

Nous allons dans ce qui suit aligner des termes d'un fichier de BNF avec d'autre de RAMEAU, et déclarer les ID des meilleurs résultats d'alignement

Nous commençons par la conversion du fichier `rameau.txt` avec un encodage en `utf-8` afin de pouvoir opérer sur son contenu. À l'aide de l'invité de commande :

```
#iconv -f macintosh -t utf-8 RAMEAU-Groupes_ethniques-2.txt > RAMEAU_utf-8.txt
```

Figure 01

Déclaration des fonctions :

Importation des différents modules à utiliser.

Ensuite, on déclare les fonctions principales à exécuter plus tard : La mesure Jaro et Levenshtein Normalisé, pour lesquelles on va déclarer des seuils minimums à respecter.

On crée des fichiers de traitement des deux sources BNF et Rameau, écrire respectivement des lignes composées d'une chaîne de caractère numériques (ID) et l'expression à aligner.

BNF On récupère la chaîne de caractère après le dernier \$ (cadre en rouge figure

RAMEAU On récupère la chaîne de caractère après le premier \$

Traitement des fichiers de source :

On réduit toutes les lettres en minuscule et on enlève tous types d'accents et caractères spéciaux

On enlève tous les stop-word manuellement

```
#liste des caractères accentués à désaccentuer
accent = ['é','è','ê','ë','à','á','É','ù','û','ú','ü','ç','ô','ó','ö','î','í','  
sans_accent = ['e','e','e','e','a','a','E','u','u','u','u','c','o','o','o','i']

#liste stop words
stopWordFr=re.compile('d\\'|l\\'| de| du|\\'|_|"| la|peuple |peuples de|civilisation ')
```

Pour les chaines de plusieurs mots on `split()` et on mesure la distance de chaque mot qui la compose avec l'autre partie

```
1 39027974
2 143 $aTraditions$mEurope$mAlbanie
3 42406496
4 143 $aTraditions$mAfrique du Nord$mAlgérie
5 42155205
6 143 $aTraditions$mAsie$mYémen$mJuifs
7 42765162
8 143 $aTraditions$mEurope$mTchécoslovaquie$mSlovaquie$mRégion de Trenčín$mMyjava
9 38603967
10 143 $aTraditions$mEurope$mBulgarie$mRhodope$mNedelino
```

Fichier de départ BNF

```
1|(39027974 albanie)
2|(42406496 algerie)
3|(42155205 juifs)
4|(42765162 myjava)
5|(38603967 nedelino)
```

BNF Après traitement

```

1 11932945      $aTswana$gpeuple d'Afrique
2 11940850      $aPortugais
3 11960038      $aAlbanais
4 11988158      $aAlgériens
5 11949880      $aBulgares
6 11974575      $aBulgares de la Volga
7 11932167      $aJuifs
8 11956425      $aSlovaques
9 12385470      $aJuifs tchécoslovaques
10 11967055      $aJuifs d'Europe de l'Est
11 11940990      $aKabyles$gpeuple berbère
12 11937529      $aBété$gpeuple de Côte d'Ivoire

```

Fichier RAMEAU de départ

```

1|(11932945 tswana)
2|(11940850 portugais)
3|(11960038 albanais)
4|(11988158 algeriens)
5|(11949880 bulgares)
6|(11974575 bulgares volga)
7|(11932167 juifs)
8|(11956425 slovaques)
9|(12385470 juifs tchecoslovaques)
10|(11967055 juifs europe est)
11|(11940990 kabyles)
12|(11937529 bete)

```

RAMEAU après traitement

Mesures de similarité :

On parcourt chaque ligne du fichier BNF en lecture, et à l'aide des expressions régulières on détermine l'ID et la chaîne de caractère à aligner avec le terme de RAMEAU

```

##### on collecte les termes à mesurer #####

mots1 = open('newBNF_utf-8.txt','r')
read_mots1= mots1.readlines()

mots2 = open('newRAMEAU_utf-8.txt','r')
read_mots2= mots2.readlines()

for ligne in read_mots1:
    res_mots1= re.search(r".*(\d{8})\s+(.*)\).*",ligne)
    if res_mots1:
        iD1 = res_mots1.group(1)
        mot1 = res_mots1.group(2)
        mot1S = mot1.split()

```

Et à l'aide des fonctions Jaro et Levenshtein Normalisé on mesure la similarité entre les termes de BNF et de RAMEAU, tout en définissant un seuil Jaro minimal >0.8 et LevN >0.75 pour les mots et Jaro 0.7 et

Cas Levenshtein > 0.7 entre chaînes de caractères

```

207         if mot1==mot2:
208             print(mot1+' iD1:'+iD1,' <> ',mot2+' iD2:'+iD2 ,' : Identity similarity',' Lev :', levenshteinN
(mot1,mot2))
209         else:
210             # Si mot1 et mot2 sont des mots uniques
211             if levenshteinN(mot1,mot2)>=0.75: #<==== ici vous mettez le seuil de distance min
212                 if jaro(mot1, mot2) > 0.8: #<==== ici vous mettez le seuil de similarité Jaro min
213
214                     #metrics(tok3.tokenize(nltk.corpus.brown.raw()),
215                     #nltk.corpus.brown.words())
216                     print(mot1+' iD1:'+iD1,' <> ',mot2+' iD2:'+iD2 ,' : ', jaro(mot1, mot2),' Lev Sim:',
levenshtein(mot1,mot2))
217             # Si mot1 et mot2 sont composés de plusieurs mots, on mesure jaro et lev entre chaque mot des deux côtés
218             else :
219                 if jaro(mot1S, mot2S) > 0.7 :#and ngram.NGram.compare(mot1,mot2,N=3)>0.5:
220                     #metrics(tok3.tokenize(nltk.corpus.brown.raw()),
221                     #nltk.corpus.brown.words())
222                     print(mot1+' iD1:'+iD1,' <> ',mot2+' iD2:'+iD2 ,' Jaro P1P2: ', jaro(mot1S, mot2S),
' Lev Sim P1P2:', levenshteinN(mot1S,mot2S))

```

Les résultats qu'on peut voir déjà à l'écran de l'invité de commande :

```
nbendjoudi@x2go3:~/TER/TER2/Etape 2$ ./lecture BNF RAMEAU_jaro.py
albanie id1:39027974 <> albanais id2:11960038 : 0.869 Lev Sim: 0.75
algerie id1:42406496 <> algeriens id2:11988158 : 0.926 Lev Sim: 0.778
juifs id1:42155205 <> juifs id2:11932167 : Identity similarity Lev : 1.0
juifs id1:42155205 <> juifs tchecoslovaques id2:12385470 Jaro P1P2: 0.833 Lev Sim P1P2: 0.5
juifs id1:42155205 <> juifs europe est id2:11967055 Jaro P1P2: 0.778 Lev Sim P1P2: 0.333
```

Cas ou Levenshtein >0.5 entre les mots

```
nbendjoudi@x2go3:~/TER/TER2/Etape 2$ ./lecture BNF RAMEAU_jaro.py
albanie id1:39027974 <> albanais id2:11960038 : 0.869 Lev Sim: 0.75
algerie id1:42406496 <> algeriens id2:11988158 : 0.926 Lev Sim: 0.778
juifs id1:42155205 <> juifs id2:11932167 : Identity similarity Lev : 1.0
```