

# Introduction à l'analyse de données textuelles

*Cours-bloc Statistique textuelle et topic models, 8-9.11.17*

Aris Xanthos

Section des sciences du langage et de l'information (SLI), Unil

# Matières

## 1. Bases conceptuelles

1. Analyse de données, modélisation, tableau de données
2. Questionnements uni-, bi- et multivariés
3. Matrice documents-termes, segmentations, annotations

## 2. Annotations

1. Définition, exemples
2. Bases du format XML

## 3. Travaux pratiques

1. Import XML et création de matrices documents-termes
2. Expressions régulières pour l'extraction de données

# 1. Bases conceptuelles

Introduction à l'analyse de données textuelles

# Le monde mis en données

- L'analyse de données implique de construire une représentation d'un phénomène:
  - Quelle est l'étendue du phénomène représenté?
  - Quels sont les objets (individus, observations, ...) représentés?
  - Par quelles propriétés (variables, caractéristiques, attributs, propriétés, métadonnées, ...) sont-ils représentés?
- Données = objets + propriétés
- Généralement organisées sous forme de tableau (lignes = objets et colonnes = propriétés)

# Exemples

- Bande dessinée

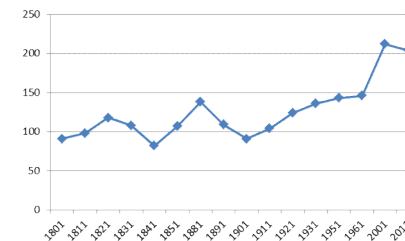
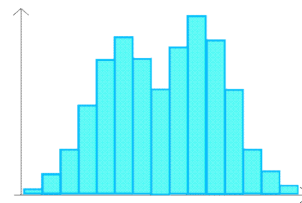
planche	titre	annee	nb_cases	haut_moy_cases	larg_moy_cases	...
1	<i>BD1</i>	1978	8	5.3	7.2	...
2	<i>BD1</i>	1978	5	5.5	9.2	...
...	...	...	...	...	...	...

- Théâtre

acte	piece	auteur	long_moy_phrases	prop_adv	...
1	<i>L'Avare</i>	<i>Molière</i>	12.3	5.6	...
2	<i>L'Avare</i>	<i>Molière</i>	11.6	3.9	...
...	...	...	...	...	...

# Questionnements uni- et bivariés

- Comment résumer l'information contenue dans une colonne de la table?
  - tendance centrale (moyenne, médiane, ...)
  - dispersion (variance, écart-type, ...)
- Quelle relation existe entre deux colonnes?
- Comment visualiser ces informations?
  - histogrammes
  - séries temporelles
  - etc.



# Questionnements multivariés

- Dans quelle mesure les profils des objets (i.e. les lignes du tableau) sont-ils similaires?
- Peut-on former des catégories d'objets sur la base de leur similarité?
- Comment peut-on visualiser cette structure de similarité?

# Matrice documents-termes (MDT)

	<i>Terme1</i>	<i>Terme2</i>	<i>Terme3</i>	...
<i>Document1</i>	Nb. d'occurrences de <i>Terme1</i> dans <i>Document1</i>	Nb. d'occurrences de <i>Terme2</i> dans <i>Document1</i>	...	...
<i>Document2</i>	Nb. d'occurrences de <i>Terme1</i> dans <i>Document2</i>	...	...	...
<i>Document3</i>	...	...	...	...
...	...	...	...	...

