

La détection des discours de haine et du langage vulgaire sur Twitter

Mohamed-Amine ROMDHANE, Amine BENTELLIS

1 Introduction

Avec l'utilisation hebdomadaire des réseaux sociaux, on se trouve des fois témoins de plusieurs discours de haine visant la race, la religion ou même certaines minorités. La détection de tels discours de discrimination a une nécessité préventive cruciale.

Ce rapport vise à mettre en place un modèle sophistiqué qui permet de détecter et prédire les discours de haine ainsi que le langage vulgaire. La langue vivante utilisé dans le traitement et l'analyse de texte est l'**anglais**.

2 Jeux de données

Les jeux de données ont été prit du site de [Hate Speech Datasets](#). On en a utilisé les suivant :

Nom de la base	Nombres d'échantillons	Étiquettes
Davidson & al.	~25000	hate speech, offensive langage, normal
Waseem & Hovy	~16000 (~5000 utilisés)	racism, sexism, neither

Le point commun de ces jeux de données est la source. Les deux contiennent du texte qui provient de Twitter. Cependant, l'étiquettage n'est pas le même. Pour cela, on a gardé l'étiquettage par défaut de la base *Davidson & al.* et on a fusionné les étiquettes de la base *Waseem & Hovy* comme suit : **racism** et **sexism** en **hate speech**, **neither** en **normal**.

Le jeu de données *Davidson & al.* contient du texte prêt à être utilisé. Néanmoins, la base *Waseem & Hovy* contient des *Tweet Status ID*. Le but sera de s'en servir de l'API de Twitter pour chercher le texte correspondant à chaque ID. Un package Python au nom de [Tweepy](#) est utilisé comme interface pour accéder à l'API.

En vu des suppressions des comptes incitant à la haine, on ne peut pas retrouver la totalité des données comprises dans la base de *Waseem & Hovy*, d'où les ~5000 échantillons utilisés.

2.1 Prétraitement

Le texte provenant des tweets bruts contient du bruit ; des emojis, des mentions, des hashtags, des liens... Le but du prétraitement est d'enlever ou de lisser ce bruit. La fonction `sanitize_tweets` de la classe `PreProcessing` permet d'accomplir ce manœuvre. Le code est le suivant :

```
1  @staticmethod
2  def sanitize_tweet(tweet, with_stemmer=True, stemmer="porter":
3      # Choix des Stemmers pour avoir la racine d'un mot.
4      if stemmer == "snowball":
5          STEMMER = SnowballStemmer("english")
6      else:
7          STEMMER = PorterStemmer()
8      # Enlève plusieurs " de suite.
9      text = re.sub("\s+", "", tweet)
10     # Enlève plusieurs espaces de suite et le remplacer par un seul.
11     text = re.sub("\s+", " ", text)
12     # Enlève RT (ancronyme des Re-Tweets)
13     text = re.sub("RT", "", text)
14     # Enlève les liens
15     text = re.sub("https?:/?/?/?[^\s]+", "", text)
16     # Enlève les caractères spéciaux codés en HTML.
17     text = re.sub("&.+;", "", text)
18     # Initialisation du TweetTokenizer
19     # @param preserve_case=False : Tout est en miniscule.
20     # @param strip_handles=True : Enlève les mentions.
21     # @param reduce_len=True : Réduit une séquence de lettres successives
22     # dans un mot à une longueur de 3 si la longueur de cette séquence
23     # est supérieur à 3.
24     tknzs = TweetTokenizer(preserve_case=False,
25                             strip_handles=True,
26                             reduce_len=True)
27     # On prend les mots constitués de lettre uniquement tout en enlevant
28     # les stopwords.
29     text = " ".join([STEMMER.stem(s)
30                      if with_stemmer
31                      else s
32                      for s in tknzs.tokenize(text)
33                      if s.isalpha() and s not in STOPWORDS])
34     return text
```

3 Analyse des données

3.1 Distribution des étiquettes

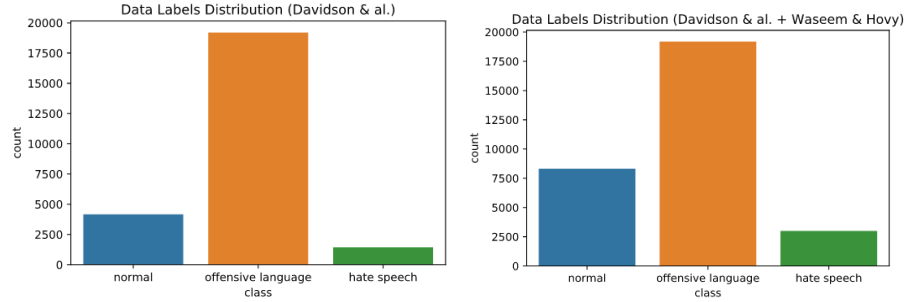
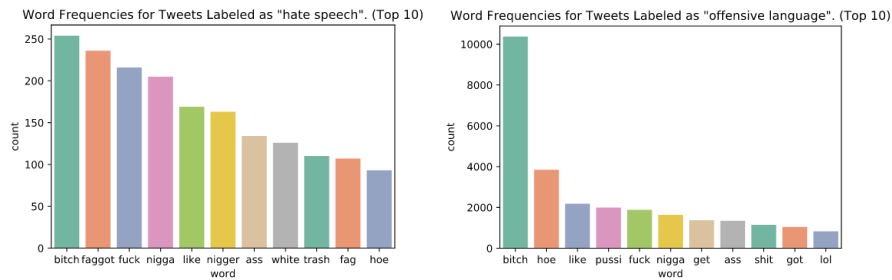


FIGURE 1 – Distribution des étiquettes sur l'ensemble des données... À droite on a *Davidson & al.*, à gauche on a le même mais en ajoutant *Waseem & Hovy*

Comme on peut voir dans ces deux figures, l'étiquetage n'est pas équilibré. *Offensive language* représente une bonne partie des données. Ceci peut poser problème s'il y'a du vocabulaire à la fois utilisé dans *offensive language* et dans *hate speech*, si c'est le cas le modèle de classification va prioriser l'étiquette qui a le plus de poids. On indique dans la partie *Classification* comment on évite un tel cas.

