



Data Science Bootcamp

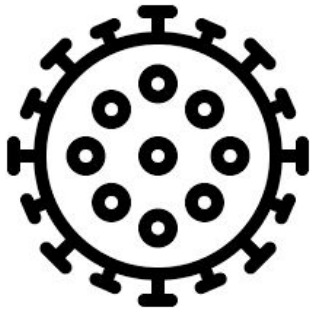
Carlo CISALE



[linkedin.com/in/carlo-cisale](https://www.linkedin.com/in/carlo-cisale)



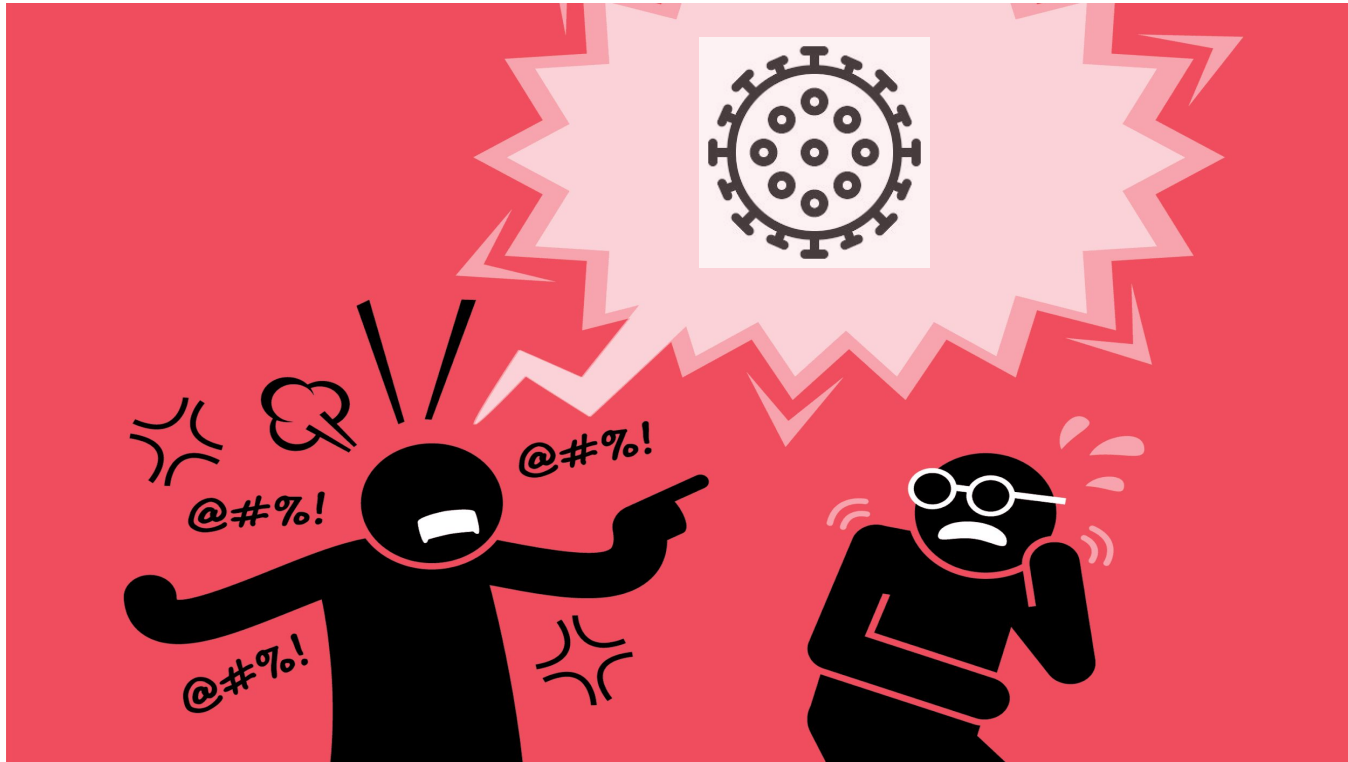
Sentiment Analysis on #coronavirus tweets



Motivation



Motivation



Steps

- Scraping tweets en anglais avec mots cles “Coronavirus” et “Italians” avec [twitterscraper](#)
- Data preparation des tweets
- Libelisation des tweets avec librarie [textblob](#)
- Entrainement modèle de ML sur un dataset de hate speech
- Test du modèle sur les tweets coronavirus

