

Sémantique Clustering of articles

Table des matières

1. Introduction.....	2
2. Choix des technologies.....	2
3. Méthodes choisis pour la classification.....	2
4. Problèmes rencontres	3
5. Pistes d'amélioration	3

1. Introduction

L'objectif principal de ce projet, c'était d'utiliser différents outils de classification pour pouvoir réaliser une étude simple sur le jeu des données des articles. Dans l'approche que j'ai choisi ce la classification sémantique des différents articles, avec l'objectif de regarder le comportement des différentes pages web et voir si on peut trouver des tendances dans la rédaction des différents types des articles.

L'idée initial était de trouver les différents groupes sémantiques qu'il existe dans le jeu de données des articles et selon leur date de publication voir quels sont topics les plus parles selon les différentes périodes des dates.

2. Choix des technologies

Pour réaliser cette étude, j'ai décidé de travailler sur des fichiers .notebook, du à leur flexibilité des au moment des lancer des taches des traitement lourdes. Comme la quantité des articles et les méthodes que j'ai utilisé ont besoin des grandes capacités des traitement, les fichier notebook permet de modifier information déjà traite plusieurs fois sans lancer tout le code. Cela permet d'essayer différentes méthodes et corriger les fautes plus facilement, surtout pour quelqu'un comme moi qui en est encore à ses premiers pas dans les applications ml.

Toutes les méthodes de clustering, réduction dimensionnel et Machine Learning utilise lors de ces études sont de la librairie python « Scikit learn ». Le choix de cette librairie est dû aux larges quantités des exemples, documentation, information utile et questions et réponses que ont peux trouver sur ces fonctions.

Pour réaliser toutes les réductions dimensionnelles, j'ai utilisé le méthode « T-distributed Stochastic Neighbor Embedding », la fonction de Sci-kit learn. J'ai décide d'utiliser T-SNE parce qu'il s'agit d'une technique non linéaire de réduction de la dimensionnalité. Dans notre cas les articles ne suivent aucune loi mathématique. Donc une approche qui est moins sensible à l'ordre des points de données se rend très utile dans notre cas.

Pour la génération des graphs j'ai opté pour la librairie « Plotly » dans python et la plateforme « Tableau ». « Plotly » parce qu'elle est bien intégrée avec l'utilisation des dataframes pandas. Tableau parce qu'on a appris à l'utiliser cette année dans le module et j'ai trouve convenaient de l'utiliser pour générer les graphs avec les résultats finales.

3. Méthodes choisis pour la classification

Pour le choix des méthodes de classification sémantique, j'ai applique deux types : K-means et Latent Dirichlet Allocation(LDA). K-means est une des méthodes les plus courants pour faires des classifications mais pas le plus approprié pour des classifications sémantiques du texte. K-means reste plus simple et plus rapide que LDA, mais selon plusieurs études LDA est efficace pour trouver les sujets caches dans la donnée.

J'ai utilisé deux méthodes de vectorisation de la librairie « SciKit Learn » : CountVectorization et TfidfVectorizer. CountVectorization est une méthode plus simple qui génère une matrice de comptage des occurrences de mots. Cependant TfidfVectorizer attribue des scores aux mots en fonction de leur importance relative dans le corpus. TF-IDF (Term Frequency-Inverse Document Frequency) prend en

compte à la fois la fréquence d'un terme dans un document particulier (Term Frequency) et son importance globale à travers tous les documents (Inverse Document Frequency).

4. Problèmes rencontrés

Le principal problème rencontré est lié à la capacité de calcul de mon pc et le manque d'accès à un centre de traitement puissant. Avec une cpu i5 (7ème gen) certains notebooks avec les 36000 articles ont pris plus de 4+ heures à tourner pour avoir des résultats. La solution qu'on a accès à Polytech est une possible solution mais l'incapacité d'ajouter des librairies limite l'utilisation. Pour cela dans certains cas les résultats affichés sur les notebooks sont avec 935 articles. Les études, la paramétrisation des vectoriseurs et classificateurs est différente selon la quantité des articles.

Autre problème que j'ai rencontré est le manque de préparation dans le sujet. L'application des modèles comme Kmeans ou LDA est possible jusqu'à certain degré. J'ai eu beaucoup de problèmes à la paramétrisation des modèles surtout en travaillant avec la totalité des articles, parce que l'utilisation de toutes les « features » génère des résultats anormaux, donc il fallait ajouter les conditions sur la vectorisation des mots. Avec les 36000 articles la quantité des « features » dépasse les millions, donc il fallait enlever les mots qui apparaissent le moins. Cela c'est une possible solution mais pas la meilleure, dû à ce qui génère la perte d'information.

5. Pistes d'amélioration

Pendant le développement de ce projet, j'avais l'idée d'utiliser un modèle de « Zero-shot classification ». J'ai récupéré le modèle de « Hugging Face » : [MoritzLaurer/DeBERTa-v3-base-mnli-fever-anli · Hugging Face](#). Ce modèle permet de donner en input un texte et un groupe de labels. En return, il t'envoie les pourcentages correspondants aux classifications du texte selon les labels. En appliquant une classification de toutes les articles publiées dans une page web, on aura trouvé quels sont les tendances de chaque page web. Le seul problème c'est qu'effectuer la classification d'un seul texte prend au tour de 5-10 secondes. Réaliser la classification avec 1000 articles aura pris un temps non négligeable. Trouver d'autres moyens de réaliser cette étude c'est comme réaliser une étude de sentiment plus complète.