

Lyricognizer

Outil d'attribution de paroles de musique à leur auteur

Aymane Hamdaoui & Charles Kayssieh

9 février 2024

Télécom Paris

Table des matières

1. Introduction

2. Fonctionnement

3. Résultats

4. Conclusion

Introduction

Problématique

Est-il possible de créer un outil qui attribue à des paroles de chansons leur auteur uniquement à l'aide d'une base de données de lyrics de chanteurs choisis préalablement et de la *Normalized compression distance* (NCD)?

Fonctionnement

Fonctionnement général du système :

1. Collecte des données
2. Traitement du texte
3. Tri de la BDD
4. Calcul de la NCD
5. Renvoi des résultats

Création d'un script *lyrics_grabber.py*

Utilisation de l'API Genius

Récupération des paroles des 100 musiques les plus populaires de chaque artiste

Artistes

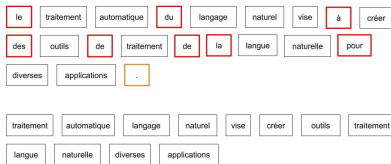
- Charles Aznavour
- B.B. Jacques
- Damso
- Drake
- Freeze Corleone
- Gazo
- Lomepal
- Mylène Farmer
- Nekfeu
- Soolking

Trois niveaux de traitement avec le script *cleaner.py* :

1. Suppression de la ponctuation, des sauts de ligne, des annotations et passage en minuscules
 2. Suppression des mots vides
 3. Lemmatisation
- } Détection de la langue

Mots vides et lemmatisation

Retrait des stop words



Fonctionnement des stop words

Lemmatisation (+ tokenisation et stopwords)

Le traitement automatique du langage naturel vise à créer des outils de traitement de la langue naturelle pour diverses applications.



Fonctionnement de la lemmatisation

Organisation de la BDD :

1. Par artiste, chaque musique est aléatoirement numérotée entre 0 et 99
 - 0-19 = utilisation pour les données test
 - 20-99 = utilisation pour les données d'entraînement
2. Traitement des 1000 musiques des 3 manières différentes

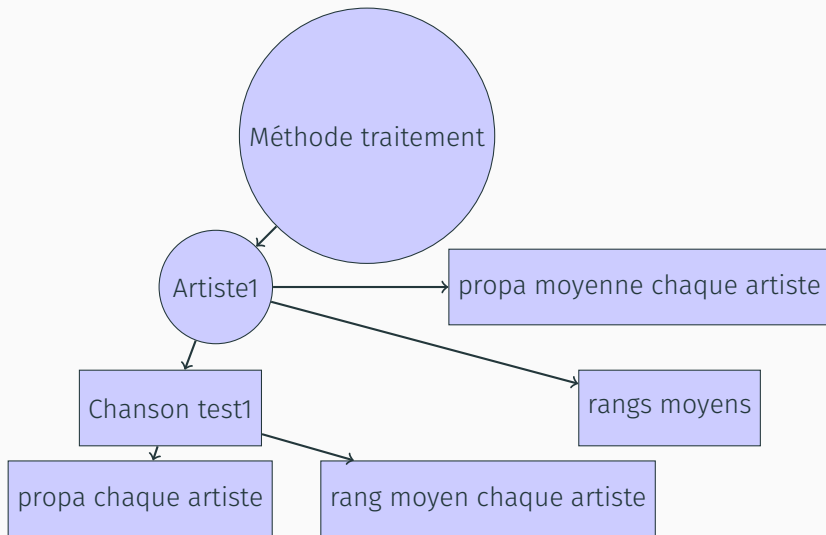
Formule théorique :

