# INTRODUCTION GÉNÉRALE SUR LES ENJEUX DU *TEXT ET DATA MINING*

PATRICE BELLOT

**8 OCTOBRE 2019**

La Science Ouverte : une révolution nécessaire
8 oct. 2019 Paris (France)

**CNRS – INS2I**

# 1
# QU'EST-CE QUE LE TDM ?

## Domaines scientifiques et difficultés

**1**

# QU'EST-CE QUE LA FOUILLE DE DONNÉES / DE TEXTES ?

**Le croisement de plusieurs domaines**

- L'analyse de données automatisée

- L'Intelligence Artificielle

- Le Traitement Automatique des Langues

**De la donnée vers l'information vers la connaissance**

**Données ?**
**Elles peuvent des nombres, des signaux, des mots, des images…**
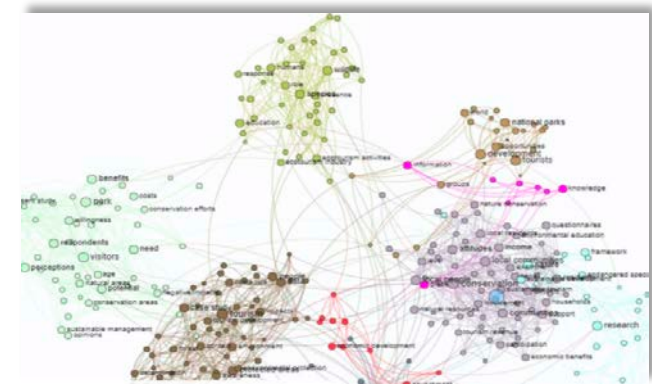**Elles peuvent être structurées ou non, liées ou non…**

**2**

## DES APPLICATIONS

- Recherche de documents, d'images, de pages Web

- Recommandation automatique de contenus, veille

- Systèmes de questions-réponses

- Résumé automatique

- Cartographie et navigation guidées

- Détection de nouveauté, analyse de sentiments, de tendances…



https://isidore.science



Gargantext

# 3 DES TRAITEMENTS GÉNÉRIQUES

- Extraction de termes, d'entités, d'informations, de relations

- Sélection de descripteurs, filtrage de données

- Calcul de similarités sémantiques

- Classification ou catégorisation automatiques

- Recherche d'information à partir de requêtes

- Structuration et segmentation

- Reconnaissance de formes, de caractères, d'images, de la parole…

# 4 QUELQUES DIFFICULTÉS...

**Quelle que soit la nature des données :**

- Structures peu normalisées, formats variés

- Les fameux V du Big Data : Volume, véracité, variabilité, valeur, vitesse

**Document, texte et langage :**

- Données hétérogènes ou multimodales, ambiguës

- Multilinguisme (lexiques, terminologies, syntaxes)

Des difficultés génériques.

Des standards nécessaires.

Des solutions à partager.

**Le Droit et les bonnes pratiques pour la mise en œuvre du TDM :**

Briefing

Requested by the JURI committee

European Parliament

The Exception for Text and Data Mining (TDM) in the Proposed Directive on Copyright in the Digital Single Market - Technical Aspects

*loi française Pour une République Numérique du 7 octobre 2016 et la directive européenne sur le droit d'auteur dans le marché unique numérique du 26 mars 2019*

**4**

# STRUCTURE LOGIQUE ET CODAGE

# ALLER AU-DELÀ DES MOTS CLEFS : PAS TOUJOURS FACILE…

# CONTEXTES ET CONNAISSANCES

**Indexation par mots, par terme, par domaine...**

# 4

## MULTILINGUISME, CITATIONS, FORMULES, RESULTATS…

OpenEdition : OpenEdition Books    OpenEdition Journals    Calenda    Hypothèses    Lettre    OpenEdition Freemium

(i)    DOI / Références    Télécharger

Accueil > Numéros > 21 > DOSSIER THÉMATIQUE Strangers at … > ἵν'ἀμνήμων τύχη γένοιτο πολλῶν δε…

