

Etape 5 : Analyse des Descripteurs Textuels et Désinformation Médicale

1. Introduction

Ce projet vise à explorer et analyser les données textuelles issues de différentes sources (articles, tweets, etc.) dans le but d'identifier des descripteurs pertinents liés à la désinformation médicale. Les méthodes employées permettent de sélectionner les termes les plus représentatifs tout en mettant en évidence leur lien avec des indices de désinformation médicale, tels que des expressions sensationnalistes ou pseudoscientifiques.

2. Méthodes Utilisées

1. Prétraitement des données textuelles :

- Nettoyage des données : suppression des stopwords, gestion des caractères spéciaux, lemmatisation.
- Représentation des textes : transformation des corpus multilingues en représentations numériques via **TF-IDF** (Term Frequency-Inverse Document Frequency).

2. Sélection des descripteurs :

- **Méthodes de filtrage :**
 - Utilisation du Chi-carré pour évaluer l'association statistique des termes avec les classes pertinentes.

```
a. Methodes de Filtrage

from sklearn.feature_selection import chi2
import numpy as np

# Exemple avec Chi-carré
X = tfidf_matrix # Matrice TF-IDF
y = np.random.randint(2, size=X.shape[0]) # Remplacez avec des labels si c
chi2_scores, p_values = chi2(X, y)

# Ajouter les scores au DataFrame TF-IDF
tfidf_df['chi2_score'] = chi2_scores # en anglais
print(tfidf_df.sort_values(by='chi2_score', ascending=False).head(10))
```

	term	score	chi2_score
9	controversial	10.692919	5.768868
91	vaccine	159.932777	3.906879
13	covid	403.455445	3.054445
58	protect	11.789267	3.037841
83	time	31.357570	2.920269
88	updated	30.588336	2.721576
64	remain	14.005084	2.683759
84	treatment	15.461314	2.603147
99	workers	17.970466	2.534372
94	virus	14.588408	2.179747

- **Méthodes intégrées :**
 - Importance des caractéristiques calculée à l'aide de modèles tels que Random Forest ou régularisation L1 (Lasso).

c. Methodes intégrées

```

### ===== Anglais

from sklearn.ensemble import RandomForestClassifier

rf_model = RandomForestClassifier()
rf_model.fit(tfidf_matrix, y)

# Importance des termes
feature_importances = rf_model.feature_importances_
important_terms = pd.DataFrame({'term': terms, 'importance': feature_importances})
print(important_terms.sort_values(by='importance', ascending=False).head(10))

```

✓ 1.3s

	term	importance
13	conseils	0.075793
52	maladie	0.045747
3	anticovid	0.031647
91	teleservice	0.031627
22	dun	0.026926
85	sanofi	0.026460
31	face	0.021901
92	temps	0.020799
10	chine	0.019716
83	sanitaire	0.017816

3. Mise en contexte :

- Identification des termes liés à la désinformation médicale, par exemple :
 - Entités : COVID-19, vaccins, traitements, maladies.
 - Motifs linguistiques : phrases pseudoscientifiques ou sensationnalistes.

3. Analyse des Résultats

1. Pertinence des termes identifiés :

- Des termes fréquents tels que "covid", "vaccine", et "disease" ont été identifiés comme des descripteurs clés.
- Les termes avec un score élevé au Chi-carré indiquent un lien fort avec des indices de désinformation.

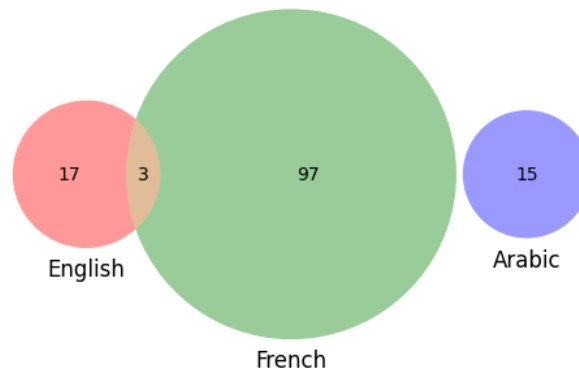
2. Observations importantes :

- Certains termes sont universels (ex. "covid") tandis que d'autres sont spécifiques à une langue ou un contexte culturel.
- Les méthodes de filtrage ont permis de réduire efficacement la dimensionnalité tout en conservant des termes pertinents.

4. Visualisations

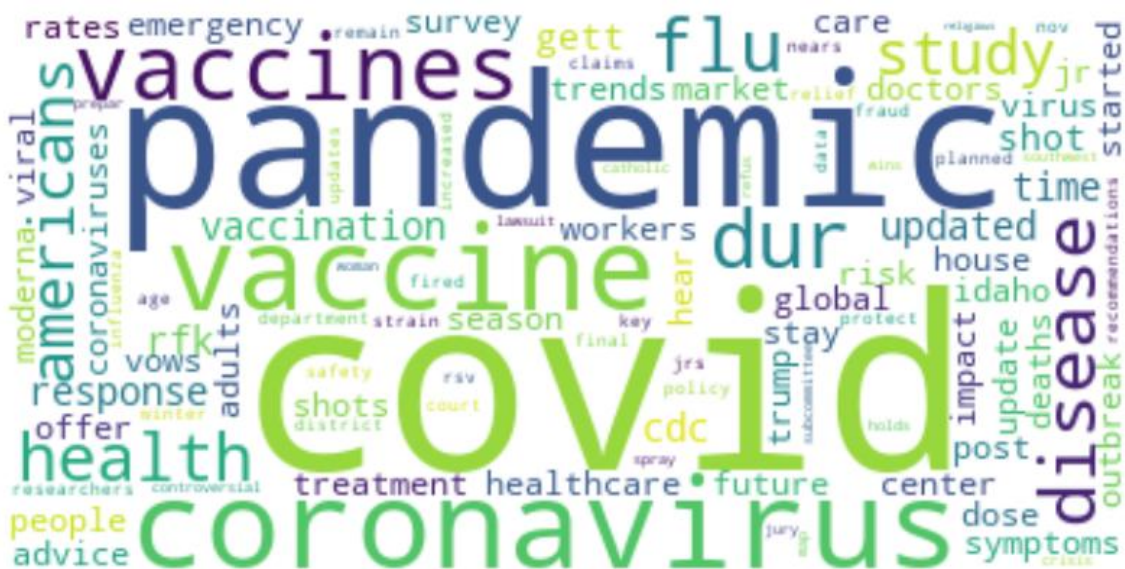
1. Diagramme de Venn :