$$\text{NCD}(x, y) = \frac{C(x + y) - \min(C(x), C(y))}{\max(C(x), C(y))}$$

Utilisation de 3 méthodes de compression

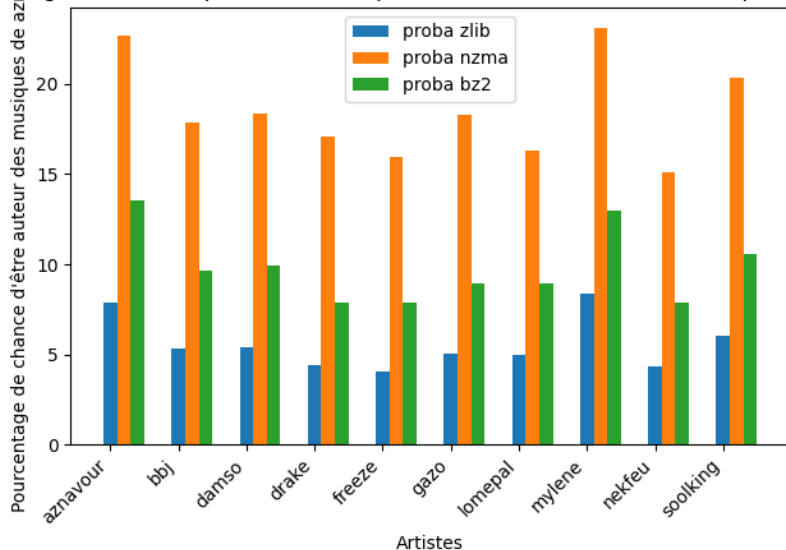
- zlib : la plus populaire (utilisée par Linux, macOS, iOS, ...) - compression par dictionnaire & encodage de Huffman
- lzma : utilisée par le programme 7zip - compression par dictionnaire
- bz2 : la plus efficace - réorganisation de donnée & encodage de Huffman

