

1 Création d'un index inversé et moteur de recherche booléen et vectoriel

Dans cette partie nous détaillerons notre approche sur les collections cacm et CS276 afin de répondre aux questions posées.

1.1 Collection cacm

1. Nous avons choisi d'utiliser la librairie NLTK pour traiter la tokenisation de la collection. Toutes les valeurs fournies sont après avoir retiré les `common_words` du corpus. Sur la collection cacm, nous avons donc obtenu un nombre de tokens de 92305.
2. La taille du vocabulaire obtenu avec ces tokens est de 10478 mots.
3. Le nombre de tokens sur la moitié de la collection est de 29916. La taille du vocabulaire associé est de 5780 mots. Nous pouvons donc calculer la valeur de k avec la formule suivante:

$$b = \log \frac{\text{vocabulaire_total}}{\text{demi_vocabulaire}} \div \log \frac{\text{tokens_total}}{\text{demi_tokens}} \quad (1)$$

Ce qui nous donne une valeur pour b de 0.53. Nous pouvons ensuite calculer la valeur de k avec la formule suivante:

$$k = \frac{\text{vocabulaire_total}}{\text{tokens_total}^b} \quad (2)$$

Ce qui nous donne une valeur de k de 25.05.

4. Maintenant que nous avons des valeurs pour k et b , nous pouvons fournir une estimation de la taille du vocabulaire pour un million de tokens avec la loi de Heaps:

$$M = k * T^b \quad (3)$$

Ce qui nous donne une estimation de l'ordre de 36865 mots de vocabulaire pour une collection ayant un million de tokens.

5. Voici les graphes de fréquence et rang pour la collection cacm, pour la fréquence nous avons limité au 200 premiers tokens afin de pouvoir mieux visualiser l'évolution:

