

Linguistique de corpus 2

Outils et méthodes de traitement de corpus 2

Patrick Paroubek

LIMSI-CNRS
Dépt. CHM - Groupe LIR
Bât. 508 Université Paris XI, 91403 Orsay Cedex
pap@limsi.fr

mercredi 01 février 2017 / L3 - Cours S2-1

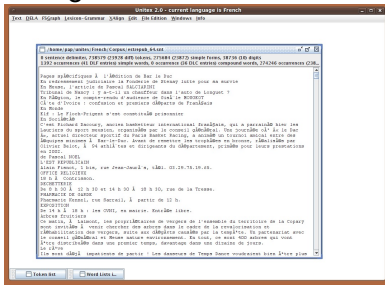
Introduction à Unitex

La boîte à outils Unitex permet d'effectuer des analyses linguistiques de documents avec une interface graphique de haut-niveau.

Unitex et NOOJ sont deux plateformes issues des mêmes travaux (Maurice Gross) et comparables en fonctionnalités. L'importation d'un texte dans Unitex propose par défaut une segmentation en mot et en phrase ainsi qu'un ensemble de pré-traitements de normalisation (par ex. encodage des caractères).

Introduction à Unix

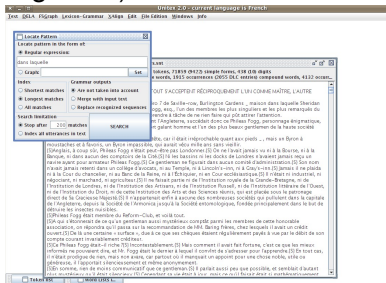
Charger un texte dans la boîte à outils :



En activant l'item de menu : Text/Open... Le chargement produit automatiquement les listes de vocabulaires des mots simples reconnus, des mots composés reconnus et des mots inconnus des dictionnaires d'Unix, accessibles depuis les onglets en bas à gauche de l'écran après importation d'un texte dans Unix.

Introduction à Unixex

La plus simple des fonctionnalités d'Unixex permet d'effectuer une recherche textuelle avec un "patron à trous" (expression régulière).



En activant l'item de menu : Text/Locate Pattern ...

- Opérateur de concaténation (le point)

```
le <A> chien  
(le <A>)chien  
le.<A>chien  
(le).<A> chien  
(le.(<A>)) (chien)
```

- Opérateur d'alternative (le plus),

```
(je+tu+il+elle+nous+vous+ils+elles) <V>  
(le brave garçon)+(le garçon)
```

- Opérateur de répétition (étoile de Kleen), de 0 à N occurrences.

```
le brave* garçon
```

Masques lexicaux

- $\langle \wedge \rangle$ saut de ligne
- $\langle \text{PNC} \rangle$ ponctuation
- $\langle \text{MOT} \rangle$ n'importe quel mot composé de lettres
- $\langle \text{NB} \rangle$ n'importe quel mot composé de chiffres
- les mots du dictionnaire $\langle \text{DIC} \rangle$
- $\langle \text{A} : \text{mp} : \text{f} \rangle$ les adjectifs masculins pluriels ou féminins.
- $\langle \text{N} \sim \text{Hum} + \text{z1} \rangle$ les noms (N) qui ne réfèrent pas à des humains ($\sim \text{Hum}$) et qui sont de niveau de langue générale ($+ \text{z1}$)

Automates d'états finis déterministes

Qu'est-ce qu'un automate d'états fini déterministe ?

Un automate d'états fini déterministe est une autre représentation d'une expression régulière.

Un automate d'états fini déterministe est défini par :

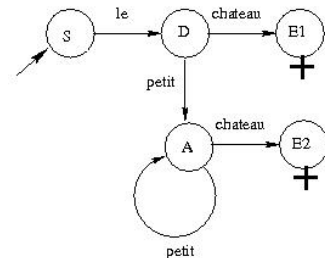
- 1 un ensemble fini d'états
- 2 un ensemble fini de transitions (arcs reliant les états)
- 3 un alphabet d'entrée (liste finie de symboles d'entrée)
- 4 un état initial
- 5 un ou plusieurs états finaux.

L'automate est déterministe si dans un état donné, pour un symbole d'entrée donné, on ne peut emprunter qu'une seule transition.

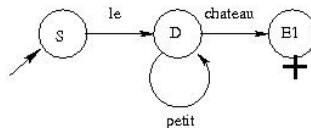
Automates d'états finis déterministes

le chateau
le petit chateau
le petit petit chateau
le petit petit petit chateau
le petit petit petit.....petit chateau

automate non minimal



automate minimal



Automates d'états finis déterministes

Un automate d'états fini déterministe est minimal, s'il n'existe aucun autre automate d'états fini déterministe avec moins d'états qui reconnaisse le même langage (automate équivalent).

Un automate d'état fini déterministe minimal n'a pas de transition vide (souvent notées ϵ).

Un automate déterministe fonctionne en “accepteur” (reconnaissance).

Automates d'états finis déterministes (suite)

Un automate d'états fini déterministe est minimal, s'il n'existe aucun autre automate d'états fini déterministe avec moins d'états qui reconnaisse le même langage (automate équivalent).

Un automate d'état fini déterministe minimal n'a pas de transition vide (souvent notées ϵ).

Un automate déterministe fonctionne en “accepteur” (reconnaissance).

Pour le faire fonctionner en “traducteur” c'est à dire en transducteur, (en anglais *transducer*), il suffit d'ajouter des symboles de sortie sur les transitions.