)'
5
6         ' ([w~-+~\.,@?^=%&:/~\+~#]*[w~-+~\.,@?^=%&:/~\+~#])?'
7         ', flags=re.MULTILINE)
8     pattern_url2 = re.compile(' (\\w+\\.\\.)+\\w+', flags=re.MULTILINE)
9     pattern_tel1 = re.compile(' \\(\\d{3}\\)\\d+~\\d+' ) # (359)1234~1000
10    pattern_tel2 = re.compile(' \\d+([ .-]?\\d+)' ) # +34 936 00 23
11
12    temp, liens['MAIL'] = re.subn(pattern_mail, ' ', texte)
13
14    temp, liens['URL'] = re.subn(pattern_url1, ' ', temp)
15    temp, nb = re.subn(pattern_url2, ' ', temp)
16    liens['URL'] += nb
17
18    temp, liens['TEL'] = re.subn(pattern_tel1, ' ', temp)
19    temp, nb = re.subn(pattern_tel2, ' ', temp)
20    liens['TEL'] += nb
21
22    return temp

```

Suppression des nombres Comme pour les liens, les nombres sont comptabilisés et retirés. Je fais la distinction entre les nombres seuls et les nombres accompagnés de sigle monétaires.

```
monnaie = '$č'
```

Nettoyage des nombres

```
1 def change_nombres(texte, liens):
2     pattern_prix1 = re.compile(f'[{monnaie}]( )?\d+([\.,]\d+)? ', flags=
    re.MULTILINE)
```



```

3     pattern_prix2 = re.compile(f' \\d+([. ,]\\d+)?( )?[{monnaie}] ', flags=
re.MULTILINE)
4     pattern_nb = re.compile('\\d+')
5
6     temp, liens['PRIX'] = re.subn(pattern_prix1, ' ', texte)
7     temp, nb = re.subn(pattern_prix2, ' ', temp)
8     liens['PRIX'] += nb
9
10    temp, liens['NOMBRE'] = re.subn(pattern_nb, ' ', temp)
11
12    return temp

```

Par module Lors du processus de nettoyage, j'utilise deux modules externes plus performants que ce que j'aurais pu faire simplement avec des expressions régulières. L'un me permet de nettoyer le code HTML, l'autre de détecter la langue du message.

Suppression du code HTML Certaines parties du corps du mail sont de type HTML. J'utilise le module *BeautifulSoup* pour parser le code et récupérer le texte affiché.

Nettoyage des nombres

```

1 from bs4 import BeautifulSoup
2
3 def clear_html(texte):
4     brut = BeautifulSoup(texte, "lxml").text
5     return brut

```

Sélection des mails en anglais Lors de mes tests, je me suis rendu compte que certains mails n'étaient pas en anglais. J'ai donc trouvé le module *langdetect* qui permet de détecter le langage utilisé dans un texte en utilisant un modèle Naïve Bayes avec une précision de 99% (voir G.2).

Je conserve dans les données à mettre en base le langage détecté avec l'idée de pouvoir traiter plusieurs langues dans le futur.

La détection de la langue est réalisée juste avant la mise en base dans l'index ElasticSearch.

Création d'un document

```

1 import langdetect
2
3 def create_document(mail, categorie):
4     corp = mail_load.extract_body(mail)
5     corp, liens = nettoyage.clear_texte_init(corp)
6     sujet, expéditeur = mail_load.extract_meta(mail)
7
8     if not corp:
9         return None
10
11    try:
12        lang = langdetect.detect(corp)
13    except langdetect.lang_detect_exception.LangDetectException:
14        return None
15

```

```

16     if lang != 'en':
17         return None
18
19     if categorie.lower() not in ['spam', 'ham']:
20         categorie = 'inconnu'
21
22     doc = {
23         'hash': hashlib.md5(corp.encode()).hexdigest(),
24         'categorie': categorie.lower(),
25         'sujet': sujet,
26         'expediteur': expediteur,
27         'message': corp,
28         'langue': lang,
29         'liens': liens
30     }
31     return doc

```

Exemple de traitement Les sections suivantes présentent des exemples de traitement de la phase 1.

Traitement initial

```

1  message = '''
2  Message dedicated to be a sample to show how the process is clearing the
3  text.
4
5  Begin reply :
6  > He once said
7  >>> that it would be great
8  End of reply.
9
10 Substitutions :
11 spamassassin-talk@example.sourceforge.net
12 https://www.inphonic.com/r.asp?r=sourceforge1&refcode1=vs3390
13 hello.foo.bar
14 between $ 25 and 25,21 $
15
16 A number is : 2588,8 588
17 Phone type a : (359)1234-1000
18 Phone type b : +34 936 00 23 23
19 Punctuation : ——## ..
20 ~~~~~
21 '''
22 text, liens = clear_texte_init(message)
23 print(liens)
24 print(text)

```

Résultat traitement initial :

```
{'URL': 2, 'MAIL': 1, 'TEL': 2, 'NOMBRE': 3, 'PRIX': 2}
```

Message dedicated to be a sample to show how the process is clearing the text.

Begin reply

End of reply.

Substitutions

between and

A number is ,
Phone type a
Phone type b
Ponctuation ..

Traitement HTML

```
1 message_html = '''
2 <!DOCTYPE html PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">
3 <html>
4 <head>
5     <title>Foobar</title>
6 </head>
7 <body>
8 I actually thought of this kind of active chat at AOL
9 bringing up ads based on what was being discussed and
10 other features
11     <pre wrap="">On 10/2/02 12:00 PM, "Mr. FoRK"
12     <a class="moz-txt-link-/rfc2396E" href="mailto:fork_
13     list@hotmail.com">&lt;fork_list@hotmail.com&gt;</a>
14     wrote: Hello There, General Kenobi !?
15 <br>
16 </body>
17 </html>
18 '''
19 print(clear_html(message_html))
20
```

Résultat traitement HTML :

Foobar

I actually thought of this kind of active chat at AOL
bringing up ads based on what was being discussed and
other features
On 10/2/02 12:00 PM, "Mr. FoRK"
<fork_list@hotmail.com>
wrote: Hello There, General Kenobi !?

```

1 message_enriched = '''
2 <smaller>I'd like to swap with someone also using Simple DNS to take
3 advantage of the trusted zone file transfer option.</smaller>
4 '''
5 print(clear_enriched(message_enriched))
6

```

Résultat traitement enriched text :

I'd like to swap with someone also using Simple DNS to take
 advantage of the trusted zone file transfer option.

1.2.4 Mise en base

Cette section détaille les éléments relatifs à la mise en base des informations récoltées. Dans ce projet, j'utilise 2 moteurs de bases de données pour stocker les données des mails.

1. un index Elasticsearch pour faire le stockage des données textuelles
2. une base PostgreSQL pour le stockage des données numériques

J'utilise des conteneurs *docker* pour héberger les services de bases de données. L'utilisation des conteneurs me permet de partager plus facilement mes configurations et limite les erreurs d'installations.

Pour chaque mail récolté le programme de la phase 1 va générer un *document* avec les informations suivantes :

- hash - signature md5 du texte nettoyé
- catégorie - Ham, Spam ou Inconnu
- sujet - correspond à l'objet du mail
- expéditeur - adresse mail
- corps - corps du mail nettoyé
- langue - la langue détectée du mail (en)
- liens - données non textuelles extraites du corps :
 - URL - liens URL
 - Mail - adresses mail
 - Téléphone - numéros de téléphone
 - Prix - nombres avec un symbole de devise
 - Nombres

Chaque document va générer une entrée dans la base Elasticsearch et une entrée dans la base PostgreSQL.

Création d'un document

```

1 def create_document(mail, categorie):
2     corp = mail_load.extract_body(mail)
3     corp, liens = nettoyage.clear_texte_init(corp)
4     sujet, expéditeur = mail_load.extract_meta(mail)
5
6     if not corp:
7         return None
8
9     try:
10         lang = langdetect.detect(corp)
11     except langdetect.lang_detect_exception.LangDetectException:

```

```

12         return None
13
14     if lang != 'en':
15         return None
16
17     if categorie.lower() not in ['spam', 'ham']:
18         categorie = 'inconnu'
19
20     doc = {
21         'hash': hashlib.md5(corp.encode()).hexdigest(),
22         'categorie': categorie.lower(),
23         'sujet': sujet,
24         'expediteur': expediteur,
25         'message': corp,
26         'langue': lang,
27         'liens': liens
28     }
29     return doc
30

```

Ci-dessous le schéma des bases de données avec les relations entre elles.

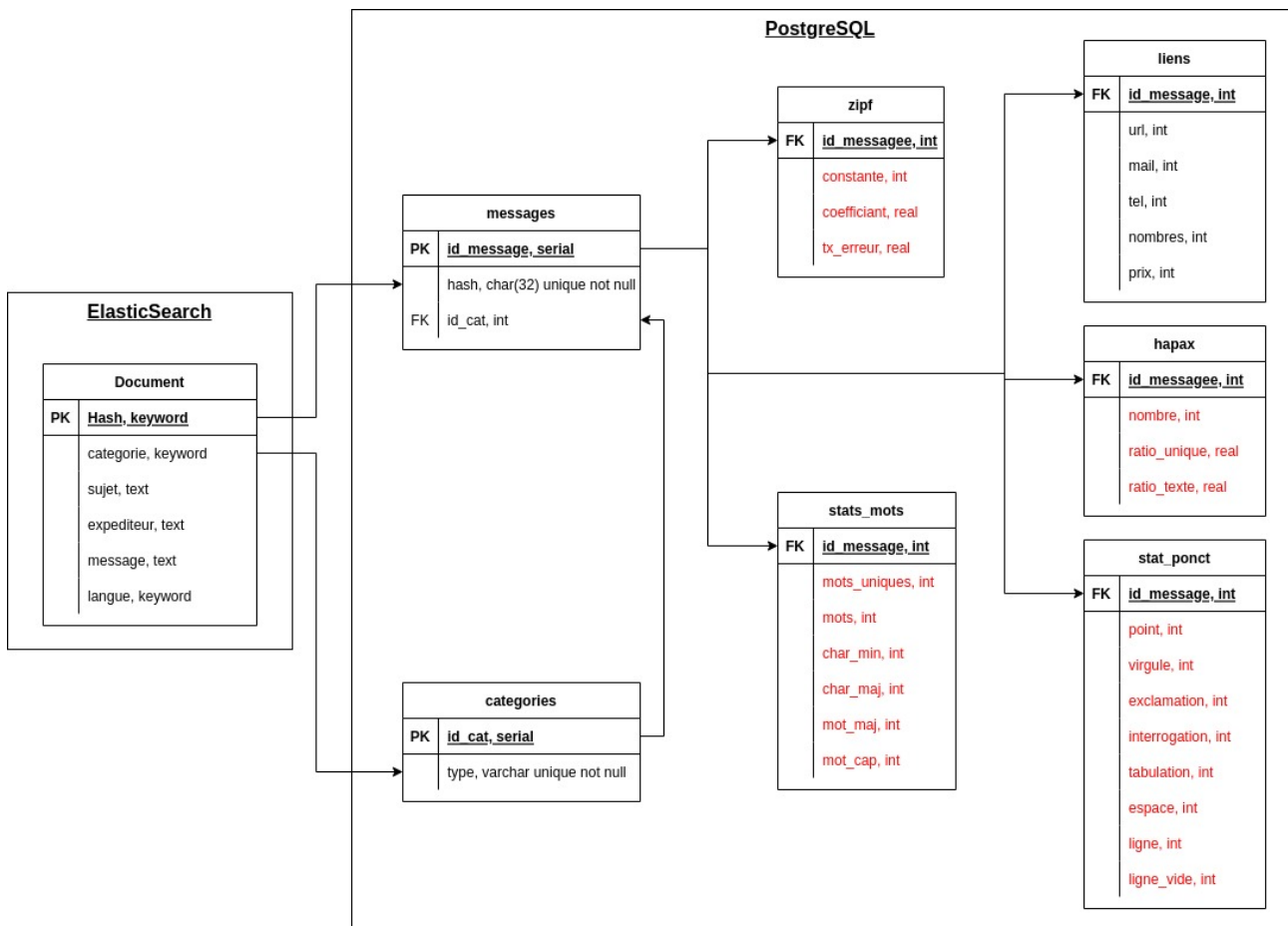


FIGURE 3 – Schéma des bases de données de l'application

Le *hash* calculé lors de la phase de traitement est l'identifiant unique du mail dans toutes les bases. La catégorie du mail est également présente dans les deux bases. Les champs en rouge sont des caractéristiques qui ne sont pas calculées lors de la phase 1.

Stockage des données : Elasticsearch Elasticsearch est un moteur de base de données NoSQL. Il intègre un moteur d'indexation des documents est assez performant pour stocker des données textuelles de taille aléatoire. Il est cependant assez compliqué de modifier le schéma d'un index une fois qu'il a été créé. Pour ces raisons, j'utilise cette technique pour ne stocker que les données textuelles après nettoyage car elles n'ont plus vocation à être modifiées.

Les corps de mail mis en base seront récupérés ultérieurement pour réaliser les opérations d'analyse. Les résultats seront stockés dans la base PostgreSQL, plus souple.

Ci-dessous les fonctions principales pour l'ajout d'un document dans la base ES.

Fonctions basique pour la liaison Elasticsearch

```
1 def es_connect(server , creds , crt):
2     """ Connexion au serveur Elasticsearch """
3     client = Elasticsearch(server , api_key=creds , ca_certs=crt)
4
5     try:
6         client.search()
7         client.indices.get(index="*")
8         return client
9
10    except (exceptions.ConnectionError , AuthenticationException ,
11    AuthorizationException) as err:
12        print("ES:conn - Informations client Elasticsearch :\n\t", err)
13        client.close()
14        return None
15
16 def es_index_doc(es_cli , index , doc):
17     """ Index un document dans la base ES """
18     id_doc = doc['hash']
19
20     if es_document_exists(es_cli , index , id_doc):
21         return 1
22
23     es_cli.index(index=index , document=doc)
24     es_cli.indices.refresh(index=index)
25     return
26
27
28 def es_document_exists(es_cli , index , hash):
29     """ Regarder dans l'index si le hash du document est deja present """
30     try:
31         resp = es_cli.search(index=index , query={"match": {"hash": hash
32 }}})
33     except elasticsearch.NotFoundError as err:
34         print("Error : hash", err , file=sys.stderr)
35         return None
36
37     return True if resp['hits']['total']['value'] == 1 else False
```

Le déploiement et les configurations de la base Elasticsearch sont disponibles dans l'annexe

B.1.

Stockage des premières informations statistiques : PostgreSQL Pour le stockage des données statiques et d'analyse, j'ai décidé de m'orienter vers le système de gestion de base de données PostgreSQL. Une base de données relationnelle est plus flexible qu'un index Elastic-Search pour l'ajout de nouvelles caractéristiques.

Durant la phase 1, j'utilise que 3 tables :

- messages - liste des mails avec le hash permettant de faire le lien avec l'index ES
- categorie - liste des catégories de mails
- liens - comptabilise le nombre d'occurrences par type de lien ou nombre

Cette base a pour but de stocker les données formatées pour l'analyse, l'exploitation et l'entraînement du modèle.

Fonctions basiques pour la liaison PostgreSQL

```
1 def connect_db(database, user, passwd, host, port):
2     """ Connexion a la base de donnees Postgres """
3     try:
4         client_psql = psycopg2.connect(database=database, user=user,
5         password=passwd, host=host, port=port)
6     except psycopg2.Error as e:
7         print("Erreur de connexion : \n{}".format(e), file=sys.stderr)
8         return None
9
10    client_psql.autocommit = True
11    return client_psql
12
13 def insert_data(client_psql, table, data):
14     """
15     Insere les donnees d'un dictionnaire dans une table de la base de
16     donnees PSQL
17     Les cles du dictionnaire doivent correspondre aux colonnes de la
18     table.
19     """
20     cols = ','.join([str(c) for c in data.keys()])
21     vals = ','.join([str(v) if (type(v) != str) else f"'{v}'" for v in
22     data.values()])
23     query = f"INSERT INTO {table}({cols}) VALUES ({vals})"
24
25     exec_query(client_psql, query)
26
27 def get_data(client_psql, table, champs, clause=None):
28     """ Recupere les donnees de la base. """
29     query = f"SELECT {'','.join(champs)} FROM {table}"
30     if clause:
31         query += f" WHERE {clause}"
32
33     result = exec_query(client_psql, query)
34     return [dict(zip(champs, ligne)) for ligne in result]
```

```

35 def insert_document_init(client_psql, data, id_cat):
36     """ Insere un nouveau document dans la base PSQL. """
37     insert_data(client_psql, 'messages', {'hash': data['hash'], 'id_cat':
38         id_cat})
39     id_message = get_data(client_psql, 'messages', ['id_message'], f"hash
40     LIKE '{data['hash']}'")[0]['id_message']
41
42     liens = data['liens']
43     liens.update({'id_message': id_message})
44     insert_data(client_psql, 'liens', liens)
45
46 def exec_query(client_psql, query):
47     """ Execute une query dans la base PSQL """
48     cursor = client_psql.cursor()
49
50     try:
51         cursor.execute(query)
52         if query.upper().find("SELECT", 0, 6) >= 0:
53             return cursor.fetchall()
54         return []
55     except psycopg2.Error as e:
56         print("Erreur d'execution de la requete : {}".format(e), file=sys
57             .stderr)
58         print("requete : {}".format(query), file=sys.stderr)
59         return []

```

Le déploiement et les configurations de la base PostgreSQL sont disponibles dans l'annexe B.2.

Stockage des données statistiques du traitement : SQLite Les données présentes dans cette base permettent de suivre l'évolution de certaines métriques lors des différentes étapes du nettoyage. A chaque grandes étapes de la phase 1 (Importation, Nettoyage, Mise en base), je calcule pour les HAM, SPAM et (HAM+SPAM) les éléments suivants :

- mails - nombre de mails
- mots - nombre de mots dans tout le corpus
- mots_uniques - nombre de mots uniques dans tout le corpus

Ces données me permettent d'estimer la quantité de données nettoyées durant cette phase.

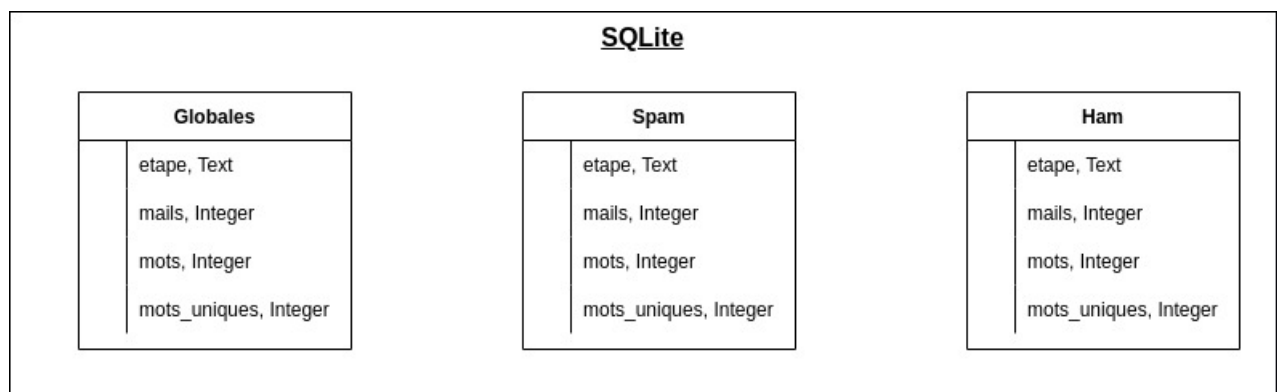


FIGURE 4 – Schéma de la base de données pour lors du traitement

1.3 Données de la phase 1

Sortie du programme : main_collecte.py