γ**GAÎA**    Revue interdisciplinaire sur la Grèce archaïque

Recherche →

Index
Auteurs
Mots-clés

Numéros en texte intégral
21 | 2018

Tous les numéros →

La revue
Présentation
Organisation scientifique
Recommandations aux auteurs

**21 | 2018**
Varia

DOSSIER THÉMATIQUE
Strangers at Home. Civilizing Immigrants between Inclusion and Exclusion in Ancient Thebes

### Ἵν'ἀμνήμων τύχη γένοιτο πολλῶν δεομένη σοφισμάτων (*Phéniciennes*, 64-65). La souveraineté brisée de la famille d'Œdipe et la crise de la parole dans le mythe tragique des *Phéniciennes*

Ἵν'ἀμνήμων τύχη γένοιτο πολλῶν δεομένη σοφισμάτων (Phoenician Women, 64–65). The Broken Sovereignty of Oedipus' Family and the Crisis of Speech in the Tragic Myth of Phoenician Women
Ἵν'ἀμνήμων τύχη γένοιτο πολλῶν δεομένη σοφισμάτων (Fenicie, 64-65). La sovranità spezzata della famiglia di Edipo e la crisi della parola nel mito tragico delle Fenicie

AVEZZÙ Guido, *Il mito sulla scena. La tragedia ad Atene*, Venise, Marsilio, 2003.

AVEZZÙ Guido, « Emulazione e antagonismo nella produzione tragica ateniese », *Dionysus ex machina*, 6 (1), 2015, p. 137-156.

BATTEZZATO Luigi, « An Introduction to Tragedy. G. Avezzù: *Il mito sulla scena. La tragedia ad Atene* », *The Classical Review*, 55 (1), 2005, p. 29.
DOI : 10.1093/clrevj/bni019

BEARZOT Cinzia, « Perdonare il traditore? La tematica amnistiale nel dibattito sul richiamo di Alcibiade », dans M. Sordi (éd.), *Amnistia, perdono e vendetta nel mondo antico*, Milan, Vita e Pensiero, 1997, p. 29-52.

BELTRAMETTI Anna, « Antigone o della questione morale. Elaborazione tragica della sovranità democratica », dans D. Ambaglio (éd.), *«Συγγραφή». Materiali e appunti per lo studio della storia e della letteratura antica*, vol. 4, Como, Edizioni New Press, 2002, p. 33-49.

# 4

Table 6. The NDCG and Hit Ratio (HT) Results of AARM and Its Variants on Five Datasets for RQ3

| Measures (%) | Movies | | CDs | | Clothings | | Cell Phones | | Beauty | |
|---|---|---|---|---|---|---|---|---|---|---|
| | NDCG | HT | NDCG | HT | NDCG | HT | NDCG | HT | NDCG | HT |
| AARM | **5.020** | **15.187** | **7.252** | **20.749** | **1.957** | **4.915** | **4.976** | **11.568** | **5.314** | **13.648** |
| A_Static | 4.376 | 13.318 | 6.794 | 19.567 | 1.898 | 4.590 | 4.728 | 11.181 | 4.918 | 12.735 |
| No-UserAtt | 4.290 | 13.104 | 6.700 | 19.108 | 1.310 | 3.217 | 4.685 | 10.786 | 4.739 | 12.297 |
| Impr A_static | 14.717 | 14.034 | 6.741 | 6.041 | 3.109 | 7.081 | 5.245 | 3.461 | 8.052 | 7.169 |
| Impr No-UserAtt | 17.016 | 15.896 | 8.239 | 8.588 | 49.389 | 52.782 | 6.211 | 7.250 | 12.133 | 10.986 |

We follow the short form convention adopted in Table 4 to name the datasets. The best performance of each measure on each dataset is highlighted in bold. The last block shows the percentage of improvements (or decrements for negative values) achieved by AARM compared with A_static (Impr A_static) and No-UserAtt (Impr No-UserAtt).

