

IFT6285 (Devoir) - Université de Montréal

COLELLA Charles

September 24, 2021

1 Introduction

Dans ce document nous allons entrainer un model de langue, de type Kneser&Ney à l'aide du 1B-word corpus¹ mis à disposition. Ce corpus d'entraînement est composé de 99 tranches différentes. Pour entrainer ce model, nous utiliserons la librairie **KenLM**³ qui permet de réduire le temps et le les coûts de mémoire.

1.1 Lissage Kneser-Ney

Le lissage de Kneser-Ney est une méthode utilisée pour calculer la distribution de probabilité de n-grammes dans un document en fonction de leur historique.

Etant donné $c : \omega_1 \in \Omega \times \omega_2 \in \Omega \rightarrow N$ une fonction définie sur l'ensemble des mots, qui exprime nombre d'occurrences du mot ω_1 suivi du mot ω_2 dans le corpus et $\lambda_{\omega_{i-1}}$ un facteur de normalisation de la probabilité p_{KN} évaluée avec $\omega_i \in \Omega$.

La probabilité d'un bigramme, s'écrit comme suit :

$$p_{KN}(\omega_i|\omega_{i-1}) = \frac{\max(c(\omega_{i-1}, \omega_i) - \delta, 0)}{\sum_{\omega'} c(\omega_{i-1}, \omega')} + \lambda_{\omega_{i-1}} p_{KN}(\omega_i) \quad (1)$$

On remarque le parametre δ , qui représente une valeur d'actualisation des fréquences des bigrammes observés en entraînement afin de correspondre aux fréquences observés lors des tests (1991 - Church & Gale).

On calcul donc d'abord la fréquence du bigramme dans le corpus (avec une constante d'actualisation) puis on ajoute la probabilité p_{KN} d'observer ω_i dans un contexte inconnu, qui est estimé comme la fréquence d'apparition après tout autre mot :

$$p_{KN}(\omega_i) = \frac{|\{\omega' : 0 < c(\omega', \omega_i)\}|}{|\{(\omega', \omega'') : 0 < c(\omega', \omega'')\}|} \quad (2)$$

L'avantage de cette méthode se trouve donc en la considération des mots possibles qui précèdent un unigramme lors de la définition de sa probabilité.

1.2 Mesure de la perplexité

Soit $(\omega_n)_{n \in N}$ une combinaison de mots appartenant à l'ensemble P des phrases. On définit la perplexité PP par :

$$\forall k \in N, \quad PP(\omega_1 \omega_2 \dots \omega_k) = p(\omega_1 \omega_2 \dots \omega_k)^{-\frac{1}{k}} \quad (3)$$

Durant la suite de ce document, nous allons utiliser la fonction **perplexity**² de **KenLM**³. Etant donné un score $S = \log_{10}(p)$ avec p la probabilité associée à une suite de mots de longueur $k \in N$, l'utilisation de cette fonction revient à calculer :

$$PP(\omega_1 \omega_2 \dots \omega_k) = (10^{S(\omega_1 \omega_2 \dots \omega_k)})^{-\frac{1}{k}} = 10^{-\frac{S(\omega_1 \omega_2 \dots \omega_k)}{k}} \quad (4)$$

¹<https://github.com/ciprian-chelba/1-billion-word-language-modeling-benchmark>

²<https://github.com/kpu/kenlm/blob/master/python/kenlm.pyx#L182>

³<https://kheafield.com/code/kenlm/>

2 Question (a)

2.1 Mise en place

KenLM est entraîné avec un model bigramme sur les 9 premières tranches du repertoire d'entraînement, prétraités au préalable, et dont le model est enregistré dans un fichier *data.arpa* puis dans *data.binary* (le format binaire se charge beaucoup plus rapidement et offre plus de flexibilité). Ensuite le model est appliqué sur les 1000 premières phrases, prétraités au préalable, de la tranche de test et les resultats enregistrés dans le fichier *data.txt*.

Le code est affiché ci-dessous :

execution.sh

```
#!/bin/bash
python search.py $training_dir -n 9 -p|$kenLM_dir/lmplz -o 2 > ./data.arpa
$kenLM_dir/build_binary data.arpa data.binary
python computing.py $test_file ./data.binary -p -n 1000 -r "data.txt"
```

Il est a noter que **search.py** et **computing.py** sont deux programmes écrits en python pour l'entraînement et le test de notre model de langue.

2.2 Résultats

Les résultats obtenus sont représentés par la moyenne, le min et le max de perplexité obtenus en utilisant la fonction **kenlm.perplexity**² présentée en introduction.

