

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.

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.

L'objectif est d'appliquer et de 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).

ECHANTILLON

| | messages | target |
|-------|---|--------|
| 0 | !vni! am sure some bashers of Pens fans are pr... | 10 |
| 1 | My brother is in the market for a high-perform... | 3 |
| 2 | !vni!vni!finally you said what you dream abou... | 17 |
| 3 | !vni!vni!vni!s the SCSI card doing the DMA L... | 3 |
| 4 | ! I have an old jasmine drive which I can... | 4 |
| 18841 | DN> From: Niyed@nrc.usc.edu (David Niyed)... | 13 |
| 18842 | !vni!t in isolated ground recepticles (usualy... | 12 |
| 18843 | !vni! installed a DX2-66 CPU in a clone moth... | 3 |
| 18844 | !vni!wouldn't this require a hyper-sphere. In... | 1 |
| 18845 | After a tip from Gary Crum (crum@com.cc.utah... | 7 |

API

- Python
- NumPy
- Pandas
- Graphviz
- Matplotlib

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 ...).

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 AVEC DISTANCE EUCLIDIENNE

In this section, state what is the purpose of your study.

KNN AVEC DISTANCE COSINUS

In this section, state what is the purpose of your study.

PERCEPTRON

In this section, state what is the purpose of your study.

ARBRE DE DÉCISION

In this section, state what is the purpose of your study.

APPROCHE NON SUPERVISÉE

CONCLUSION

To wrap up your poster, present two to three key findings. You can also add a brief explanation or narrative to these that can encourage conversation or dialogue with the audience. These findings can be actionable items that can lead to implementation, policy creation or further study.

Related literature

References can take up a lot of space, so cite only the key references used in the study.

CHARTS