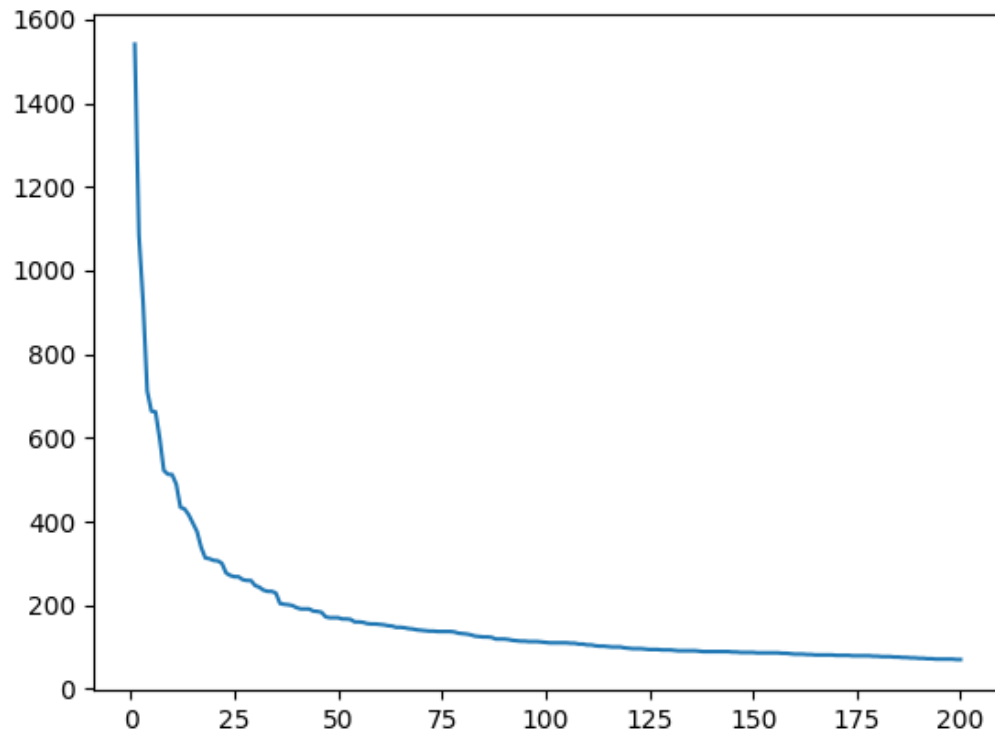


Figure 1: freq vs rang

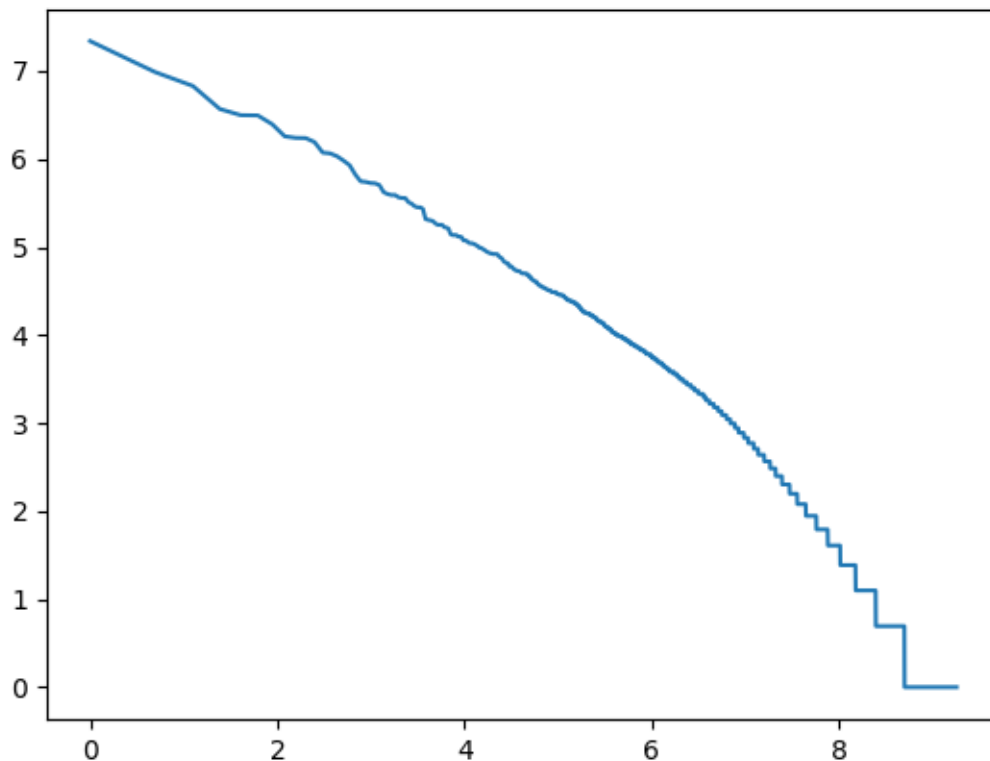


Figure 2: $\log(\text{freq})$ vs $\log(\text{rank})$

1.2 Collection CS276

1. Pour cette collection, aucune liste de common_words n'est fournie avec. Nous avons donc utilisé celle de nltk pour l'anglais. Ce traitement nous donne un nombre de tokens de 18 918 184.
2. La taille du vocabulaire associé est de 335 167 mots.
3. En utilisant l'équation que nous avons utilisé pour la collection cacm, nous trouvons des valeurs pour k et b de 0.71 pour b et 2.29 pour k.
4. Le nombre de tokens sur la moitié de la collection est de 8 905 232. La taille du vocabulaire associé est de 196 338 mots. La taille du vocabulaire estimé pour une collection avec un million de tokens en suivant la loi de Heaps est de 41590 mots.

5. Voici les graphes de fréquence et rang pour la collection CS276, comme la collection contient beaucoup de tokens, nous avons représenté ces graphes uniquement pour les 200 premiers tokens:

