

# Détecter les BadBuzz grâce au Deep Learning

---



**Air Paradis**

P7 - Sylvia Bankowska - Septembre 2022

# Sommaire



1. Enjeux et objectifs
  2. Données
  3. Prétraitement
  4. Word Embeddings
  5. Modèle sur mesure simple
  6. Modèles sur mesure avancés
  7. Résultats
  8. Mise en production
-

# Enjeux et objectifs

# Enjeux

- Anticiper des tweets potentiellement **nuisibles**
- Prédire le **sentiment** associé à un tweet

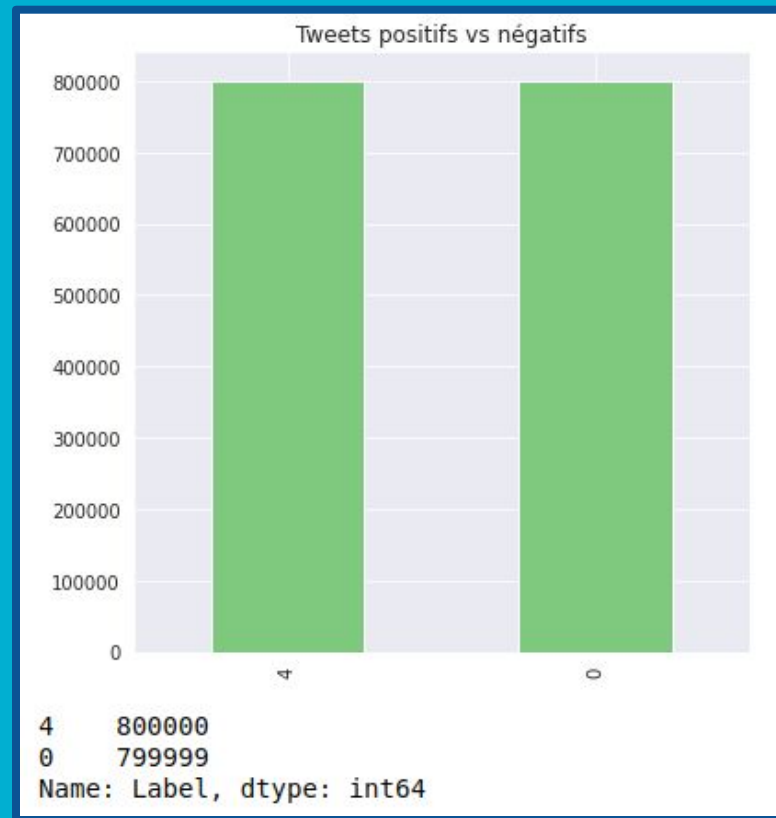
# Objectifs

- Comparer plusieurs modèles de réseaux de neurones
  - Créer le prototype d'un produit IA
-

# Données

Analyse du dataset

- **Jeu de données :**
  - 1 600 000 tweets
  - propre
- **Echantillon :**
  - contenu d'un tweet + label
  - 8000 pos + 8000 neg
- **Target :**
  - variable binaire
  - classes très équilibrées
- **Métrique :**
  - score F1 + AUROC