- Illustre les chevauchements des descripteurs sélectionnés entre trois langues.
- Exemples :
 - Termes communs : "vaccine", "covid".
 - Termes spécifiques : "miracle" (anglais), "remède" (français), "لقاح" (arabe).



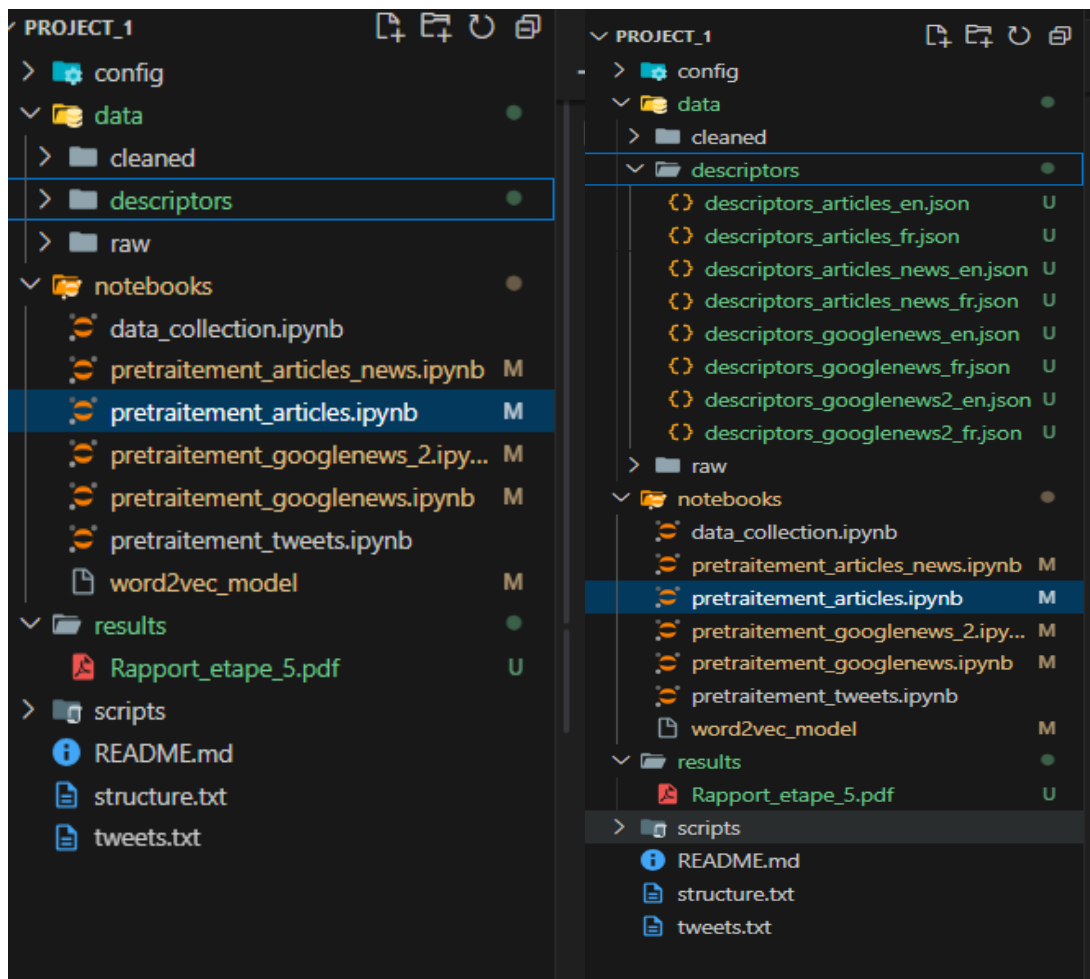
2. Nuages de mots pondérés :

- Représentation visuelle des termes avec leur poids respectif (TF-IDF).
- Les termes fortement associés à la désinformation (ex. "remède miracle") sont mis en avant.



Sauvegarde des résultats

Les descripteurs sélectionnés, accompagnés de leurs scores et méthodes associées, ont été sauvegardés dans le répertoire suivant :



5. Conclusion

Les méthodes employées permettent d'identifier efficacement les termes pertinents pour analyser la désinformation médicale. Les résultats mettent en évidence des similitudes et des spécificités entre les langues, ainsi que des termes particulièrement liés à la désinformation. Ces observations peuvent être utilisées pour développer des systèmes de détection automatique ou sensibiliser à la désinformation médicale.