```
=== Vérification des conteneurs Docker ===
* Conteneur 'docker-es01-1'... OK
* Conteneur 'docker-kibana-1'... OK
* Conteneur 'docker_pgadmin_1'... OK
* Conteneur 'docker_pgdb_1'... OK
=== Phase 1 : collecte & mise en base ===
== Création de la base SQLITE
SQLITE table globales : CREATED
SQLITE table spam : CREATED
SQLITE table ham : CREATED
== Recolte ==
-- Création de la liste des fichiers... OK
--- Process de statistiques après la récolte
-- Stats - étape : Récolte spam... OK
-- Stats - étape : Récolte ham... OK
--- Sauvegarde des stats de l'étape: recolte... OK
Données stats de l'étape: recolte:
HAM, mails: 4153 mots: 1863674 mots uniques: 178812
SPAM, mails: 1898 mots: 918256 mots uniques: 129601
GLOBALES, mails: 6051 mots: 2781930 mots uniques: 287637
== Création document ==
-- Importation - Création spam... OK
-- Importation - Création ham... OK
--- Process de statistiques après la création de document
-- Stats - étape : création documents spam... OK
-- Stats - étape : création documents ham... OK
--- Sauvegarde des stats de l'étape: creation document... OK
Données stats de l'étape: creation document:
HAM, mails: 4124 mots: 844151 mots uniques: 73725
SPAM, mails: 1784 mots: 584012 mots uniques: 40099
GLOBALES, mails: 5908 mots: 1428163 mots uniques: 97235
== Mise en base des documents ==
-- Création de l'index ElasticSearch... OK
-- Création de la base PostgreSQL
User 'data' déjà existant
-- Création des tables PostgreSQL... OK
-- Mise en base ES & PSQL des spam... OK (256 doublons)
-- Mise en base ES & PSQL des ham... OK (78 doublons)
-- Récupération des spam... OK
-- Récupération des ham... OK
--- Process de statistiques après la mise en base
-- Stats - étape : Mise en base spam... OK
-- Stats - étape : Mise en base ham... OK
--- Sauvegarde des stats de l'étape: mise en base... OK
Données stats de l'étape: mise en base:
HAM, mails: 4046 mots: 839368 mots uniques: 73725
SPAM, mails: 1528 mots: 529121 mots uniques: 40099
```

GLOBALES, mails: 5574 mots: 1368489 mots uniques: 97235
 == FIN ==

Évolution lors du traitement Les données recueillis pendant le traitement permettent de constater les comportements suivants :

- La diminution des documents spam est plus importante que celle des ham
- Le nombre de mots uniques ne diminue plus après la création de document
- La diminution du nombre de mots est plus importante dans les ham que dans les spam.
- le nombre de mot uniques est plus important dans les ham que dans les spam

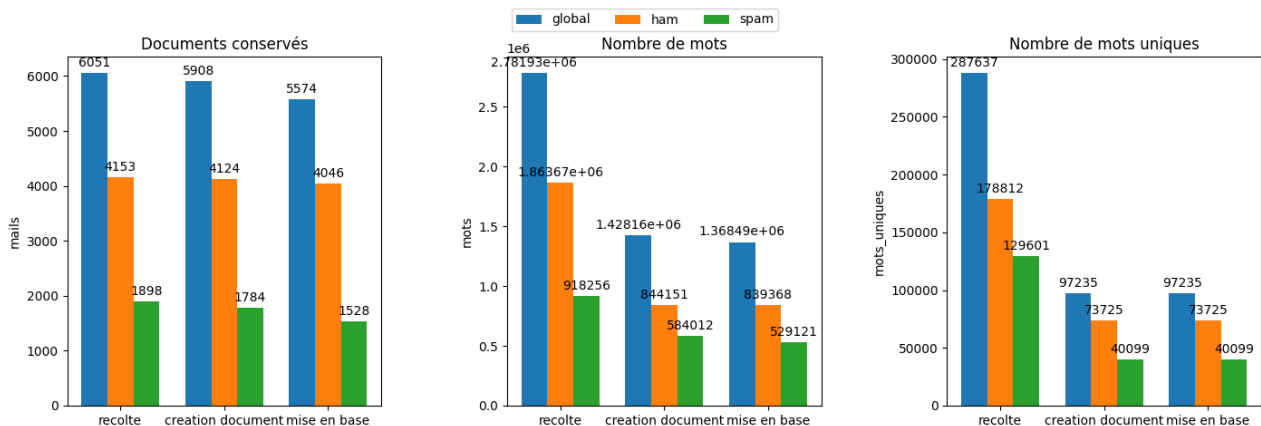


FIGURE 5 – Évolution des principaux marqueurs lors du traitement

Conclusion de la récolte A l'issue de la phase de récolte, on remarque qu'il y a un nombre plus important de document en double dans la catégorie *spam*. Il y a également une réduction importante du nombre de mots ainsi que du nombre de mots uniques dans les *ham*. Cette réduction peut s'expliquer par le nettoyage du corps des mails :

- Retrait des réponses
- Retrait de certaines ponctuations
- Suppression des balises HTML et enriched text
- Retrait des liens, et nombres

Enfin on peut voir que les *ham* utilisent plus de mots uniques que les *spam*. Il est donc possible que le vocabulaire des *spam* soit plus restreint.

Premières données statistiques Les documents conservés à l'issue de cette phase de récupération se constitue à 72,6% de Ham et à 27.4% de Spam.

Répartition des ham/spam

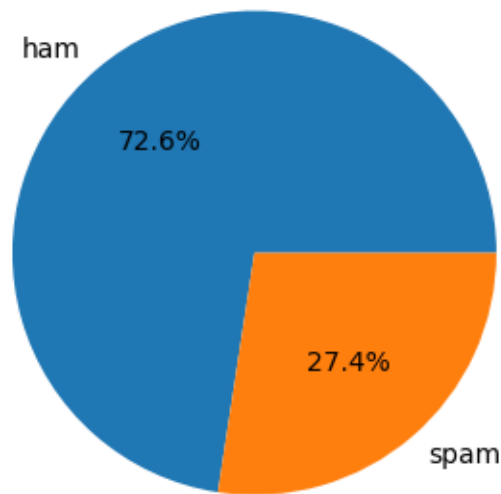


FIGURE 6 – Répartition des ham/spam dans le dataset

Grâce aux informations et des pré-traitements réalisés lors de cette première phase, il est possible de présenter quelques données statistiques. Notamment sur la présence et l'utilisation de liens (url, mail, téléphonique) et des données numériques (nombres et prix).

TABLE 1 – Statistiques sur les liens

	url			mail			tel		
	Globale	Ham	Spam	Globale	Ham	Spam	Globale	Ham	Spam
moyenne	5.47	5.54	5.26	1.16	1.21	1.01	0.16	0.20	0.64
écart-type	49.74	57.77	13.71	2.50	1.72	3.88	1.06	0.63	1.73
minimum	0			0			0		
25%	1			0			0		
médiane	2	2	3	1	1	0	0		
75%	4	4	5	2	2	1	0		
maximum	3557	3557	295	69	34	69	67	25	67

A partir des données des liens présentées dans le tableau ci-dessus il est possible d'émettre les hypothèses suivantes.

On trouve en moyenne environ 5 url dans les deux types de mails. Cependant la dispersion est notablement plus faible pour les spams. On voit également que 75% des spam contiennent 5 liens ou moins.

La présence d'adresses mail est assez marginale environ une adresse mail par corps de texte. L'écart-type est plus faible dans les ham. On voit également que 75% des ham contiennent deux adresse mails ou moins contre 1 pour les spam.

Enfin la présence de numéro de téléphone est quasiment nulle pour les deux catégories.

TABLE 2 – Statistiques sur les données numériques

	nombre			prix		
	Globale	Ham	Spam	Globale	Ham	Spam
moyenne	15.14	9.38	30.38	0.28	0.08	0.81
écart-type	219.16	49.35	410.53	1.54	0.66	2.66
minimum	0			0		
25%	1	1	2	0		
médiane	5	5	7	0		
75%	9	8	14	0		
maximum	15801	2794	15801	29	17	29

Les données numériques du dataset montrent les hypothèses suivantes.

Les nombres simples semblent être les données les plus présentes dans le corps du mail, environ 15 en moyenne par mail, cependant la dispersion est élevée. Cela étant 75% des spam ont 14 nombres ou moins contre 8 ou moins pour les ham.

Enfin la présence de prix est assez marginale par rapport aux nombres. Il y a cependant en moyenne 10 fois plus de prix dans les spam que les ham.

Distribution des mails en fonction du nombre de liens

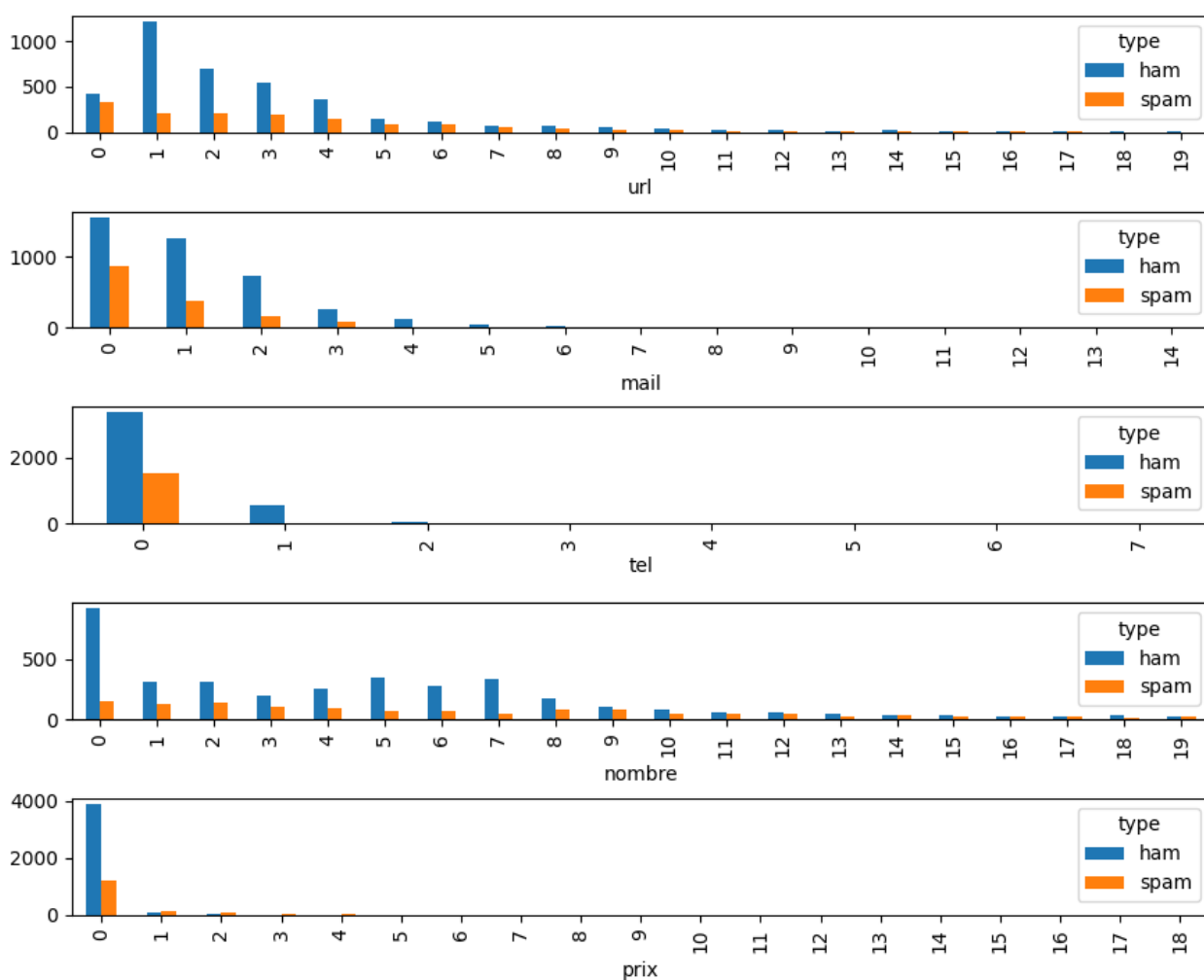


FIGURE 7 – Distribution du nombre de liens selon la catégorie du mail

Au vu des analyses précédents, on peut dire qu'il est statistiquement plus probable de trouver dans les spam plus de liens url, de nombre et de prix que dans les ham. La présence d'adresse mail est plus probable dans les ham. Il est enfin possible d'exclure la présence de numéro de téléphone qui est trop marginale.

Les codes sources de cette partie sont disponibles dans le fichier `./analyse/analyse_p1.py` ou dans la section ??

2 Phase 2 : Traitement des données

Les opérations de cette phase ont pour but d'extraire les caractéristiques des textes collectés à la phase précédente. Puis, Nous allons rechercher des informations numériques sur la forme des messages (nombre de mots, de ponctuations,...). Enfin, nous appliqueront des techniques de traitement du langage naturel pour travailler sur le fond des corps de mail.

2.1 Génération des données statistiques

Nous verrons ici le code utilisé pour récupérer ou générer les données des caractéristiques pour chaque message ainsi que la distribution de Zipf de manière globale.

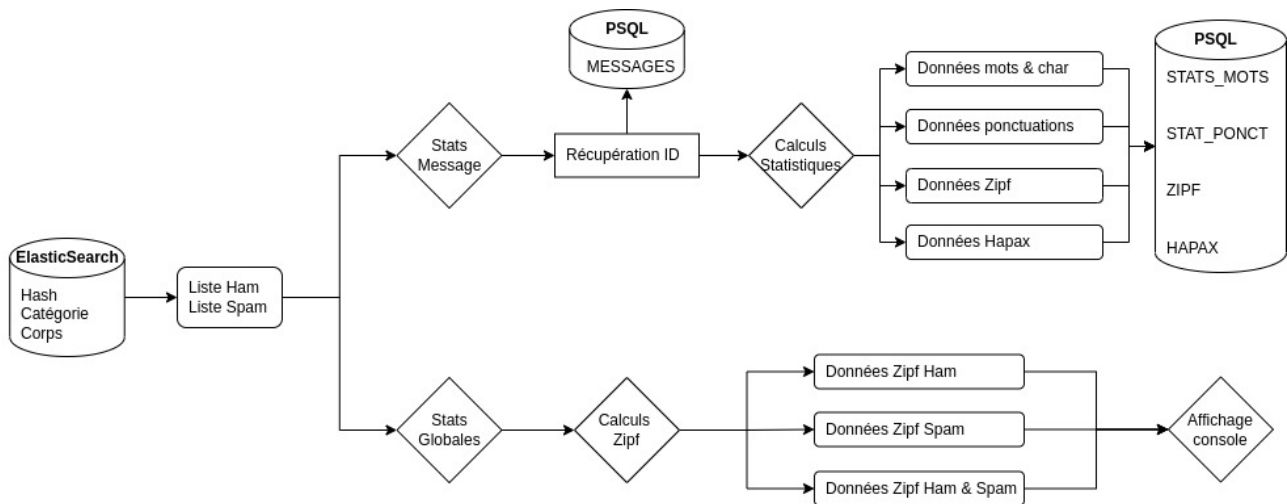


FIGURE 8 – Schéma des étapes de récupération des données

Les mails sont récupérés dans la base ElasticSearch. Les Ham et les Spam sont séparés dans deux listes distinctes.

Récupération des messages dans la base ES

```

1 def recup_mails(es_index):
2     es_cli = elastic_cmd.es_connect(es_secrets.serveur, (es_secrets.apiid
3     , es_secrets.apikey), es_secrets.ca_cert)
4     data = {}
5     for cat in ['spam', 'ham']:
6         print("— Recuperation des {}".format(cat), end=' ')
7         data[cat] = {entry.get('_source').get('hash'): entry.get('_source
8         ').get('message') for entry in elastic_cmd.es_get_all(es_cli, es_index
9         , sort={'hash': 'asc'}, query={'match': {'categorie': cat}})}
10        print('OK')
11    return data
  
```

Chaque message va ensuite passer par une série de traitement pour générer les données des caractéristiques avant d'être mise en base dans PostgreSQL.

Pipeline pour chaque message

```

1 def stats_pipe_message(hash_message, message, psql_cli):
2     resp = psql_cmd.get_data(psql_cli, 'messages', ['id_message'], f"hash
3     LIKE '{hash_message}'")
4     if not resp:
  
```

```

4         print(f"No id_message found for {hash_message}", file=sys.stderr)
5         return
6     id_mess = resp[0].get('id_message')
7
8     for table, stat_func in {'stat_ponct': stats_ponct,
9                             'stats_mots': stats_mot,
10                            'zipf': stats_zipf,
11                            'hapax': stats_hapax}.items():
12         resp = psql_cmd.get_data(psql_cli, table, ['id_message'], f"
id_message={id_mess}")
13         if resp:
14             continue
15         psql_cmd.insert_data(psql_cli, table, stat_func(id_mess, message))
16
17 def stats_ponct(id_text, texte):
18     return {
19         "id_message": id_text,
20         "point": texte.count('.'),
21         "virgule": texte.count(','),
22         "exclamation": texte.count('!'),
23         "interrogation": texte.count('?'),
24         "tabulation": texte.count('\t'),
25         "espace": texte.count(' '),
26         "ligne": texte.count('\n') + 1,
27         "ligne_vide": len(re.findall(r'^\s*$', texte, re.MULTILINE))
28     }
29
30
31 def stats_mot(id_text, texte):
32     tokens = re.findall(r'\w+', texte, re.MULTILINE).
33     return {
34         'id_message': id_text,
35         'char_min': len(re.findall(r'[a-z]', texte, re.MULTILINE)),
36         'char_maj': len(re.findall(r'[A-Z]', texte, re.MULTILINE)),
37         'mots': len(tokens),
38         'mots_uniques': len(set(tokens)),
39         'mot_maj': sum(mot.isupper() for mot in tokens),
40         'mot_cap': sum(bool(re.match(r'[A-Z][a-z]+', mot)) for mot in
tokens)
41     }
42
43
44 def stats_zipf(id_text, texte):
45     tokens = re.findall(r'\w+', texte, re.MULTILINE)
46     data = stats.zipf_process(tokens)
47     data['id_message'] = id_text
48     data['constante'] = data.pop('const_moy')
49     data['coefficient'] = data.pop('coef_min')
50     data['tx_erreur'] = data.pop('cout_min')
51     return data
52
53
54 def stats_hapax(id_text, texte):

```

```

55     tokens = re.findall(r'\w+', texte, re.MULTILINE)
56     data = stats.hapax(tokens)
57     data['id_message'] = id_text
58     data['nombre'] = data.pop('nombres')
59     data['ratio_unique'] = data.pop('ratio_mots_uniques')
60     return data
61

```

J'ai également réaliser les calculs de la distribution de Zipf sur l'ensemble du corpus ainsi que sur le regroupement des Ham et des Spam pour voir si des différences majeures pouvaient être mises en lumière. Les résultats de cette recherche sont affiché directement dans la console.

Pipeline pour les données globales

```

1  def stats_pipe_globale(data):
2      ls_ham = [mess for mess in data['ham'].values()]
3      ls_spam = [mess for mess in data['spam'].values()]
4      tokens = []
5      for mess in ls_ham + ls_spam:
6          tokens += re.findall(r'\w+', mess, re.MULTILINE)
7      data = stats.zipf_process(tokens, True)
8      print("Donnees Zipf ham+spam:", data)
9      data = stats.hapax(tokens)
10     print("Donnees Hapax ham+spam:", data)
11
12     tokens = []
13     for mess in ls_ham:
14         tokens += re.findall(r'\w+', mess, re.MULTILINE)
15     data = stats.zipf_process(tokens, True)
16     print("Donnees Zipf ham:", data)
17     data = stats.hapax(tokens)
18     print("Donnees Hapax ham:", data)
19
20     tokens = []
21     for mess in ls_spam:
22         tokens += re.findall(r'\w+', mess, re.MULTILINE)
23     data = stats.zipf_process(tokens, True)
24     print("Donnees Zipf spam:", data)
25     data = stats.hapax(tokens)
26     print("Donnees Hapax spam:", data)
27

```

Les fonctions pour les instructions pour la distribution de Zipf sont donnés dans l'annexe ??.

2.2 Sortie de la partie 2 : Stats d'exploitation

```

=== Phase 2 : Stats Exploitation ===
-- Récupération des spam... OK
-- Récupération des ham... OK
-- Récupération des statistiques par message
Stats Spam... OK
Stats Ham... OK
-- Récupérations des statistiques globales

```



