

Analyse et représentation de la langue

Damien Nouvel



Plan

1. Analyses lexicales et morphologiques
2. Analyses syntaxiques
3. Vers la sémantique

Correction orthographique

- ▶ **Recherche dans une base lexicale**
 - Liste des mots pour une langue donnée
 - Comparaison mots du texte vs lexique
 - Disposer de toutes les formes possibles : **couverture**
 - ⇒ Complexe : par exemple $1000 * 50K = 50M \dots$
- ▶ Optimisation de la recherche
 - Utilisation d'automates
 - Structure de préfixes (TRIE)
- ▶ Besoin d'informations contextuelles

Couverture du lexique

- ▶ Difficultés de couverture
 - Morphologie, dont **flexions** (déclinaisons / conjugaisons)
 - ⇒ Utilisation de **paradigmes** de flexion
 - Variantes d'écritures
 - Noms propres (classe **ouverte**) acronymes, etc.
 - ⇒ Automates dédiés aux reconnaissances partielles

Applications

- ▶ Éditeurs de texte
 - ▶ Saisie prédictive / autocorrect (T9)
 - Handicap (aveugle, moteurs)
 - Grand public
- ⇒ Logiciels **hors-ligne** (contrairement aux reconnaissances)
- ▶ Détection de nouveaux mots (néologismes, noms propres, etc.)

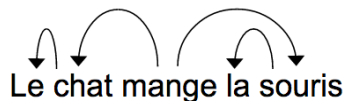
Plan

1. Analyses lexicales et morphologiques
2. Analyses syntaxiques
3. Vers la sémantique

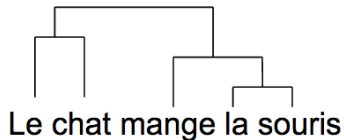
Lier les mots dans les énoncés

- ▶ Niveau de représentation
 - Exploite les approches morphologiques / lexicales
 - ⇒ Importance des **catégories morpho-syntaxiques** (POS)
 - Détermine les liens / relations / **dépendances** entre les mots
 - Identifie des **syntagmes**
- ▶ Deux catégories d'approches
 - **Constituants** : arbres syntaxiques (projectivité)
 - **Dépendances** : relations de dépendances (Tesnière, 1940)
 - ⇒ Possibilité de convertir de l'un à l'autre
 - ⇒ Approches : LFG, GPSG, TAG, HPSG ...

Exemple



(Bourigault / SYNTEX)



Quelques difficultés

- ▶ Comment traiter
 - Le rattachement prépositionnel
- ⇒ « Il voit le chef de la bande » / « Il voit le chef de sa fenêtre »
 - Les propositions (relatives, subordonnées, etc.) ?
 - La coordination ?

Quelques difficultés

- ▶ Quelques essais avec **FRMG** (de la Clergerie)
 - Portail : <http://alpage.inria.fr/frmgwiki/>
 - Analyseur : http://alpage.inria.fr/frmgwiki/frmg_main/frmg_server
 - Embeddings : <http://alpage.inria.fr/depglove/process.pl>

Plan

1. Analyses lexicales et morphologiques
2. Analyses syntaxiques
3. Vers la sémantique

Profondeur de la syntaxe

- ▶ Analyse syntaxique
 - Morpho-syntaxe
 - Syntaxe de surface / parenthésage (syntagmes)
 - Syntaxe « grammaticale »
- ▶ Syntaxe « profonde »
 - Sujet du verbe *vs* agent de l'action (passif)
 - Objet : patient, récepteur, etc.
 - Sujet distant (anaphores / coréférences ...)

⇒ Interface syntaxe / sémantique

Les entités nommées

- ▶ Pas de définition stable
 - « Noms propres et quantités d'intérêt » (MUC)
 - « Entités du monde concret qui ont un nom » (ESTER)
 - « Expression qui réfère à une entité unique » (Ehrmann)
 - « Objets mentaux pour la logique » (Nouvel)
 - ...
- ▶ Deux grandes catégories linguistiques
 - **Noms propres**
 - ⇒ Limites : « Superman », « l'Histoire », etc.
 - **Descriptions définies** (syntagmes nominaux définis)
 - ⇒ Limites : « cette voiture », « le roi des forains », « l'occident »
- ⇒ Tenir compte des **synonymes**, **homonymes**, **métonymies**
- ⇒ Beaucoup d'applications ...prétraitement ?

Structurer les textes en connaissances

▶ Des **formats** plus ou moins **structurés**

- Texte brut : sauts de lignes
- PDF : graphiquement stable, mais peu structuré
- HTML : sections, styles, liens, micro-formats
- Wikipedia : articles et info-boxes
- BabelNet : <http://babelnet.org/synset?word=Inalco&lang=FR>

⇒ Structuration : du texte aux **graphes** (valués)

▶ Quelques **réseaux sémantiques**

- WordNet
- FrameNet
- Graphes conceptuels
- Web sémantique (RDF/OWL)

⇒ Interactions **homme-machine** et **machine-machine**