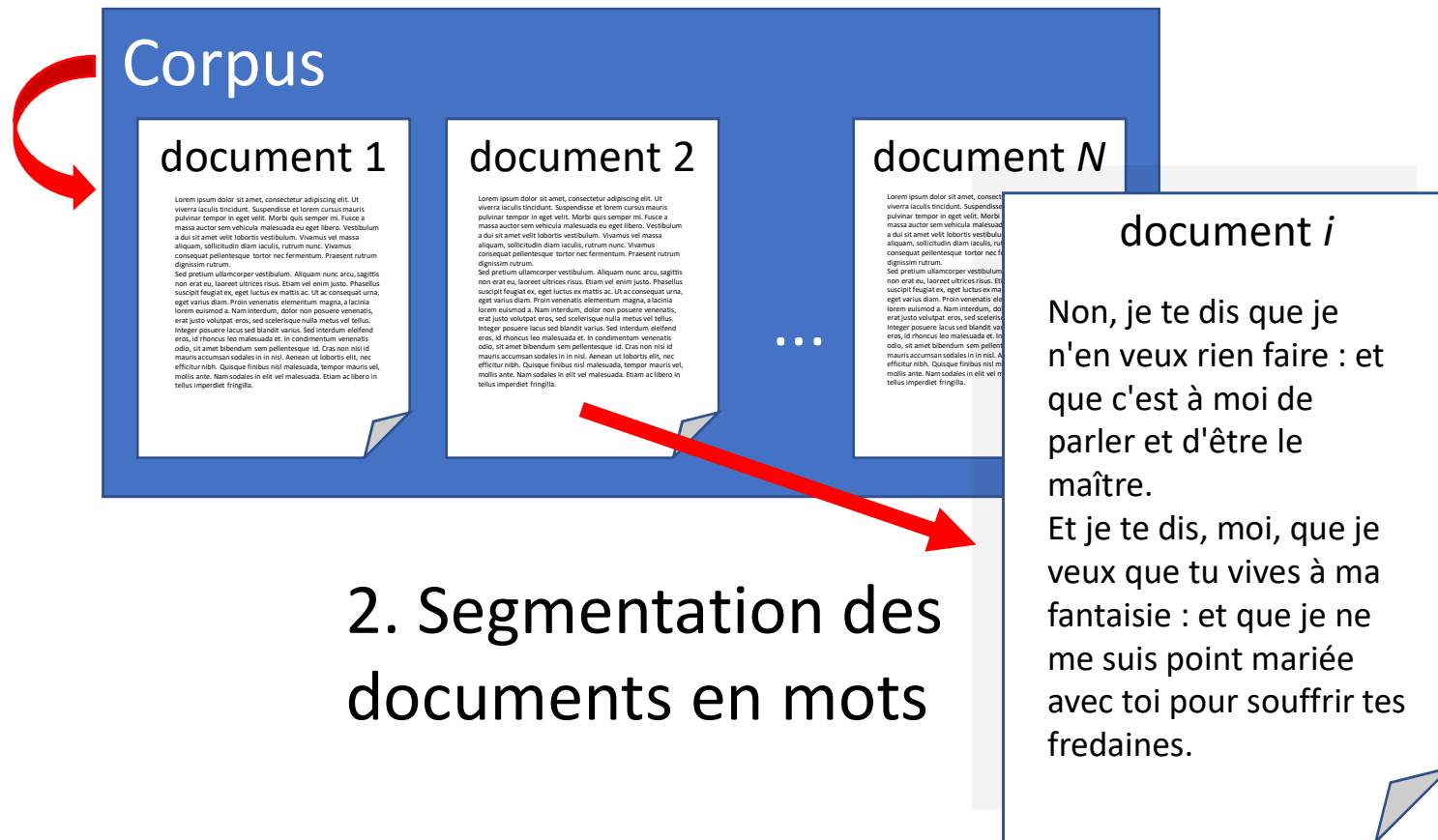


# Un exemple de MDT

	<i>A</i>	<i>Abandonnez</i>	<i>Acaste</i>	<i>Acceptez</i>	<i>...</i>	<i>total</i>
<i>L'Avare</i>	4	1	0	0	...	20 870
<i>L'Ecole des femmes</i>	4	0	0	0	...	16 494
<i>Le Médecin malgré lui</i>	0	0	0	0	...	9 219
<i>Le Misanthrope</i>	2	0	2	1	...	16 988
<i>total</i>	10	1	2	1	...	63 571

# Données doublement segmentées

## 1. Segmentation du corpus en documents



## 2. Segmentation des documents en mots

# Contextes et unités

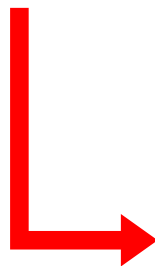
- Contextes (contenants):
  - typiquement documents, mais aussi groupe de documents ou parties de documents
- Unités (contenus):
  - typiquement mots, mais aussi séquences de mots, parties de mots, lettres individuelles, etc...
- Exemples:
  - fréquence des désinences verbales dans les pièces groupées par genre
  - fréquence des syntagmes nominaux dans les actes
  - etc.

# Contenu textuel vs annotations

- Identité des segments traditionnellement basée sur leur contenu textuel, p.ex.:
  - une seule unité *le* dans *le bateau arrive* et *il le faut*
  - *vais* et *irai* considérés comme deux unités sans relation particulière
- Il est possible et parfois nécessaire de fonder l'identité des segments sur des *annotations* qui leur sont associées, p.ex.:
  - deux unités distinctes *le/art* vs. *le/pro*
  - un seul lemme (lexème/vocable) *ALLER* derrière la flexion

# Une MDT peut en cacher une autre

	<i>A</i>	<i>Abandonnez</i>	<i>Acaste</i>	<i>Acceptez</i>	<i>...</i>
<i>L'Avare</i>	4	1	0	0	<i>...</i>
<i>L'Ecole des femmes</i>	4	0	0	0	<i>...</i>
<i>Le Médecin malgré lui</i>	0	0	0	0	<i>...</i>
<i>Le Misanthrope</i>	2	0	2	1	<i>...</i>

