

Exploiter les
données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

Exploiter les données

Introduction à la textométrie avec TXM

Jean-Baptiste Camps & Simon Gabay

Univ. de Neuchâtel

Formation en philologie numérique :
encoder, exploiter, diffuser
12-16 février 2018

De la production des données à l'exploitation

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

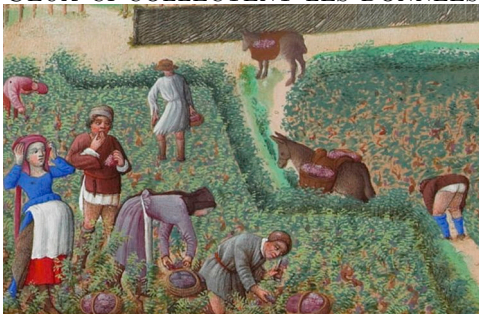
Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

CEUX-CI COLLECTENT LES DONNÉES



CEUX-LÀ LES EXPLOITENT



Objectifs

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

- S'initier à la **textométrie**,
 - créer des corpus ;
 - les interroger ;
 - faire (un peu) d'analyse quantitative.
- dans le cadre d'un logiciel "tout en un" et convivial, **TXM** [Heiden et al., 2010],
- sur des cas tirés de la littérature du XVII^e siècle.

La textométrie selon [Pincemin et al., 2008]

La textométrie développe les possibilités de consultation et d'analyse de corpus textuels en faisant appel à des décomptes et des modélisations statistiques et en combinant aux possibilités de repérage d'occurrences des calculs de tri, de sélection et de réorganisation statistique.

TXM : un logiciel de textométrie



`http://textometrie.ens-lyon.fr/`
Base de français médiéval :
`http://txm.bfm-corpus.org/.`

- Logiciel libre et multiplateforme ;
- développé à l'ÉNS-LSH de Lyon ;
- dévoué à la textométrie ;
- repose sur des technologies de référence :
 - XML/TEI pour les données ;
 - R pour l'analyse statistique ;
 - CQP pour l'interrogation de corpus ;
 - TreeTagger pour l'annotation.

Plan

Exploiter les
données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

- 1 Mise en jambe : *Andromaque*
- 2 Importer des données et créer un corpus
- 3 Interroger les données
- 4 Quelques notions d'analyse quantitative

Plan

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

1 Mise en jambe : *Andromaque*

2 Importer des données et créer un corpus

3 Interroger les données

4 Quelques notions d'analyse quantitative

Création d'un corpus à partir de notre édition d'*Andromaque*

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

- ❶ Fichier, importer, import XML/W + CSV ;
- ❷ sélectionner le dossier avec les sources et remplir les paramètres du corpus ;
- ❸ lancer la création du corpus.

Premières fonctionnalités

Exploiter les
données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

- ❶ Consulter la description du corpus ;
- ❷ parcourir l'édition ;
- ❸ regarder le lexique ;
- ❹ ouvrir l'index, y chercher les occurrences de 'Seigneur' ;
 - ❶ clic-droit, envoyer vers les concordances ;
 - ❶ double-clic sur une occurrence pour aller au texte ;
 - ❷ clic-droit, envoyer vers les cooccurents ;
 - ❶ aller d'un cooccurrent aux concordances, puis au texte
 - ❸ clic-droit, envoyer vers la progression ;

Partitions et quelques éléments descriptifs

Exploiter les
données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

- ❶ créer une partition, en sélectionnant la structure `sp` et l'attribut `who` ;
- ❷ consulter les dimensions ;
- ❸ créer une table lexicale, expérimenter avec les tris, la fusion ou suppression des colonnes, etc.

Statistiques de base

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

À partir de la table lexicale créée,

- ➊ calculer les spécificités ;
- ➋ quels sont les mots les plus spécifiques d'Oreste ? de Pylade ?
- ➌ en sélectionner quelques uns qui sont pertinents ;
- ➍ calculer le diagramme en bâton des lignes sélectionnées.

Qu'en déduire ?

Exploiter les données

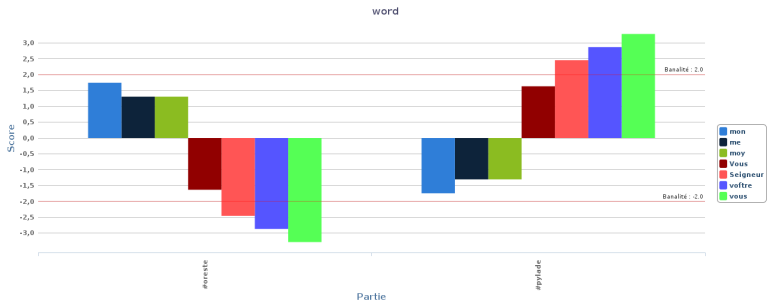
Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative



Exploiter les
données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

Prêts à passer aux choses sérieuses ?

Plan

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

1 Mise en jambe : *Andromaque*

2 Importer des données et créer un corpus

3 Interroger les données

4 Quelques notions d'analyse quantitative

Différents modes d'import

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

flemmards import (avec ou sans métadonnées complémentaires) depuis

- presse-papier ;
- des fichiers txt ;
- traitement de texte.

XML (TEI ou autre) ;
formats de logiciels de textométrie.

Corpus du jour : un peu de théâtre du XVII^e siècle

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

Corpus constitué pour le cours d'aujourd'hui :

- source : Paul Fièvre,
<http://www.theatre-classique.fr/>;
- documents encodés en XML/TEI (ou dans plusieurs XML/TEI) ;
- corpus de 36 pièces de théâtre en vers du XVII^e siècle,
- appartenant à trois genres principaux :
 - comédie,
 - tragédie,
 - tragi-comédie.
- sélectionnées un peu au hasard,
 - mais en essayant de conserver un équilibre entre les genres principaux (comédie, tragédie, tragi-comédie),
 - d'avoir des pièces de longueur similaire (entre 1250 et 2000 vers),
 - un équilibre entre les auteurs (4 pièces par auteur),
 - et entre les générations (4 auteurs par génération).

