

Corpus et méthode expérimentale

module X9IT080

Florian Boudin

Département informatique, Université de Nantes

Révision 1 du 30 juillet 2012

Présentation du module

- ▶ Notions abordées dans ce module
 - ▶ ...
- ▶ Volume horaire : 15 séances (20h)
 - ▶ 6 CM (8h) - 6 TD (8h) - 3 TP (4h)
- ▶ Notation

Plan

La notion de corpus

Les normes de caractérisation et d'annotation

Démarche expérimentale

Références

Définitions I

- ▶ un **corpus** est un ensemble de documents (textes, images, vidéos) regroupés dans une optique précise
- ▶ dans le cadre de ce cours : corpus \rightarrow textes
- ▶ caractéristiques d'un corpus :
 - ▶ la taille
 - ▶ le langage
 - ▶ le temps couvert par les textes du corpus
 - ▶ le registre de langage

Définitions II

- ▶ le corpus doit atteindre une **taille** critique pour permettre des traitements statistiques fiables
 - ▶ 10K+ phrases pour entraîner un *POS-tagger*
 - ▶ 20M+ mots pour construire un modèle de langue
- ▶ un corpus **monolingue** couvre une seule langue, et une seule déclinaison de cette langue
 - ▶ e.g. français de France et français du Québec

Définitions III

- ▶ la **période couverte** par les textes doit être courte
 - ▶ français d'aujourd'hui \neq français parlé au 19^e siècle
 - ▶ français d'aujourd'hui \neq français de 1990 (néologismes)
- ▶ les **registres de langage** doivent être séparés
 - ▶ un corpus de textes scientifiques ne peut être utilisé pour extraire des informations sur les textes vulgarisés
 - ▶ un corpus de textes scientifiques et vulgarisés ne permettra pas de tirer de conclusion sur ces deux registres

Corpus parallèle

- ▶ un **corpus parallèle** est un ensemble de paires de textes tel que, pour une paire, un des textes est la traduction de l'autre
- ▶ un **alignement** des unités textuelles est nécessaire :
 - ▶ mettre en correspondance des unités textuelles en langue source avec celles de la langue cible
- ▶ un alignement peut être manuel ou automatique
- ▶ la granularité de l'alignement dépend de l'utilisation
 - ▶ traduction automatique, dictionnaires bilingues, etc.

Alignement de phrases

Il s'agit de l'un des sauropodes les plus connus.

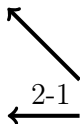
C'était un très grand quadrupède au long cou, avec une longue queue en forme de fouet.

Ses membres antérieurs étaient légèrement plus courts que ses membres postérieurs, ce qui lui donnait une posture horizontale.

One of the best-known sauropods, Diplodocus was a very large long-necked quadrupedal animal, with a long, whip-like tail.


Its forelimbs were slightly shorter than its hind limbs, resulting in a largely horizontal posture.

It is the longest dinosaur known from a complete skeleton.




Alignement de mots

Je suis Français .
I am French .



Je m' appelle Paul .
My name is Paul .



Corpora parallèles disponibles

- ▶ Europarl [Koe05]
 - ▶ délibérations du Parlement européen disponibles en 21 langues et alignées au niveau de la phrase.
- ▶ OPUS [Tie09]
 - ▶ *European Medicines Agency documents, KDE4 localization files, Europarl, OpenSubtitles, etc.*
 - ▶ <http://opus.lingfil.uu.se/>
- ▶ livres numériques librement disponibles
 - ▶ <http://www.gutenberg.org/>
- ▶ Wikipédia : plus comparable que parallèle

Corpus comparable I

- ▶ les corpora parallèles sont très coûteux et ils ne sont disponibles que dans peu de langues/domaines
- ▶ les corpus dits **comparables** sont plus répandus
- ▶ d'après Déjean & Gaussier [DG02] :
 - ▶ « Deux corpus de deux langues l_1 et l_2 sont dits comparables s'il existe une sous-partie non négligeable du vocabulaire du corpus de langue l_1 , respectivement l_2 , dont la traduction se trouve dans le corpus de langue l_2 , respectivement l_1 . »

Corpus comparable II

- ▶ exemple de corpus comparable :
 - ▶ un ensemble d'articles de journaux dans différentes langues, traitant d'une même actualité et à la même époque
- ▶ applications
 - ▶ extraction de phrases parallèles [SQT10]
 - ▶ constitution de dictionnaires bilingues [Rap99]
- ▶ peu (ou pas ?) de corpora comparables disponibles
 - ▶ Wikipédia, le web multilingue, etc.

La constitution de corpus

1. définir les caractéristiques du corpus :

- ▶ quels sont les phénomènes que l'on souhaite observer ?
- ▶ quelle est la tâche que l'on souhaite réaliser ?
- ▶ le corpus doit-il pouvoir être distribué ?