Table 7. The Corresponding Precision and Recall Results of AARM
and Its Variants on Five Datasets for RQ3

Figure 3 shows an exemplary 6-step route for an intermediate of a drug candidate synthesis reported in 2015, which has been found by our algorithm in 5.4 s. It matches the published route.[45]

Figure 3: An e spect to product $v$, the interactions between $a_i$ and all the aspects i
published first d up:
autonomously v

$$\mathbf{x}_{v,i} = \mathbf{g}_v \odot \mathbf{c}_i,$$
$$\mathbf{g}_v = \sum_{j \in A_v} \mathbf{c}_j.$$

interaction between two aspects represents their similarity, $\mathbf{x}_{v,i}$ rep between the aspect $a_i$ and the product $v$. To measure the importa $\mathbf{x}_{v,i}$ is used as aspect $a_i$'s input to the user-level attention layer l as

$$\hat{\alpha}_{u,v,i} = \mathbf{w}_{att_2}{}^T \mathbf{x}_{v,i},$$
$$\alpha_{u,v,i} = \frac{\exp(\hat{\alpha}_{u,v,i})}{\sum_{j \in A_u} \exp(\hat{\alpha}_{u,v,j})}.$$

$_{att_2} \in \mathbb{R}^{d_a}$ is a learnable vector, and $\alpha_{u,v,i}$ represents the importan s preferences with regard to product $v$. This attention layer is differe

## 5 CONCLUSION AND FUTURE WORK

In this article, we presented an AARM, which carefully captures the interactions betv extracted from reviews for recommendation. AARM first calculates the interactions between pect embeddings to estimate how a product fits a user's requirements on each aspect, and then timates the user's overall satisfaction on the product by synthesizing the product's performan on each aspect. To deal with the problem that the number of shared aspects between a user an product is often limited, AARM takes the interactions between different aspects into considerati With a well-designed aspect-level attention module, not only the shared aspects but also other lated aspect pairs can be selected and assigned higher attention values. In addition, we hold assumption that a user's interests toward aspects are varied when examining different products. achieve the goal, an attention module which simultaneously considers user and product informa tion is designed in AARM. In the experiments on five real-world datasets, AARM outperforms the state-of-the-art methods on the top-N recommendation task. In particular, compared with multi-modal (textual reviews, product images, and numerical ratings) methods JRL and eJRL, AARM can still achieve better results in all datasets. To demonstrate the effectiveness of each component in AARM, a lot of quantitative experiments and qualitative case studies are conducted.

**2
L'EXISTANT**

De nombreux acteurs

# 1 DE NOMBREUX ACTEURS INDUSTRIELS
*(et de très nombreux laboratoires publics)*

# 2

## DES PROJETS OUVERTS

# DES LOGICIELS OUVERTS POUR L'IST : EXEMPLE *BILBO*

## 2



https://lab.hypotheses.org/category/bilbo-bibliographical-robot

**2**

## UNE INITIATIVE PUBLIQUE : ISTEX

https://www.istex.fr/

**Un moteur de recherche en texte intégral, une chaîne de traitements, des enrichissements, des APIs et des services**

**23 millions de documents**
**9 279 revues**





APERÇU DU GRAPHE DES JEUX DE DONNÉES



https://data.istex.fr

# 2 DES INITIATIVES POUR UN COUPLE IST / TDM OUVERT



https://www.openaire.eu

Home page | **Project mining** | Data citation mining | Document Classification | Software mining | Interactive project mining | Citation matching | Document similarity

### Project mining details

Provide your **UTF-8** encoded text, on the current URL using the HTTP POST method.
You may also choose among the available mining processes.

HTTP POST parameters:

- **document**: UTF-8 encoded text
- **projects**: Project processing (**on**/off)
- **datacitations**: Data citation processing (on/**off**)
- **classification**: Classification processing (on/**off**)

The service will return a **JSON** encoded result containing the following fields:

- **funding_info**: Result category
- **fund**: Funder name (e.g., *FP7*, *Wellcome Trust*)
- **acronym**: The project acronym (only for FP7 projects)
- **grantid**: The project grant identifier
- **confidence**: Confidence weight

Result example:

```
{"funding_info": [{"fund": "fp7",
"acronym": "CORONET", "grantid":
"269459", "confidence": 0.96}]}
```

Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint funded by fp7 project OpenAIREplus occaecat cupidatat non proident,fp7 project with grant agreement number 318338 sunt in culpa qui officia deserunt mollit anim id est laborum.

Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. In posuere felis nec tortor. Pellentesque faucibus. Ut accumsan ultricies elit. Maecenas at justo id velit placerat molestie. Donec dictum lectus non odio. Cras a ante vitae enim iaculis aliquam. Mauris nunc quam, venenatis nec, euismod sit amet, egestas placerat, est. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Cras id elit. Integer quis urna. Ut ante enim, dapibus malesuada, fringilla eu, condimentum quis, tellus. Aenean porttitor eros vel dolor. Donec convallis pede venenatis nibh. Duis quam. Nam eget lacus. Aliquam erat volutpat. Quisque dignissim congue leo.

Mauris vel lacus vitae felis vestibulum volutpat. Etiam est nunc, venenatis in, tristique eu, imperdiet ac, nisl. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. In iaculis facilisis massa. Etiam eu urna. Sed porta. Suspendisse quam leo, molestie sed, luctus quis, feugiat in, pede. Fusce tellus. Sed metus augue, convallis et, vehicula ut, pulvinar eu, ante. Integer orci tellus, tristique vitae, consequat nec, porta vel, lectus. Nulla sit amet diam. Duis non nunc. Nulla rhoncus dictum metus. Curabitur tristique mi condimentum orci. Phasellus pellentesque aliquam enim. Proin dui lectus, cursus eu, mattis laoreet, viverra sit amet, quam. Curabitur vel dolor ultrices ipsum dictum tristique. Praesent vitae lacus. Ut velit enim, vestibulum non, fermentum nec, hendrerit quis, leo. Pellentesque rutrum malesuada neque.

References:
Sieger, Rainer; (2012): PanGet - downloads multiple data sets from PANGAEA; PANGAEA - Data Publisher for Earth & Environmental Science.
Grobe, Hannes; (2005): Description and user manual of the information system PANGAEA; PANGAEA - Data Publisher for Earth & Environmental Science. http://dx.doi.org/10.1594/PANGAEA.319947

Appendix:
Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. In posuere felis nec tortor. Pellentesque faucibus. Ut accumsan ultricies elit. Maecenas at justo id velit placerat molestie. Donec dictum lectus non odio. Cras a ante vitae enim iaculis aliquam. Mauris nunc quam, venenatis nec, euismod sit amet, egestas placerat, est. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Cras id elit. Integer quis urna. Ut ante enim, dapibus malesuada, fringilla eu, condimentum quis, tellus. Aenean porttitor eros vel dolor. Donec convallis pede venenatis nibh. Duis quam. Nam eget lacus. Aliquam erat volutpat. Quisque dignissim congue leo.

☑ Projects ☑ Data citations ☑ Classification

[ Submit ]   [ Insert example text ]

{"funding_info": [{"fund": "EC", "acronym": "OPTIQUE", "grantid": "318338", "confidence": 0.92}], "classification_info": [], "datacitation_info": [{"organization": "datacite", "related_doi": "10.1594/PANGAEA.319947", "confidence": 1, "resource_type": "Dataset"}, {"organization": "datacite", "related_doi": "10.1594/PANGAEA.804562", "confidence": 0.69, "resource_type": "Dataset"}]}

# DES INITIATIVES EUROPEENNES : L'ERIC *CLARIN*

CLARIN

About ▾   Participants   Services   Knowledge Base ▾   Events   News   Contact        Applications   Intranet login

## CLARIN - European Research Infrastructure for Language Resources and Technology

CLARIN makes digital language resources available to scholars, researchers, students and citizen-scientists from all disciplines, especially in the humanities and social sciences, through single sign-on access. CLARIN offers long-term solutions and technology services for deploying, connecting, analyzing and sustaining digital language data and tools. CLARIN supports scholars who want to engage in cutting edge data-driven research, contributing to a truly multilingual European Research Area. Read more...

**CLARIN**

Common Language Resources and Technology Infrastructure

**Search** 🔍

**CLARIN Newsflash September 2019 is out!**

**CLARIN NEWSFLASH**

Read the most recent CLARIN Newsflash: September 2019 here

### Participating Organizations

The operations, services and centres of the CLARIN infrastructure are provided and funded by the national consortia in the countries that have joined CLARIN ERIC or by associated centres.

