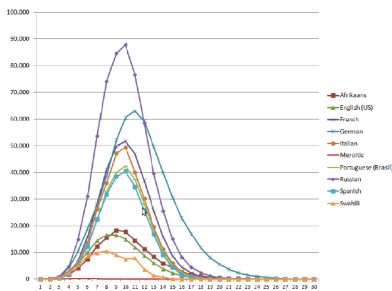


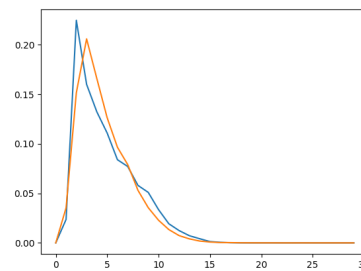
### Objectifs

- Comprendre la factorisation
- Améliorer un code et un rendu
- Procéder à une analyse de données

### Exercice 1 : Rappel de l'attendu



(a) La figure recherchée



(b) Là où nous nous sommes arrêtés

C'est la dernière partie de notre code qu'il faut améliorer :

```
import matplotlib.pyplot as pyplot
for chemin in ["13846-0.txt", "4300-0.txt"]:
    chaine = lire_fichier(chemin)
    liste_mots = decouper_en_mots(chaine)
    dic_longueurs = get_effectifs(liste_mots)
    liste_effectifs = vecteur_longueurs(dic_longueurs)
    pyplot.plot(liste_effectifs)
pyplot.savefig("frequences.png")
pyplot.show()
```

La première ligne qu'il faut améliorer c'est la sélection des fichiers, en effet si l'on doit traiter 1.000 fichiers on ne va pas écrire leurs noms à la main. Ici on va faire deux choses:

- regrouper nos textes dans un dossier "data"
- utiliser une librairie qui va nous permettre de lister les textes de ce dossier :

```
import glob #la bonne librairie
print(glob.glob("./*")) #tout ce qu'il y a dans le dossier courant
print(glob.glob("data/*")) #tout ce qu'il y a dans le dossier "data"
```

Nous allons modifier notre code de la façon suivante :

---

```
import matplotlib.pyplot as pyplot
import glob
for chemin in glob.glob("data/*"):
    print(chemin)#pour voir si tout se passe bien
    ....
```

---

Maintenant nous allons ajouter des étiquettes aux axes (*labels*), un titre (*title*), et une légende (*legend*). Ce qui va nous donner la Figure 2.

---

```
...
for chemin in glob.glob("../data/*"):
    print(chemin)
    nom_pour_legende = chemin
    ....
    pyplot.plot(liste_effectifs, label = nom_pour_legende)
pyplot.legend()
pyplot.title("Une magnifique Courbe")
pyplot.xlabel("Longueur des Mots")
pyplot.ylabel("Fréquence")
pyplot.savefig("frequences2.png")
pyplot.show()
```

---

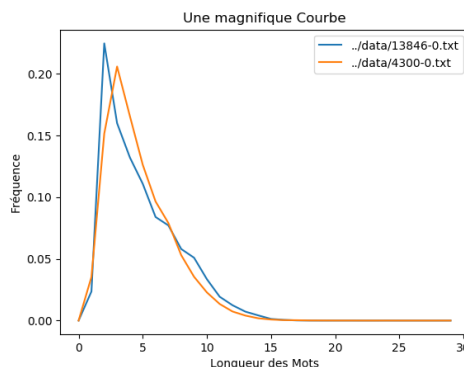


Figure 2: Avec la légende

Pour rendre la légende de la figure plus explicite, on va modifier le nom des fichiers (à la main) pour faire apparaître la langue sous la forme d'un code à deux lettres <sup>1</sup>. Nos fichiers se nomment donc désormais :

- en\_4300-0.txt
- fr\_13846-0.txt

Puis nous allons formater le nom du fichier pour que seul le nom de la langue apparaisse dans la légende au moyen d'une expression régulière (attention à la librairie):

---

```
nom_fichier = re.split("/",chemin)[-1]#Windows: re.split("\\\\",...)
nom_pour_legende = nom_fichier[:2]
```

---

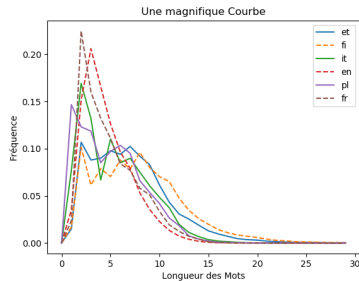
Pour avoir une idée du rendu, regardez la légende de la Figure 3a.