2. assembler les unités textuelles

- ▶ est-ce que le processus peut être automatisé ?
 - ▶ e.g. moissonnage à partir du web (*web scraping*)

→ garder en mémoire la méthodologie utilisée

3. annotation du corpus (optionnel)

- ▶ annoter les unités textuelles (e.g. *Part-Of-Speech*)
- ▶ créer un référenciel pour l'évaluation (e.g. termes-clés)

→ définir des *guidelines* à joindre au corpus

Étude de cas : extraction de termes clés

- ▶ un corpus pour entraîner et évaluer un système d'extraction de termes-clés
 - nature des documents : **articles**, blogs, tweets, etc.
 - langue(s) des documents : anglais, **français**, etc.
 - source(s) des documents : lemonde.fr, **wikinews**, etc.
 - nombre de documents : 10, 20, 50, **100**, 1000, etc.
- ▶ récupération des documents
 - ▶ automatiser le processus à partir d'un *dump* de wikinews
 - ▶ définir un format pour les documents (XML, txt, etc.)
- ▶ création d'un référentiel pour l'évaluation
 - ▶ tâche subjective → plusieurs annotations par document

Étude de cas : analyse en dépendances

- ▶ un corpus de phrases annotées en dépendances entraîner un *parser*
 - langue(s) des phrases : _____
 - source(s) des phrases : _____
 - nombre de phrases : _____
- ▶ récupération des phrases
 - automatisée ou manuelle ?
- ▶ annotation des phrases
 - ▶ recruter des spécialistes ? définir des *guidelines* ?

L'annotation de corpus I

- ▶ différents niveaux d'annotation :
 - ▶ collection de documents, e.g. regroupement par sujets
 - ▶ document, e.g. jugement de pertinence en RI
 - ▶ paragraphe, e.g. découpage thématique
 - ▶ phrases, e.g. segmentation en phrases
 - ▶ multi-mots, e.g. détection d'entités nommées
 - ▶ mots, e.g. tokenization
 - ▶ caractères, e.g. analyse morphologique
- ▶ l'annotation d'un point de vue technique
 - ▶ séparation des annotations du texte
 - ▶ format (e.g. XML) et encodage standard

L'annotation de corpus II

- ▶ critères de sélection des annotateurs
 - ▶ spécialistes, e.g. médecin pour de la RI médicale
 - ▶ natif de la langue, e.g. évaluation de la grammaticalité
 - ▶ extérieur au projet, e.g. annotation par les auteurs...
 - ▶ attention au *crowdsourcing* !
- ▶ variabilité des annotations (*consistency*)
 - ▶ résoudre les désaccords par consensus ?
 - ▶ considérer plusieurs annotations ?

Étude de cas : annotation de tweets

- ▶ détecter la **polarité** d'un tweet

- ▶ annoter un sentiment négatif, neutre ou positif

1. « Décès de Whitney Houston : la vilaine bourde de Sony Music »
2. « Interview des #DaftPunk sur France Inter 1 mois après la sortie de l'album »
3. « Seuls 126% des marseillais exagèrent. »

Étude de cas : annotation de tweets

- ▶ détecter la **polarité** d'un tweet

- ▶ annoter un sentiment négatif, neutre ou positif

1. « Décès de Whitney Houston : la vilaine bourde de Sony Music » → négatif
2. « Interview des #DaftPunk sur France Inter 1 mois après la sortie de l'album » → neutre
3. « Seuls 126% des marseillais exagèrent. » → ?

Un point sur l'encodage

- ▶ la lecture/écriture de fichiers textes suppose l'utilisation d'un encodage des caractères
 - ▶ anglais : ASCII
 - ▶ français : ISO-8859-1
 - ▶ japonais : ISO-2022-JP
 - ▶ etc.

Un point sur l'encodage

- ▶ la lecture/écriture de fichiers textes suppose l'utilisation d'un encodage des caractères
 - ▶ anglais : ASCII
 - ▶ français : ISO-8859-1
 - ▶ japonais : ISO-2022-JP
 - ▶ etc.

→ UTF-8

Un point sur les méta-données

- ▶ une **méta-donnée** est une donnée servant à définir ou décrire une autre donnée
 - ▶ e.g. associer à une donnée la date de création, à une photo les coordonnées GPS du lieu où elle a été prise
- ▶ le *Dublin Core* est la principale initiative visant à la convergence des éléments de méta-données à utiliser
- ▶ 15 éléments de description :
 - ▶ formels (titre, créateur, éditeur)
 - ▶ intellectuels (sujet, description, langue, etc.)
 - ▶ relatifs à la propriété intellectuelle (droits)

Plan

La notion de corpus

Les normes de caractérisation et d'annotation

Démarche expérimentale

Références

Les normes de caractérisation et d'annotation

- ▶ Partie de Béatrice

Plan

La notion de corpus

Les normes de caractérisation et d'annotation

Démarche expérimentale

Références

Introduction

- ▶ **démarche expérimentale** → valider une hypothèse par des expériences
 - ▶ e.g. expérience du cerf-volant de Benjamin Franklin
- ▶ **évaluer** un système de TAL consiste à vérifier qu'il produit le résultat pour lequel il a été conçu
- ▶ le TAL regroupe

Introduction

Plan

La notion de corpus

Les normes de caractérisation et d'annotation

Démarche expérimentale

Références

References I



Hervé Déjean and Eric Gaussier.

Une nouvelle approche de l'extraction de lexiques bilingues à partir de corpus comparables.

[Lexicometrica, Alignement lexical dans les corpus multilingues](#), pages 1–22, 2002.



Philipp Koehn.

Europarl : A parallel corpus for statistical machine translation.

In [MT summit](#), volume 5, 2005.



Reinhard Rapp.

Automatic identification of word translations from unrelated english and german corpora.

In [Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics](#), pages 519–526. Association for Computational Linguistics, 1999.

References II



Jason R Smith, Chris Quirk, and Kristina Toutanova.

Extracting parallel sentences from comparable corpora using document level alignment.

In Human Language Technologies : The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics, pages 403–411. Association for Computational Linguistics, 2010.



Jörg Tiedemann.

News from opus-a collection of multilingual parallel corpora with tools and interfaces.

In Recent Advances in Natural Language Processing, volume 5, pages 237–248, 2009.