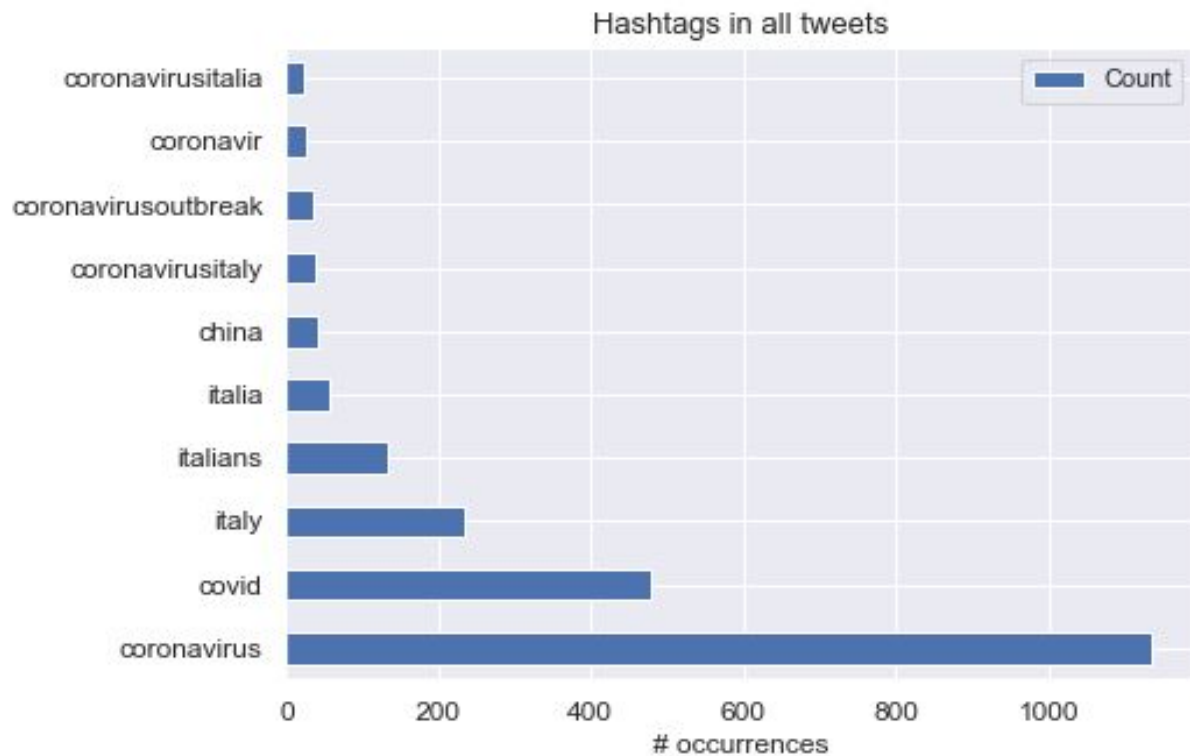
Difficulté avec la data preparation

- Format des tweets “scrapé”: Json avec beaucoup d’info
- Doublons
- Suppression des liens: `http...`, `pic.twitter...`

Word Cloud



Hashtags les plus populaires



Comment convertir données textuelles en numériques?

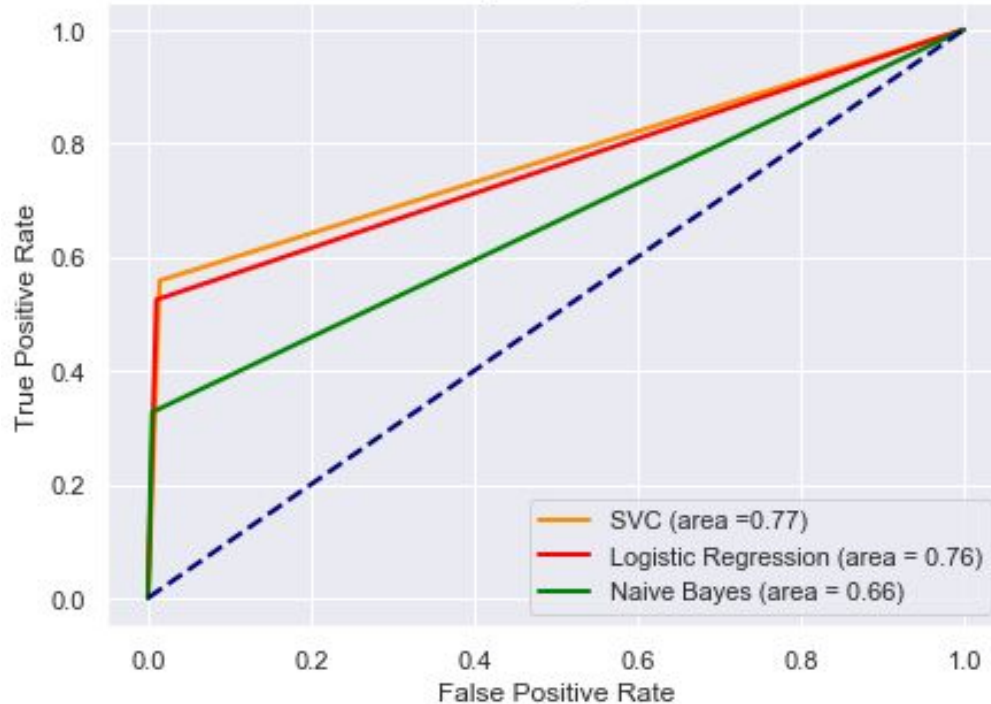
— TF*IDF: Term Frequency – Inverse Document Frequency

$$TF(t, d) = \frac{\text{number of times term}(t) \text{ appears in document}(d)}{\text{total number of terms in document}(d)}$$

$$IDF(t, D) = \log \left(\frac{\text{total number of documents}(D)}{\text{number of documents with the term}(t) \text{ in it}} \right)$$