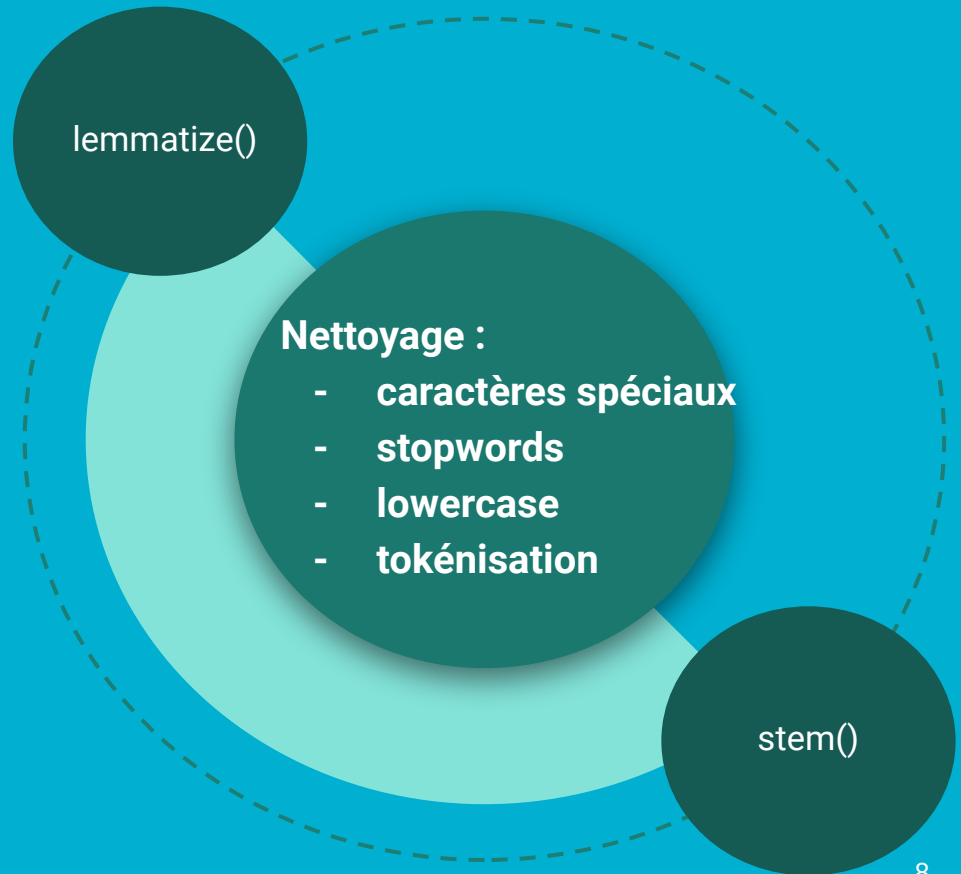


# Prétraitement

Transformer le texte en chiffres

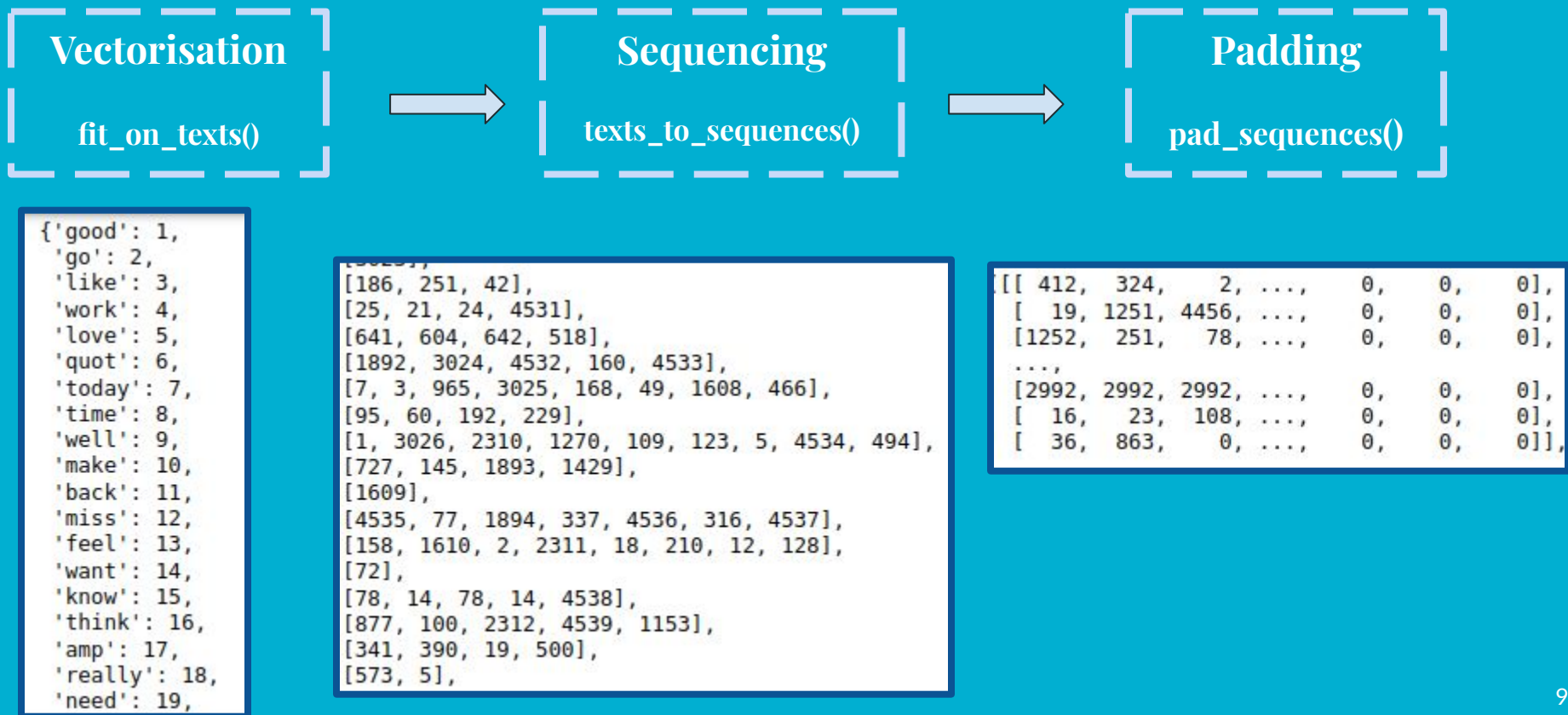
---

Stemming	Lemmatisation
Forme tronquée : <ul style="list-style-type: none"><li>➤ suppression des suffixes</li><li>➤ suppression des flexions</li></ul>	Forme canonique : <ul style="list-style-type: none"><li>➤ analyse morpho</li></ul>





# Transformation du texte en chiffres

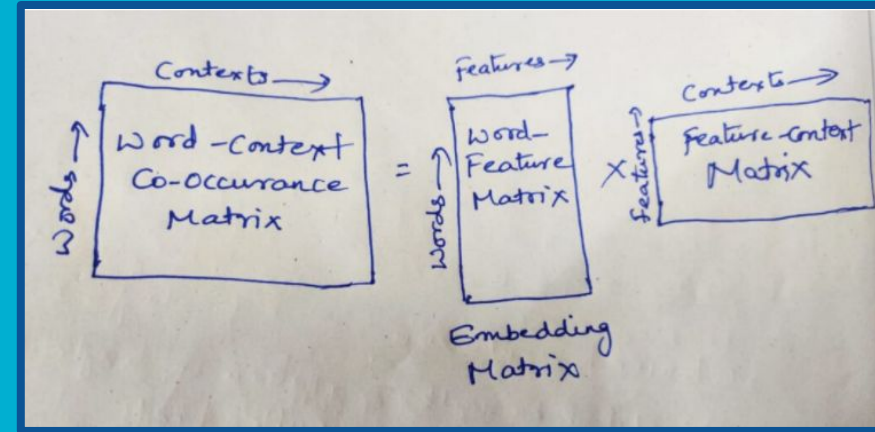
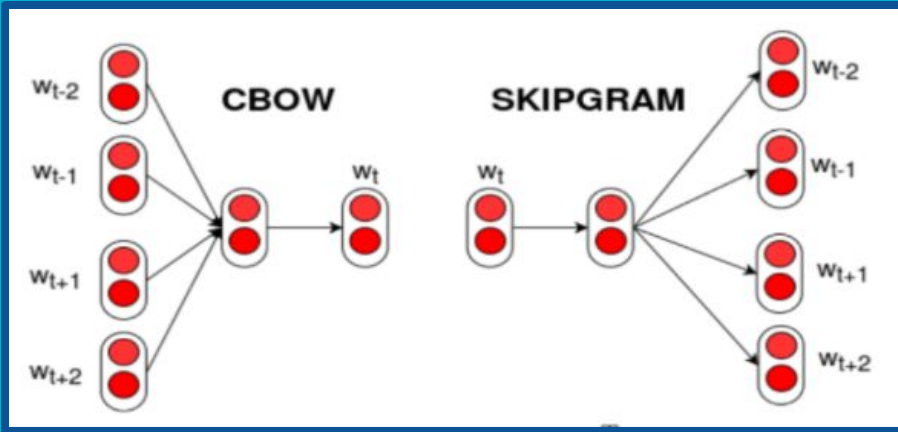


