

NLPvs

Outil de classification de documents au format PDF

NLP^{VS}

par Vincent DUBOIS

Dossier contenant les fichiers à classer

Sélectionner un dossier

TEST

Sélection du modèle de classification

Sélectionner un dossier

MODEL/DEFAULT

Entraîner un nouveau modèle

☒ oui

☐ non

☐ Accélération CUDA (GPU Nvidia requis)

☒ CPU

Nombre de clusters

Dossier contenant les fichiers d'entraînement

Sélectionner un dossier

PDF

Documents clusterisés

Rapport_et_conclusions_Modif_2_PLU_St_Denis.pdf affecté au cluster N° 1
LM.pdf affecté au cluster N° 0
DUBOIS_Vincent_Mémoire_18.12.pdf affecté au cluster N° 2
Vincent_Dubois_Rapport_Licence.pdf affecté au cluster N° 0
Uni de tous les savoirs_4.10.2000_materiaux_intelligents.pdf affecté au cluster N° 2
Contest Task - Version française.pdf affecté au cluster N° 1

fichier GEXF enregistré sous : GEXF/2020-05-11 20:24:10.gexf avec succès

Graph exporté sous :GRAPH/2020-05-11 20:24:10.png

Montrer le graph

Clusters

Cluster N° 0:[' bim', ' vincent', ' dubois', ' monsieur']
Cluster N° 1:[' saintdenis', ' pleyel', ' local', ' modification']
Cluster N° 2:[' beton', ' materiaux', ' intelligences', ' capables']

Temps d'exécution

effectué en 15.45 secondes.

Lancer la classification

Quitter

ABSTRACT

NLPvs est un outil de classification non-supervisé de documents au format PDF.

Pour un corpus de documents à classer défini par l'utilisateur, l'outil se charge de les répartir au sein de clusters selon le modèle de Machine Learning pré-entraîné choisi, et de restituer cette classification sous forme graphique montrant la proximité des documents entre eux ainsi que leurs répartition au sein des clusters.

Il est également possible d'entraîner son propre modèle, en spécifiant un corpus de documents d'entraînement ainsi qu'un nombre de clusters (avec une possibilité de l'accélérer grâce aux coeurs CUDA sur un GPU Nvidia). Après l'entraînement, ce modèle est automatiquement sauvegardé (pour permettre une réutilisation pour un autre jeu de documents) et procède à la classification du corpus de documents à classer.

Pour chaque corpus (entraînement ou classification) de PDF, un traitement initial est requis:

- **"Parsing" des fichiers PDF** : il s'agit ici d'extraire le contenu textuel de chaque fichier (sous forme de chaîne de caractères) et d'en supprimer toute mise en page.
- **Uniformisation et suppression de la ponctuation de chaque texte.**
- **"Tokenisation"** : découpage d'un texte en termes (mots) afin de créer une entrée de vocabulaire (un "jeton")
- **Suppression des articles et "stop words"** (verbes être et avoir, etc...) : élimination des termes les plus communs entre chaque texte afin de les alléger pour les étapes suivantes.
- **"Racinisation"** : réduction de chaque terme à sa "racine" en enlevant conjugaison, suffixe ou préfixe.
- **Mesure du "Text Frequency - Inverted Frequency Index" (TF-IDF)** : Estimation de la fréquence d'un terme (ou groupe de termes) à travers un corpus, sur la base de sa rareté au sein du corpus (moins un terme est fréquent, plus il a de chances d'être porteur de sens pour pouvoir différencier les textes entre eux).
- **Vectorisation de chaque texte** : Associe un vecteur à chaque texte sur la base des indices de fréquence.

Un tel corpus peut ensuite servir de base pour entraîner un modèle (basé sur l'algorithme de clusterisation Kmeans) afin d'en tirer des clusters propres à ce modèle puisque lié aux vocabulaire de termes du corpus d'entraînement (les features de ces clusters seront les termes les plus significatifs qui ont servi au modèle pour différencier les textes entre eux).

Le modèle est ensuite chargé de "prédire" auquel de ses clusters appartient chaque texte du corpus à classer. Enfin, la distance cosinus entre chaque texte vectorisé sera calculée, dans le but de générer un graphique légendé qui traduira la clusterisation et cette distance en un seul document au format PNG. Un document au format GEXF sera également produit de manière systématique.