	Sans prétraitement	Avec prétraitement
Moyenne	412.933	339.687
Maximum	21410.839	15752.395
Minimum	28.735	18.588

Le prétraitement consiste à mettre, autant pour l'entraînement que pour le test, le texte en minuscule, remplacer les URL par `__URL__`, les nombres par `__NUMBER__`, les emails par `__E_MAIL__` et les numéros de téléphone par `__PHONE__`. Cette tâche à pu être effectuée avec la librairie **cleantext**⁴.

3 Question (b)

3.1 Mise en place

Pour analyser les performances du model de langue, un entraînement sur un nombre de tranches $i \in [1, 99]$ du 1B-word corpus a été établi. Chaque model entraîné sur les i premières tranche est ensuite confronté au texte final pour être comparé aux autres.

script.compare.sh

```
#!/bin/bash
for i in $(eval echo {1..$number_training_files})
do
    train $i # on entraine notre model sur les i premieres tranches
    compute # on l'applique a la tranche de test
done
```

⁴<https://pypi.org/project/clean-text/>

3.2 Résultats

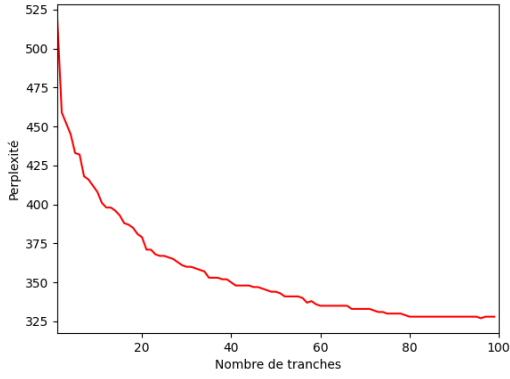


Figure 1: Moyenne perplexité

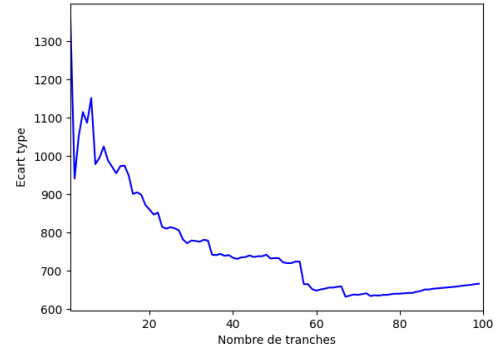


Figure 2: Ecart-type perplexité

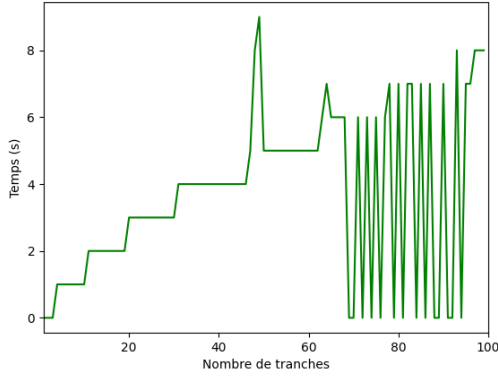


Figure 3: Temps de traitement de la tranche de test

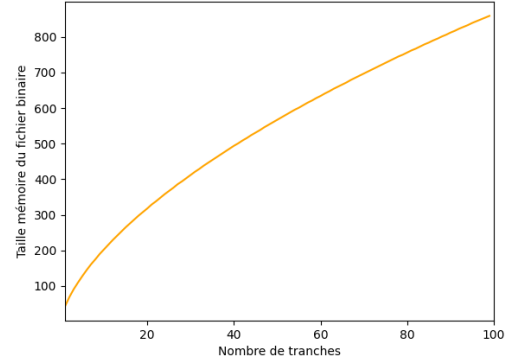


Figure 4: Taille du fichier binaire contenant le modèle

La perplexité d'une phrase prise indépendamment étant très sensible à la distribution de probabilités, la figure 2 met en avant quelques disfluences de l'écart type par rapport aux nombres de tranches d'entraînement (n'excédant pas 15%).

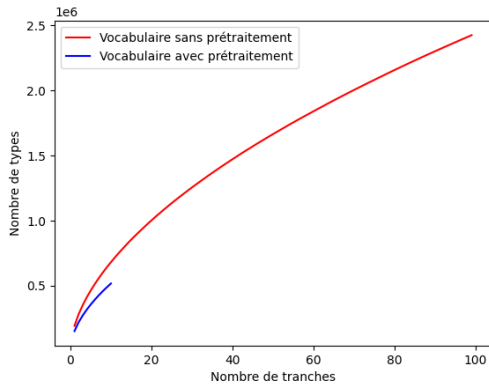


Figure 5: Nombre de types ($\times 10^6$) du corpus d'entraînement

On constate sur la figure 1 une diminution de plus en plus lente de la perplexité en fonction du nombre de tranches étudié lors de l'entraînement. Les figures 1 et 4 semblent fortement corrélées à la figure 5 représentant le nombre de types étudié.