# Word Embeddings

# Word2Vec

vs

# GloVe



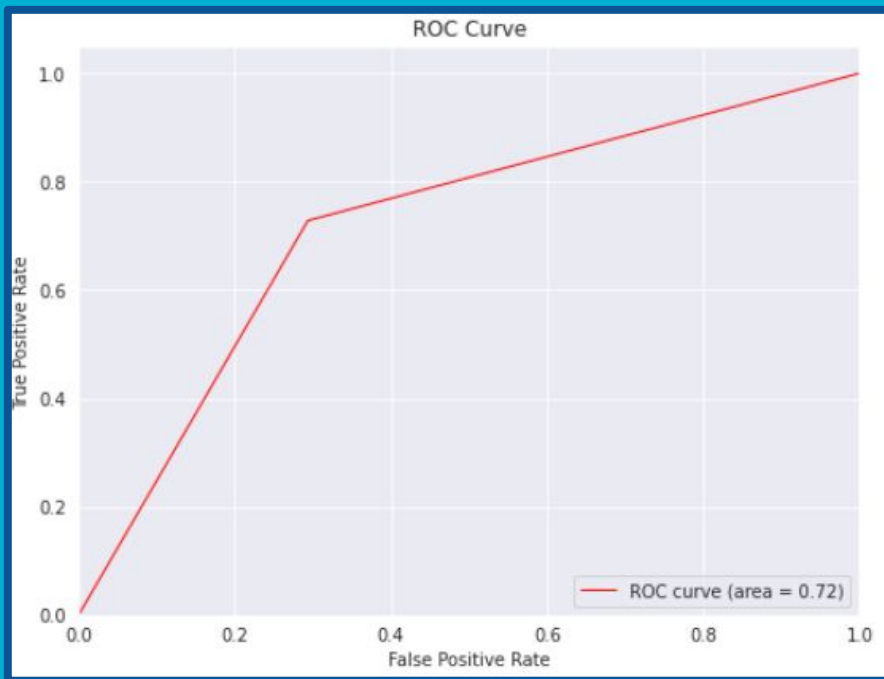
- réseau de neurones feed forward
- capture le contexte similaire

- calcul basé sur les cooccurrences de mots sur l'ensemble du corpus
- capture la probabilité que 2 mots apparaissent ensemble

# Modèle sur mesure simple

# Régression Logistique

La **FONCTION SIGMOIDE** permet de mesurer si une entrée a dépassé le seuil de classification.