```
Données Zipf ham+spam: {'const_moy': 73746.97464716957, 'cout_min': 6.37798252533617,
'coef_min': 0.93}
Données Hapax ham+spam: {'nombres': 33666, 'ratio_mots_uniques': 0.47656526478207323,
'ratio_texte': 0.024792183385275213}
Données Zipf ham: {'const_moy': 52712.404742906976, 'cout_min': 4.193202011347324,
'coef_min': 0.93}
Données Hapax ham: {'nombres': 26231, 'ratio_mots_uniques': 0.48674175650850793,
'ratio_texte': 0.03150367568904908}
Données Zipf spam: {'const_moy': 35485.13912209956, 'cout_min': 4.848725944493132,
'coef_min': 0.93}
Données Hapax spam: {'nombres': 12396, 'ratio_mots_uniques': 0.4097309446684736,
'ratio_texte': 0.023598168648092978}
```

2.3 Analyses statistiques

Dans cette partie nous nous intéressons aux données statistiques des caractéristiques récupérées lors de la phase précédente. Nous verrons dans un premier temps les données des mots :

- nombre de mots
- nombre de mots uniques
- nombre de caractère en minuscule
- nombre de caractère en majuscule
- nombre de mots écrits complètement en majuscule
- nombre de mots écrits complètement avec la première lettre en capitale

A cela nous ajouterons une recherche de caractéristique sur les données de ponctuation :

- nombre de points '.'
- nombre de virgules ','
- nombre de points d'exclamation
- nombre de points d'interrogation
- nombre de tabulations
- nombre d'espaces
- nombre de lignes totales
- nombre de lignes vides

Enfin nous appliquerons les méthodes de distribution de Zipf sur chaque message¹. Nous récupérerons alors les données suivantes

- constante estimée
- coefficient estimés
- taux d'erreur absolu moyen entre les fréquences réelles et théoriques de chaque mot
- nombre de mots n'apparaissant qu'une seule fois dans le texte (hapax)
- nombre d'hapax par rapport au vocabulaire du texte
- nombre d'hapax par rapport au total des mots du texte

Statistiques des mots et caractères Le tableau ??, sur les données des mots, montre que les Spam ont généralement plus de mots, 343 contre 205 en moyenne pour les Hams. Cela est également vrai pour les mots uniques par messages, les spams en contiennent globalement plus. Concernant la forme des mots, on voit que les mots avec majuscules sont plus souvent utilisés dans les Spam. Enfin, il est important de noter la forte dispersion des données de toutes ces métriques quelque soit la catégorie.

1. J'ai conscience que cette distribution est plus pertinente quand le corpus est grand

TABLE 3 – Statistiques sur les données des mots

	nombre de mots			mots uniques			mots capitalisés			mots majuscules		
	Global	Ham	Spam	Global	Ham	Spam	Global	Ham	Spam	Global	Ham	Spam
moyenne	243	205	343	130	116	167	41	34	59	14	8	29
écart-type	586	531	701	187	180	198	114	111	119	49	39	66
minimum	1	3	1	1	3	1	0	0	0	0	0	0
25%	53	44	94	46	39	73	9	8	15	2	2	3
médiane	105	85	161	80	68	110	17	14	32	4	3	10
75%	215	174	326	139	120	191	35	25	60	10	7	24
maximum	14127	14127	11274	4322	4322	2479	4656	4656	2192	1969	1969	763

TABLE 4 – Statistiques sur les données des caractères

	caractères en minuscule			caractères en majuscule		
	Globale	Ham	Spam	Globale	Ham	Spam
moyenne	1084	923	1511	111	66	232
écart-type	2816	2518	3448	736	222	1351
minimum	0	3	0	0	0	0
25%	223	187	380	15	13	43
médiane	442	367	675	31	23	84
75%	917	755	1410	76	44	173
maximum	68082	68082	54406	50543	6450	50543

Le tableau ?? qui montre les données des caractères minuscules et majuscules montre une taille moyenne des Spam plus grande que celle des Ham. Il est également intéressant de noter que l'utilisation des majuscules est plus prononcée dans les Spam.

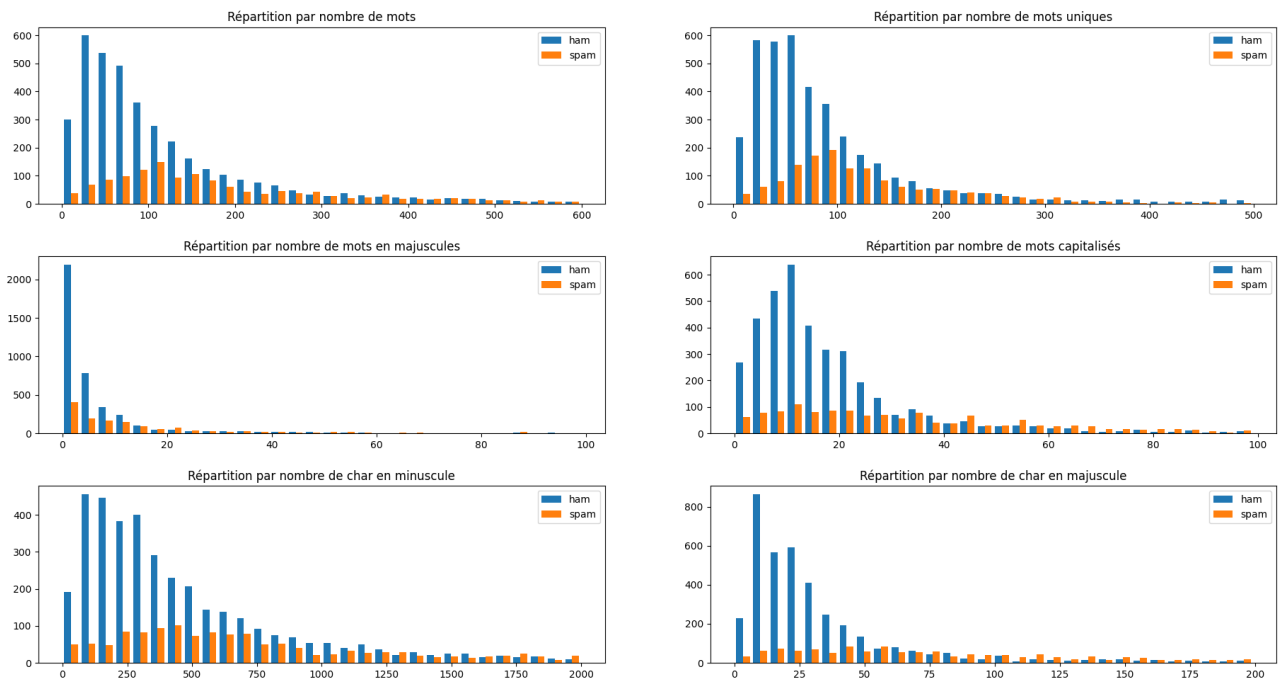


FIGURE 9 – Distribution des données de mots selon la catégorie du mail

Concernant les distributions des mails visibles dans la Figure ?? on remarque que les Ham se concentrent généralement vers des métriques basses puis l'effectif se réduit à mesure que

la valeur de la métrique augmente. Les Spam quand à eux ont plus tendance à suivre une distribution de type normale ou uniforme.

TABLE 5 – Statistiques sur les ponctuations

	points			virgules			exclamations			interrogations		
	Global	Ham	Spam	Global	Ham	Spam	Global	Ham	Spam	Global	Ham	Spam
moyenne	14	12	20	12	11	16	2	0.7	5.9	1	0.9	1.3
écart-type	35	30	46	37	31	50	7	5	10	2	2	3
minimum	0	0	0	0	0	0	0	0	0	0	0	0
25%	3	3	3	2	2	2	0	0	1	0	0	0
médiane	6	5	9	5	4	6	0	0	3	0	0	0
75%	13	11	18	10	9	14	2	1	7	1	1	1
maximum	848	848	643	1268	890	1268	343	343	92	79	79	58

Statistiques des ponctuations Le tableau ?? des données de ponctuation donne l'information que les textes des Spam utilisent plus de signes de ponctuation que les Ham. Il faut relever que l'utilisation des points d'interrogation est anecdotique pour les deux catégories.

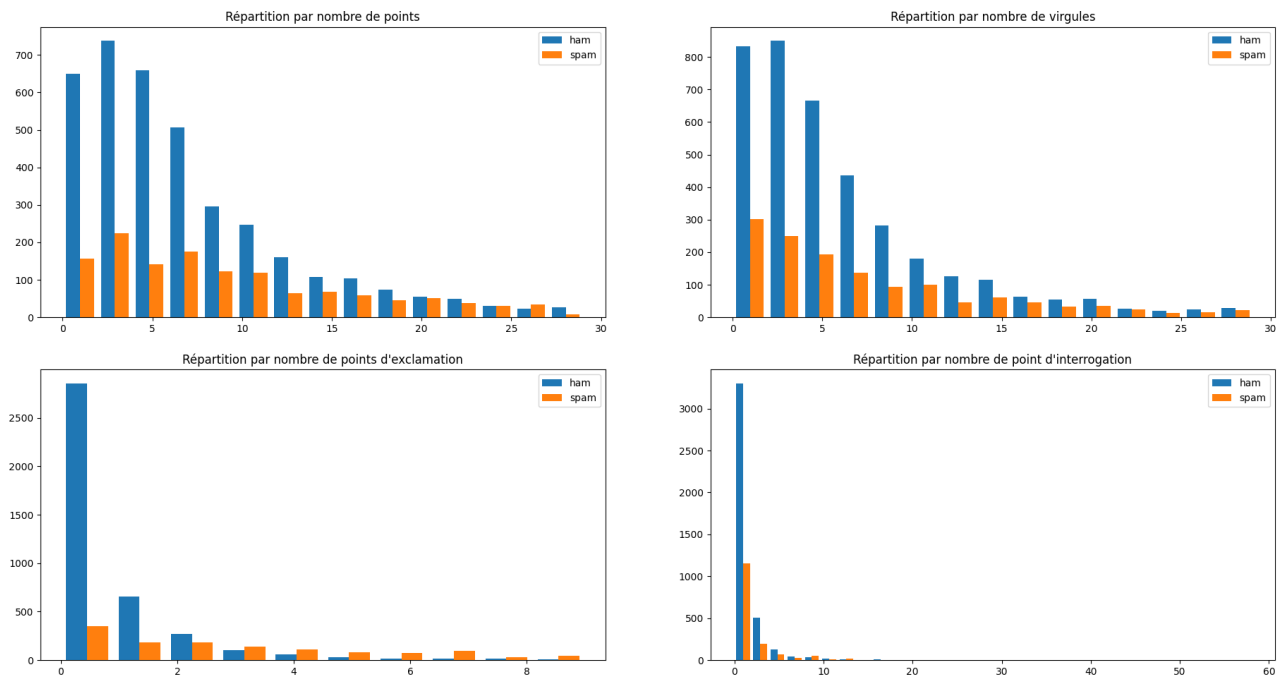


FIGURE 10 – Distribution des données de la ponctuation selon la catégorie du mail

Les schéma de la Figure ?? montrent un comportement similaire pour les deux catégories avec le nombre de points, de virgule et de point d'interrogation. Seule la répartition des effectifs de spam pour le nombre de point d'exclamation semble ne pas suivre le même cheminement et semble rester uniforme.

Statistiques des espaces et des lignes Le tableau ?? regroupe les données des espace vide et le calcul du nombre de lignes. On remarque que les tabulations ne sont pas utilisées, sauf en grande quantité. Il pourrait être intéressant de voir si cet élément peut être utiliser pour détecter les documents générant des valeurs aberrantes. Pour ce qui est du nombre d'espaces, de lignes et de lignes vide, la moyenne des Spam est approximativement le double de celle des Ham.

TABLE 6 – Statistiques sur les espaces et les lignes

	tabulations			espaces			lignes			lignes vides		
	Globale	Ham	Spam	Globale	Ham	Spam	Globale	Ham	Spam	Globale	Ham	Spam
moyenne	2	2	3	287	240	411	65	57	86	15	13	19
écart-type	29	32	19	874	917	733	137	144	113	44	50	22
minimum	0	0	0	1	2	1	1	1	1	0	0	0
25%	0	0	0	55	45	99	21	19	28	6	6	8
médiane	0	0	0	114	95	196	34	30	52	9	8	13
75%	0	0	0	255	196	435	61	48	97	14	12	23
maximum	1118	1118	379	36350	36350	11436	6320	6320	1695	2882	2882	291

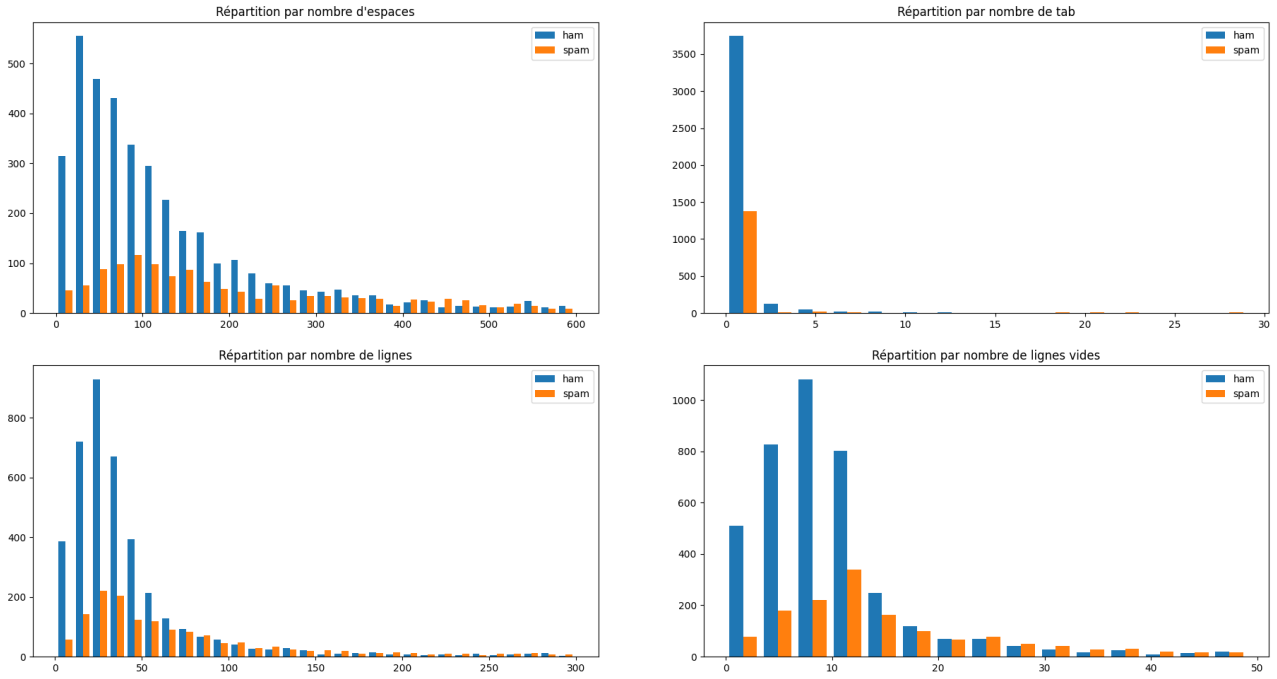
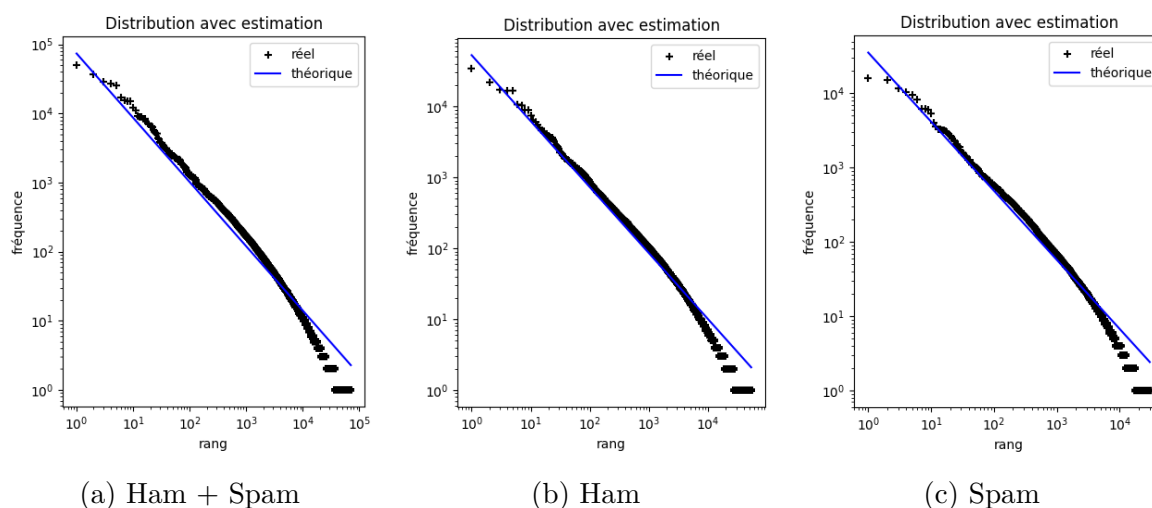


FIGURE 11 – Distribution des données de lignes et espaces selon la catégorie du mail

On peut voir dans les histogrammes de la Figure ?? que les mails se répartissent de manière relativement normale pour le comptage du nombre de lignes et de lignes vides pour des valeurs faibles. Puis à mesure que la valeur augmente, l'effectif semble se réduire de manière linéaire. La répartition du nombre d'espace ressemble fortement à la répartition du nombre de mots. Il est possible que le nombre d'espaces soit corrélé au nombre de mots.

Distribution de Zipf Cette distribution est une loi empirique qui permet de décrire un classement des mots selon leur fréquence d'apparition. Les fonctions utilisées pour la mise en œuvre de cette distribution ont été réalisées par mes soins. L'annexe ?? détaille le développement de cette fonctionnalité.

Les données suivantes montre la représentation de cette distribution pour l'ensemble du corpus.



On voit dans les graphiques ci dessus que les corpus Ham, Spam et l'association des deux respectent une distribution zipfienne.

	constante	coefficient	erreur moyenne	hapax	hapax/vocab	hapax/total
Ham + Spam	73746.97	0.93	6.37	33666	0.47	0.02
Ham	52712.40	0.93	4.19	26231	0.48	0.03
Spam	35485.13	0.93	4.84	12396	0.40	0.02

On remarque à partir du tableau ci dessus que la constante estimée est plus importante dans les corpus Ham que dans le corpus Spam. La constante des Ham (52712) est assez proche de la constante moyenne découverte avec le corpus de Brown lors du développement de cette méthode (56525). Il est possible d'émettre l'hypothèse que le vocabulaire des Ham est plus fournis. On remarque également que le taux d'hapax dans le vocabulaire des Ham est de 48% alors que cette proportion est de 40% pour les Spam.

Distribution de Zipf par message Le processus appliqué pour le calcul de la distribution de Zipf sur les différent corpus a été utilisé sur chaque mail individuellement. Il a alors été possible d'ajouter ces données (constante, coefficient, taux d'erreur, nombre d'hapax et ratio d'hapax) pour chaque mail. Le résultat de ce traitement est détaillé dans le Tableau ?? et dans la Figure ??.

TABLE 7 – Statistiques sur les données des calculées de Zipf et Hapax

	constante		coefficient		taux d'erreur		hapax		ratio vocabulaire		ratio texte	
	Ham	Spam	Ham	Spam	Ham	Spam	Ham	Spam	Ham	Spam	Ham	Spam
moyenne	66	102	1.21	1.19	1.37	1.54	87.7	111.53	0.83	0.72	0.67	0.52
écart-type	115	141	0.08	0.07	0.24	0.41	123.8	123.4	0.12	0.21	0.20	0.22
minimum	2	1	0.86	0.86	0.44	0	0	0	0	0	0	0
25%	21	39	1.16	1.14	1.26	1.37	35	51	0.77	0.71	0.52	0.40
médiane	36	63	1.22	1.19	1.37	1.46	57	80	0.84	0.77	0.67	0.54
75%	65	112	1.29	1.26	1.46	1.57	93	136	0.91	0.83	0.83	0.67
maximum	2311	1866	1.30	1.30	3.49	3.49	3552	1690	1	1	1	1

Le tableau ?? montre que les constantes estimées sont globalement plus importantes dans les spams. Le coefficient est plus élevé dans les Ham. Cependant il est à noté que la moyenne des coefficient pour les Ham (1.21) et pour les Spam (1.19) est assez supérieure au coefficient

estimé pour tout le corpus (0.93). On remarque également pour le coefficient que la valeur se rapproche rapidement de 1.30 qui est la limite haute fixée pour l'estimation du coefficient.

Le taux d'erreur, qui correspond à la moyenne de l'écart absolu entre la valeur réelle et estimée de la constante, est plus élevé dans les spam (1.54). Mais cette valeur reste inférieure à celle observée pour l'ensemble du corpus. Cela peut s'expliquer par un écart de fréquence moins important entre les mots les plus présents et les moins présents.

Le calcul des hapax semble plus parlant mail par mail que dans la globalité du corpus. On remarque effectivement un écart de 15 points entre la moyenne des hapax dans les textes des Ham (67%) que dans ceux des Spam (62%). En comparaison cet écart n'est que de 1 point sur l'ensemble de chaque corpus.

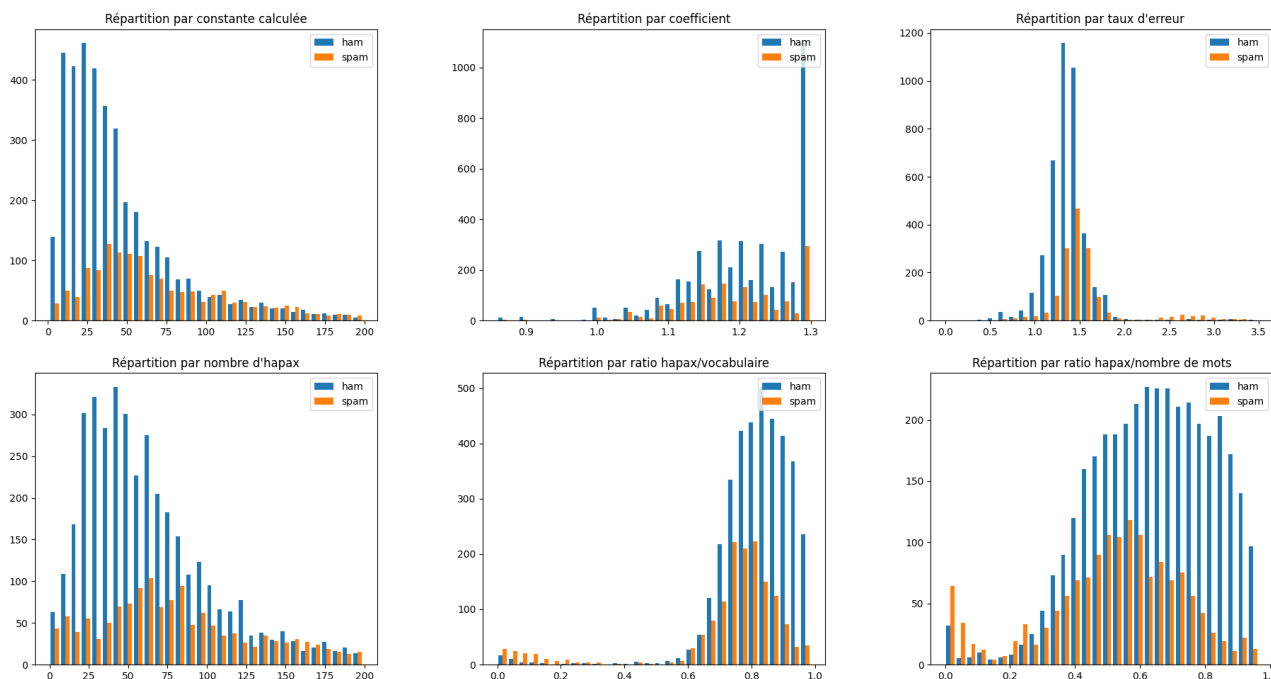


FIGURE 13 – Distribution des données de Zipf selon la catégorie du mail

On voit dans les schéma de la Figure ?? que les Ham se concentre également dans les valeurs basses de la constante. Les Spams semblent plus dispersés pour cette métrique.

Le nombre d'hapax lui semble suivre une distribution normale pour les Ham. Pour cette métrique la répartition est plus diffuse.

Les ratio du nombre d'hapax par rapport au vocabulaire et à l'ensemble du corps du mail ont une distribution normale sauf pour des valeurs marginales entre 0% et 20%. Il est intéressant de remarquer que les spam sont majoritaires dans cette tranche alors que le nombre de Spam est nettement inférieur au nombre de Ham dans le data set (70/30).

Le pic du taux d'erreur des Ham (1.4) est effectivement inférieur à celui des Spam (1.5). On remarque également un regroupement de Spam avec un taux d'erreur compris entre (2.3 et 3).

Enfin il semblerait que le coefficient de chaque type gravite autour de la médiane, avec base assez large. Le pic en bout de graphique peut être ignoré car il correspond à la limite haute de la recherche de coefficient bornée entre 0.8 et 1.30. augmenter cette limite impliquerait d'augmenter les temps que calcul qui ne semble pas justifié dans ce cas. Moins de 25% des documents ont été classés avec un coefficient supérieur à 1.29.

Matrice de corrélation La corrélation entre 2 variables permet de voir l'interdépendance entre elles. Il est alors possible de voir si leurs augmentations ou diminutions sont de nature à être liées. Cette relation peut dans certain cas montrer une relation de cause à effet.

Le calcul des coefficients présentés dans la matrice de la Figure ?? utilise la méthode de corrélation linéaire de Pearson. Il s'agit ici de calculer le ratio (r) de la covariance entre 2 variables par le produit de leurs écart-types.

$$r_{X,Y} = \frac{cov(X,Y)}{\sigma_X \sigma_Y}$$

$$r_{X,Y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Une corrélation forte sera proche des bornes 1 et -1 . Une corrélation faible sera proche de 0.

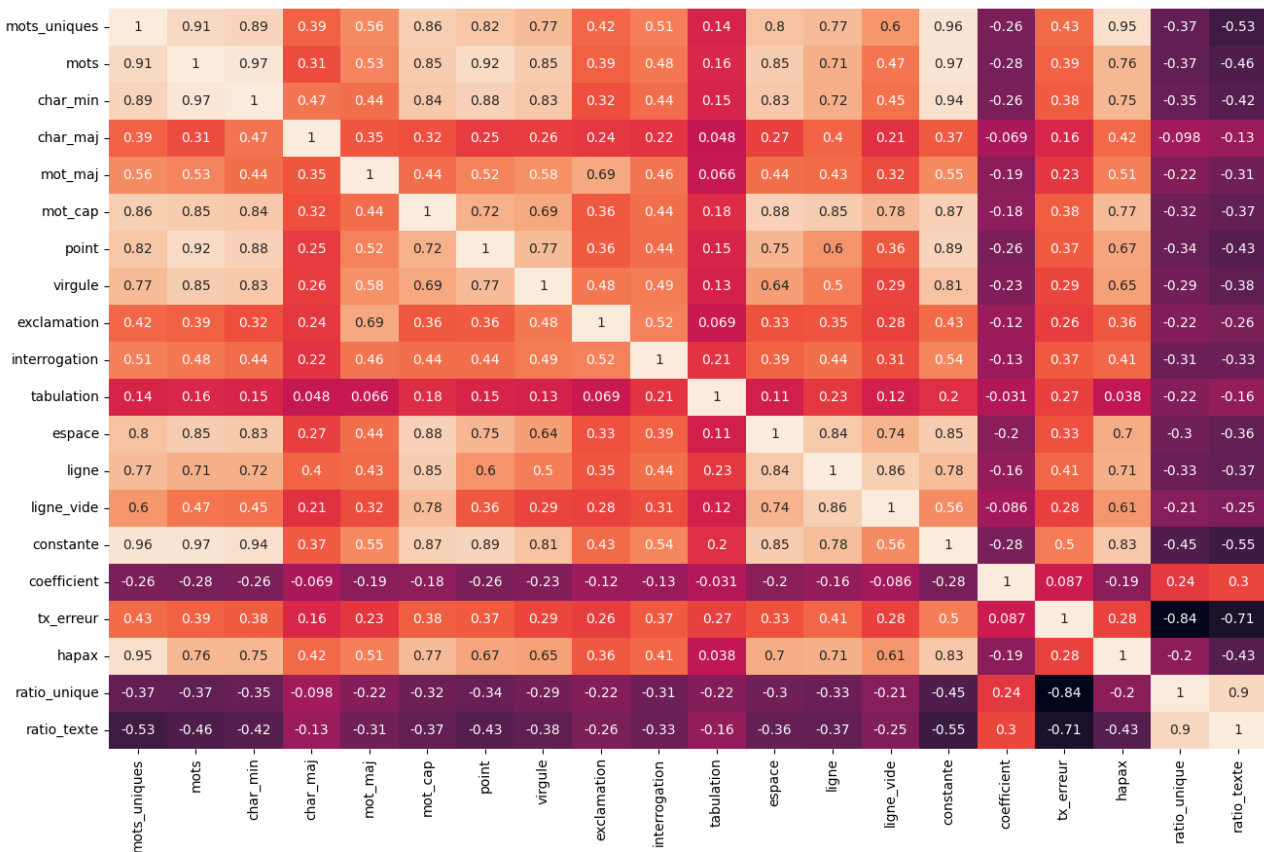


FIGURE 14 – Matrice de corrélation des métriques

On peut remarquer que les nombre de mots, de caractères en minuscule, de points et de virgules sont fortement corrélés. On voit également que certaines métriques sont quasiment indépendante de la plus part des autres variables :

- nombre de caractères en majuscules
- nombre de mots en majuscules
- nombre de point d'exclamation
- nombre de point d'interrogation
- nombre de tabulation
- coefficient Zipf déterminé
- ratio d'hapax dans le vocabulaire
- ratio d'hapax dans tout le texte

Les variables dénombrant les espaces et les lignes (vides ou non) sont faiblement liées aux signes de punctuations.

Résumé L’analyse de l’utilisation des données des mots permet de mettre en lumière les points suivants :

- Les Spam sont généralement plus long que les Ham
- On retrouve plus de caractères en majuscules dans les Spam que dans les Ham
- Les nombres de mots et de mots uniques sont fortement corrélés

Les données sur l’utilisation des ponctuations ne montre pas d’écart majeur entre les Ham et les Spam sur l’utilisation des points et des virgule. Cependant l’utilisation des points d’exclamation est nettement plus forte dans les Spam.

L’utilisation de point d’exclamation et de majuscule est susceptible d’induire une notion d’urgence chez le lecteur.

Les données des lignes et des espaces ne permettent de faire une différence franche entre les Ham et les Spams au regards des autres données.

Globalement, les Ham et les Spam respectent les principes de la distribution de Zipf. Le coefficient est identique pour chacun des deux corpus. La constante et le nombre d’Hapax plus importantes pour les Ham s’expliquent par un nombre de document largement supérieur.

En regardant les données de la distribution de Zipf appliquée à chaque mail on remarque un coefficient moyen plus élevé et qui semble plus précis pour les Ham. Il est également notable que les mots se répètent moins dans les messages Ham.

En conclusion, on peut supposer les faits suivants :

- Le vocabulaire d’un Ham est plus fourni que celui d’un Spam.
- Il y a plus de mots dans un Spam.
- Un Spam aura probablement une construction induisant une notion d’urgence.

2.4 Traitement du langage

Dans cette section nous aborderons des traitements du langage naturel qui vont être employés afin de tenter de détecter des similitudes dans les thèmes abordés dans les mails de chaque catégorie. Pour cela nous allons transformer les messages en un ensemble de clé/valeur avec le mot ainsi que sa fréquence dans le document. Les données récoltées seront utilisées pour vectoriser les document avec la méthode TF-IDF.

Schéma d’exécution Les corps de mail brut vont être récupérés dans la base de données ElasticSearch.

La première étape vise à vérifier que le message est déjà passé par l’étape de traitement de la phase 1 et qu’il est présent dans la base de données PSQl. Pour le moment, l’absence du message dans la base PSQl termine le traitement NLP du message. Une évolution possible serait de déclencher automatiquement les traitements précédents.

La deuxième étape permet d’éviter que le traitement d’un même message s’effectue plusieurs fois. cela est possible en vérifiant la présence de l’identifiant de message dans la table *nlp_status*. Dans l’étape de lemmatisation nous allons utiliser le moteur de *StanfordNLP* [?] [citeq2020stanza]. Les processus utilisés sont décrits plus bas.

Le calcul de la fréquence des mots utilise le code développé lors la recherche de la distribution de Zipf. A l’issue de cette phase les données de fréquence pour chaque mots et messages sont mis en base.

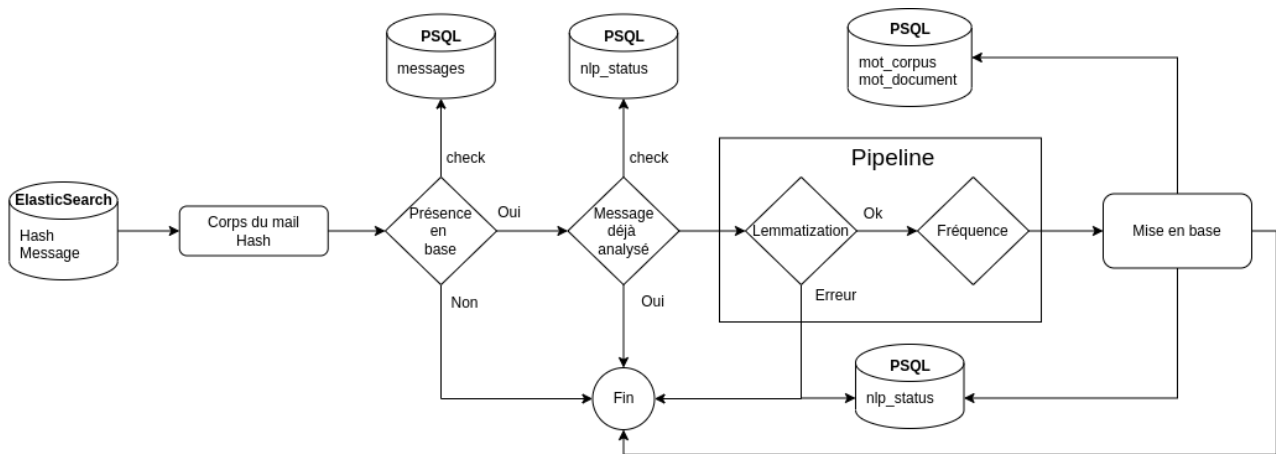


FIGURE 15 – Schéma du déroulement des instructions

Modification de la base de données Afin de stocker les informations générées pendant cette phase 3 nouvelles tables vont être ajoutée à la base de données PSQL existante :

1. **mot_corpus** - Cette table regroupe les mots du corpus
 - id_mot - clé primaire du mot dans la base
 - mot - valeur du mot
 - freq_corpus - fréquence du mot dans le corpus
 - freq_doc_all - nombre total de document avec le mot
 - freq_doc_spam - nombre de spam avec le mot
 - freq_doc_ham - nombre de ham avec le mot
2. **mots_document** - Table faisant la liaison entre les mot et les messages
 - id_message - identifiant unique du message dans la base
 - id_mot - identifiant unique du mot dans la base
 - occurence - nombre d'occurrence de mot dans le message
3. **nlp_status** - enregistrement de résultat du statut du traitement NLP
 - id_message - identifiant unique du message dans la base
 - success - réussite ou échec
 - raison - raison de l'échec

La table **nlp_status** va permettre de ne pas traiter plusieurs fois le même message. Le module de traitement NLP standford est susceptible de généré des erreurs lors du traitement de certains messages. Ces messages seront exclus du traitement.

Les erreurs rencontrées pour la fonction lemmatise sont :

- TypeError - Provient de la présence de données non textuelle dans le corps du mail (ex MIME)
- OutOfMemoryError - Présence de mot beaucoup trop long pour être parsé.

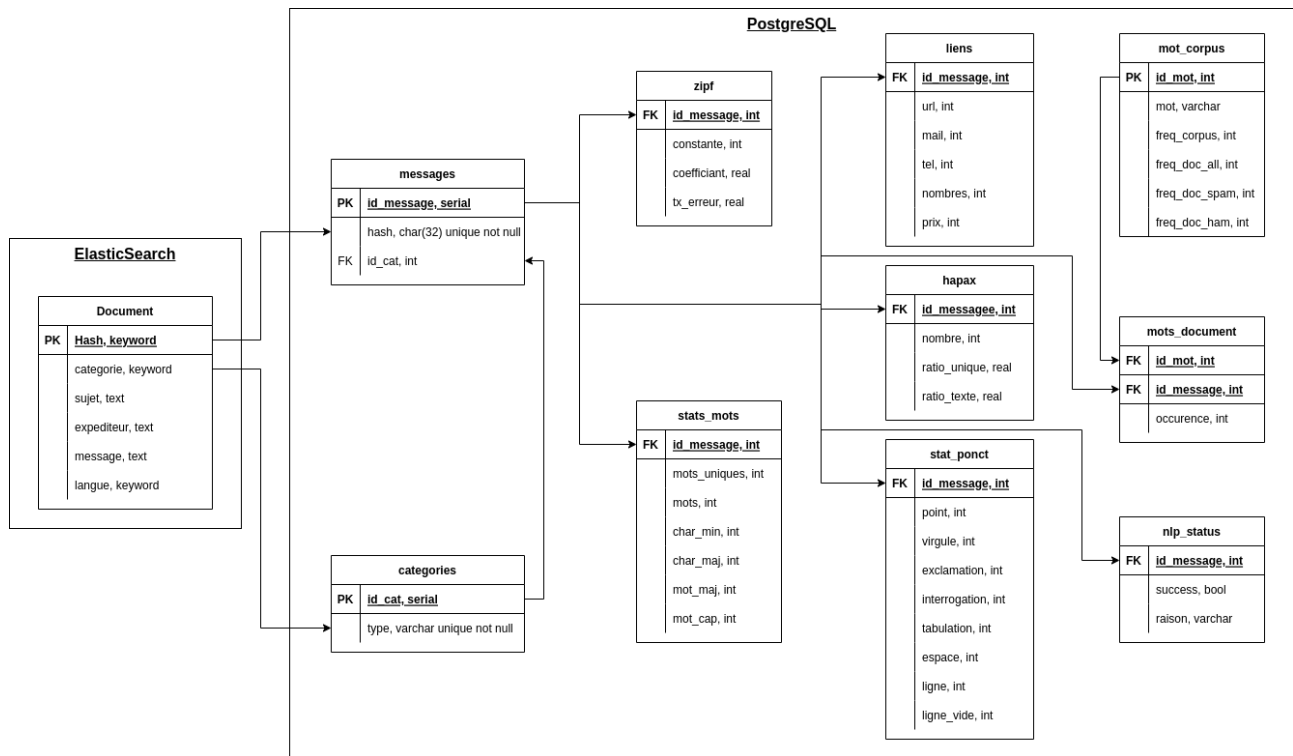


FIGURE 16 – Schéma des bases de données pour accueillir les données NLP

Les nouvelles tables ont été ajoutées en utilisant les fonctions présentées dans l'annexe B.2 avec le mapping suivant :

Mapping nouvelles tables