Learn more

### Services, Tools and Data

CLARIN provides a wide variety of Services and Data sets

Learn more

● ● ● ● ●

**ANNUAL CONFERENCE 2019**
Leipzig, Germany

Conference programme

Twitter: #CLARIN2019

**HN Huma-Num**
la TGIR des humanités numériques

| Members | National Consortium (NC) | Leading NC partner |
|---|---|---|
| Austria | Digital Humanities Austria | ACDH-OEAW |
| Bulgaria | CLaDA-BG | Bulgarian Academy of Sciences |
| Croatia | HR-CLARIN | University of Zagreb |
| Cyprus | CLARIN-CY | Digital Heritage Research Lab (Cyprus University of Technology) |
| Czech Republic | LINDAT/CLARIN | Charles University Prague |
| Denmark | CLARIN-DK | University of Copenhagen |
| Estonia | CLARIN Estonia | Center of Estonian Language Resources |
| Finland | FIN-CLARIN | University of Helsinki |
| Germany | CLARIN-D | University of Tuebingen |
| Greece | clarin:el | ILSP-ATHENA Research Center |
| Hungary | HunCLARIN | Research Institute for Linguistics, Hungarian Academy of Sciences |
| Italy | CLARIN-IT | Institute for Computational Linguistics A. Zampolli, Italian National Research Council |
| Latvia | CLARIN-LV | Institute of Mathematics and Computer Science, University of Latvia |
| Lithuania | CLARIN-LT | Vytautas Magnus University |
| The Netherlands | CLARIAH-NL | Utrecht University |
| Norway | CLARINO | University of Bergen |
| Poland | CLARIN PL | Wroclaw University of Technology |
| Portugal | PORTULAN CLARIN | University of Lisbon |
| Slovenia | CLARIN.SI | Jožef Stefan Institute |
| Sweden | SWE-CLARIN | Språkbanken |
| **Observer** | **National Consortium (NC)** | **Leading partner NC** |
| France | Huma-Num | the National Center for Scientific Research (CNRS) |
| Iceland | CLARIN Iceland | The Árni Magnússon Institute for Icelandic Studies |
| South Africa | SADiLaR | North-West University |
| United Kingdom | CLARIN-UK | Oxford University |

https://www.clarin.eu

# DES INITIATIVES EUROPÉENNES : L'INFRASTRUCTURE *OPENMINTED* 2015-2018

**| Resource Type**

☑ Applications (41)

**| Refine**

**Licence**
☐ Affero General Public License v1.0 (1)
☐ BSD-2-Clause (Simplified) License (1)
☐ Creative Commons Attribution Share Alike 3.0 Unported (1)
☐ GNU Lesser General Public License v3.0 (1)
☐ Non standard Licence or Terms of use (2)
View more...

**Rights Statement**
☐ Open Access (41)

**Language**
☐ Bulgarian (1)
☐ Catalan; Valencian (1)
☐ Czech (1)
☐ Danish (1)
☐ German (1)
☐ Spanish; Castilian (1)
☐ Irish (1)
☐ Galician (1)
☐ Croatian (1)
☐ Latvian (1)
View more...

**Function**
☐ Extraction of funding information

---

Showing 1 - 10 of 41 results

‹ Previous
PAGE 1 OF 5

### GeoPolitical Extractor App

Extracts geopolitical terms.

### IXA pipes for Basque for PDF files

IXA pipes for Basque with PDF reader. eu-ixa-pipes-omtd prov
tokenizer, POS tagger, lemmatizer, NER tagger, Chunking and D
Classification. It reads from an input folder containing XMI doc
and outputs the added annotations in XMI format to an outpu
directory.

