

Étape 1

- Écrire un premier script Python nommé « **extraire.py** » permettant d'extraire les entités médicales de type noms de médicaments par substance active de A à Z, à partir du contenu des 26 pages HTML du dossier « VIDAL » que je vous ai mis en pièce-jointe.
- Le script Python « **extraire.py** » doit générer en sortie un dictionnaire au format **.dic** (format DELAF vu en cours 4). Ce dictionnaire DELAF doit être encodé en **UTF-16 LE avec BOM** (UCS-2 LE BOM).
- Ce dictionnaire DELAF doit s'appeler « **subst.dic** » et doit donc contenir les noms de médicaments par substance active des 26 pages HTML du dossier « VIDAL ».
- Chaque entrée lexicale de ce dictionnaire doit être suivie par les informations (codes) **„N+subst**
- L'information **N** est de type grammatical et l'information **subst** est de type sémantique.
- Vous devez donc obtenir une sortie ayant le format **DELAF** d'UNITEX, comme illustré dans l'exemple suivant :

```
abacavir,.N+subst
abatacept,.N+subst
abciximab,.N+subst
abiratérone,.N+subst
.....
```

Remarque : L'encodage UTF-8 sans BOM des pages HTML du dossier « VIDAL » ne doit pas être modifié.

- Générer un fichier nommé « **infos1.txt** » contenant :
 - le nombre d'entités médicales de type noms de médicaments par substance active du dictionnaire DELAF « **subst.dic** » généré préalablement, pour chaque lettre de l'alphabet ;
 - et le nombre total d'entités médicales de type noms de médicaments par substance active du dictionnaire DELAF « **subst.dic** ».

Remarque : Ce premier script python « **extraire.py** » doit impérativement avoir 1 argument : le nom du dossier « VIDAL ». Ce script doit également afficher sur l'invite de commandes le nombre d'entités médicales de type noms de médicaments par substance active du dictionnaire DELAF « **subst.dic** » généré préalablement, pour chaque lettre de l'alphabet.

Étape 2

- Après avoir extrait les entités médicales de type noms de médicaments par substance active à partir du dossier « VIDAL » et généré le dictionnaire « **subst.dic** », vous devrez écrire un deuxième script Python « **enrichir.py** », permettant d'alimenter et d'enrichir le dictionnaire « **subst.dic** » (généré dans l'étape précédente) avec de nouvelles entités médicales de type noms de médicaments par nom commercial ou par substance active, à partir du fichier « **corpus-medical.txt** » donné en argument.
- L'encodage UTF-8 sans BOM du fichier du corpus médical ne doit pas être modifié et le dictionnaire « **subst.dic** » après enrichissement doit conserver son encodage de départ, à savoir l'« **UTF-16 LE avec BOM** » (UCS-2 LE BOM).
- Après l'étape d'enrichissement, le dictionnaire « **subst.dic** » ne doit pas contenir de doublons et doit être trié par ordre croissant (a-z). Il contiendra donc toutes les entités médicales de type noms de médicaments par substance active issues du dossier « VIDAL » + les nouveaux noms de médicaments issus du corpus médical « **corpus-medical.txt** ».
- Le script d'enrichissement « **enrichir.py** » doit garder une trace des noms de médicaments trouvés dans le fichier « corpus-medical.txt », en les stockant dans un autre fichier qui doit s'appeler « **subst_corpus.dic** », en mettant ses entrées lexicales en minuscules. Cependant, ce dictionnaire ne doit subir ni tri, ni suppression de doublons et doit être encodé en « **UTF-16 LE avec BOM** » (UCS-2 LE BOM).
- Le script « **enrichir.py** » doit également générer un fichier nommé « **infos2.txt** » sans doublons contenant :
 - le nombre de médicaments issus du corpus pour chaque lettre de l'alphabet ;
 - et le nombre total de médicaments issus du corpus.
- Le script « **enrichir.py** » doit générer un autre fichier nommé « **infos3.txt** » sans doublons contenant :
 - le nombre de médicaments conservés pour l'enrichissement pour chaque lettre de l'alphabet ;
 - et le nombre total de médicaments conservés pour l'enrichissement.

Étape 3

- Construire un graphe d'extraction (.grf) sous UNITEX, qui se base **impérativement** sur l'étiquette **<N+subst>** du dictionnaire « **subst.dic** », afin d'**extraire les occurrences de** « **posologies** » à partir du fichier « corpus-medical.txt ». Le graphe d'extraction **doit s'appeler** « **posologie.grf** ». Le résultat de cette extraction sera placé par UNITEX dans le fichier « **concord.html** », qui se trouve dans le dossier « corpus-medical_snt » généré par UNITEX.

Remarque : Une « posologie » contient généralement le nom du médicament, le dosage du médicament (**50 mg, 20 mg, 10 mg, 500, 400, 0,4 ml, 0.4 ml, 0,4, 4000 UI**, etc.), le rythme ou fréquence d'administration (**2 fois par jour, 3 fois par jour, 4 fois par jour, 1 le matin et 1 le soir (donc 2 fois par jour)**, etc.), l'heure-moment de prise du médicament (**à 8 heures, à 20h00, le soir, le matin, trois fois par jour (donc le matin, le midi et le soir)** etc.) et la durée de traitement (**pendant un mois, pendant encore 21 jours, de J1 à J7**, etc.).