```

1 {
2   "mail_features_prod": {
3     "mot_corpus": {
4       "id_mot": ["SERIAL", "PRIMARY KEY"],
5       "mot": ["VARCHAR", "UNIQUE", "NOT NULL"],
6       "freq_corpus": ["INT"],
7       "freq_doc_all": ["INT"],
8       "freq_doc_spam": ["INT"],
9       "freq_doc_ham": ["INT"]
10    },
11    "mots_document": {
12      "id_message": ["INT"],
13      "id_mot": ["INT"],
14      "occurrence": ["INT"],
15      "pk": ["id_message", "id_mot"],
16      "fk": {
17        "fk_message": ["id_message", "messages(id_message)", "SET NULL"],
18        "fk_mot": ["id_mot", "mot_corpus(id_mot)", "SET NULL"]
19      }
20    },
21    "nlp_status": {
22      "id_message": ["INT"],
23      "success": ["BOOL"],
24      "raison": ["VARCHAR"],
25      "pk": ["id_message"],
26      "fk": {

```

```

27         "fk_message": ["id_message", "messages(id_message)", "SET NULL"]
28     }
29 }
30 }
31 }
32

```

2.4.1 Lemmatisation

La lemmatisation d'un texte vise à réduire la taille d'un texte en ramenant chaque mot à la forme du mot présente dans le dictionnaire. Ce traitement permet donc de compter le nombre d'occurrence d'un mot sans se soucier de sa forme ou des modifications grammaticales appliquées lors de la rédaction du texte. A l'issue de ce traitement nous serons en mesure de connaître les mots les plus utilisés dans les corps des mails et tenter de déterminer les thèmes les plus récurrents. Dans ce projet, nous allons utiliser le moteur de lemmatisation de *StanfordNLP*[?][4] Les traitements nécessaires pour arriver à une lemmatisation sont :

1. Tokenisation : Séparer les phrases en token. Ici chaque token correspond à un mot ou à une ponctuation.
2. Multi-Word Token Expansion : Ce traitement n'est pas nécessaire en anglais selon la documentation *Stanza*. Ce traitement permet d'étendre un token s'il correspond à une contraction de plusieurs termes. Par exemple en français le terme *du* sera transformé en *de le*.
3. Part of Speech Tagging : Cette opération vise à marquer chaque mot du texte avec sa position grammaticale dans la phrase qui le contient et va permettre d'effectuer le transfert vers le lexème correspondant.

L'initialisation de la pipeline de lemmatisation de Stanza se fait en appelant la classe *Pipeline* avec les arguments de langage et de processus dans l'ordre d'utilisation. A cette étape j'ajoute 2 filtres. Le premier va me permettre de ne conserver que les mots avec des lettres. La deuxième étape permet de retirer les stopwords en anglais qui ne permettent pas d'analyser le sens du message.

La dernière étape avant la mise en base permet de compter la fréquence de chaque mot dans le texte

Pipeline de lemmatisation

```

1  import nltk
2  import stanza
3  import stats
4
5  nltk.download("stopwords")
6  en_stopwd = set(stopwords.words('english'))
7  pipe = stanza.Pipeline(lang='en', processors='tokenize,mwt,pos,lemma')
8  pat = re.compile(r'\w+')
9
10 def lemmatise(message, stopwds, pipeline, pattern):
11     doc = pipeline(message)
12     lemma = [mot.lemma for phrase in doc.sentences for mot in phrase.
13              words]
14     return [lem.lower() for lem in lemma if re.match(pattern, lem) and
15             lem.lower() not in stopwds]

```

15 stats.frequence_mot(lemmatise(value, en_stopwd, pipe, pat))

16

Ci dessous un exemple des étapes du traitement sur un message de la base

```
>>> print(value)
```

```
Heres the hottest thing in DVDs. Now you can make a personal backup
copy of a DVD right onto CDR. Our Hot new software easily takes you through
the steps to make a copy of your own DVDs.
```

```
NOW INCLUDED FOR FREE! Copy PLAYSTATION, MUSICMPs and all Software.
```

```
Step by Step Interactive Instructions
```

```
All Software Tools Included On CD
```

```
No DVD Burner Required
```

```
FREE Live Technical Support
```

```
Day Risk Free Trial Available
```

```
FREE Dvd Movie of your choice LIMITED TIME OFFER!
```

```
We have All the software you need to COPY your own DVD Movies.
```

```
This email has been screened and filtered by our in house OPTOUT system in
compliance with state laws. If you wish to OPTOUT from this mailing as well
as the lists of thousands of other email providers please visit
```

```
ZKZmblanzxaDBmNTTicorikgl
```

```
>>> lemmatise(value, en_stopwd, pipe, pat)
```

```
['hot', 'thing', 'dvd', 'make', 'personal', 'backup', 'copy', 'dvd', 'right', 'onto',
'cdr', 'hot', 'new', 'software', 'easily', 'take', 'step', 'make', 'copy', 'dvd',
'include', 'free', 'copy', 'playstation', 'musicmps', 'software', 'step', 'step',
'interactive', 'instruction', 'software', 'tool', 'include', 'cd', 'dvd', 'burner',
'require', 'free', 'live', 'technical', 'support', 'day', 'risk', 'free', 'trial',
'available', 'free', 'dvd', 'movie', 'choice', 'limited', 'time', 'offer',
'software', 'need', 'copy', 'dvd', 'movie', 'email', 'screen', 'filter', 'house',
'optout', 'system', 'compliance', 'state', 'law', 'wish', 'optout', 'mailing',
'well', 'list', 'thousand', 'email', 'provider', 'please', 'visit',
'zkzmbblanzxadbmntticorikgl']
```

```
>>> stats.frequence_mot(lemmatise(value, en_stopwd, pipe, pat))
```

```
{'hot': 2, 'thing': 1, 'dvd': 6, 'make': 2, 'personal': 1, 'backup': 1, 'copy': 4,
'right': 1, 'onto': 1, 'cdr': 1, 'new': 1, 'software': 4, 'easily': 1, 'take': 1,
'step': 3, 'include': 2, 'free': 4, 'playstation': 1, 'musicmps': 1,
'interactive': 1, 'instruction': 1, 'tool': 1, 'cd': 1, 'burner': 1,
'require': 1, 'live': 1, 'technical': 1, 'support': 1, 'day': 1, 'risk': 1,
'trial': 1, 'available': 1, 'movie': 2, 'choice': 1, 'limited': 1, 'time': 1,
'offer': 1, 'need': 1, 'email': 2, 'screen': 1, 'filter': 1, 'house': 1,
'optout': 2, 'system': 1, 'compliance': 1, 'state': 1, 'law': 1, 'wish': 1,
'mailing': 1, 'well': 1, 'list': 1, 'thousand': 1, 'provider': 1, 'please': 1,
'visit': 1, 'zkzmbblanzxadbmntticorikgl': 1}
```

Après cette phase de traitement nous avons les données numériques suivantes :

Le nombre de mail par catégorie

```
SELECT c.type, COUNT(*) FROM nlp_status as n
JOIN messages as m ON n.id_message = m.id_message
```

```
JOIN categories as c ON m.id_cat = c.id_cat
WHERE n.success = true
GROUP BY c.type;
```

type	count
spam	1453
ham	3918

Le nombre de mots uniques et total de mots dans le corpus

```
SELECT COUNT(mot) as mots_uniques, SUM(freq_corpus) as mots FROM mot_corpus;
```

mots_uniques	mots
39392	687121

Le nombres de mots uniques et globalement par catégorie

```
SELECT c.type, COUNT(m.id_mot) as mots_uniques, SUM(m.occurrence) as mots
FROM mots_document as m
JOIN messages as mess ON m.id_message = mess.id_message
JOIN categories as c ON mess.id_cat = c.id_cat
GROUP BY c.type;
```

type	mots_uniques	mots
ham	285831	424261
spam	149430	262860

Afin de préparer l'étape suivante qui est celle de la vectorisation, nous allons récupérer une liste de mots répondant aux caractéristiques suivantes :

- X mots les plus présents dans le corpus
- X mots les plus présents dans chaque catégorie
- X mots moins présents dans une catégorie, mais très présent dans l'autre
- X mots au moins deux fois plus présent dans une catégorie que dans l'autre
- X mots avec un ratio d'apparition fort dans une catégorie par rapport à l'autre

Pour cela nous allons effectuer les requêtes suivantes dans la base PSQL.

