# POSTER PROJET LU3IN026

Dans ce projet, nous avons étudié les performances de plusieurs algorithmes d'apprentissage automatique sur les données textuelles du forum 20 Newsgroups, qui contient des messages issus de 20 forums de discussion différents.

L'objectif est d'appliquer et comparer différents algorithmes d'apprentissage supervisé et non-supervisé sur des données textuelles. Étudier l'influence des représentations et des classifieurs sur la performance (accuracy, temps, robustesse).

#### Nous avons appliqué à la fois des approches supervisées en approche supervisée et non-supervisées (clustering), en expérimentant différentes représentations vectorielles du texte : bag-of-words, fréquences, TF-IDF, avec ou sans suppression des stopwords.

#### **ECHANTILLON**

- API
- Python
- NumPy Pandas
- Graphviz Matplotlib
- **GUILLAUME DUPART** YANIS TOUTAIN GROUPE 2

LU31N026 - 2024 2025

**SCIENCES** 

## **PRÉTRAITEMENT**

### **NETTOYAGE DU TEXTE**

- Nettoyage → Met le texte en minuscule, remplace la ponctuation par un espace
- text2vect → Supprime les mots inutiles (conjonction, verbes trop communs ... ).  $\rightarrow$  S-W

## REPRÉSENTATION VECTORIELLES

- Bag-of-words binaire (1 si au moins une fois dans le document, 0 sinon).
- Bag-of-words par comptage (compteur de mots par document)
- Fréquence relative (compteur mot / total)
- TF-IDF (l'importance d'un mot dans un document donné, tout en réduisant le poids des mots trop fréquents dans l'ensemble du corpus).

## DONNÉES

- 18846 lignes
- Train (avec stop-words): 173 lignes • Test (avec stop-words): 932 lignes
- Train (sans stop-words): 173 lignes
- Test (sans stop-words): 932 lignes

## APPROCHE SUPERVISÉE

#### KNN DISTANCE EUCLIDIENNE

Avec BOW binaire

- Avec filtre S-W : 5.7% → 105 sec
- Sans filtre S-W: 5.7% → 84 sec Avec BOW comptage
- Avec filtre S-W : 6.3% → 108 sec • Sans filtre S-W: 8.8% → 79 sec
- Avec Fréquence relative
- Avec filtre S-W: 5% → 74 sec • Sans filtre S-W: 8.4% → 21.7 sec
- Avec TF-IDF:
- Avec filtre S-W: 5.3% → 73 sec • Sans filtre S-W: 7% → 22 sec

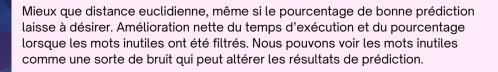


pourcentage de bonne prédiction et le temps d'exécution.

## KNN DISTANCE COSINUS

#### Avec BOW binaire

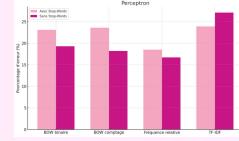
- Avec filtre S-W : 22.6% → 5.14 sec • Sans filtre S-W: 5.7% → 41.9 sec
- Avec BOW comptage
- Avec filtre S-W : 24.7% → 5.63 sec • Sans filtre S-W: 8.8% → 44.3 sec Avec Fréquence relative
- Avec filtre S-W : 24.6% → 5.34 sec • Sans filtre S-W: 8.4% → 23.31 sec
- Avec TF-IDF: • Avec filtre S-W : 33.5% → 5.11 sec
- Sans filtre S-W: 7% → 24.3 sec



#### **PERCEPTRON**

Avec BOW binaire

- Avec filtre S-W : 23.1% → 0.05 sec • Sans filtre S-W : 19.3% → 0.05 sec
- Avec BOW comptage
- Avec filtre S-W : 23.6% → 0.05 sec • Sans filtre S-W : 18.2% → 0.05 sec Avec Fréquence relative
- Avec filtre S-W : 18.5% → 0.02 sec • Sans filtre S-W : 16.7% → 0.02 sec
- Avec filtre S-W: 23.9% → 0.02 sec
- Sans filtre S-W : 27.1% → 0.02 sec



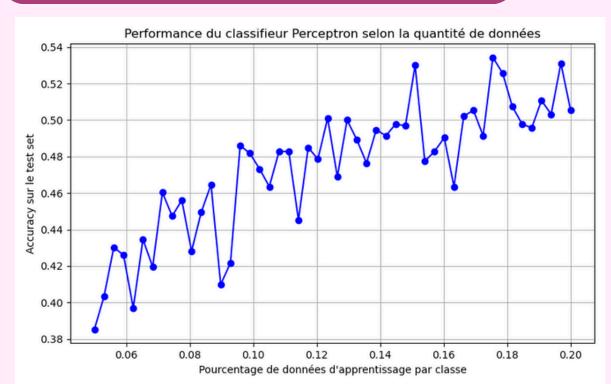
Avec ou sans stop words, le pourcentage de précision de prédiction stagne entre 16 et 24 %, le temps d'exécution est très efficace.

## ARBRE DE DÉCISION



Les expérimentations montrent que les arbres de décision fonctionnent mieux avec des représentations simples comme le comptage ou le binaire, surtout lorsqu'un filtrage des stopwords est appliqué. Cela suggère que la qualité de la représentation des données textuelles joue un rôle crucial dans les performances des classifieurs traditionnels. Les méthodes plus complexes comme TF-IDF n'apportent pas de bénéfices ici.

# PERFORMANCE DU CLASSIFIEUR PERCEPTRON



Le perceptron bénéficie clairement d'une augmentation de la quantité de données d'apprentissage. Même si ses performances restent limitées dans l'absolu (autour de 50 % au mieux), il montre une progression régulière. Cela suggère que ce type de classifieur linéaire a besoin de suffisamment d'exemples pour stabiliser son apprentissage, mais il reste peu adapté à des problèmes complexes et non linéaires comme ici (texte multiclasse).

## APPROCHE NON SUPERVISÉE

#### REMARQUE

Nous n'avons pas eu le temps de faire cette partie. En revanche tout les codes concernant cette partie sont disponibles dans notre dossier "iads".

## CONCLUSION

Nos résultats ne sont pas assez représentatif, car notre jeux de données est limité par la capacité de nos ordinateurs (PPTI compris).

Néanmoins, pour ce qui est de l'apprentissage supervisé, dans la plupart des cas, la méthode TF-IDF est la plus approprié pour le problème de classification de mots. Nous remarquons que sans l'application des tris sur les stop-words, les performances sont biaisés par tous les mots inutiles, ce qui baissent les performances de bonne prédiction.

Concernant les classifieurs, le Perceptron de Rosenblatt est le classifieur ayant obtenu les meilleurs résultats. Pour les arbres de décision, la représentation en fréquence et TF-IDF se découpent en branches plus homogènes, comparé aux autres qui s'enchainent.