

Datamining : TP3
M2 Informatique, Université Paris 7
Anne-Claire Haury

Exercice 1

1. Ecrire un script qui prend en entrée le fichier `SMSSpamCollection.txt`, lit les messages et leurs labels ('spam' ou 'ham') et les range dans une liste `messages` et une liste `labels`.
2. Créer une liste nommée `spams` contenant tous les spams et une liste nommée `hams` contenant tous les messages légitimes.
3. Ecrire une fonction `frequence_mots` qui prend en entrée une liste de messages et renvoie un dictionnaire permettant, pour chaque mot, d'accéder au nombre de fois qu'il apparaît dans l'ensemble des messages.
4. Explorer ces données en regardant par exemple quels mots apparaissent le plus souvent et en dessinant un histogramme. Faire de même pour les spams d'une part et les messages légitimes de l'autre.

Exercice 2

Ecrire une fonction `nettoyage_texte` qui prend en entrée une chaîne de caractères et renvoie le texte nettoyé de la ponctuation, des chiffres, des majuscules, etc. Si vous avez le temps, nettoyez également les urls, les accents, les mots trop courts, etc.

Pour cela, vous aurez sans doute besoin de la librairie `re` de python (expressions régulières).

Exercice 3

Ecrire une classe `Tfidf` qui traite un corpus de documents et calcule les valeurs `tf-idf` de chacun des mots. Un objet de cette classe doit être capable d'updater ces valeurs lorsqu'un nouveau document est lu.

Exercice 4

En utilisant la fonction `scipy.stats.ttest_ind` qui effectue un test de différence de deux moyennes comme vu en cours (cours 2), dire si la fréquence d'apparition du mot 'love' est significativement différente entre les spams et les non-spams. Idem pour le mot 'call'. Faites ce test pour chacun des mots et faites apparaître les mots dont la fréquence d'apparition est significativement différente entre les deux types de messages.