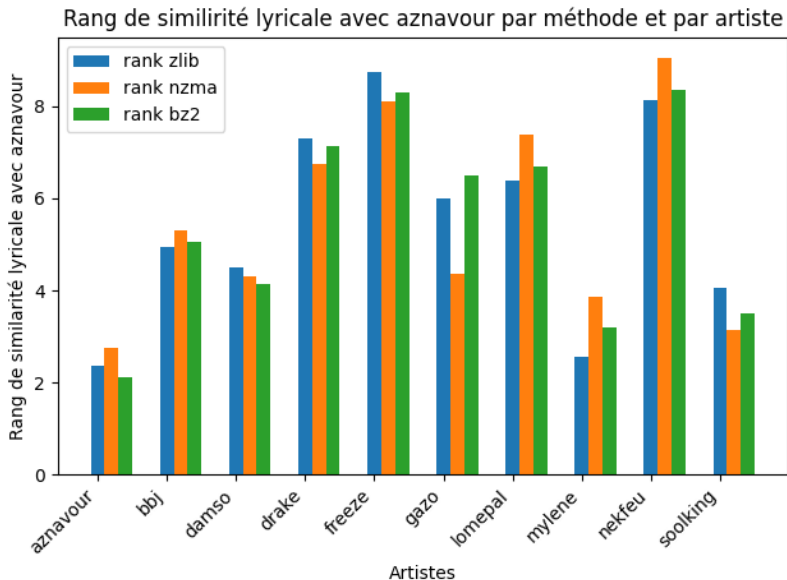
Structure des résultats

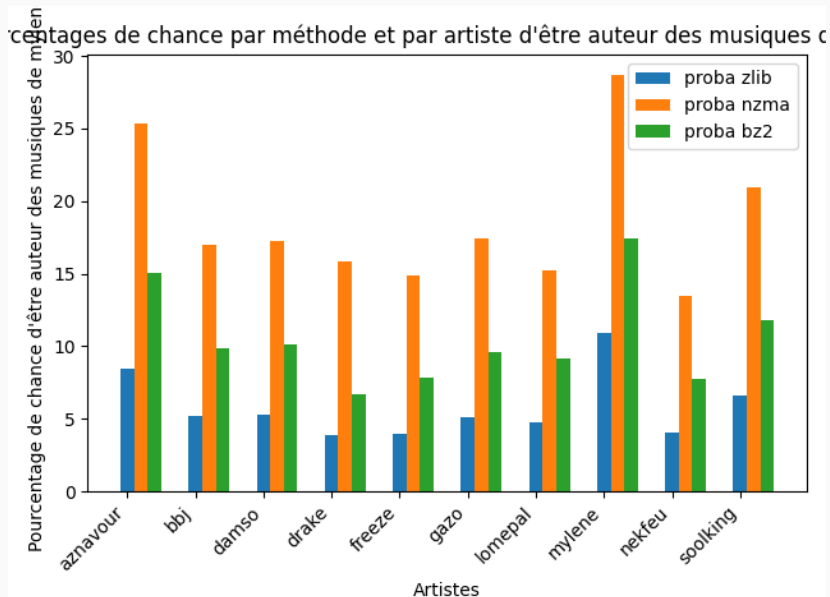


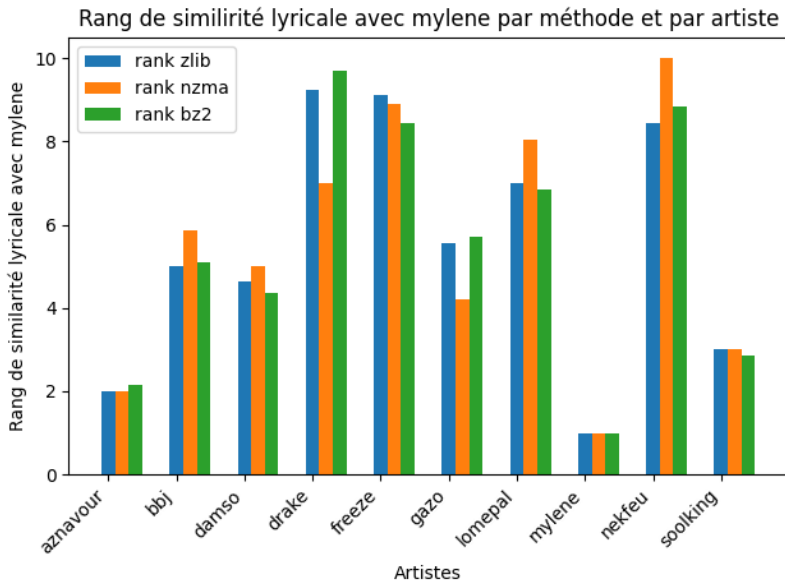
Résultats

pourcentages de chance par méthode et par artiste d'être auteur des musiques de aznavour