9 auteurs, 2 générations

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

G1, c. 1630-1650

- Pierre Du Ryer
(fl. 1628-1655) ;
- Georges de Scudéry
(fl. 1631-1643).
- Jean de Rotrou
(fl. 1635-1649) ;
- Paul Scarron
(fl. 1648-1660) ;

G2, c. 1650-1690

- Claude Boyer
(fl. 1646-1697) ;
- Thomas Corneille
(fl. 1651-1696)
- Molière
(fl. 1655-1673)
- Jean Racine
(fl. 1664-1691) ;

Et un monstre sacré, **Pierre Corneille** (fl. 1629-1675).

Sources et pré-traitements

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

- graphies modernisées (dommage... mais va faire notre affaire dans ce cas précis) ;
- Faut-il supprimer la distinction majuscule/minuscule ?
 - Pour : suppression de biais éditoriaux.
 - Contre : les majuscules peuvent conserver de l'information syntaxique.
- Veut-on garder tout ce qui est extérieur aux répliques (liste des personnages, page de titre, etc.) ?
 - Peut être retiré grâce à la structuration XML/TEI.
- pas de lemmatisation : possibilité de lemmatiser automatiquement.

Transformations avant l'import

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

Dossier xsl, feuille source_to_txt.xsl, deux sorties :

- ❶ fichier metadata.csv : métadonnées des documents extraites automatiquement des fichiers TEI (teiHeader et page de titre, docDate) ;
- ❷ dossier txt : transformation en txt des pièces :
 - passage en bas de casse ;
 - suppression du teiHeader ;
 - suppression du castList et des mentions de personnage, speaker ;
 - suppression du front, des docTitle, docDate, docAuthor, docImprint, printer, performance, div[\@type='dedicace'] ;
 - suppression des titres, notes.

Import txt + csv

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

- 1 sélectionner le répertoire de sources `txm_import1_txt` (contenant fichiers `txt` et métadonnées `csv`, corpus tronqué pour gagner un peu de temps) ;
- 2 paramétrer l'import (nommer le corpus `THEATRENEUCHTXT`) ;
- 3 demander la lemmatisation ;
- 4 vérifier que les métadonnées sont bien comprises ;
- 5 lancer l'import ;
- 6 (regarder le log de l'import) ;
- 7 une fois l'import réussi, jeter un œil à la description, à l'édition, etc.

Import XML amélioré

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

- ➊ sélectionner le répertoire de sources `txm_import2_xml` (contenant fichiers `xml` et métadonnées `csv`, complet) ;
- ➋ paramétrer l'import (nommer le corpus THEATRENEUCHXML) ;
- ➌ demander la lemmatisation ;
- ➍ vérifier que les métadonnées sont bien comprises ;
- ➎ lancer l'import ;
- ➏ (regarder le log de l'import) ;
- ➐ une fois l'import réussi, jeter un œil à la description, à l'édition, etc.

Plan

Exploiter les
données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

1 Mise en jambe : *Andromaque*

2 Importer des données et créer un corpus

3 Interroger les données

4 Quelques notions d'analyse quantitative

Plan

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

1 Mise en jambe : *Andromaque*

2 Importer des données et créer un corpus

3 Interroger les données

4 Quelques notions d'analyse quantitative

Approches principales : la forme (linguistique) ou le fond (sémantique) ?

Exploiter les données

Jean-Baptiste
Camps &
Simon Gabay

Mise en
jambe :
Andromaque

Importer des
données et
créer un
corpus

Interroger les
données

Quelques
notions
d'analyse
quantitative

Sans prétention à l'exhaustivité, on peut distinguer deux approches principales qui font emploi des méthodes que nous allons présenter :

- **approche stylométrique**, qui se préoccupe d'attribution, datation, localisation des textes ; pour ce type d'approche, on va s'intéresser à l'information graphique (variantes d'orthographe), flexionnelle (terminaisons verbales, ...), etc. On va aussi tendre à s'intéresser seulement aux **mots les plus fréquents** (mots-outils, mots-vides), les moins sensibles aux variations intentionnelles de leurs auteurs (genre, sujet, etc.). Le **sens des textes nous indiffère** (ou presque).
- **approche sémantique**, lexicométrique, etc., qui est à peu près l'inverse de la précédente : intérêt pour les mots en tant que lemmes et leurs rapports entre eux (cooccurrences), pour les réseaux de sens, les thèmes des

Bibliographie

Exploiter les
données

Jean-Baptiste
Camps &
Simon Gabay



Heiden, S., Magué, J-P., et Pincemin, B., « TXM : Une plateforme logicielle open-source pour la textométrie – conception et développement », dans *Proc. of 10th International Conference on the Statistical Analysis of Textual Data - JADT 2010*, éd. Sergio Bolasco, Isabella Chiari, Luca Giuliano, Rome, 2010, t. 2, p. 1021-1032, <https://halshs.archives-ouvertes.fr/halshs-00549779/fr/>.



Pincemin, Bénédicte, Céline Guillot, Serge Heiden, Alexei Lavrentiev, et Christiane Marchello-Nizia, « Usages linguistiques de la textométrie : analyse qualitative de la consultation de la Base de Français Médiéval via le logiciel Weblex », *Syntaxe et Sémantique*, 9 (2008), p. 87–110, <https://halshs.archives-ouvertes.fr/halshs-00355461>.