```
SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
ORDER by freq_corpus DESC
LIMIT 200;
```

```
SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
ORDER by freq_doc_spam DESC
LIMIT 200;
```

```
SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
ORDER by freq_doc_ham DESC
LIMIT 200;
```

```
SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
ORDER by freq_doc_ham, freq_doc_spam DESC
LIMIT 200;
```

```
SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
ORDER by freq_doc_spam, freq_doc_ham DESC
LIMIT 200;
```

```
SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
WHERE freq_doc_ham >= 2*freq_doc_spam
ORDER BY freq_doc_ham DESC
LIMIT 200;
```

```
SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
WHERE freq_doc_spam >= 2*freq_doc_ham
ORDER BY freq_doc_spam DESC
LIMIT 200;
```

```
SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham,
freq_doc_ham/freq_doc_spam as "ratio ham/spam"
FROM mot_corpus
WHERE freq_doc_ham > 0 AND freq_doc_spam > 0
ORDER BY "ratio ham/spam" DESC
LIMIT 200;
```

```
SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham,
freq_doc_spam/freq_doc_ham as "ratio spam/ham"
FROM mot_corpus
WHERE freq_doc_ham > 0 AND freq_doc_spam > 0
ORDER BY "ratio spam/ham" DESC
LIMIT 200;
```

Le résultat de ces commandes détaillé dans l'annexe ??.

En limitant les résultats à 200 lignes pour chaque requête, j'arrive à obtenir une liste de plus de 1000 mots uniques. J'ai pris la liberté de retirer des mots très spécifiques au dataset (spamassinsightings, deathspamdeathspamdeathspam) ainsi que tous les mots des moins de 2 caractères. La liste complète des mots retenus est également disponible dans l'annexe ??.

2.4.2 Vectorisation TF-IDF

La vectorisation d'un texte doit permettre de transposer les données textuelles en données numériques afin de pouvoir les utiliser dans des calculs statistiques ou dans les modèles d'apprentissage automatiques.

J'ai choisi d'utiliser une vectorisation TF-IDF^[1] (Term Frequency-Inverse Document Frequency). Cette méthode se rapproche de la distribution de Zipf détaillée précédemment. En effet le score

d'un terme est dépendant de la fréquence dans le document et de la fréquence de ce terme dans l'ensemble du corpus. La formule utilisée est la suivante :

$$W_{i,j} = tf_{i,j} * \log\left(\frac{N}{df_i}\right)$$

Avec $tf_{i,j}$ la fréquence d'apparition du mot i dans le document j , N le nombre de document dans le corpus et df_i le nombre de document dans le corpus contenant le terme i .

Cette méthode permet d'attribuer des score plus élevé aux mots apparaissant plus dans un document que dans les autres documents du corpus. Nous pouvons ainsi conserver une forme de contexte d'utilisation.

Lors de la phase précédente de lemmatisation, nous avons commencé à préparer les données pour ce calcul. En effet dans la table *mots_document* nous avons stocké la fréquence de chaque mot dans chaque document et dans la table *mot_corpus* se trouve le nombre de document contenant chaque mot.

Le vecteur de chaque document sera conservé dans une table *tfidf_vector* créée à la volée avec un nombre de colonnes variable selon le nombre de mots retenus pour à l'étape précédente. Une table annexe (*tfidf_assoc* permet de conserver une trace de l'association entre les labels des colonnes des vecteur et les identifiants des mots dans la base.

Le code ci-dessous présente le code de création des ces tables :

Préparation de la base vectorielle

```

1  tfidf_assoc_field = {"id_mot": ['INT', 'UNIQUE', 'NOT NULL'],
2                        "label": ['VARCHAR'],
3                        "fk": {"fk_mot": ['id_mot', 'mot_corpus(id_mot)', '
    CASCADE']}]}}
4  psql_cmd.create_table(psql_cli, 'tfidf_assoc', tfidf_assoc_field)
5
6  tfidf_vector_fields = {"id_message": ['INT', 'UNIQUE', 'NOT NULL'],
7                        "fk": {"fk_message": ['id_message', 'messages(
    id_message)', 'CASCADE']}]}}
8
9  n_label = 0
10 for mot in mots:
11     result = psql_cmd.get_data(psql_cli, 'mot_corpus', ['id_mot'], f"mot
    LIKE '{mot}'")
12     if not result:
13         print(f"id_mot non recupere pour {mot}", file=sys.stderr)
14         continue
15
16     label = f"feat_mot_{n_label}"
17     data = {'id_mot': result[0]['id_mot'], 'label': label}
18     psql_cmd.insert_data(psql_cli, 'tfidf_assoc', data)
19
20     tfidf_vector_fields[label] = ['DECIMAL']
21     n_label += 1
22
23 psql_cmd.create_table(psql_cli, 'tfidf_vector', tfidf_vector_fields
24
```

Le processus de vectorisation est réalisé avec le code ci-dessous

```

1 def tfidf_vectorise(client_psql, id_message, nb_documents):
2     vector = {'id_message': id_message}
3     associations = psql_cmd.get_data(client_psql, 'tfidf_assoc', ['id_mot', 'label'])
4
5     for entry in associations:
6         id_mot = entry.get('id_mot')
7         col_label = entry.get('label')
8
9         res = psql_cmd.get_data(client_psql, 'mots_document', ['occurrence'],
10                                f' id_message = {id_message} AND id_mot = {id_mot}'))
11         term_freq = res[0].get('occurrence', 0) if res else 0
12
13         res = psql_cmd.get_data(client_psql, 'mot_corpus', ['freq_doc_all'],
14                                f' id_mot = {id_mot}'))
15         doc_freq = res[0].get('freq_doc_all', None) if res else None
16
17         cell_value = term_freq * math.log(nb_documents/doc_freq) if doc_freq else 0
18         vector[col_label] = cell_value
19
20     return vector
21
22 # — IN MAIN —
23 result = psql_cmd.exec_query(psql_cli, "SELECT COUNT(id_message) FROM nlp_status WHERE "
24                                "success = true")
25 n_docs = result[0][0]
26
27 result = psql_cmd.get_data(psql_cli, 'nlp_status', ['id_message'], 'success = true')
28 for row in tqdm.tqdm(result, desc='— Vectorisation', leave=False, file=sys.stdout, ascii=True):
29     id_mess = row.get('id_message')
30     vecteur = tfidf_vectorise(psql_cli, id_mess, n_docs)
31     psql_cmd.insert_data(psql_cli, 'tfidf_vector', vecteur)
32

```

Analyse de la vectorisation A l'issue de la vectorisation nous avons 7 messages (6 hams et 1 spam) qui ont un vecteur null, c'est à dire que la somme de toutes les colonnes est égale à zéro.

Le tableau ci dessous présente les caractéristiques des messages avec un vecteur null.

id	type	point	exclamation	espace	ligne	mots	char min	char maj	nombre
974	spam	0	70	10	73	108	1267	1346	46
1594	ham	1	0	19	5	20	88	7	0
3698	ham	1	0	6	6	7	39	1	0
3779	ham	1	0	20	9	21	84	5	1
4470	ham	4	0	19	9	18	87	2	0
4592	ham	0	0	12	5	10	40	2	2
1677	ham	2	0	7	21	9	27	2	0

Le nombre de message étant assez petits, ils ne seront pas sortis pour l'entraînement et l'évaluation du futur modèle.

Les mots utilisés dans ces messages sont :

id message	mot	freq corpus	freq doc	freq spam	freq ham
1677	sex	162	91	44	47
1677	mr	71	52	20	32
1677	woman	298	150	62	88
4592	kill	122	99	13	86
4592	pid	6	5	0	5
4592	etchttpd	2	1	0	1
4470	become	461	293	109	184
4470	fuck	18	12	4	8
4470	bush	196	44	3	41
4470	rest	159	139	49	90
4470	nation	257	74	22	52
4470	cheney	38	9	0	9
4470	suddenly	30	27	6	21
4470	eventuality	1	1	0	1
4470	unavailable	11	10	3	7
4470	gg	13	13	0	13
3779	c	529	295	69	226
3779	matter	336	225	105	120
3779	x	452	235	59	176
3779	hardware	222	149	10	139
3779	transform	29	15	0	15
3779	java	209	69	7	62
3779	xml	176	34	1	33
3779	slow	106	74	9	65
3779	parsing	19	15	0	15
3779	outofthebox	2	2	0	2
3698	artist	80	59	9	50
3698	draw	110	78	28	50
3698	numerous	43	37	7	30
3698	character	100	54	1	53
3698	cliff	7	7	0	7
1594	national	257	148	77	71
1594	similar	235	185	49	136
1594	review	382	189	43	146
1594	rag	6	6	0	6
1594	usual	56	48	4	44
1594	crud	4	4	0	4
1594	moron	8	7	0	7
1594	rant	18	16	0	16
1594	beat	86	72	31	41
1594	chest	11	10	2	8
1594	merit	24	13	5	8
1594	forking	3	3	0	3
974	hrefshaved	1	1	1	0
974	sweetiesa	1	1	1	0

3 Phase 3

A Développement visualisation distribution de Zipf

Présentation La loi de distribution de Zipf est une loi empirique (basée sur l'observation) qui veut que le mot le plus fréquent est, à peu de chose près, 2 fois plus fréquent que le 2^{ème}, 3 fois plus fréquent que le 3^{ème} etc.

La formulation finale de la 1^{ère} loi de Zipf est la suivante :

$$|mot| = constante \times rang(mot)^{k \approx 1}$$

avec $|mot|$ la fréquence d'apparition d'un mot, *constante* une valeur propre à chaque texte, $rang(mot)$ la place du mot dans le tri décroissant par fréquence d'apparition et k un coefficient proche de 1.

Développement Afin de pouvoir utiliser les résultats de cette distribution dans ce projet, j'ai développé un ensemble de fonctions sur un corpus "*reconnu*". Mon choix s'est porté sur le corpus *Brown* (voir G.1) présent dans la librairie *nltk*. Ce corpus contient environ 500 documents contenant 1 millions de mot en anglais.

Le processus d'analyse se fait sur 2 versions de ce corpus.

— la première version contient tous les mots sans modifications

— la seconde version contient tous les mots sans les *stopwords*

Les *stopwords* sont des mots qui n'ont pas ou peu de signification dans un texte. Ces mots sont retirés dans la 2^e version pour voir l'effet d'une réduction sur la distribution de Zipf.

Les paragraphes ci-dessous détaillent les étapes du développement :

Étape 1 - Ordonner les mots La première étape est de compter les occurrences de tous les mots des 2 corpus et de les ranger en fonction de leur nombre d'occurrence.

Triage des mots

```
1 def frequence_mot(bag, freq=None):
2     """
3     Calcule la frequence de chaque mot dans un sac de mot
4     :param bag: <list> - liste de tous les mots d'un texte
5     :param freq: <dict> - dictionnaire avec {<str> mot: <int> frequence}
6     :return: <dict> - dictionnaire avec la frequence par mot {mot:
7     frequence}
8     """
9     if freq is None:
10         freq = {}
11     for mot in bag:
12         freq[mot] = freq.get(mot, 0) + 1
13     return freq
14
15 def classement_zipf(dico):
16     """
17     Trie un dictionnaire de mots : occurrence et leur assigne un rang en
18     fonction du nombre d'occurrence
19     :param dico: <dict> dictionnaire de mot: occurrences
20     :return: <list> {"rang": <int>, "mot": <str>, "frequence": <int>}
```

```

19     """
20     ranked = []
21     for rang, couple in enumerate(sorted(dico.items(), key=lambda item:
22         item[1], reverse=True), start=1):
23         ranked.append({"rang": rang,
24             "mot": couple[0],
25             "frequence": couple[1]})
26     return ranked

```

On obtient les représentations suivantes :

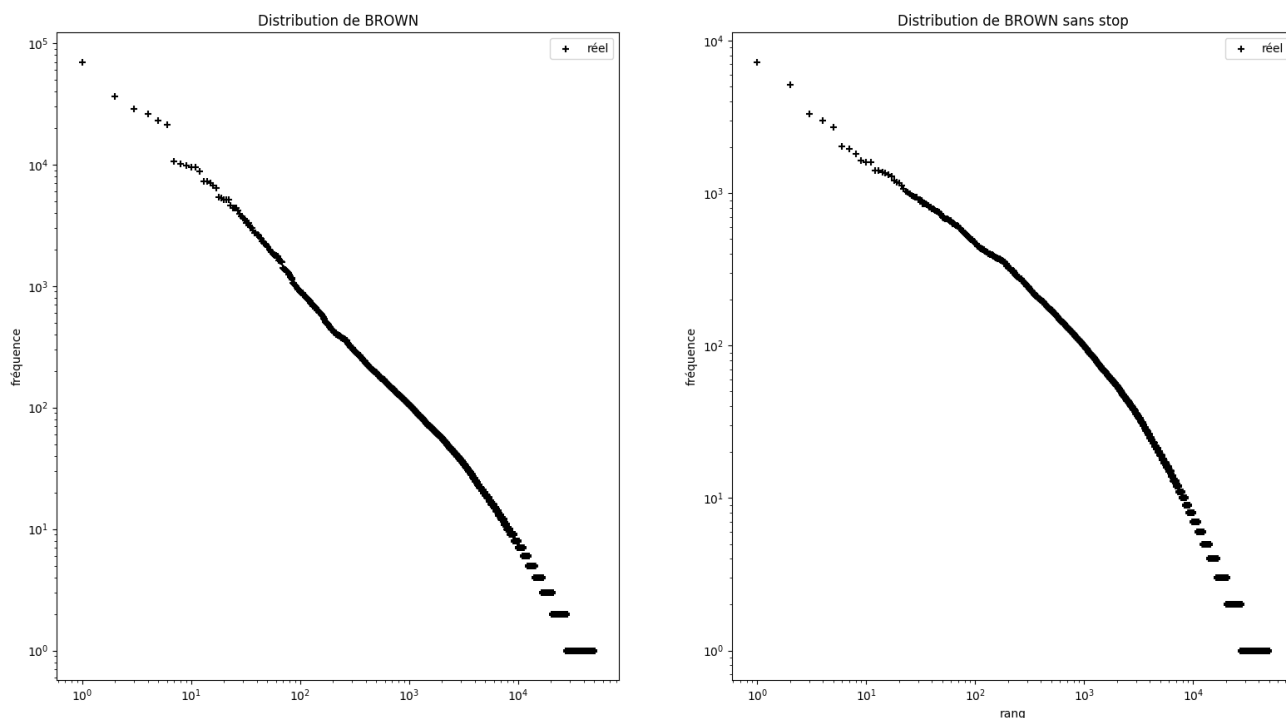


FIGURE 17 – Distribution de Zipf pour les deux corpus

- Nombre de mots dans brown : mots : 49398 occurrences : 1012528
- Nombre de mots dans brown stop : mots : 49383 occurrences : 578837

La distribution de la version complète du corpus semble à première vue plus fidèle à la représentation classique de la distribution de Zipf.

Etape 2 - calcul de la constante Le premier paramètre qu'il faut déterminer est la *constante*. Pour ce faire j'effectue le calcul suivant pour tous les mots :

$$constante = |mot| \times rang(mot)$$

On obtient une liste de toutes les constantes théoriques pour chaque mot selon son rang. De cette liste, nous allons extraire la moyenne et la médiane.

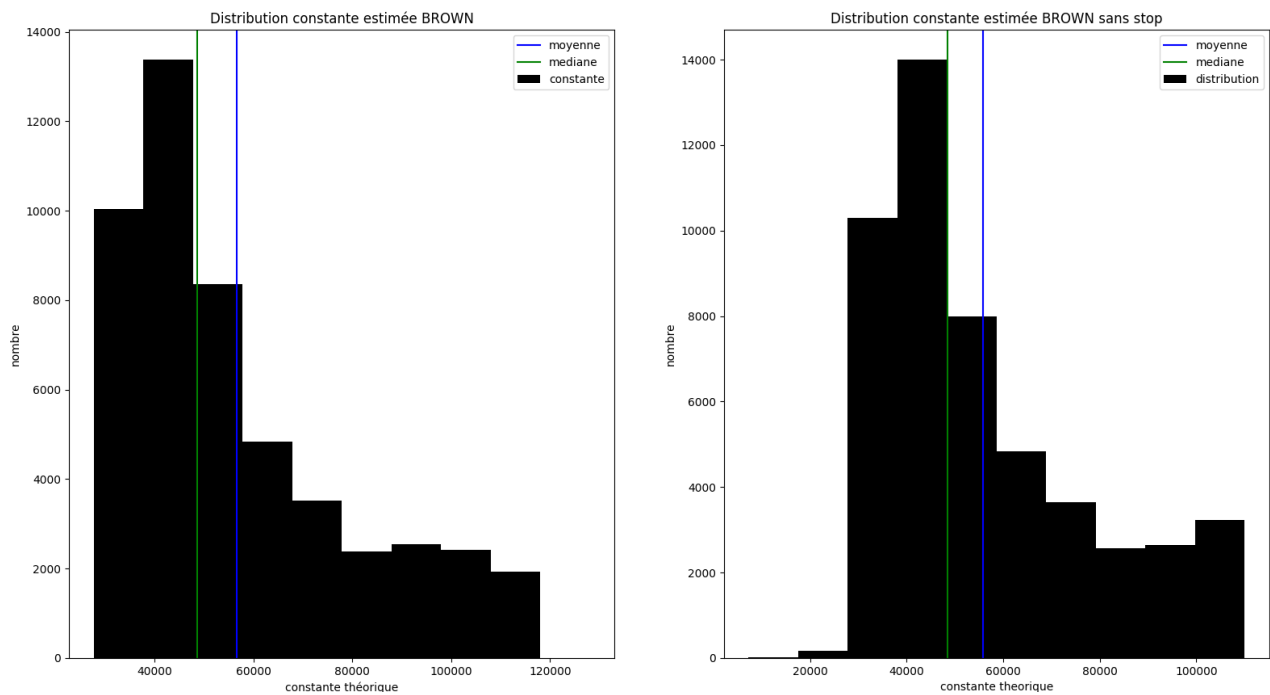


FIGURE 18 – Distribution des constantes théoriques pour les deux corpus

On voit qu'il y a une majorité de mots donnant une constante brute comprise entre 20.000 et 60.000. Dans les deux corpus La différence entre les moyennes et médianes des deux corpus n'est pas flagrante :

- Brown moyenne : 56525.81, médiane : 48601.50
- Brown (- stopwords) moyenne : 55809.97, médiane : 48494.00

Etape 3 - recherche du coefficient Le coefficient k permet d'ajuster le résultat, et pourra éventuellement donner une indication de complexité. La recherche de k se fera sur les deux corpus avec utilisant les moyennes et médianes.

Pour ce faire nous allons :

1. Faire la liste de tous les coefficients possibles dans l'intervalle $[0.86, 1.3]$ avec un pas de 0.01^2 .
2. Calculer toutes la fréquences théoriques de tous les rangs avec tous les coefficients possibles en utilisant les constantes moyenne et médiane de chaque corpus.
3. Calculer la moyenne des coûts absolus entre les fréquences théoriques par coefficient avec la fréquence réelle observée pour chaque corpus.

Le couple coefficient/constante avec le coup minimal sera retenu pour l'utilisation dans la phase de *feature engineering*.

Fonctions utilisées dans la recherche du coefficient

```

1 def zipf_freq_theorique(constante, rang, coef):
2     """
3     Calcul la frequence theorique d'un mot selon son rang, la constante
    du texte et un coeficiant d'ajustement

```

2. les bornes et le pas sont totalement arbitraire afin d'obtenir un graphique présentable

```

4      :param constante: <int> constante determinee par la distribution de
      Zipf
5      :param rang: <int> rang du mot selon sa frequence
6      :param coef: <float> variable d'ajustement
7      :return: <float> frequence theorique zipfienne
8      """
9      return constante / (rang ** coef)
10
11 def cout(l1, l2, methode):
12     """
13     Calcul le cout de l'ecart entre les elements de l1 et le l2, place
14     par place
15     :param l1: <list> liste d'entier
16     :param l2: <liste> liste d'entier
17     :param methode: <str> methode de calcul du cout
18     :return: <float> cout selon methode
19     """
20     if len(l1) != len(l2):
21         print("Erreur, fonction cout: l1 & l2 de taille differente", file
22             =sys.stderr)
23         return None
24
25     if len(l1) == 0:
26         print("Erreur, fonction cout: liste vide", file=sys.stderr)
27
28     if methode.lower() not in ['absolue', 'carre', 'racine']:
29         print("Erreur, fonction cout - methode '{}' inconnue".format(
30             methode), file=sys.stderr)
31         return None
32
33     if methode.lower() == 'absolue':
34         return np.mean([abs(x-y) for x, y in zip(l1, l2)])
35
36     if methode.lower() == 'carre':
37         return np.mean([(x-y)**2 for x, y in zip(l1, l2)])
38
39     if methode.lower() == 'racine':
40         return np.sqrt(np.mean([(x-y)**2 for x, y in zip(l1, l2)]))
41
42     return None

```

Calcul des fréquences par coefficient

```

1  ls_coef = list(np.arange(0.86, 1.3, 0.01))
2  zbmo_th = {coef: [stats.zipf_freq_theorique(zb_const_moyen, r, coef)
3  for r in zb_rang] for coef in ls_coef}
4  zbme_th = {coef: [stats.zipf_freq_theorique(zb_const_median, r, coef)
5  for r in zb_rang] for coef in ls_coef}
6  zbmoth_cmoy = [stats.cout(zb_freq, zbmo_th[coef], 'absolue') for coef
7  in ls_coef]
8  zbmeth_cmoy = [stats.cout(zb_freq, zbme_th[coef], 'absolue') for coef
9  in ls_coef]
10
11 zbsmo_th = {coef: [stats.zipf_freq_theorique(zbs_const_moyen, r, coef)

```