4 Question (c)

4.1 Mise en place

Afin de mettre en place le model le plus adapté à notre tranche de test, trois entrainements différents ont été testés sur les 10 premieres tranches du repertoire d'entrainement :

- Entrainement sur un texte **sans prétraitement** en utilisant des **bigrammes**
- Entrainement sur un texte **avec prétraitement** en utilisant des **bigrammes**
- Entrainement sur un texte **sans prétraitement** en utilisant des **trigrammes**

L'entrainement sur un texte prétraité avec des trigrammes n'a pas pu être réalisé du fait du manque de performances de la machine utilisé.

4.2 Résultats

Les même données que précédemment ont été récoltés. Cependant, par soucis de performances⁵, seules les dix premieres tranches ont pu être comparées.

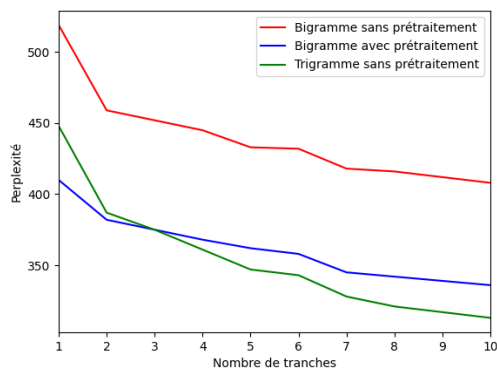


Figure 6: Moyenne perplexité

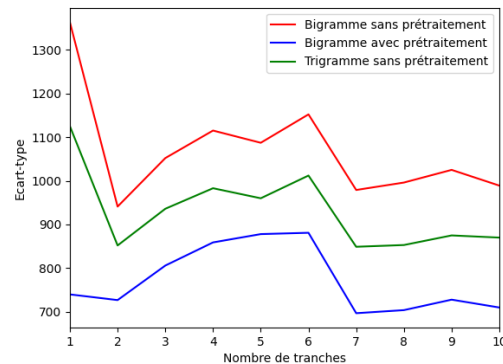


Figure 7: Ecart-type perplexité

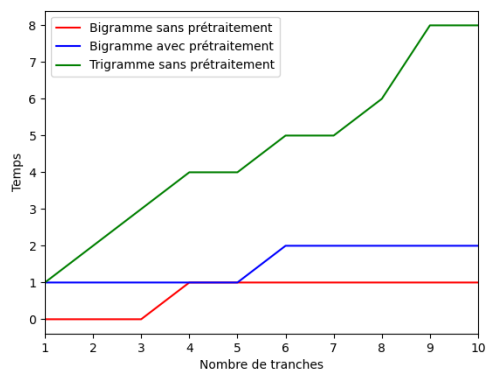


Figure 8: Temps de traitement de la tranche de test

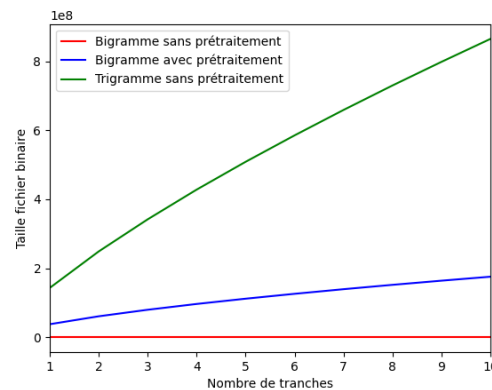


Figure 9: Taille du fichier binaire contenant le modèle

Sans grosse surprise, le minimum de perplexité est atteint pour le modèle utilisant des trigrammes sur la figure 6. On peut supposer que la perplexité diminuerait encore en utilisant un modèle d'ordre supérieur. La figure 7 met en évidence l'impact du prétraitement, qui engendre une normalisation des mots, sur l'écart-type de la perplexité. Il toutefois important de noter, relativement aux figures 8 et 9 qu'en augmentant l'ordre du modèle, la taille du fichier en mémoire ainsi que le temps de calcul augmente drastiquement ce qui rend le traitement beaucoup plus compliqué à effectuer.

⁵Les programmes ont été exécutés sur une machine du DIRO de l'Université de Montréal