F1 score: 0.7197274698048932

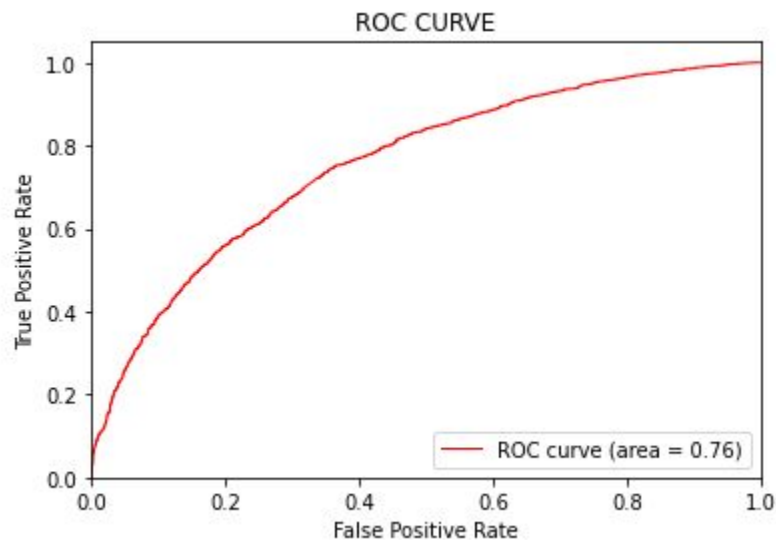
## Classification Report

	precision	recall	f1-score	support
0	0.72	0.71	0.71	1604
1	0.71	0.73	0.72	1596
accuracy			0.72	3200
macro avg	0.72	0.72	0.72	3200
weighted avg	0.72	0.72	0.72	3200

# Modèles sur mesure avancés

Keras + Embedding / RNN & LSTM / CNN / BERT

# Couche Embedding de Keras

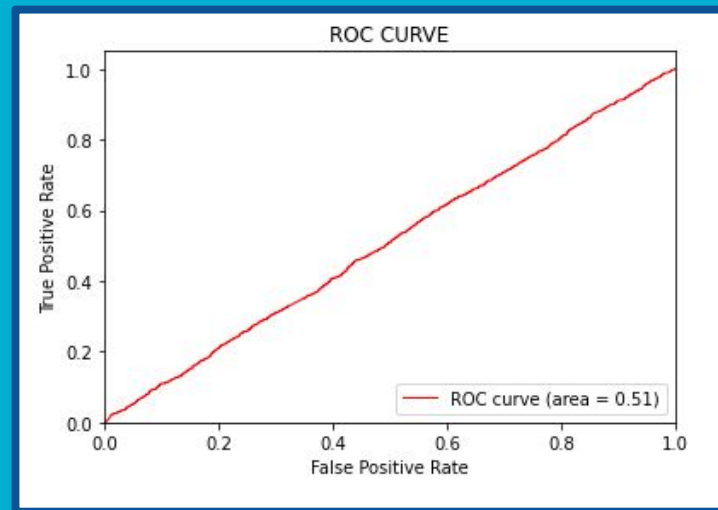


	precision	recall	f1-score	support
0	0.52	0.45	0.48	2452
1	0.50	0.57	0.53	2348
accuracy			0.51	4800
macro avg	0.51	0.51	0.50	4800
weighted avg	0.51	0.51	0.50	4800

- optimisation d'hyperparamètres
- tuner Hyperband de Keras Tuner
- tuning de l'hypermodèle :
  - best learning rate
  - best epochs
  - nb de neurones couche Dense

# RNN & LSTM

- Entraînement :
  - avec Word2Vec et GloVe
  - sur un corpus lemmatisé et stemmé
- Résultats : médiocres
- **Le meilleur score F1** => corpus prétraité avec le stemmeur :
  - Word2Vec : 0.54
  - GloVe : 0.51