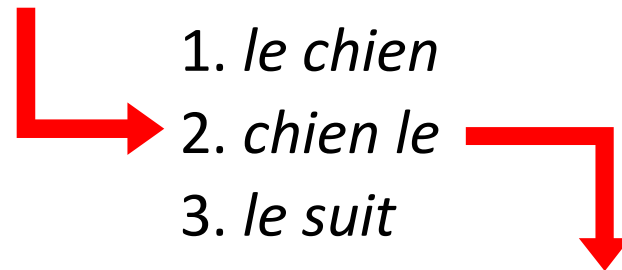


	<i>Consonne</i>	<i>Voyelle</i>
<i>Prose</i>	68 206	54 748
<i>Vers</i>	74 123	55 737

# Cas particulier: fenêtre coulissante

- Construction de MDT à partir d'une seule segmentation, p.ex.:

*le chien le suit*, fenêtre de 2 mots



	<i>le</i>	<i>chien</i>	<i>suit</i>
<i>1</i>	1	1	0
<i>2</i>	1	1	0
<i>3</i>	1	0	1

# Formats de MDT

	<i>t1</i>	<i>t2</i>	<i>t3</i>
<i>d1</i>	1	2	0
<i>d2</i>	0	1	1

tabulaire

<i>document</i>	<i>terme</i>
<i>d1</i>	<i>t1</i>
<i>d1</i>	<i>t2</i>
<i>d1</i>	<i>t2</i>
<i>d2</i>	<i>t2</i>
<i>d2</i>	<i>t3</i>

creux

<i>document</i>	<i>terme</i>	<i>occurrences</i>
<i>d1</i>	<i>t1</i>	1
<i>d1</i>	<i>t2</i>	2
<i>d2</i>	<i>t2</i>	1
<i>d2</i>	<i>t3</i>	1

creux pondéré

# MDT, et après?

- Questionnements typiquement multivariés:
  - peut-on identifier des catégories de contextes ou d'unités?
    - clustering, topic models
  - peut-on visualiser les similarités entre contextes ou entre unités, ainsi que les attractions/répulsions entre contextes et unités?
    - analyse factorielle des correspondances
- Ici: mise en œuvre avec les logiciels Iramuteq et R



# 2. Annotations

Introduction à l'analyse de données textuelles

# Définition

« L'annotation [...] est la pratique consistant à ajouter des informations *interprétatives* [...] à un corpus existant [...], par le biais d'une forme de marquage associé [...] à la *représentation* électronique des données »

Geoffrey Leech (1993). Corpus Annotation Schemes. *Literary and Linguistic Computing*, 8(4): 275-281 (ma traduction)

# Exemple

Catégories morpho-syntaxiques et lemmes:

<b>Nous</b>	<b>PRO:PER</b>	<b>nous</b>
<b>misons</b>	<b>VER:pres</b>	<b>miser</b>
<b>sur</b>	<b>PRP</b>	<b>sur</b>
<b>les</b>	<b>DET:ART</b>	<b>le</b>
<b>transports</b>	<b>NOM</b>	<b>transport</b>
<b>publics</b>	<b>ADJ</b>	<b>public</b>
<b>et</b>	<b>KON</b>	<b>et</b>
<b>sur</b>	<b>PRP</b>	<b>sur</b>
<b>...</b>		

# Exemple (2)

Annotation selon la convention [EmotionML](#)  
(2 émotions catégorisées et quantifiées):

```
<emotion>  
  <category name="Disgust" value="0.82"/>  
  'Come, there's no use in crying like that!'  
</emotion>  
  said Alice to herself rather sharply;  
<emotion>  
  <category name="Anger" value="0.57"/>  
  'I advise you to leave off this minute!'  
</emotion>
```

# Formats d'annotation

- Formats ad hoc vs standards
- XML (eXtensible Markup Language) :
  - norme d'encodage des données textuelles
  - standard
  - extensibilité
  - intégration (navigateurs, lang. de prog.)
  - ensemble de technologies (XSLT, XPath, ...)

# XML: principes de base

- Balises :

```
<texte>un exemple simple</texte>
```

```
<fin_de_page></fin_de_page> → <fin_de_page/>
```

- Emboîtement (strict) :

```
<texte>
```

```
  <mot>un</mot>
```

```
  <mot>exemple</mot>
```

```
</texte>
```

- Attributs :

```
<mot lemme="exemple" cms="NOM">exemple</mot>
```

# XML: document complet

- Déclaration XML
- Un et un seul élément racine (ici `<texte>`)

```
<?xml version="1.0" encoding="utf-8"?>
<texte>
  <mot type="art">un</mot>
  <mot type="nom">exemple</mot>
  <mot type="adj">simple</mot>
  et un fragment non annoté...
</texte>
```

# 3. Travaux pratiques

Introduction à l'analyse de données textuelles