Comparison algorithms (hate speech data)

Receiver operating characteristic

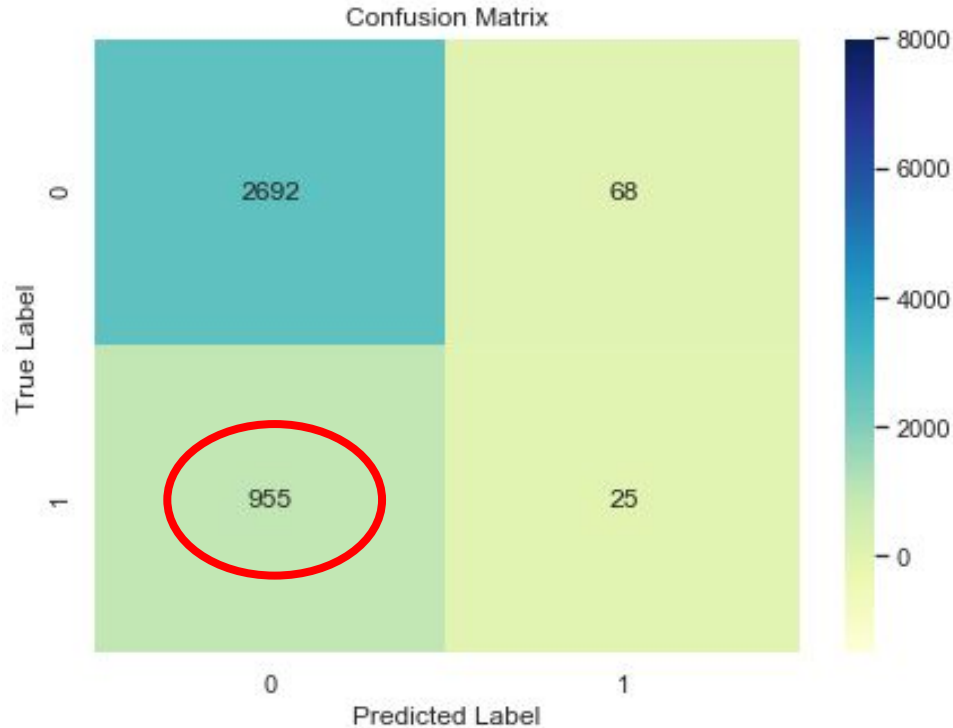


Algo	F-1 score
SVC	0.64
Logistic	0.64
NB	0.53

Matrice de confusion

	Predicted label class 1	Predicted label class 2
True label class 1	correct true positive for class 1	wrong false positive for class 2
True label class 2	wrong false positive for class 1	correct true positive for class 2

Matrice de confusion (coronavirus data)



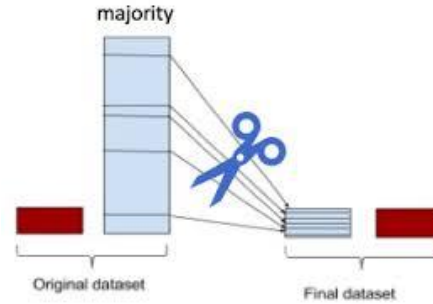
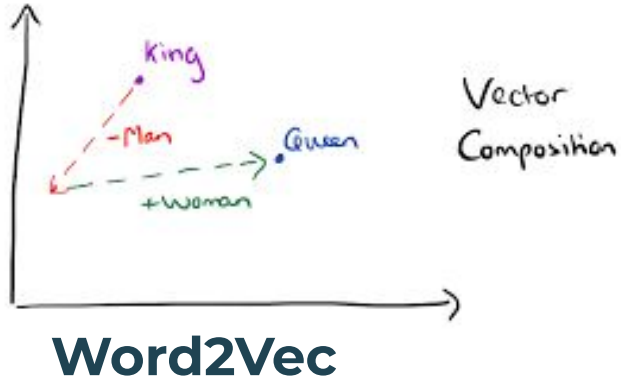
Matrice de confusion (coronavirus data)



Limitations. Pourquoi?

1. Mauvaise classification de TextBlob des tweets négatifs (des tweets normaux étaient classifié comme péjoratifs).
2. Les deux datasets (hate speech et coronavirus) n'avaient pas la même taille.
3. Le modèle prends en considération “high signal words”. Mais qu'est ce que il se passe avec des nouveaux mots?
4. Le dataset initial de hate speech (mais aussi du coronavirus) sont imbalanced.

Conclusions et améliorations



**Downsampling/
Upsampling**



Human classification





Data Science Bootcamp

Des questions ?