```

) for r in zbs_rang] for coef in ls_coef}
8   zbsme_th = {coef: [stats.zipf_freq_theorique(zbs_const_mediane, r,
coef) for r in zbs_rang] for coef in ls_coef}
9   zbsmoth_cmoy = [stats.cout(zbs_freq, zbsmo_th[coef], 'absolue') for
coef in ls_coef]
10  zbsmeth_cmoy = [stats.cout(zbs_freq, zbsme_th[coef], 'absolue') for
coef in ls_coef]

```

La recherche du coefficient nous retourne les éléments suivants :

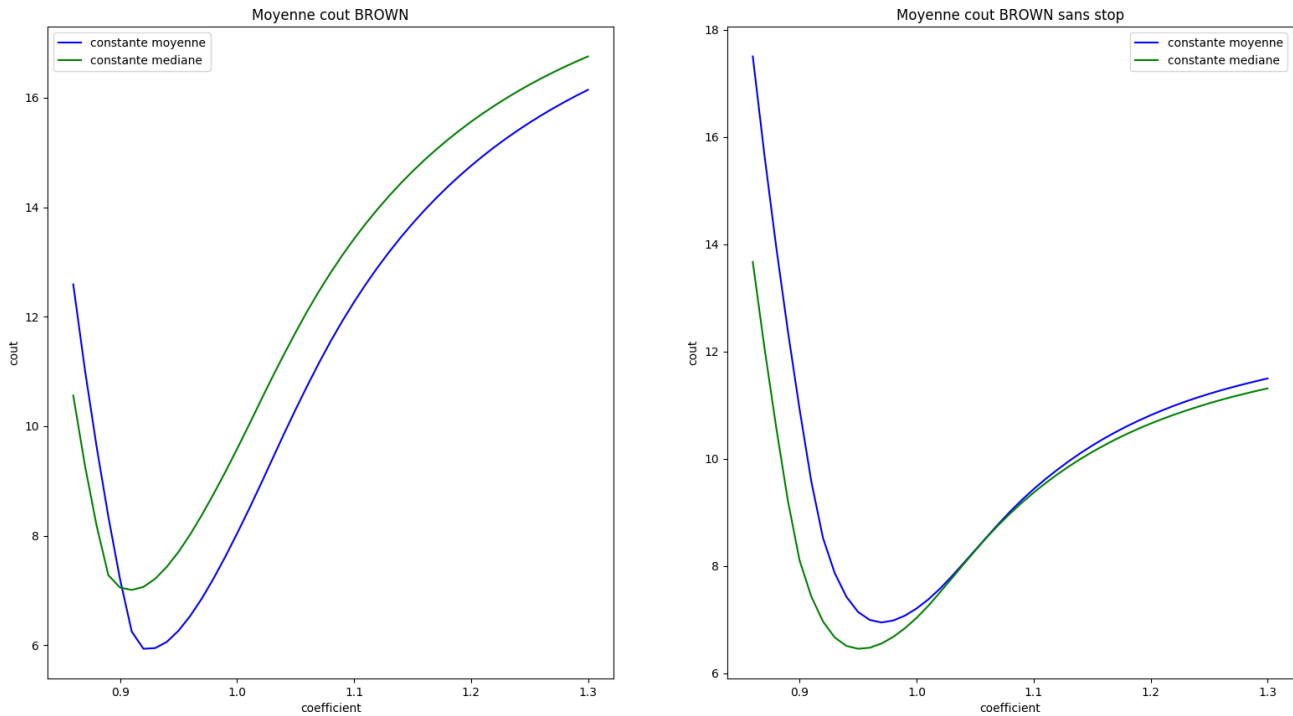


FIGURE 19 – Coût absolu moyen par coefficient

- Coût min brown moyenne : 5.93, median : 7.01
- Coût min brown (- stopwords) moyenne : 6.95, median : 6.46
- Coefficient min brown moyenne : 0.92, median : 0.91
- Coefficient min brown (- stopwords) moyenne : 0.97, median : 0.95

Résultats Le tableaux ci dessous rappelle les données récupérées au long de la recherche :

	BROWN avec stopwords	BROWN sans stopwords
nombre de mots uniques	49398	49383
nombre de mots total	1012528	578837
Constante moyenne	56525.81	55809.97
Constante médiane	48601.50	48494.00
Coefficient avec moyenne	0.92	0.97
Cout du coefficient moyenne	5.93	6.95
Coefficient avec médiane	0.91	0.95
Cout du coefficient médiane	7.01	6.46

D'après les données il est possible de dire que l'on obtient de meilleurs résultats si on conserve tous les mots du corpus. Dans ce cas l'utilisation de la moyenne des constantes génère un taux d'erreur plus faible que la médiane.

Ci-dessous la représentation des fréquences théoriques avec le coefficient optimal pour chaque corpus et chaque méthode. On voit que la courbe de la constante moyenne sur le corpus brute est celle qui suit le mieux les données réelles.

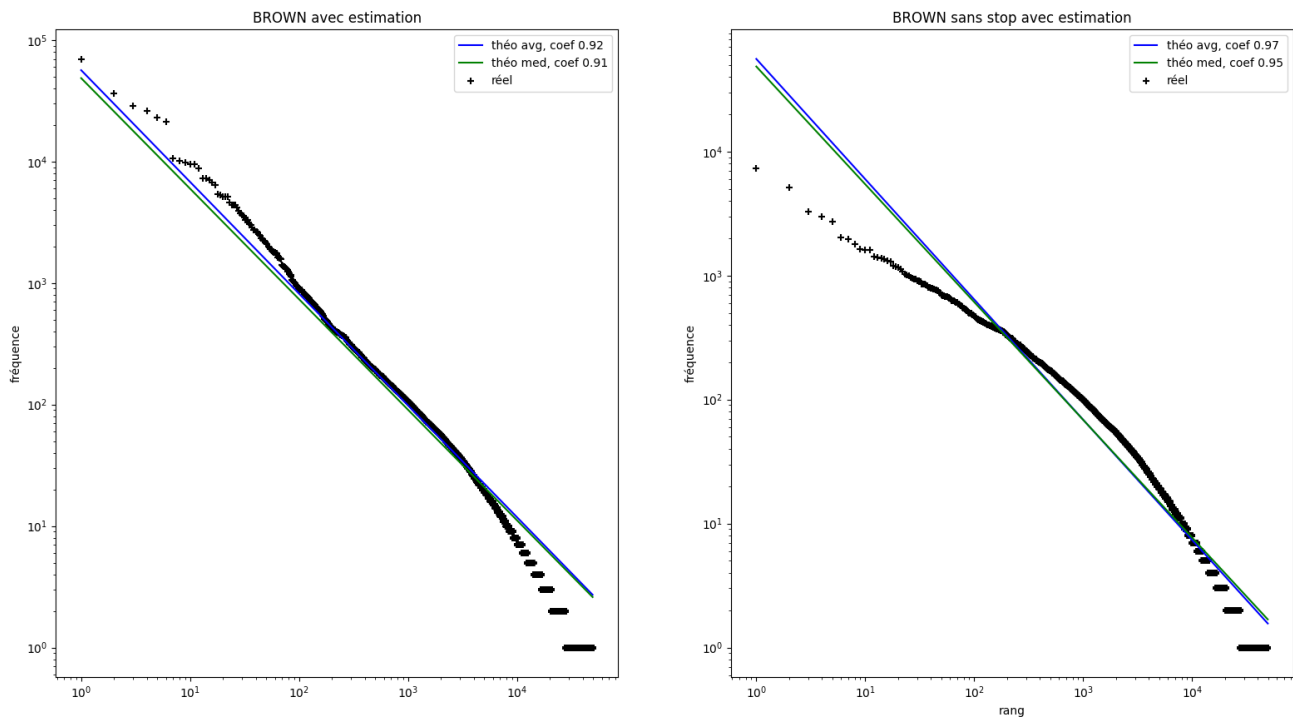


FIGURE 20 – Distribution de Zipf avec les estimations

En conclusion, j'utiliserais la moyenne des constantes sur un document complet afin de déterminer le coefficient dans ma recherche de spam.

Notes : L'ensemble des codes sources pour cette partie est disponible dans les fichiers :

./analyse/rech_zipf.py
./traitement/stats.py

B Déploiement des bases de données

Cette annexe détaille la mise en place de l'infrastructure de base de données pour le projet. J'utilise 2 environnements de base de données en version conteneur (docker version 23.0.4) :

- Elasticsearch
 - 1 Node Elastic, pour le service de base de données
 - 1 Service Kibana, pour la visualisation des index
 - 1 Service de certificat, pour sécuriser les échanges
- PostgreSQL
 - 1 service PostgreSQL, pour la base de données
 - 1 service PgAdmin, pour la visualisation de la base

Chaque environnement est indépendant, et possède une ouverture sur le PC hôte.

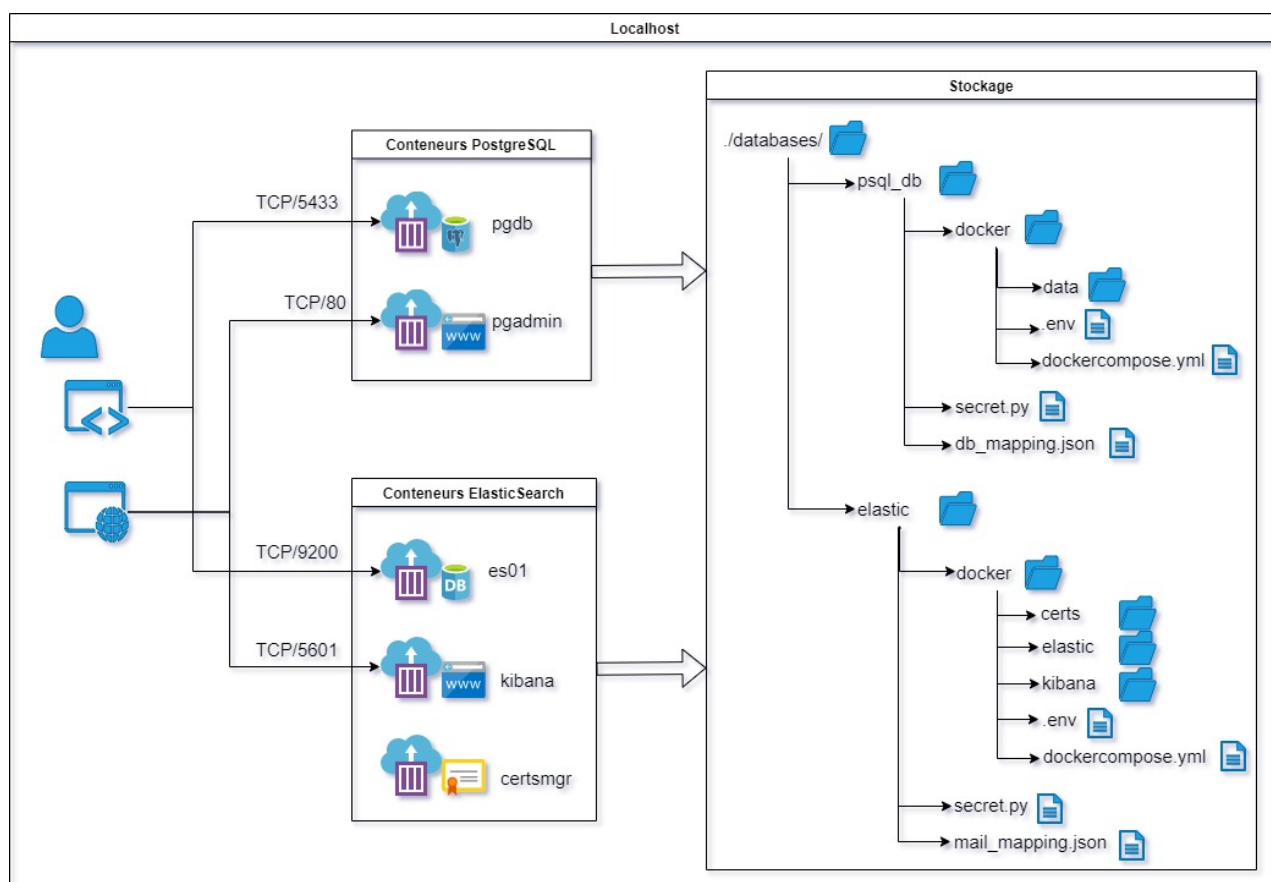


FIGURE 21 – Schéma de l'architecture Docker

B.1 Elasticsearch

Cet environnement se lance avec la commande suivante :

```
docker compose -f ./databases/elastic/docker/docker-compose.yml up -d
```

B.1.1 Conteneurisation

DockerCompose

```
1 version: "3.2"
2
3 services:
```

```

4      certsmgr:
5          image: elasticsearch:${VERSION}
6          volumes:
7              - ./certs:/usr/share/elasticsearch/config/certs
8          user: "0"
9          command: >
10             bash -c '
11                 if [ x${ELASTIC_PASSWORD} == x ]; then
12                     echo "Set the ELASTIC_PASSWORD
environment variable in the .env file";
13                     exit 1;
14                 elif [ x${KIBANA_PASSWORD} == x ]; then
15                     echo "Set the KIBANA_PASSWORD
environment variable in the .env file";
16                     exit 1;
17                 fi;
18                 if [ ! -f certs/ca.zip ]; then
19                     echo "Creating CA";
20                     bin/elasticsearch-certutil ca —
silent —pem —out config/certs/ca.zip;
21                     unzip config/certs/ca.zip -d
config/certs;
22                 fi;
23                 if [ ! -f certs/certs.zip ]; then
24                     echo "Creating certs";
25                     echo -ne \
26                         "instances:\n"\
27                         "  — name: es01\n"\
28                         "      dns:\n"\
29                         "          — es01\n"\
30                         "          — localhost\n"\
31                         "      ip:\n"\
32                         "          — 127.0.0.1\n"\
33                         > config/certs/instances.
yaml;
34                     bin/elasticsearch-certutil cert
—silent —pem —out config/certs/certs.zip —in config/certs/instances
.yml —ca-cert config/certs/ca/ca.crt —ca-key config/certs/ca/ca.key;
35                     unzip config/certs/certs.zip -d
config/certs;
36                 fi;
37                 # echo "Setting file permissions"
38                 echo "chown -R root:root config/certs";
39                 echo "find . -type d -exec chmod 750 \{\}
\;";
40                 echo "find . -type f -exec chmod 640 \{\}
\;";
41                 echo "Waiting for Elasticsearch
availability";
42                 until curl -s —cacert config/certs/ca/ca
.crt https://es01:9200 | grep -q "missing authentication credentials";
do sleep 30; done;
43                 echo "Setting kibana_system password";

```

```

44         until curl -s -X POST --cacert config/
certs/ca/ca.crt -u elastic:${ELASTIC_PASSWORD} -H "Content-Type:
application/json" https://es01:9200/_security/user/kibana_system/
_password -d "{\"password\":\"${KIBANA_PASSWORD}\"}" | grep -q "^{}";
do sleep 10; done;

45         echo "All done!";
46     ,
47     healthcheck:
48         test: ["CMD-SHELL", "[ -f config/certs/es01/es01.
crt ]"]
49         interval: 1s
50         timeout: 5s
51         retries: 120
52
53     es01:
54         depends_on:
55             certsmgr:
56                 condition: service_healthy
57         image: elasticsearch:${VERSION}
58         volumes:
59             - ./certs:/usr/share/elasticsearch/config/certs
60             - ./elastic/data:/usr/share/elasticsearch/data
61         ports:
62             - ${ES_PORT}:9200
63         environment:
64             - discovery.type=single-node
65             - ELASTIC_PASSWORD=${ELASTIC_PASSWORD}
66             - bootstrap.memory_lock=true
67             - xpack.security.enabled=true
68             - xpack.security.http.ssl.enabled=true
69             - xpack.security.http.ssl.key=certs/es01/es01.key
70             - xpack.security.http.ssl.certificate=certs/es01/
71             es01.crt
72             - xpack.security.http.ssl.certificate_authorities
=certs/ca/ca.crt
73             - xpack.security.http.ssl.verification_mode=
certificate
74             - xpack.security.transport.ssl.enabled=true
75             - xpack.security.transport.ssl.key=certs/es01/
es01.key
76             - xpack.security.transport.ssl.certificate=certs/
es01/es01.crt
77             - xpack.security.transport.ssl.
certificate_authorities=certs/ca/ca.crt
78             - xpack.security.transport.ssl.verification_mode=
certificate
79             - xpack.license.self_generated.type=${LICENSE}
80         mem_limit: ${MEM_LIMIT}
81         ulimits:
82             memlock:
83                 soft: -1
84                 hard: -1

```

```

85         healthcheck:
86             test:
87                 [
88                     "CMD-SHELL",
89                     "curl -s --cacert config/certs/ca
/ca.crt https://localhost:9200 | grep -q 'missing authentication
credentials'",
90                 ]
91             interval: 10s
92             timeout: 10s
93             retries: 120
94
95     kibana:
96         depends_on:
97             es01:
98                 condition: service_healthy
99         image: kibana:${VERSION}
100        volumes:
101        - ./certs:/usr/share/kibana/config/certs
102        - ./kibana/data:/usr/share/kibana/data
103        ports:
104        - ${KIBANA_PORT}:5601
105        environment:
106            - SERVERNAME=kibana
107            - ELASTICSEARCH_HOSTS=https://es01:9200
108            - ELASTICSEARCH_USERNAME=kibana_system
109            - ELASTICSEARCH_PASSWORD=${KIBANA_PASSWORD}
110            - ELASTICSEARCH_SSL_CERTIFICATEAUTHORITIES=config
/certs/ca/ca.crt
111        mem_limit: ${MEM_LIMIT}
112        healthcheck:
113            test:
114                [
115                    "CMD-SHELL",
116                    "curl -s -I http://localhost:5601
| grep -q 'HTTP/1.1 302 Found'",
117                ]
118            interval: 10s
119            timeout: 10s
120            retries: 120
121

```

Fichier d'environnement

```

1  # Password for elastic user
2  ELASTIC_PASSWORD=XXXXXXXXXX
3  # Password for kibana_system user
4  KIBANA_PASSWORD=XXXXXXXXXX
5  # Version elastic product
6  VERSION=8.1.2
7  # Licence to use
8  LICENSE=basic
9  # Port for elastic HTTP API
10 ES_PORT=9200

```

```

11 #ES_PORT=127.0.0.1:9200
12 # Port for kibana access
13 KIBANA_PORT=5601
14 # Memory available (bytes)
15 MEM_LIMIT=322122547
16

```

B.1.2 Initialisation de l'index

Exemples de secrets

```

1 serveur = "https://localhost:9200"
2 apiid = "XXXXXXXXXXXXXXXXXXXXXXXXXXXX"
3 apikey = "XXXXXXXXXXXXXXXXXXXXXXXXXXXX"
4 ca_cert = "databases/elastic/docker/certs/ca/ca.crt"
5

```

Il est possible de générer une clé API via l'interface Kibana > Stack Management > Sécurité > API keys > Create API Key > JSON.

Mapping de l'index

```

1 {
2   "properties": {
3     "hash": {"type": "keyword"},
4     "categorie": {"type": "keyword"},
5     "sujet": {"type": "text"},
6     "expediteur": {"type": "text"},
7     "message": {"type": "text"},
8     "langue": {"type": "keyword"}
9   }
10 }
11

```

Fonctions Python utiles

```

1 def es_connect(server, creds, crt):
2     """ Connexion au serveur Elasticsearch """
3     client = Elasticsearch(server, api_key=creds, ca_certs=crt)
4
5     try:
6         client.search()
7         client.indices.get(index="*")
8         return client
9
10    except (exceptions.ConnectionError, AuthenticationException,
11            AuthorizationException) as err:
12        print("ES:conn - Informations client Elasticsearch :\n\t", err)
13        client.close()
14        return None
15
16 def es_create_indice(es_cli, index, mapping):
17     """ Créer un indice s'il n'existe pas déjà """
18     indices = es_cli.indices.get(index='*')

```

```

19     if indices and index in indices:
20         print("Warning: Indice {} deja present".format(index), end=' ')
21         return
22
23     try:
24         res = es_cli.indices.create(index=index, mappings=mapping)
25     except elasticsearch.ApiError as err:
26         print(err)
27         return
28
29     if not res['acknowledged']:
30         print("Error : Echec de la creation de l'indice {}".format(index)
31 )

```

B.2 PostgreSQL

Cet environnement se lance avec la commande suivante :

```
docker compose -f ./databases/psql_db/docker/docker-compose.yml up -d
```

Le programme se charge automatiquement de créer le compte utilisé pour se connecter à la base.

B.2.1 Conteneurisation

DockerCompose

```

1  version: "3.3"
2  services:
3      pgdb:
4          image: postgres:${PS_VERSION}
5          restart: always
6          environment:
7              POSTGRES_PASSWORD: ${PS_PASSWORD}
8
9          volumes:
10             - ./data:/var/lib/postgresql/data
11
12         ports:
13             - ${PS_PORT}:5432
14
15     pgadmin:
16         image: dpage/pgadmin4:${PS_VERSION}
17         environment:
18             PGADMIN_DEFAULT_EMAIL: ${PGA_MAIL}
19             PGADMIN_DEFAULT_PASSWORD: ${PGA_PASSWORD}
20
21         ports:
22             - ${PGA_PORT}:80
23
24         depends_on:
25             - pgdb

```

```
1 PS_VERSION=latest
2
3 PS_PASSWORD=XXXXXXXXXX
4 PS_PORT=5433
5
6 PGA_MAIL=data@data.org
7 PGA_PASSWORD=XXXXXX
8 PGA_PORT=80
9
```

B.2.2 Initialisation de la base de données

Exemples de secrets.py

```
1 owner = "XXX"
2 owner_pw = "XXX"
3 admin = "XXXX"
4 admin_pw = "XXXX"
5 host = "localhost"
6 port = "5432"
7
```

Mapping

```
1 {
2   "mail_features": {
3     "categories": {
4       "id_cat": ["SERIAL", "PRIMARY KEY"],
5       "type": ["VARCHAR", "UNIQUE", "NOT NULL"]
6     },
7     "messages": {
8       "id_message": ["SERIAL", "PRIMARY KEY"],
9       "hash": ["CHAR(32)", "UNIQUE", "NOT NULL"],
10      "id_cat": ["INT", "NOT NULL"],
11      "fk": ["fk_message", "id_cat", "categories(id_cat)", "SET NULL"]
12    },
13    "liens": {
14      "id_message": ["INT"],
15      "url": ["INT"],
16      "mail": ["INT"],
17      "tel": ["INT"],
18      "nombre": ["INT"],
19      "prix": ["INT"],
20      "fk": ["fk_liens", "id_message", "messages(id_message)", "CASCADE"]
21    },
22    "stats_mots": {
23      "id_message": ["INT"],
24      "mots_uniques": ["INT"],
25      "mots": ["INT"],
26      "char": ["INT"],
27      "char_maj": ["INT"],
28      "mot_maj": ["INT"],
```



```

29         "mot_cap": ["INT"],
30         "fk": ["fk_stats_mot", "id_message", "messages(id_message)", "
CASCADE"]
31     },
32     "stat_ponct": {
33         "id_message": ["INT"],
34         "point": ["INT"],
35         "virgule": ["INT"],
36         "exclamation": ["INT"],
37         "interrogation": ["INT"],
38         "espace": ["INT"],
39         "tabulation": ["INT"],
40         "ligne": ["INT"],
41         "ligne_vide": ["INT"],
42         "fk": ["fk_stats_ponct", "id_message", "messages(id_message)", "
CASCADE"]
43     },
44     "zipf": {
45         "id_message": ["INT"],
46         "constante": ["INT"],
47         "coefficient": ["REAL"],
48         "tx_erreur": ["REAL"],
49         "fk": ["fk_zipf", "id_message", "messages(id_message)", "CASCADE"]
50     },
51     "hapax": {
52         "id_message": ["INT"],
53         "ratio_unique": ["REAL"],
54         "ratio_texte": ["REAL"],
55         "fk": ["fk_hapax", "id_message", "messages(id_message)", "CASCADE"]
56     }
57 }
58 }
59

```

Fonctions utiles

```

1 def create_db(nom, owner, user, passwd, host, port):
2     """ Créer une nouvelle base de donnees """
3     client_psql = psycopg2.connect(user=user, password=passwd, host=host,
4     port=port)
5     client_psql.autocommit = True
6     cursor = client_psql.cursor()
7
8     cursor.execute("DROP DATABASE IF EXISTS {};".format(nom))
9     cursor.execute("CREATE DATABASE {};".format(nom))
10    cursor.execute("ALTER DATABASE {} OWNER TO {};".format(nom, owner))
11    client_psql.close()
12
13 def connect_db(database, user, passwd, host, port):
14     """ Connexion a la base de donnees Postgres. Penser a fermer la
15     connexion """
16     try:
17         client_psql = psycopg2.connect(database=database, user=user,

```