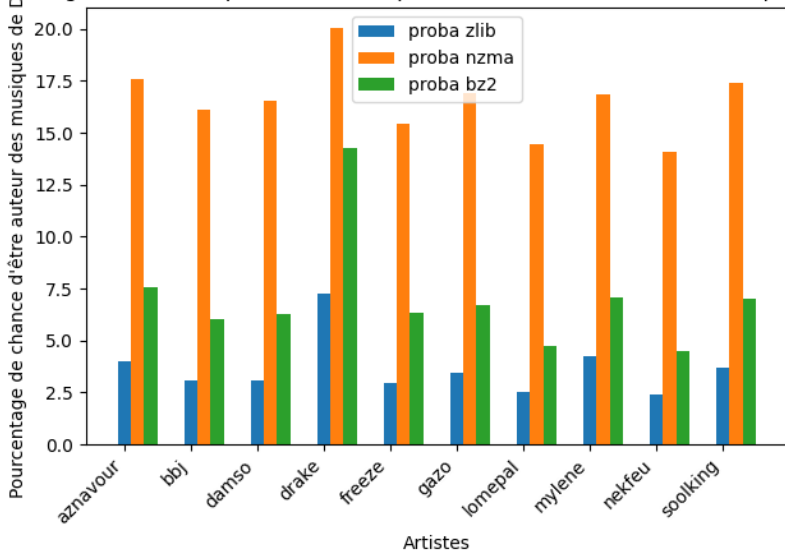


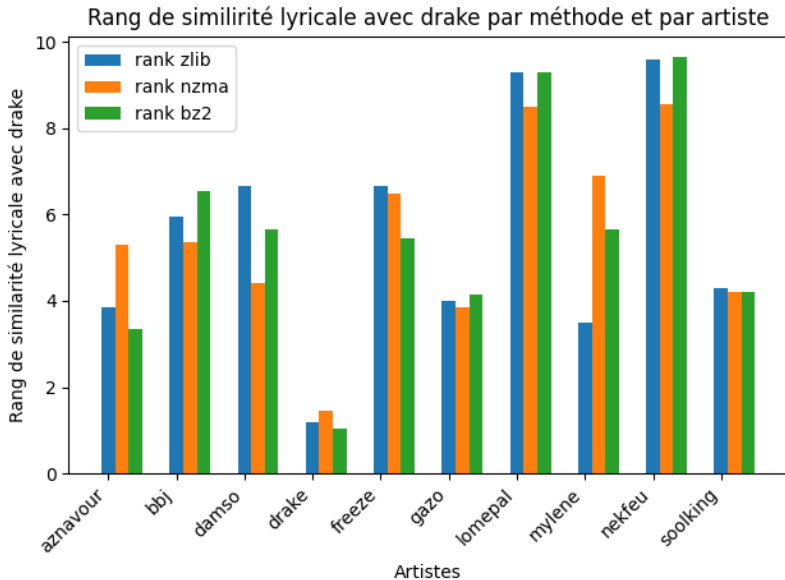




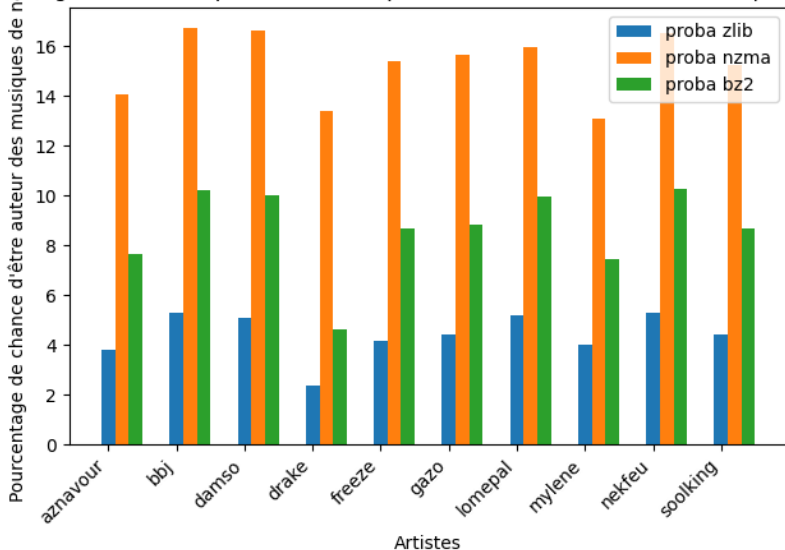


Percentages de chance par méthode et par artiste d'être auteur des musiques de Drake