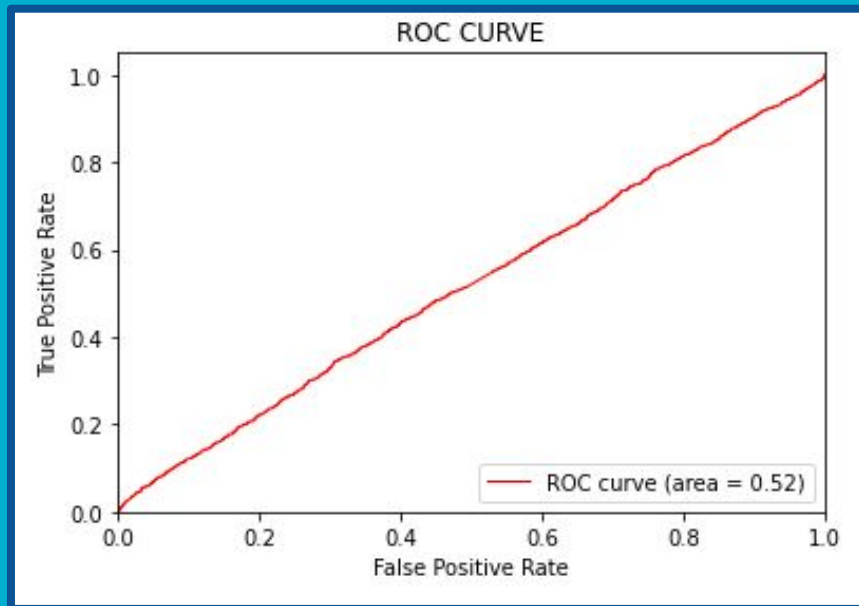
```
predict ("I don't like cats", model_rnn)
1/1 [=====] - 1s 1s/step
('Label: ', 'POSITIVE', 'Score: 0.60', 'Elapsed_time: 1.46')

predict ("I love cats", model_rnn)
1/1 [=====] - 0s 24ms/step
('Label: ', 'POSITIVE', 'Score: 0.60', 'Elapsed_time: 0.06')
```



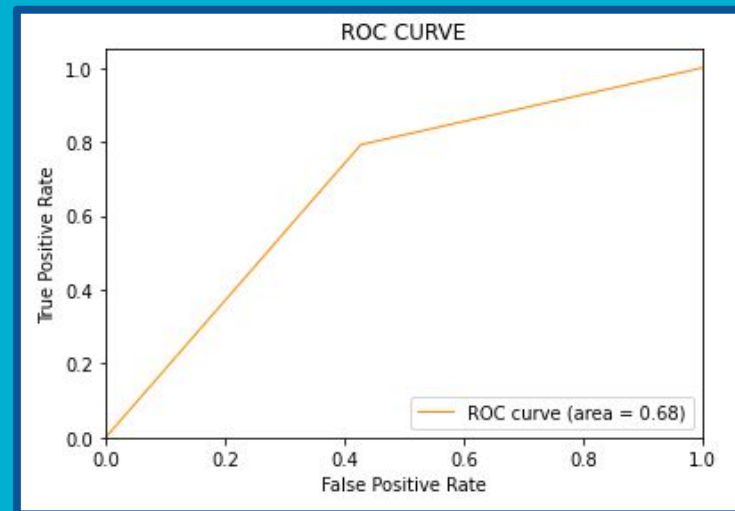
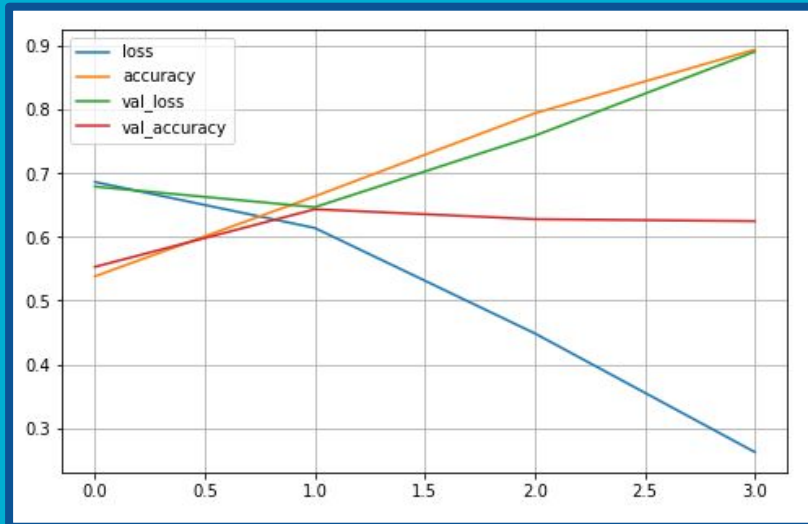
# CNN