3.2 Fréquence des mots

Il est intéressant de voir les mots les plus utilisés pour les étiquettes qu'on a. Les figures suivantes nous montrent un aperçu des top 10 pour chaque classe :



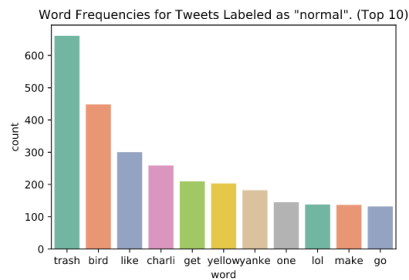


FIGURE 2 – Fréquence des mots pour chaque étiquette

Comme on peut le constater, *offensive language* et *hate speech* partagent beaucoup de mots. On peut conclure qu'il faut pas seulement avoir un vocabulaire rempli que de mots mais aussi de bigrammes ou/et même des trigrammes, sinon la précision du modèle sera très basse.

4 Classification

On utilise la librairie Scikit-Learn pour l'extraction des caractéristiques de notre corpus ainsi que son processus de classification.

La classification contient 3 modules :

- **Count Vectorizer** : Construit une matrice de décomptes des tokens appartenant à un vocabulaire tiré à partir d'une collection de documents. Ce vocabulaire peut contenir des mots simples ainsi que de n-grammes.
- **TF-IDF Transformer** : Transforme une matrice de décomptes à une représentation normalisée de TF-IDF (Term-Frequency - Inverse Document-Frequency). Cette représentation permet de minimiser l'impact des termes qui sont très fréquents dans le corpus tout en valorisant l'impact des zones de corpus qui semblent peu informatives.
- **SGD Classifier** : Classificateur linéaire avec un entraînement qui utilise l'algorithme de *Stochastic Gradient Descent*. Ce classificateur peut être paramétré pour utiliser plusieurs modèles.

On peut s'amuser à observer l'impacte des n-grammes sur la précision du modèle. On peut à la suite conclure que l'utilisation d'un vocabulaire consti-

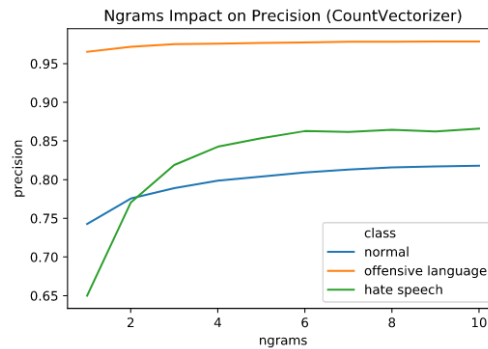


FIGURE 3 – Impacte des n-grammes sur la précision du modèle

tué de mots uniques, bigrammes, trigrammes, quadrigrammes, pentagrammes et hexagrammes améliore considérablement la précision de classification des étiquettes de type *hate speech*.

Le code utilisé pour installer ce système de classification est le suivant :

```

1  ...
2  pipe = Pipeline([
3      # Matrice de décomptes avec un vocabulaire de mots uniques,
4      # bigrammes, trigrammes... jusqu'au hexagrammes.
5      ('count', CountVectorizer(ngram_range=(1, 6))),
6      # Matrice de transformation TF-IDF
7      ('tfidf', TfidfTransformer()),
8      # SGD Classifier: (Support Vector Machines par défaut)
9      # @param class_weight="balanced" : paramètre qui tient en compte
10     # les poids des étiquettes. Ce mode règle les poids
11     # inversement proportionnels aux fréquences des étiquettes.
12     # @param shuffle=True : permet de mélanger les documents de
13     # corpus.
14     ('clf', SGDClassifier(class_weight="balanced", shuffle=True)),
15 ])
16
17 # On lance la classification
18 # data["tweet"] est la colonne qui contient les 30000 tweets.
19 # data["class"] est la colonne qui contient les étiquettes
20 # respectives de ces tweets.
21 pipe = pipe.fit(data["tweet"], data["class"])
22 ...

```

Quand on vérifie la performance du modèle après l'exécution du code précédent, on obtient ces statiques.

	precision	recall	f1-score	support
hate speech	0.86	0.84	0.85	3004
normal	0.81	0.98	0.89	8319
offensive language	0.98	0.89	0.93	19190
accuracy			0.91	30513
macro avg	0.88	0.90	0.89	30513
weighted avg	0.92	0.91	0.91	30513

On obtient bien des precisions respectables pour chacune des étiquettes... Les F1-scores sont aussi considérablement élevés. Cela indique que les precisions qu'on a représentent bien du texte qui a été correctement identifié.

5 Conclusion

Bien qu'on peut bien classifier les discours de haine et le langage vulgaire avec des precisions très proches de 1, un jeu de données avec 30000 de tweets reste très peu pour couvrir la totalité des termes qu'on peut voir dans un cas réel. Comme vu dans la partie de fréquences de mots, il y'a des mots qui sont très répondus et d'autres pas, même s'ils ont le même sens ainsi que le même impacte sur la personne.