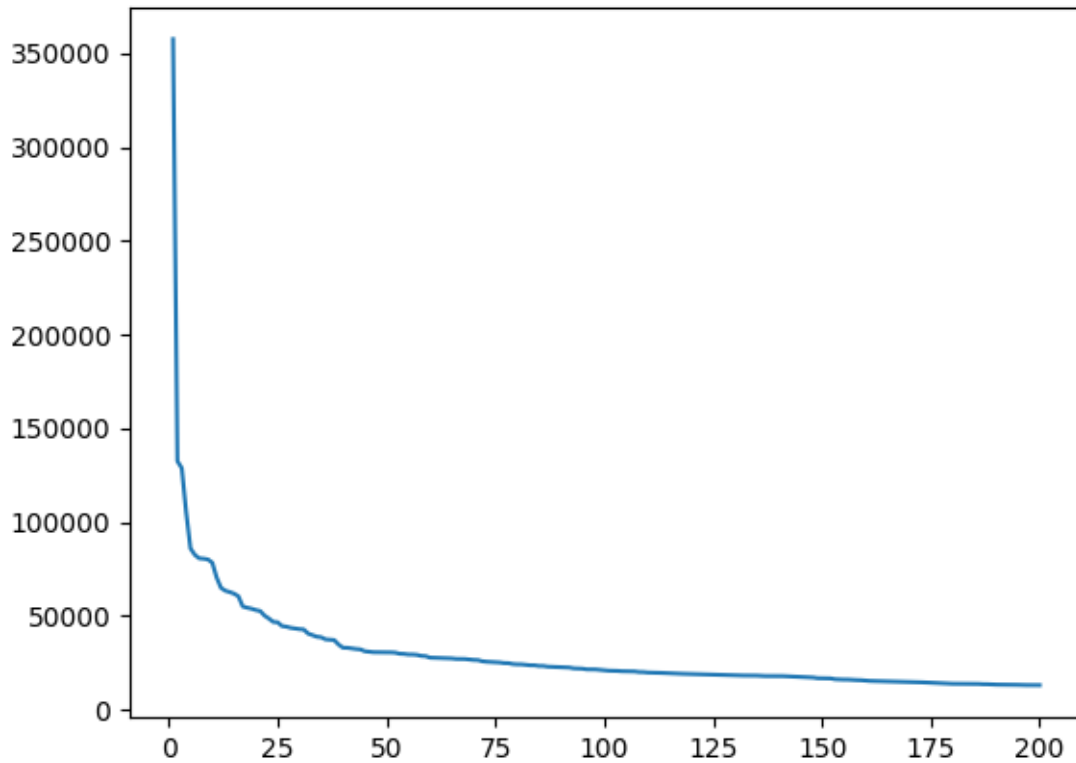


Figure 3: freq vs rank

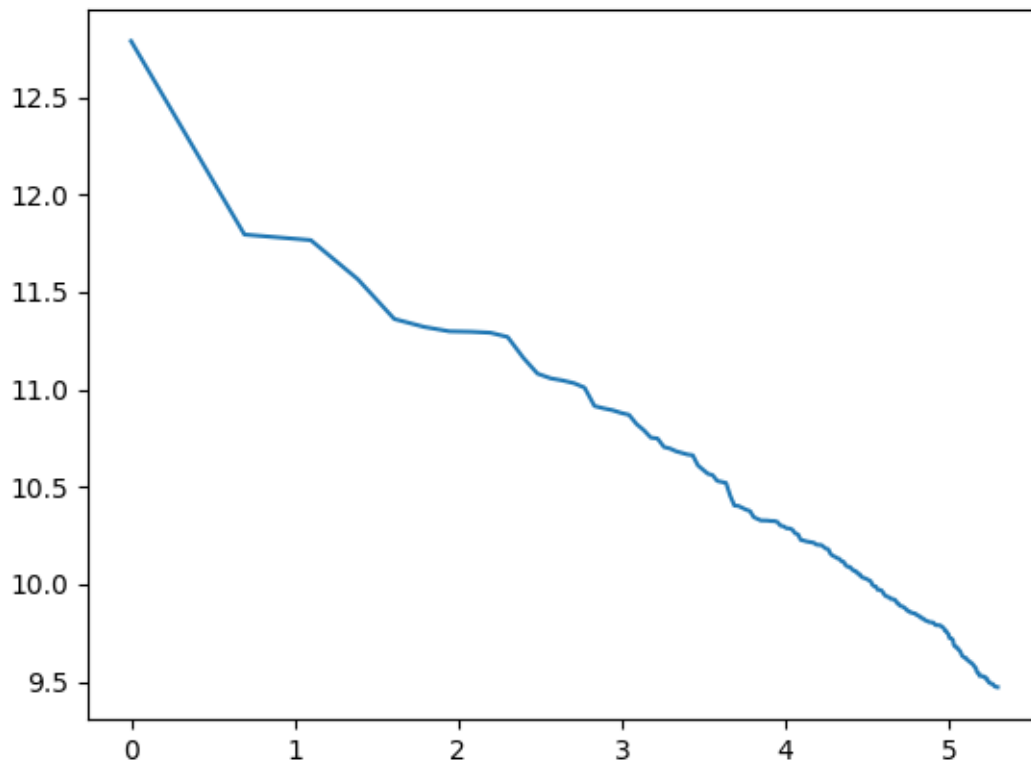


Figure 4: $\log(\text{freq})$ vs $\log(\text{rank})$

Indexation et évaluation pour la collection CACM

1.3 Indexes et recherches simples pour les deux collections

Pour la collection CACM, comme la collection est petite et que ce n'est pas très coûteux de générer l'index, nous n'avons pas implémenté de fonction capable de la sauvegarder dans un fichier. Pour la collection CS276, celle-ci étant beaucoup plus volumineuse, l'index ainsi que les documents parsés sont sauvegardés dans des fichiers.

Afin de pouvoir observer les performances de chaque index, nous avons créé plusieurs fichiers pour faciliter ceci. Leur utilisation est détaillée dans le fichier readme du projet github.

1.4 Mesures de performance pour la collection CACAM

1. Les mesures de performance sont accessibles en exécutant les recherches de test grace aux fichiers préparés.
2. Les resultats détaillés pour l'évaluation des performance sont accessibles et facilement modifiables (modification du nombre de documents retournés en sortie du fichier `pre-pare_cacm.py`. La F-mesure de notre systeme est de: 0,12. La E-mesure de notre systeme est de: 0,82. La prcision moyenne est de 0,3.