---

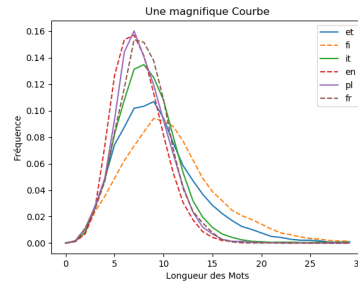
<sup>1</sup>Norme ISO : [https://fr.wikipedia.org/wiki/Liste\\_des\\_codes\\_ISO\\_639-1](https://fr.wikipedia.org/wiki/Liste_des_codes_ISO_639-1)

## Exercice 2 : Devoir à faire et à déposer en TD

Nous allons ajouter de nouveaux textes, récupérés sur [gutenberg.org](http://gutenberg.org) par exemple, pour avoir une figure avec plusieurs langues. Attention, il faut avoir des données suffisamment grands pour être significatifs, au moins 100.000 mots soit 500 Ko. Au besoin, concaténez (fusionnez) plusieurs textes pour une même langue.



(a) Nouvelle courbe avec plus de langues



(b) Changement de la sortie de `decouper_mots`

Efforcez vous d'avoir des données pour plusieurs grandes familles de langues <sup>(2)</sup> et plusieurs types de langues <sup>3)</sup>. Votre corpus devra contenir au moins 10 langues avec au moins 3 familles et 3 types différents. Vous produirez une autre figure avec les mêmes données mais où vous aurez simplement changé la sortie de la fonction `decouper_mots` de la façon suivante :

---

```
return set(liste_mots)
```

---

Ce qui devrait vous donner quelque chose qui ressemble à la Figure 3b.

## Devoir

Vous produirez un document d'une page (à l'aide du traitement de textes de votre choix) environ incluant :

- des statistiques sur votre corpus (nombre de caractères et de mots pour chaque langue)
- les figures que vous avez généré
- des commentaires sur ces figures et leurs différences
- quelques phrases de conclusion (qu'est-ce qui était attendu, qu'est-ce qui est inattendu...)

Vous déposerez sur Moodle une archive zip nommée NUMETU.zip (où NUMETU est votre numéro d'étudiant) et contenant :

- Votre code exporté au format Python `.py` (et pas ipynb)
- le PDF du document que vous avez produit

**Date limite : 12/10/2020 23h55**

[La suite page suivante...](#)

---

<sup>2</sup>[https://fr.wikipedia.org/wiki/Langues\\_par\\_famille](https://fr.wikipedia.org/wiki/Langues_par_famille)

<sup>3</sup>[https://fr.wikipedia.org/wiki/Typologie\\_morphologique\\_des\\_langues](https://fr.wikipedia.org/wiki/Typologie_morphologique_des_langues)

**Exercice bonus** : une fois que le travail est terminé et déposé, essayez de trouver un moyen d'améliorer la fonction `decouper_mots` en tenant compte de la langue. Il s'agit de chercher un tokenizer (découpeur en mots) adapté pour chaque langue. Considérez l'exemple suivant (il faut au préalable installer la librairie `MosesTokenizer`<sup>4</sup> voir par exemple la commande `pip install mosestokenizer` dans le terminal. Ce qui nécessite d'installer `pip` si vous ne l'avez pas déjà<sup>5</sup>.

**Créez la variable suivante:** `phrase = "L'élision est l'effacement d'une voyelle en fin de mot devant la voyelle commençant le mot suivant."`

Et testez ce code :

---

```
print(phrase.split())#tokenisation de base
from mosestokenizer import *
tokenize = MosesTokenizer(lang='fr')
mots = tokenize(phrase)
print(mots)#tokenisation plus correcte
```

---

Attention, ça ne marchera peut être pas pour toutes les langues comme ça !

---

<sup>4</sup><https://pypi.org/project/mosestokenizer/>

<sup>5</sup><https://pip.pypa.io/en/stable/installing/>