### Leica Model Annotation App

The application annotates a corpus with Leica Microsystems p

### MADIS FUNDING MINING

The Funding Mining application mines the fulltext of publicatic
extracts links to projects. Currently, projects from EC (FP7/H2(
(National Science Foundation, USA), NIH (National Institute of
USA), Wellcome Trust, FCT (Fundação para a Ciência e a Tecno
Portugal), ARC (Australian Research Council), NHMRC (National
and Medical Research Council, Australia), CSF/HRZZ (Hrvatska
Za Znanost, Croatia), MSES-MZOS (Ministarstvo Znanosti, Obra
športa, Croatia), SFI (Science foundation Ireland), NWO (Nederl
Organisatie voor Wetenschappelijk Onderzoek, Netherlands) ar
supported, but new funders are added regularly.

---

**Function**
☐ Analyzer (1)
☐ Annotator of semantic role labels (1)
☐ Chunker (1)
☐ Co-reference annotator (1)
☐ Constituency parser (1)
☐ Information extraction (1)
☐ Parser (1)
☐ Variables dectector (1)
☐ Document classifier (2)
View more...

**Component Distribution Medium**
☐ Executable Code (28)
☐ Docker Image (35)

**Processing Resource Type**
☐ Lexical Conceptual Resource (6)
☐ Document (28)
☐ Corpus (29)

**Data Format**
☐ Binary CAS (1)
☐ Binary format (1)
☐ CoNLL-U (1)
☐ CSV (1)
☐ GATE format (1)
☐ TSV (1)
☐ UIMA CAS format (1)
☐ ALVIS Enriched Document format (3)
☐ XML (4)
View more...

**Annotation Type**
☐ Citation (1)
☐ Constituent (1)
☐ Coreference (1)
☐ Dependency tree (1)
☐ Discourse annotation type (1)
☐ Event (1)
☐ Funding (1)
☐ Grape variety (1)

---

...Extractor

Recognizes phenotypes, genes, markers, and wheat-related taxa. It categorizes the phenotypes with the Wheat Trait Ontology.

### Alvis Arabidopsis Gene Regulation Extractor

Recognizes Gene, Protein and RNA of Arabidopsis thaliana. It normalizes them with Gene Locus and identifies interacts_with relationships between Gene and Protein.

### Lancaster Stemmer (DKPro Core)

This Paice/Husk Lancaster stemmer implementation only works with the English language so far.

### Variable Disambiguator

Assign variable IDs to sentences based on calculating the similarity between the sentence text and the description of the variable.

### ClearNLP Segmenter (DKPro Core)

Tokenizer using Clear NLP.

### de-ixa-pipes-omtd

IXA pipes for German. It provides tokenizer, POS tagger, lemmatizer and NER tagger. It reads from an input folder containing XMI documents and outputs the added annotations in XMI format to an output directory.

## CONCLUSION : LE TDM…

**Nécessite :**

- un corpus cible, des ressources de spécialité

- d'intégrer différents composants logiciels

- un scénario et une référence pour évaluer la chaîne de traitements

**Faisable si :**

- les composants sont interopérables, les métadonnées compatibles

- une infrastructure ouverte de services de fouille de textes est disponible…

- l'exception TDM sur la directive du droit d'auteur / copyright est soutenue

**Concerne et impacte :**

- La recherche scientifique dans son ensemble

- La société au travers d'applications du quotidien



https://www.jisc.ac.uk/reports/value-and-benefits-of-text-mining

# MERCI DE VOTRE ATTENTION



## Visa TM Day le 15 novembre : vers une infrastructure de services avancés en text-mining

Dans le cadre du projet Visa TM du Comité pour la Science Ouverte, un « **Visa TM Day** » sera organisé **vendredi 15 novembre 2019** au ministère de l'Enseignement supérieur, de la Recherche et de l'Innovation à **Paris**.

Le projet Visa TM a pour objectif l'étude d'une e-infrastructure de recherche pour la création d'une offre de service en fouille de textes pour la recherche, basée sur l'analyse sémantique et s'appuyant sur le potentiel de combinaison et d'adaptation offert par la plateforme européenne OpenMinTeD.

Autour de conférences et d'ateliers prospectifs, cette journée est destinée à dresser un état des lieux et discuter des perspectives concrètes ouvertes par les résultats du projet.

**Les inscriptions sont ouvertes jusqu'au 15 octobre.**

Programme et inscriptions sur https://journees.inra.fr/visa-tm-day/