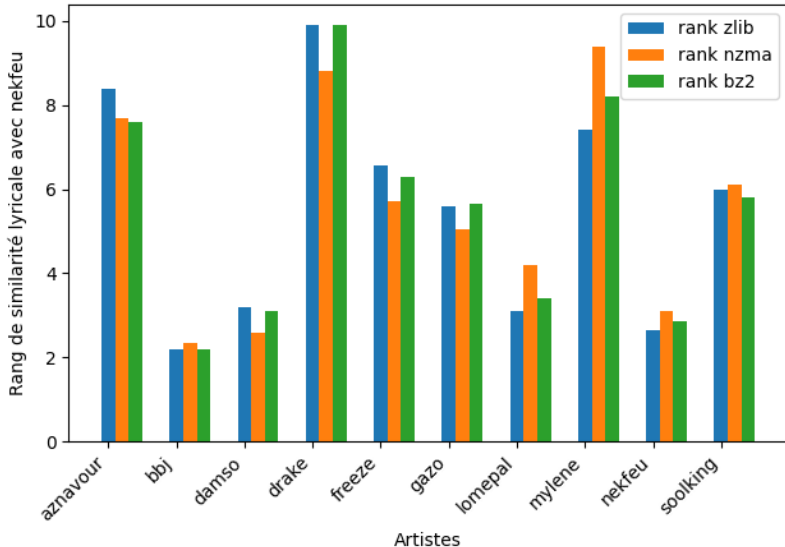




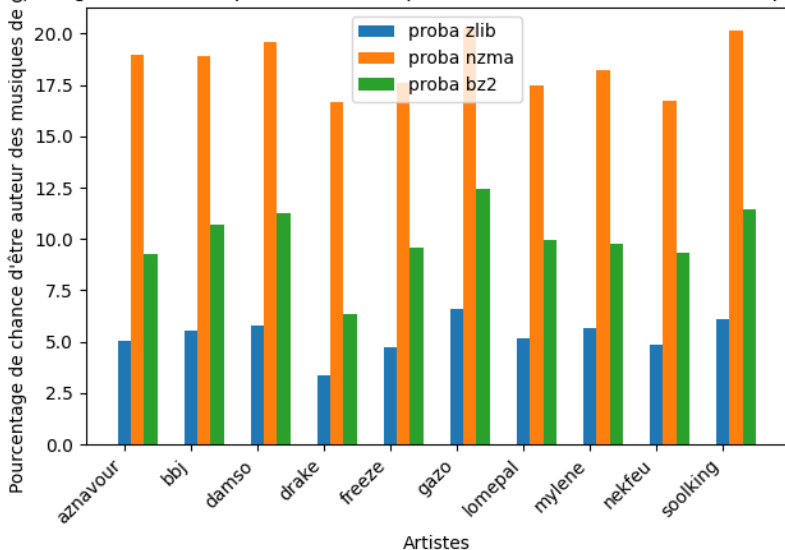
pourcentages de chance par méthode et par artiste d'être auteur des musiques de Nekfeu



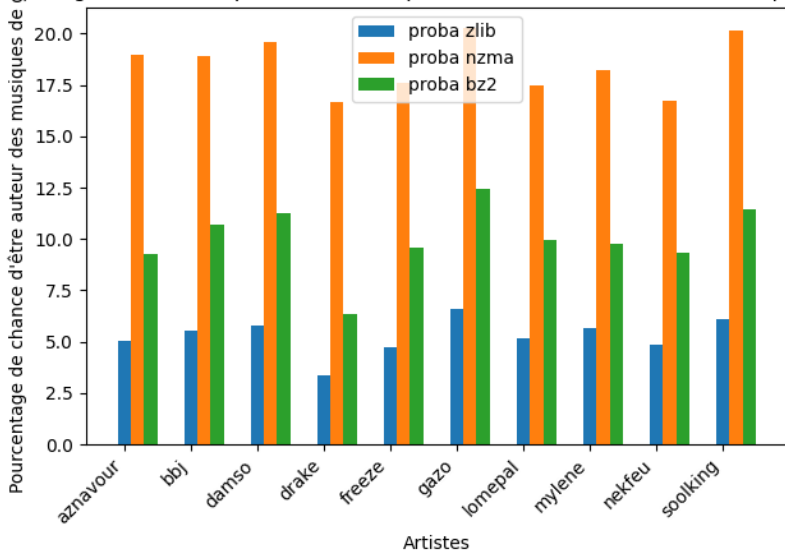
Rang de similarité lyrique avec nekfeu par méthode et par artiste



pourcentages de chance par méthode et par artiste d'être auteur des musiques de gazo



pourcentages de chance par méthode et par artiste d'être auteur des musiques de gazo



Conclusion

Conclusion et ouverture

- Utilisation de la NGD

$$NGD(A, B) = \frac{\max(\log(f(A)), \log(f(B))) - \log(f(A \cap B))}{\log(N) - \min(\log(f(A)), \log(f(B)))}$$

- Prévision des résultats possibles ?

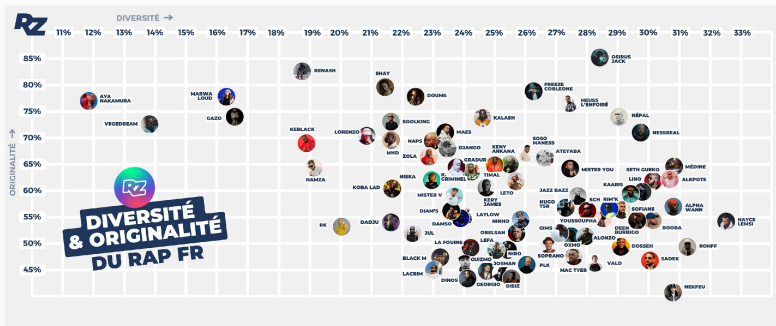


Figure 1 : Graphique diversité & originalité des auteurs de rap français

Questions?