- Résultats :
  - médiocres
  - meilleurs que RNN
- **Le meilleur score F1 =>**  
corpus prétraité avec le  
lemmatiseur :
  - Word2Vec : 0.54
  - GloVe : 0.54



# BERT

- Lenteurs à l'exécution
- Résultats :
  - meilleurs que RNN ou CNN
- Le score  $F_1$  :
  - corpus lemmatisé : 0.70
  - corpus stemmé : 0.62



# Résultats

Comparaison des scores

# Corpus

Stemming

Lemmatisation

	Loss	Acc	AUROC	Time	F1
Logistic Regression	/	0.72	0.72	1.19	0.72
Keras Emb STEM	0.6	0.69	0.76	6.88	0.47
W2V RNN / LSTM STEM	0.91	0.5	0.5	123.71	0.54
W2V CNN STEM	0.87	0.51	0.51	21.72	0.52
Glove RNN / LSTM STEM	0.9	0.5	0.51	47.84	0.51
Glove CNN STEM	0.8	0.51	0.52	7.64	0.53
BERT STEM	0.73	0.49	0.62	2016.06	0.62

	Loss	Acc	AUROC	Time	F1
Logistic Regression	/	0.72	0.72	1.19	0.72
Keras Emb LEMMA	0.6	0.69	0.76	5.18	0.53
W2V RNN / LSTM LEMMA	0.81	0.51	0.51	49.0	0.51
W2V CNN LEMMA	0.81	0.5	0.5	11.88	0.54
Glove RNN / LSTM LEMMA	0.89	0.5	0.5	37.26	0.48
Glove CNN LEMMA	0.84	0.51	0.51	7.63	0.54
BERT LEMMA	0.67	0.6	0.68	1940.3	0.7

# Mise en production

Fast API + GitHub + Heroku



Curl

```
curl -X 'POST' \  
  'https://fastapi-projet.herokuapp.com/prediction' \  
  -H 'accept: application/json' \  
  -H 'Content-Type: application/json' \  
  -d '{  
    "review": "I enjoyed the flight, thank you."  
  }'
```

Request URL

```
https://fastapi-projet.herokuapp.com/prediction
```

Server response

Code

Details

200

Response body

```
{  
  "prediction": "The sentiment is positive :-)"  
}
```

Curl

```
curl -X 'POST' \  
  'https://fastapi-projet.herokuapp.com/prediction' \  
  -H 'accept: application/json' \  
  -H 'Content-Type: application/json' \  
  -d '{  
    "review": "This flight was absolutely awful."  
  }'
```

Request URL

```
https://fastapi-projet.herokuapp.com/prediction
```

Server response

Code

Details

200

Response body

```
{  
  "prediction": "The sentiment is negative :-("  
}
```

# Conclusions

Envisager des actions supplémentaires pour obtenir de meilleurs résultats :

- augmenter la volumétrie de l'échantillon
- mieux prétraiter les données textuelles
- une recherche d'hyper paramètres plus poussée

---

A photograph taken from the perspective of a passenger looking out of an airplane window. The wing of the aircraft is visible in the upper left corner, extending towards the center. Below the wing, a vast expanse of white, fluffy clouds stretches across the horizon. The sky is a mix of light blue and warm orange, suggesting the sun is low on the horizon. The overall mood is serene and expansive.

**Avez-vous des questions ?**