```

17 password=password, host=host, port=port)
18 except psycopg2.Error as e:
19     print("Erreur de connexion : \n{}".format(e), file=sys.stderr)
20     return None
21
22 client_psycopg2.autocommit = True
23 return client_psycopg2
24
25 def create_table(client_psycopg2, nom, champs):
26     """ Créer une nouvelle table dans la base de données """
27     fk = ""
28     fields = []
29
30     curseur = client_psycopg2.cursor()
31     curseur.execute(f"DROP TABLE IF EXISTS {nom}")
32
33     if 'fk' in champs.keys():
34         ls = champs.pop('fk')
35         fk = f"CONSTRAINT {ls.pop(0)} FOREIGN KEY({ls.pop(0)}) REFERENCES
36         {ls.pop(0)}"
37         if ls:
38             fk += f" ON DELETE {ls.pop(0)}"
39
40     for key, value in champs.items():
41         fields.append(f"{key} {value}")
42
43     query = f"CREATE TABLE {nom} ({', '.join(fields)})"
44     if fk:
45         query = query[:-1] + f", {fk}"
46     curseur.execute(query)
47
48 def create_index(client_psycopg2, nom, table, colonne):
49     """ Index sur une colonne """
50     query = f"CREATE UNIQUE INDEX {colonne} ON {table}({colonne})".format(nom, table, colonne)
51     exec_query(client_psycopg2, query)
52

```

C Sorties PSQL et la liste de mots

```
SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
ORDER by freq_corpus DESC
LIMIT 50;
```

mot	freq_corpus	freq_doc_spam	freq_doc_ham
email	5326	933	1009
get	3904	613	1266
use	3868	433	1347
list	3818	719	1797
one	3268	511	1028
make	3186	441	826
free	2803	644	450
new	2689	484	782
time	2550	487	783
send	2485	565	593
go	2338	407	845
say	2327	153	792
information	2310	569	777
would	2288	332	851
good	2242	455	795
work	2226	348	779
receive	2216	653	312
people	2100	260	510
mailing	2077	535	1062
message	2057	363	727
like	2039	355	843
write	1969	116	1520
click	1963	815	148
please	1954	781	296
business	1915	400	234
address	1909	438	319
year	1761	363	427
company	1754	395	238
find	1745	285	663
report	1733	148	257
want	1732	438	538
money	1706	354	153
order	1697	283	162
see	1696	315	746
need	1671	402	559
also	1655	281	617
take	1649	370	574
name	1617	366	224
know	1568	277	697
user	1546	124	908
system	1544	194	527
group	1527	116	924

change		1494		160		507
file		1491		91		469
day		1454		360		376
program		1423		228		182
remove		1407		741		181
include		1406		395		361
may		1402		332		448
look		1357		240		584

```

SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
ORDER by freq_doc_spam DESC
LIMIT 50;

```

mot		freq_corpus		freq_doc_spam		freq_doc_ham
-----+-----+-----+-----						
email		5326		933		1009
click		1963		815		148
please		1954		781		296
remove		1407		741		181
list		3818		719		1797
receive		2216		653		312
free		2803		644		450
get		3904		613		1266
information		2310		569		777
send		2485		565		593
mailing		2077		535		1062
one		3268		511		1028
time		2550		487		783
new		2689		484		782
offer		1320		457		171
good		2242		455		795
make		3186		441		826
address		1909		438		319
want		1732		438		538
use		3868		433		1347
go		2338		407		845
need		1671		402		559
business		1915		400		234
include		1406		395		361
company		1754		395		238
wish		643		373		97
take		1649		370		574
name		1617		366		224
message		2057		363		727
year		1761		363		427
day		1454		360		376
like		2039		355		843
money		1706		354		153
today		826		352		207
work		2226		348		779
service		1344		345		285

state		1313		332		202
would		2288		332		851
may		1402		332		448
form		869		331		148
reply		625		323		98
future		679		321		179
call		1187		317		344
see		1696		315		746
help		1165		313		363
home		1229		303		304
contact		631		302		118
subject		883		294		403
phone		947		286		299
find		1745		285		663

```

SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
ORDER by freq_doc_ham DESC
LIMIT 50;

```

mot		freq_corpus		freq_doc_spam		freq_doc_ham
-----+-----+-----+-----						
list		3818		719		1797
write		1969		116		1520
use		3868		433		1347
get		3904		613		1266
mailing		2077		535		1062
one		3268		511		1028
email		5326		933		1009
group		1527		116		924
user		1546		124		908
would		2288		332		851
go		2338		407		845
like		2039		355		843
make		3186		441		826
good		2242		455		795
say		2327		153		792
time		2550		487		783
new		2689		484		782
work		2226		348		779
information		2310		569		777
date		955		71		747
see		1696		315		746
linux		1302		71		738
message		2057		363		727
know		1568		277		697
think		1354		155		685
im		1133		105		667
find		1745		285		663
url		899		48		632
also		1655		281		617
send		2485		565		593

well		1236		225		587
look		1357		240		584
irish		732		68		581
take		1649		370		574
maintainer		655		65		571
way		1295		234		569
unsubscription		650		66		566
need		1671		402		559
try		1076		135		551
want		1732		438		538
could		1205		175		535
system		1544		194		527
sponsor		717		131		518
give		1226		245		516
people		2100		260		510
change		1494		160		507
run		963		68		502
thing		900		112		498
even		1163		247		484
first		1261		222		472

```

SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
ORDER by freq_doc_spam, freq_doc_ham DESC
LIMIT 50;

```

mot		freq_corpus		freq_doc_spam		freq_doc_ham
-----+-----+-----+-----						
rpmlist		365		0		363
razoruser		242		0		191
cheer		187		0		184
matthias		264		0		182
lawrence		159		0		149
saou		133		0		131
razor		256		0		125
kernel		208		0		123
exmhuser		108		0		98
exmh		183		0		92
rpm		123		0		82
picasso		81		0		81
teledynamics		80		0		80
hettinga		104		0		74
pudge		107		0		73
beberg		74		0		70
garrigue		72		0		69
zdnet		167		0		68
redhat		97		0		67
gmt		71		0		64
hmm		68		0		64
vendor		106		0		62
spammer		122		0		62
barcelona		61		0		61

norte		61		0		61
planta		61		0		61
edificio		61		0		61
pablo		60		0		58
evildoer		58		0		58
exmhworker		59		0		57
mason		60		0		57
argue		74		0		57
bearer		66		0		56
wrongdoer		56		0		56
viricio		56		0		56
farquhar		58		0		55
blog		62		0		54
repository		63		0		54
pgp		154		0		54
gibbon		55		0		52
usefulness		54		0		51
agreeable		53		0		50
cvs		59		0		49
spamassassin-devel		48		0		48
cio		56		0		48
procmail		71		0		47
weird		51		0		47
owen		51		0		46
damn		47		0		45
img		57		0		44

```

SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
ORDER by freq_doc_ham, freq_doc_spam DESC
LIMIT 50;

```

mot		freq_corpus		freq_doc_spam		freq_doc_ham
-----+-----+-----+-----						
annuity		117		34		0
testimonial		56		28		0
multilevel		47		26		0
folloring		23		22		0
moneyback		26		22		0
faithfully		22		21		0
oprah		20		20		0
featurepack		36		20		0
assuming		18		18		0
okidata		17		17		0
tonersgo		34		17		0
ez		39		17		0
libido		17		17		0
hazardous		17		17		0
seen		21		17		0
paperwork		26		17		0
penile		18		17		0
replenishable		34		17		0

jody		67		17		0
emem		17		17		0
merciless		17		17		0
hgh		26		16		0
easytoremember		16		16		0
dieting		33		16		0
spout		16		16		0
smoker		39		16		0
otcbb		24		16		0
cam		23		15		0
craving		15		15		0
suza		15		15		0
diploma		35		15		0
deductible		21		15		0
sondstrom		15		15		0
professionally		16		15		0
lodge		19		15		0
hedland		15		15		0
friendsrelative		15		15		0
wsuperior		15		15		0
downline		31		15		0
ravage		19		15		0
payable		18		15		0
states		18		15		0
refinancing		15		14		0
westport		14		14		0
dellaca		14		14		0
pharmacy		16		14		0
act.		13		13		0
propecia		15		13		0
rated		21		13		0
xenical		17		13		0

```

SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
WHERE freq_doc_ham >= 2*freq_doc_spam
ORDER BY freq_doc_ham DESC
LIMIT 50;

```

mot		freq_corpus		freq_doc_spam		freq_doc_ham
-----+-----+-----+-----						
list		3818		719		1797
write		1969		116		1520
use		3868		433		1347
get		3904		613		1266
one		3268		511		1028
group		1527		116		924
user		1546		124		908
would		2288		332		851
go		2338		407		845
like		2039		355		843
say		2327		153		792

work		2226		348		779
date		955		71		747
see		1696		315		746
linux		1302		71		738
message		2057		363		727
know		1568		277		697
think		1354		155		685
im		1133		105		667
find		1745		285		663
url		899		48		632
also		1655		281		617
well		1236		225		587
look		1357		240		584
irish		732		68		581
maintainer		655		65		571
way		1295		234		569
unsubscription		650		66		566
try		1076		135		551
could		1205		175		535
system		1544		194		527
sponsor		717		131		518
give		1226		245		516
change		1494		160		507
run		963		68		502
thing		900		112		498
first		1261		222		472
problem		925		130		470
file		1491		91		469
world		1133		185		461
ive		640		42		421
seem		588		30		420
still		726		120		416
add		847		121		410
many		1079		189		387
something		588		59		386
really		627		78		382
thanks		571		99		380
set		758		85		379
pm		502		26		376

```

SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham
FROM mot_corpus
WHERE freq_doc_spam >= 2*freq_doc_ham
ORDER BY freq_doc_spam DESC
LIMIT 50;

```

mot		freq_corpus		freq_doc_spam		freq_doc_ham
click		1963		815		148
please		1954		781		296
remove		1407		741		181
receive		2216		653		312

offer		1320		457		171
wish		643		373		97
money		1706		354		153
form		869		331		148
reply		625		323		98
contact		631		302		118
hour		597		279		124
visit		551		274		121
special		553		246		109
low		580		241		109
within		644		238		106
request		452		234		85
credit		659		234		44
fill		392		225		50
guarantee		463		224		36
marketing		551		219		34
rate		676		215		103
dollar		497		207		41
simply		411		205		102
purchase		419		203		50
opportunity		502		194		74
interest		533		193		92
website		431		193		87
thousand		419		189		80
hundred		329		189		59
minute		383		187		92
dear		250		179		55
thank		284		179		53
legal		535		174		77
family		409		161		75
sale		467		161		75
financial		392		159		47
cash		484		155		33
income		391		150		24
risk		350		148		65
removal		213		145		12
investment		421		145		28
insurance		512		142		19
instruction		381		136		58
secret		312		126		39
absolutely		230		125		30
ad		304		121		49
mortgage		247		117		4
prove		193		114		54
payment		219		112		21
obligation		164		112		13

```

SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham,
freq_doc_spam/freq_doc_ham as "ratio spam/ham"
FROM mot_corpus
WHERE freq_doc_ham > 0 AND freq_doc_spam > 0

```

ORDER BY "ratio spam/ham" DESC
LIMIT 50;

mot	freq_corpus	freq_doc_spam	freq_doc_ham	ratio
-----	-----	-----	-----	-----
spamassassinsightings	108	107	1	
deathtosпамdeathtosпамdeathtosпам	101	100	1	
refinance	106	72	1	
lender	85	62	1	
optout	57	49	1	
attain	52	42	1	
guaranteed	84	75	2	
tollfree	41	35	1	
mortgage	247	117	4	
mastercard	38	28	1	
honesty	32	27	1	
inkjet	54	26	1	
medication	44	25	1	
herbal	49	24	1	
seeker	31	22	1	
erection	25	21	1	
postal	57	42	2	
susan	24	21	1	
reap	22	21	1	
confidentiality	57	40	2	
originator	21	20	1	
laserjet	56	19	1	
recruiting	26	19	1	
thermal	23	19	1	
utmost	23	19	1	
systemworks	20	19	1	
untitled	21	19	1	
mlm	85	39	2	
sincere	21	18	1	
pam	18	17	1	
icann	18	17	1	
furnish	19	17	1	
wi	25	17	1	
rela	18	17	1	
residual	19	16	1	
ofcourse	17	16	1	
grumble	17	16	1	
resell	41	32	2	
homeowner	41	32	2	
consolidate	33	31	2	
creditor	26	15	1	
attn	17	15	1	
astonishment	16	15	1	
mayle	16	15	1	
casino	44	15	1	

aging		86		31		2	
financially		54		45		3	
worst		16		15		1	
subjectremove		15		14		1	
unclaimed		20		14		1	