Remarque : Il est à noter que dans certains cas, le dosage de médicament n'est pas présent, par exemple, "METOPROLOL : ½ le matin, ½ le soir". Dans cet exemple, le dosage du METOPROLOL n'est pas précisé. Pourtant, il en existe différents dosages, comme le "METOPROLOL 100 mg", employé dans "METOPROLOL 100 mg : ½ le matin, ½ le soir" ou le "METOPROLOL 50 mg", employé dans "METOPROLOL 50 : 1/jour".

Exemples d'extraction de « posologies » à partir du corpus médical « corpus-medical.txt » :

TOPALGIC 100 mg 1 amp, 3 fois par jour, pendant 5 jours

INNOHEP 3 500 unités : 1 injection par jour pendant encore 21 jours

SIMVASTATINE 20 mg : 1 cp/j à 8 heures pendant un mois

CYTARABINE 100 mg/m² de J1 à J7

PLAVIX 75 mg : 1 cp/jour

ZOLPIDEM 10 mg 1 cp au coucher

METFORMINE 850 mg 3 fois par jour

SPECIAFOLDINE 5 mg : 1 cp matin – 1 cp soir pendant un mois

ALADACTONE 25 mg : 1 cp/jour le midi

INEXIUM 40 1 cp par jour le soir

TEGRETOL 200 mg : 1 cp 2 fois par jour

PAROXETINE 20 mg : 1 fois par jour

EQUANIL 400 : 3 fois / jour

KEPPRA 500 : 2/jour

CRESTOR 10 mg : 1 comprimé par jour le soir

LOVENOX 0,4 : 20h

LOVENOX 4000 UI : 1/jour

LOVENOX 0,4 : 19h

LOVENOX 0,4 ml : 1 injection/jour le soir

LOVENOX 0.4 1 inj/jour à midi

Remarque : Dans certains cas, l'heure-moment de prise du médicament et la durée de traitement ne sont pas précisés, comme dans la posologie suivante :

PLAVIX 75 mg : 1 cp/jour

- Écrire un troisième script nommé « **unitex.py** » permettant d'appeler UNITEX pour exploiter votre graphe, à partir de l'emplacement **C:\.....\Unitex-GramLab\App>**
 - a. **Pour appeler UNITEX, vous devrez utiliser le script du cours dédié au lancement d'UNITEX à partir d'un script Python.** Ce troisième script Python « **unitex.py** » **doit exploiter** les ressources suivantes :
 - I. le dossier « **corpus-medical_snt** » créé automatiquement à chaque lancement du script « **unitex.py** » ;
 - II. le fichier : « **corpus-medical.txt** » ;
 - III. le fichier : « **corpus-medical.snt** » ;
 - IV. le fichier : « **Norm.txt** » ;
 - V. le fichier : « **Alphabet.txt** » (préciser dans quelle phase du script « **unitex.py** » ce fichier « **Alphabet.txt** » doit être utilisé et expliquer à quoi sert ce fichier TXT, en donnant des exemples précis. Cette réponse doit être écrite sous forme de commentaires dans le script « **unitex.py** ».) ;
 - VI. le fichier : « **subst.dic** » ;
 - VII. le fichier : « **subst.bin** » ;

- VIII. le fichier : « **Dela_fr.bin** » ;
- IX. le fichier : « **Dela_fr.inf** » ;
- X. le fichier : « **posologie.grf** » ;
- XI. le fichier : « **posologie.fst2** » ;
- XII. le fichier : « **concord.ind** » du dossier « **corpus-medical_snt** ».

Remarque : Lors de la phase d'extraction, il est **nécessaire** d'utiliser comme ressource supplémentaire le dictionnaire système « **Dela_fr.bin** » fourni par UNITEX, afin de pouvoir exploiter les masques lexicaux comme <PREP>, <DET> ou <PREPDET>, etc. **Vérifiez aussi que vous avez bien « Dela fr.inf » à côté du « Dela fr.bin », afin que ce dernier puisse être exploité.**

- Pour lancer votre application d'extraction d'information, placez vos 3 scripts (**extraire.py**, **enrichir.py** et **unitex.py**) dans l'emplacement **C:\.....\Unitex-GramLab\App>**

Pour l'évaluation de votre travail, vous **devrez m'envoyer par mail** :

- **Le script d'extraction** : « **extraire.py** » doit générer « **subst.dic** » et « **infos1.txt** ». Ce script prend 1 argument : le dossier « VIDAL ».
- **Le script d'enrichissement** : « **enrichir.py** » doit enrichir le DELAF « **subst.dic** » à partir du fichier « **corpus-medical.txt** » donné en argument. Ce script doit générer 4 fichiers :
 - I. « **subst.dic** » (dictionnaire enrichi à partir du fichier « **corpus-medical.txt** ») ;
 - II. « **subst_corpus.dic** » ;
 - III. « **infos2.txt** » ;
 - IV. « **infos3.txt** ».
- **Le script Python qui appelle UNITEX** : « **unitex.py** » doit exploiter les ressources citées ci-dessus, telles que le graphe « **posologie.grf** » et le DELAF « **subst.dic** ».
- **Le graphe d'extraction** : « **posologie.grf** » doit extraire à partir du fichier « **corpus-medical.txt** » les posologies, en s'appuyant sur les DELAF « **Dela_fr.bin** » et « **subst.bin** ». Le résultat doit contenir au minimum **1000 extractions correctes**.

Consignes du projet « Extraction d'information »

Pour résumer, vous devrez m'envoyer **5 fichiers** :

- les 3 scripts Python ;
- Le fichier « concord.html » ;
- et le graphe d'extraction au format **.grf**.