```

SELECT mot, freq_corpus, freq_doc_spam, freq_doc_ham,
freq_doc_ham/freq_doc_spam as "ratio ham/spam"
FROM mot_corpus
WHERE freq_doc_ham > 0 AND freq_doc_spam > 0
ORDER BY "ratio ham/spam" DESC
LIMIT 50;

```

mot		freq_corpus		freq_doc_spam		freq_doc_ham		ratio ham/spam
-----+-----+-----+-----+-----								
gary		167		1		157		157
tue		157		1		153		153
aug		382		2		304		152
spamassassintalk		156		1		147		147
murphy		137		1		130		130
perl		439		1		121		121
osdn		125		1		120		120
folder		235		1		104		104
jul		243		2		184		92
cnet		401		1		89		89
argument		110		1		78		78
bug		322		2		153		76
score		128		1		72		72
object		127		1		68		68
apt		96		1		61		61
yesterday		66		1		60		60
van		67		1		55		55
hack		72		1		55		55
character		100		1		53		53
sep		343		5		265		53
default		210		2		107		53
stupid		84		1		52		52
icq		58		1		52		52
useless		105		2		100		50
rebuild		74		1		50		50
thread		87		1		50		50
kevin		124		2		101		50
judge		107		1		50		50
bunch		55		1		48		48
ip		78		1		48		48
output		74		1		48		48
kate		49		1		48		48
spamassassin		155		2		97		48
adam		117		2		92		46
innovation		60		1		45		45
odd		52		1		45		45
fairly		50		1		44		44

chris		198		3		130		43
spec		54		1		43		43
gnome		84		1		43		43
beta		90		1		43		43
binary		71		1		42		42
warn		45		1		42		42
scientist		78		1		42		42
ilug		48		1		41		41
complain		48		1		41		41
universe		130		1		41		41
speech		87		1		41		41
evil		59		1		41		41
alan		42		1		40		40

Liste des mots retenus

email, get, use, list, one, make, free, new, time, send,
say, information, would, good, work, receive, people, mailing, message, like,
write, click, please, business, address, year, company, find, report, want,
money, order, see, need, also, take, name, know, user, system,
group, change, file, day, program, remove, include, may, look, think,
service, offer, state, linux, way, internet, first, software, well, home,
give, could, site, call, web, help, product, even, mail, come,
world, much, many, try, check, link, price, start, government, read,
number, run, date, million, phone, week, problem, server, computer, month,
thing, follow, url, subject, form, back, available, grant, long, add,
line, every, life, network, right, today, spam, provide, security, online,
easy, window, pay, support, two, set, news, unsubscribe, irish, late,
code, still, keep, since, sponsor, last, great, show, version, without,
future, rate, case, market, page, end, credit, issue, maintainer, technology,
unsubscription, let, tell, within, wish, ive, never, contact, old, next,
save, place, really, reply, sell, point, another, card, sure, release,
base, search, anyone, part, full, build, hour, access, per, put,
something, seem, professional, low, cost, result, thanks, question, communication,
lot, country, special, fast, marketing, visit, buy, open, high, application, might,
package, source, must, legal, person, allow, customer, request, fill, guarantee,
dollar, simply, purchase, require, opportunity, website, interest, thousand, hundred,
minute, simple, dear, thank, ever, family, sale, account, financial, increase, cash,
income, friend, risk, less, lose, personal, investment, top, removal, error,
deal, complete, insurance, type, fax, notice, easily, learn, instruction, plan,
accept, experience, city, real, benefit, large, current, secret, fact, rpmlist,
cant, aug, post, though, actually, wednesday, sep, around, rights, mean,
monday, copyright, probably, different, big, install, bit, however, welcome, supply,
original, idea, geek, ill, thursday, didnt, create, august, bad, someone,
example, stop, update, reason, enough, yahoo, little, job, ask, yes,
friday, discuss, heaven, isnt, yet, reserve, anything, hard, annuity, testimonial,
multilevel, folloring, moneyback, faithfully, oprah, featurepack, assuming, tonersgo,
libido, hazardous, replenishable, merciless, emem, okidata, jody, paperwork, penile,
seen, otcbb, easytoremember, spout, hgh, smoker, dieting, friendsrelative, ravage,
states, downline, suza, hedland, craving, sondstrom, diploma, cam, payable,
professionally, deductible, lodge, wsuperior, dellaca, westport, refinancing,

pharmacy, propecia, rated, climax, act., xenical, mega, mailto?subjectremove, classified, sent, hrefclick, metabolism, turnkey, lagos, congo, cosigner, abacha, phentermine, enlarge, walaa, sirmadam, nursery, hormones, zyban, mailto?, fastest, disclaimerwe, horny, repose, erotic, alimony, customized, leasing, homebased, ltc, cordial, diete, nationally, lowest, lien, cloak, kabila, worksheet, prizemama, energetic, committment, erectile, homeownership, extracurricular, forwardlooking, utilitiesall, towing, tables, roadside, optedin, grantwriting, remittance, mustdo, nofrills, grantswriting, feasibility, redtape, payer, lowinter, mustknow, onehalf, usmillion, taxable, nonsmoker, onestopmoneyshop, emailing, indulgence, abidjan, massand, divinorum, rd., salvia, valtrex, hrefclickherea, advisement, highrisk, sani, recewveng, recruiter, mitsubishi, pleasureincreased, reg., staminaincreased, nospam, deliverable, sensationincreased, marvelously, supplemental, tonereverse, chavana, ferocity, commencement, growthwrinkle, weightbuild, fige, untraceable, leone, herba, prosaka, afree, enlargement, skinnew, prostate, vjestika, payperview, viripotent, calea, herpes, highratio, resale, seizure, angelica, dagga, hookah, anheuserbusch, aphrodisia, cellulite, capillaris, ffa, employment, pcinternet, refining, cbyi, geenergy, contemplativeness, inexpensively, nc., services, nonephedrine, calbay, sativa, chevron, mrsa, dreaming, nonmahuang, ok,oh, noncaffeine, flowertop, amalgamation, uplifting, dosage, carton, crocus, prash, drops, aligncenterfont, signify, alkaloid, razoruser, cheer, matthias, lawrence, saou, razor, kernel, exmhuser, exmh, rpm, picasso, teledynamics, hettinga, pudge, beberg, garrigue, zdnet, redhat, gmt, hmm, spammer, vendor, planta, barcelona, edificio, norte, pablo, evildoer, mason, argue, exmhworker, vircio, bearer, wrongdoer, farquhar, pgp, repository, blog, gibbon, usefulness, agreeable, cvs, cio, weird, procmail, owen, damn, img, config, niall, duncan, rohit, guido, btw, ham, funny, context, hash, valhalla, unseen, angle, geege, suse, functionality, debate, corpus, meatspace, brent, aptget, debian, disable, partition, vipul, drunken, warming, skeptic, bitbitch, architecture, sober, whitelist, matthew, rah, looney, credulity, mathematics, mailtoon, powershot, ssh, rival, welch, hal, somewhat, sha, gibson, expression, python, geometry, scalable, compiler, crucial, alsa, buffer, loop, fiction, experiment, gpl, compromise, freebsd, crap, ved, fault, spamd, prakash, gcc, distro, router, megapixel, sendmail, peoplesoft, elias, voyage, zealot, vice, dispatch, trivial, vulnerability, eugen, apache, barrera, byte, exploit, lyda, razorcheck, kre, technomad, rossum, fahrlander, moen, gnupg, revision, evansville, wifi, ziggy, huckleberry, iirc, finn, bounce, aprils, constantly, bullet, behaviour, debug, harley, pioneering, dimension, phil, agile, collapse, activist, gregory, presumably, neary, compatibility, diff, logic, pointer, bbc, syntax, turpin, panasas, merge, spokesman, cnets, vulnerable, forwardedby, assignee, jesse, digest, pdraig, leitzl, classifier, configurator, invoke, notion, zdnet, ought, andrew, filesystem, attacker, tokenizer, interaction, sysadmin, accuse, widespread, stuff, leave, course, happen, far, datum, bythinkgeek, least, test, feature, either, inc, turn, john, talk, answer, thats, man, story, maybe, else, mailto, jul, september, rather, fix, log, center, quite, cause, move, design, play, lead, cell, manage, comment, early, hope, hand, july, ago, gary, key, anyway, break, understand, exist, public, bug, tue, fall, kind, begin, pretty, content, often, subscribe, appear, drive, several, mark, newsletter, absolutely, mortgage, prove, payment, obligation, immediately, shipping, bank, earn, loan, fund, sincerely, independent, debt,

commission, loss, optin, limited, fee, contract, promotion, assistance, advertisement, successful, invest, health, retail, approve, telephone, expense, subscriber, guaranteed, weight, transaction, fortune, obtain, refinance, compliance, bonus, monthly, confidential, affiliate, unsolicited, valuable, improvement, assist, sum, immediate, savings, lender, profitable, discount, exciting, deposit, bulk, adult, strongly, qualify, blank, brand, sir, apologize, drug, assure, optout, opt, medical, venture, forever, correspondence, proposal, unlimited, resident, incredible, financially, membership, employment, zip, boss, estate, qualified, auto, attain, minimum, legally, postal, urgent, exercise, premium, exclude, cooperation, confidentiality, supplies, prescription, affordable, prohibit, consultation, broker, equity, convenience, guideline, presentation, utilize, mlm, carefully, tremendous, nigeria, handling, paragraph, inconvenience, visa, anytime, consolidation, classify, unwanted, living, approval, repair, instantly, tollfree, sexual, outstanding, son, specialist, currency, strictly, lifetime, shall, comprehensive, africa, kindly, buyer, insider, substantial, homeowner, safely, maximum, resell, sleep, consolidate, extract, warranty, aging, cartridge, receipt, ebook, exceed, refund, trading, inform, mastercard, honesty, inkjet, medication, herbal, seeker, susan, erection, reap, originator, utmost, thermal, untitled, laserjet, systemworks, recruiting, sincere, rela, furnish, icann, pam, grumble, ofcourse, residual, worst, attn, mayle, astonishment, creditor, casino, subjectremove, toner, unclaimed, lowcost, procurement, strenuous, intimidate, kin, modality, wrinkle, tenant, searcher, vault, idaho, ads, madam, foreigner, unlisted, debtor, hidden, resources, beneficiary, athlete, montana, stringent, hewlett, lottery, reconnect, housekeeper, mammoth, stepbystep, investigative, urgently, factual, dysfunction, gasoline, reclaim, laurent, cram, wiretap, enhancer, hispanic, missouri, dreams, tobacco, arkansas, nbc, daycare, etc, tutoring, childcare, skyrocket, productsservice, cock, dental, intensify, solicite, botanical, ext., repaid, multibillion, mitchell, profession, glossary, repay, counseling, ministry, hardcore, lauderdale, physician, epson, noninvasive, cbs, ent, exceedingly, alberta, gen, serenity, growing, mrs, children, thanx, facearial, christian, banker, toptoftheline, ave., formulate, lust, resin, rates, omaha, amalgamate, barrister, sold, comprise, cholesterol, rewarding, downsizing, cordially, bud, booster, remit, ingredient, reputable, potency, quotation, consumable, costeffective, ordering, getter, hobby, promotional, mentor, nevada, recieve, fedex, diagnostics, illegality, legality, doctorate, herb, cote, communion

D Tableau des choix technologiques

Élément	Retenu	Raisons	Observations
Datasets			
Mail de la compagnie Enron	Non	Mails non classés	Non retenu pour la phase de développement car pas de moyen fiable de contrôler la sortie automatiquement
Mail du projet SpamAssassin	Oui	Mails déjà pré-triés	Mails principalement en Anglais déjà pré-trié en catégorie Spam et Ham
Brown dataset (nlk)	Oui	Corpus d'un million de mots en Anglais publié en 1961	Dataset utilisé pour le développement de la visualisation de la distribution de Zipf
Stopwords (nlk)	Oui	Corpus de mots commun non significatif dans un texte	Utilisation dans le développement de la visualisation de la distribution de Zipf
Langage et Modules			
Python	Oui	Langage polyvalent pour le traitement des données	
Module email	Oui	Module natif pour le traitement des mails	Grande flexibilité pour la lecture des mails
Virtualisation			
Docker	Oui	configuration et environnement dans les fichiers. Cela facilite le portage vers d'autre environnement	Lors du développement j'ai pu utilisé mon PC personnel sous Linux et un autre PC sous Windows. L'utilisation de Docker m'a évité de nombres configurations sur Windows
Bases de données			
ElasticSearch	Oui	Technologie utilisée dans mon entreprise. Présence d'une interface de visualisation des données Kibana.	Application dockerisée. Basculement possible pour une base MongoDB plus souple et moins gourmande en ressource
PostgreSQL	Oui	Moteur de base de données relationnelle plus facilement scalable que ElasticSearch pour l'ajout de nouvelle catégorie de données. Il n'est pas nécessaire de ré-indexer toute la base pour ajouter des champs	Application dockerisée
SQLite	Oui	Base de données légère pour stocker uniquement les données statistiques des étapes de la phase 1	Rapide à mettre en place et déjà intégrée

E Modèles

E.1 Naïves Bayes

Ce type de modèle est utilisé par le module *langdetect* qui me sert pour la détection des langues.

Introduction Les modèles Naïves Bayes se basent sur le théorème de probabilité de Bayes. Il permet de déterminer la probabilité conditionnelle d'apparition d'un événement A sachant qu'un événement B s'est produit. Le terme naïf fait référence au fait que l'on présuppose que les événements A et B ne sont pas corrélés.

Ces techniques sont utilisées pour des modèles de classification en apprentissage supervisé.

La formule mathématique de ce théorème est la suivante :

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (1)$$

On recherche ici $P(A|B)$, c'est à dire la probabilité d'apparition d'un événement A sachant que l'évènement B s'est produit.

Pour ce faire nous avons besoin des données suivantes :

- $P(B|A)$ est la probabilité que l'évènement B s'est produit sachant que l'évènement A s'est produit
- $P(A)$ est la probabilité d'apparition de l'évènement A
- $P(B)$ est la probabilité d'apparition de l'évènement B

Exemples d'utilisation Les exemples ci dessous vont permettre d'illustrer l'utilisation de cette technique. D'abord manuellement sur un petit jeu de données puis à l'aide d'un code pré-existant sur un autre jeu de données plus important.

Manuel Dans cet exemple nous allons déterminer la probabilité qu'a un joueur d'aller sur le terrain selon les conditions météorologiques. Cette probabilité sera calculée en fonction des données récupérées lors des matchs précédents.³

On recherchera ainsi la probabilité de présence sur le terrain d'un joueur selon la météo $P(A|B)$. Pour ce faire nous auront besoin de :

- $P(A)$ Probabilité de jouer quelque soit le temps
- $P(B)$ Probabilité de l'évènement météorologique
- $P(B|A)$ Probabilité de l'évènement sachant que le joueur a été sur le terrain

TABLE 8 – Données de présence sur le terrain

météo	soleil	soleil	couvert	pluie	pluie	pluie	couvert
présent	non	non	oui	oui	oui	non	oui
météo	soleil	soleil	pluie	soleil	couvert	couvert	pluie
présent	non	oui	oui	oui	oui	oui	non

3. Les données présentées sont inventées

TABLE 9 – Synthèse et probabilité simple $P(A)$ et $P(B)$

météo	oui	non	$P(B)$
couvert	4	0	$4/14$
soleil	2	3	$5/14$
pluie	3	2	$5/14$
$P(A)$	$9/14$	$5/14$	

On peut déterminer les probabilités de chaque météo en fonction de la présence du joueur sur le terrain $P(B|A)$. Pour ce faire on divise le nombre d'évènements de présence du joueur lors d'un évènement météo par le nombre total d'évènements de présence du joueur

TABLE 10 – Probabilité météo selon présence du joueur

météo	$P(B oui)$	$P(B non)$
couvert	$4/9$	$0/5$
soleil	$2/9$	$3/5$
pluie	$3/9$	$2/5$

On va maintenant calculer la probabilité qu'à un joueur d'être sur le terrain si le temps est couvert.

On commence par la probabilité du oui :

$$\begin{aligned}
 P(A|B) &= \frac{P(B|A)P(A)}{P(B)} \\
 P(A|B) &= \frac{\frac{4}{9} \cdot \frac{9}{14}}{\frac{4}{14}} \\
 P(A|B) &= \frac{\frac{4}{14}}{\frac{4}{14}} \\
 P(A|B) &= \frac{4}{14} \cdot \frac{14}{4} \\
 P(A|B) &= 1
 \end{aligned}$$

On enchaîne sur la probabilité de ne pas jouer si le temps est couvert

$$\begin{aligned}
 P(A|B) &= \frac{P(B|A)P(A)}{P(B)} \\
 P(A|B) &= \frac{\frac{0}{5} \cdot \frac{5}{14}}{\frac{4}{14}} \\
 P(A|B) &= 0 \cdot \frac{14}{4} \\
 P(A|B) &= 0
 \end{aligned}$$

On peut dire que si le temps est couvert le joueur très probablement sur le terrain On peut également déterminer la probabilité de jouer pour chaque évènement météo

TABLE 11 – Probabilité présence du joueur selon la météo

météo	oui	non	plus probable
couvert	1	0	oui
soleil	$2/5$	$3/5$	non
pluie	$3/5$	$2/5$	oui

Cas polynomial : Il est possible de déterminer la probabilité d'un évènement par rapport à plus autres. Dans ce cas, il faudra multiplier entre elles les probabilités de ces évènements selon l'apparition de l'évènement voulu.

Calcul pour un évènement (A) selon 2 autres évènements (B et C)

$$P(A|BC) = \frac{P(B|A)P(C|A)P(A)}{P(B)P(C)}$$

En code Dans cet exemple nous allons utiliser un code existant dans la librairie python `scikit-learn`[3]. Ce moteur Naïves Bayes va nous permettre cette fois-ci de catégoriser des variétés d'iris selon la longueur et la largeur des pétales et des sépales. Les données proviennent cette fois-ci d'un dataset également disponible dans `scikit-learn`.

Nous allons utilisé le modèle *GaussianNB* de `scikit-learn` qui est adapté lorsque les données utilisées suivent une distribution normale. Ce qui semble être le cas pour les longueurs et largeur des sépale.

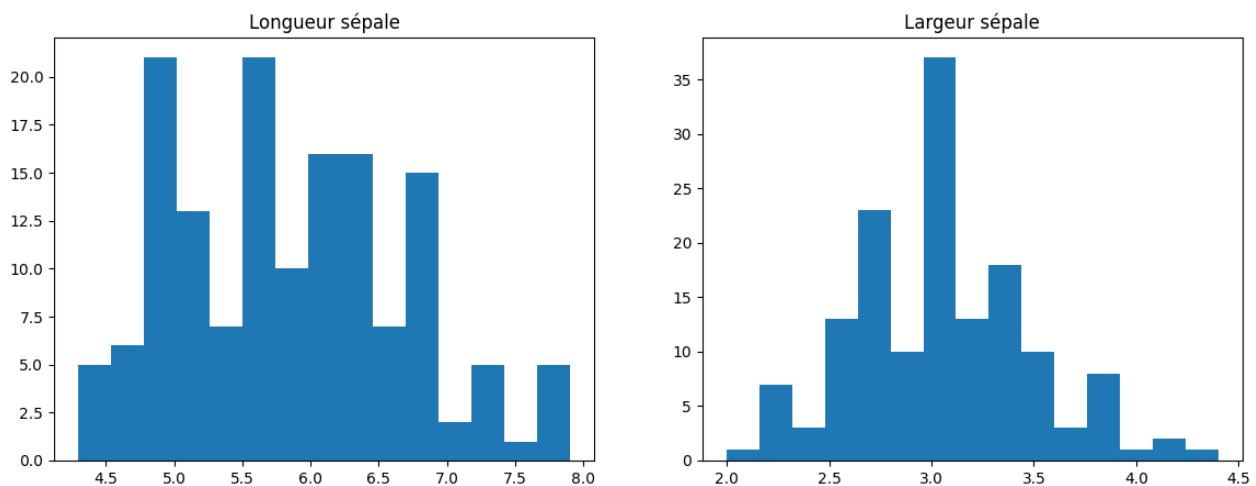


FIGURE 22 – Distribution des longueurs et largeurs des sépales

Progamme complet

```

1 from sklearn.datasets import load_iris
2 from sklearn.model_selection import train_test_split
3 from sklearn.naive_bayes import GaussianNB
4 from sklearn.metrics import accuracy_score, confusion_matrix,
   ConfusionMatrixDisplay, f1_score, \
5     recall_score
6
7 import matplotlib.pyplot as plt
8
9 X, y = load_iris(return_X_y=True)
10
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33,
12     random_state=0)
13 model = GaussianNB()
14 model.fit(X_train, y_train)
15
16 y_pred = model.predict(X_test)

```

```

16 precision = accuracy_score(y_pred, y_test)
17 recall = recall_score(y_test, y_pred, average="weighted")
18 f1 = f1_score(y_pred, y_test, average="weighted")
19
20 print("Precision:", precision)
21 print("Rappel:", recall)
22 print("Score F1:", f1)
23
24 plt.figure('Donnees du modele', figsize=(14, 5))
25 plt.subplot(1, 3, 1, title='Donnees du train set')
26 plt.scatter(X_train[:, 0], X_train[:, 1], c=y_train)
27 plt.xlabel('Sepale long.')
28 plt.ylabel('Sepale larg.')
29 plt.subplot(1, 3, 2, title='Donnees du test set')
30 plt.scatter(X_test[:, 0], X_test[:, 1], c=y_test)
31 plt.xlabel('Sepale long.')
32 plt.subplot(1, 3, 3, title='Donnees test apres evaluation')
33 plt.scatter(X_test[:, 0], X_test[:, 1], c=y_pred)
34 plt.xlabel('Sepale long.')
35 plt.show()
36
37 cm = confusion_matrix(y_test, y_pred, labels=[0, 1, 2])
38 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=[0, 1,
39                               2])
40 disp.ax_.set_title('Matrice de confusion')
41 disp.plot()
42 plt.show()
43
44 plt.figure('Distribution des donnees Iris', figsize=(14, 5))
45 plt.subplot(1, 2, 1, title='Longueur sepale')
46 plt.hist(X[:, 0], bins=15)
47 plt.subplot(1, 2, 2, title='Largeur sepale')
48 plt.hist(X[:, 1], bins=15)
49 plt.show()

```

Les données du dataset ont été séparés en 2 jeux, un pour l'entraînement du modèle et un pour le test. On obtient alors la représentation suivantes après entraînement et test du modèle

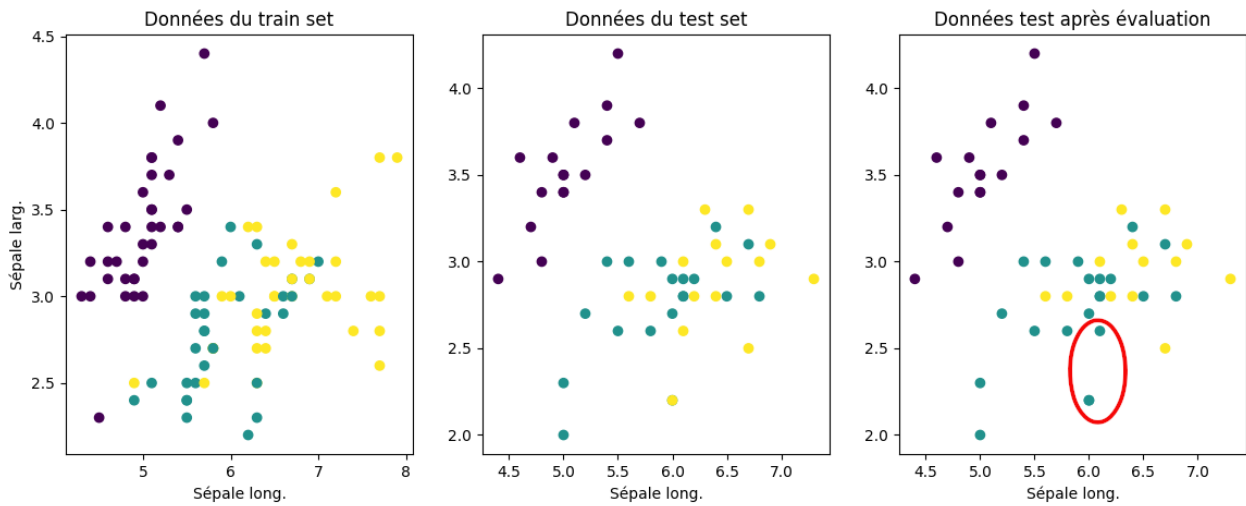


FIGURE 23 – Représentation des données

Dans les données de test nous avons 2 catégorisations qui n'ont pas été réalisées correctement. On obtient les scores suivants :

- Précision : 0.96⁴
- Rappel : 0.96⁵
- Score F1 : 0.9604285714285714⁶

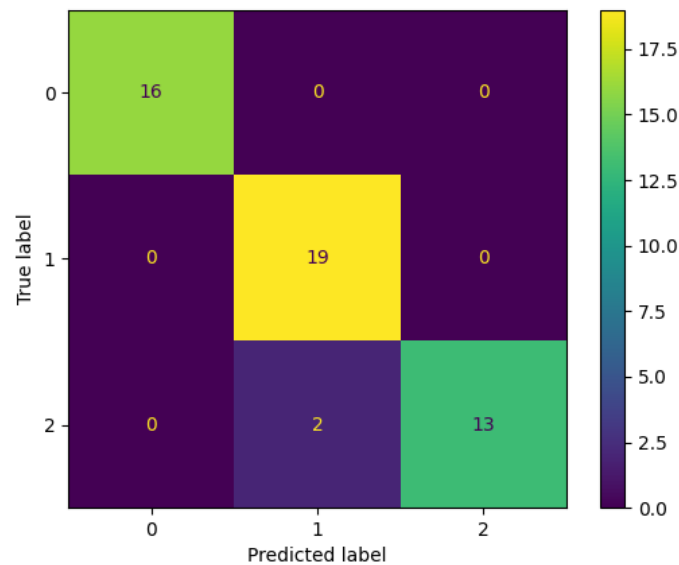


FIGURE 24 – Matrice de confusion

A l'aide de ce modèle nous devrions avoir une 96% de chance de déterminer la bonne variété d'iris en se basant sur la longueur et la largeur des sépales.

Avantages et inconvénients Le modèle Naïve Bayes est un modèle simple et rapide qui ne nécessite pas de grande capacités de calcul. De ce fait il permet de traiter une grande quantité

4. La précision est la proportion des éléments correctement identifiés sur l'ensemble des éléments prédit

5. Le rappel est la proportion des éléments correctement identifiés sur l'ensemble des éléments de la catégorie

6. Le Score F1 est la moyenne harmonique calculée de la manière suivante $2 * (precision * rappel) / (precision + rappel)$

de données.

Cependant, les données qui lui sont fournies ne doivent pas être corrélées ce qui est rarement le cas dans les problèmes du monde réel. Ce type de modèle est limité à des problèmes de classification supervisée. Si on se fie à l'équation (1) la probabilité d'apparition de l'évènement $B : P(B)$ ne peut pas être nulle.

E.2 Random tree forest

Cette section a pour but de détailler les principes généraux de la forêt d'arbres de décisions aléatoires utilisée pour faire la classification des mails.

Arbre de décision

Random forest

F Bibliographie

Références

- [1] Madjid Khichane. *Le Machine Learning avec Python*. 2021.
- [2] Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard, and David McClosky. The Stanford CoreNLP natural language processing toolkit. In *Association for Computational Linguistics (ACL) System Demonstrations*, pages 55–60, 2014.
- [3] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn : Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [4] Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. Stanza : A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics : System Demonstrations*, 2020.

G Sitotec

G.1 Corpus

- Enron company mails, fichier CSV contenant l'ensemble des mails d'une entreprise ayant fermée ses portes (33.834.245 mails) [en ligne], <https://www.kaggle.com/wcukierski/enron-email-dataset> (consulté le 27/01/2022)
- Mails project SpamAssassin, projet opensource de détection de spam (6065 fichiers email déjà trier en ham et spam) [en ligne], <https://spamassassin.apache.org/old/publiccorpus/> (consulté le 27/01/2022)
- Brown corpus, ensemble de texte en anglais publié en 1961 qui contient plus d'un million de mots <https://www.nltk.org/book/ch02.html> (consulté le 20/08/2022)

G.2 Modules

Module langdetect

- Page Github du projet *langdetect* capable de différencier 49 langages avec une précision de 99%, [en ligne] <https://github.com/Mimino666/langdetect> (consulté le 04/12/2022)
- Language Detection Library, présentation du module (anglais) [en ligne] <https://www.slideshare.net/shuyo/language-detection-library-for-java> (consulté le 04/12/2022)

G.3 Modèles

Naïves Bayes Le modèle Naïves Bayes est employé dans le module langdetect (G.2)

- Les algorithmes de Naïves Bayes, Explication sommaire du principe de ces type d'algorithme, [en ligne] <https://brightcape.co/les-algorithmes-de-naives-bayes/> (consulté le 26/03/2023)
- Naive Bayes Classification Tutorial using Scikit-learn, exemple d'utilisation de ce type de modèle avec python (anglais) [en ligne] <https://www.datacamp.com/tutorial/naive-bayes-scikit-learn> (consulté le 26/03/2023)
- Scikit learn Naive Bayes, description des types d'algorithme disponibles dans le module Scikitlearn en python (anglais) [en ligne] https://scikit-learn.org/stable/modules/naive_bayes.html (consulté le 26/03/2023)

H Codes sources

H.1 Github

Le lien vers l'ensemble des sources est disponible en publique via le lien ci dessous : <https://github.com/peredur0/mercury>

H.2 Analyse statistiques de la phase 1

Code pour l'affichage des statistiques de la phase 1

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3
4 from databases.psql_cmd import connect_db, exec_query
5 from databases.psql_db import secrets as psql_secrets
6
7
8 def get_p1_data():
9     psql_cli = connect_db('mail_features_prod', psql_secrets.owner,
10                          psql_secrets.owner_pw,
11                          psql_secrets.host, psql_secrets.port)
12
13     column = ['id_message', 'type', 'url', 'mail', 'tel', 'nombre', 'prix']
14     query = """select m.id_message, c.type, l.url, l.mail, l.tel, l.
15 nombre, l.prix from messages
16 as m
17 join categories as c on m.id_cat = c.id_cat
18 join liens as l on m.id_message = l.id_message"""
```



```

17
18     df = pd.DataFrame(exec_query(psql_cli, query), columns=column)
19     psql_cli.close()
20
21     return df.set_index('id_message')
22
23
24 def set_bar_graph(data, feat, subplot, pos):
25     df = data[data[feat] < 20].groupby(['type', data[feat]]).size()
26     df.unstack(0).plot(kind='bar', ax=subplot[pos])
27
28
29 if __name__ == '__main__':
30     df_all = get_p1_data()
31     df_spam = df_all[df_all['type'] == 'spam']
32     df_ham = df_all[df_all['type'] == 'ham']
33
34     d_pie = df_all.groupby(['type']).size()
35     fig, ax = plt.subplots()
36     fig.suptitle('Repartition des ham/spam')
37     ax.pie(d_pie, labels=d_pie.index, autopct='%1.1f%%')
38     plt.show()
39
40     print("Statistiques Liens")
41     print("Globales: \n", df_all.describe())
42     print("Ham: \n", df_ham.describe())
43     print("Spam: \n", df_spam.describe())
44
45     fig, ax = plt.subplots(nrows=5, ncols=1)
46     fig.suptitle("Distribution des mails en fonction du nombre de liens")
47     fig.tight_layout(pad=0.5)
48     position = 0
49     for feat in ['url', 'mail', 'tel', 'nombre', 'prix']:
50         set_bar_graph(df_all, feat, ax, position)
51         position += 1
52     plt.show()
53

```
