

BETWEEN HISTORY AND NATURAL LANGUAGE PROCESSING : STUDY, ENRICHMENT AND ONLINE PUBLICATION OF FRENCH PARLIAMENTARY DEBATES OF THE EARLY THIRD REPUBLIC (1881–1899)

The AGODA project

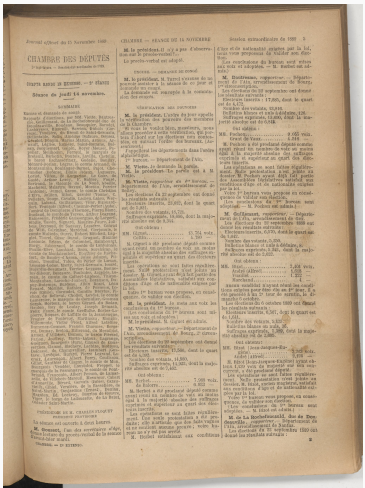
Marie Puren, Aurélien Pellet, Nicolas Bourgeois, Pierre Vernus, Fanny Lebreton

20 June 2022

ParlaCLARIN III at LREC2022 - Workshop on Creating, Enriching and Using Parliamentary Corpora, Marseille

- AGODA : **A**nalyse sémantique et **G**raphes relationnels pour l'**O**uverture et l'étude des **D**ébats à l'**A**ssemblée nationale
- Funded by the Bibliothèque nationale de France for a period of one year
- One of the 5 pilot projects supported by the **DataLab**
- Collaboration between Epitech (MNSHS), Inria (ALMAAnaCH) and the University Lumière Lyon 2 (LARHRA)

PARLIAMENTARY DEBATES DURING THE THIRD REPUBLIC



- Debates in the lower house of French Parliament (“Chambre des députés”) carefully recorded in the Journal officiel de la République française. Débats parlementaires (1881-1940)
- Available online on the digital library of the National Library of France (Gallica)
- Still difficult to work on this corpus : better if you already know what you’re looking for!

Figure – Parliamentary sitting - 14 Nov. 1889

Working on a sub-part of the corpus : legislature 1889-1893, i.e. **10418 images to be processed**

- Give easier access to ancient transcriptions of parliamentary debates
- Facilitate research in this corpus
- Offer new ways of visualizing the documents



Figure – Workflow

OCERISE THE DEBATES

THE CASE OF PARLIAMENTARY DEBATES

- Retrieval of ocerised texts via Gallica's **API Document** => uneven quality of the OCR
- Errors due to :
 - Document quality : stains and overprinting
 - **Curvature of the page** at the level of the binding

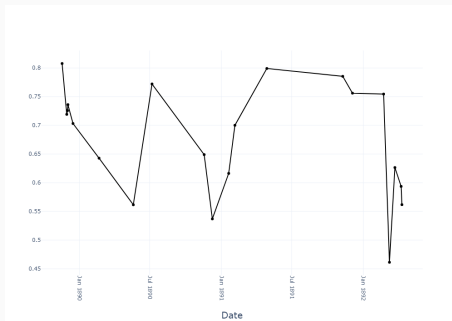


Figure – Quality assessment of the OCR provided by Gallica

EFFECT OF CURVATURE ON OCR RESULTS

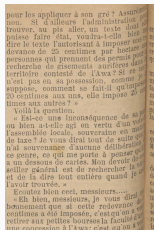


Figure – 14
May 1890
(p.786)

pour les appliquer à son gré ? Assure.

non. Si d'ailleurs l'administration fj, trouver, au pis aller, un texte (1011t en puisse faire état, voudra-t-elle dire le texte l'autorisant à imposer u devance de 25 centimes par hectar , personnes qui prennent des permis Pourl recherche de gisements aurifères a teSI territoire contesté de l'Awa ? Si ce tes

n'est pas en sa possession, cornas le suppose, comment se fait-il quilIWjs3j 20 centimes aux uns, elle impose 20 cet times aux autres ? » , Voilà la question. é

« Est-ce une inconséquence de sa te ou bien a-t-elle agi en vertu d'un v l'assemblée locale, souveraine en OlqallC de taxe? Je vous dirai tout de suite Il n'ai souverance d'aucune délibératl"lIV ce genre, ce qui me porte à penser à un dessous de cartes. Mon devoir seiller général est de rechercher la cJ"l" et de la dire tout entière quand Je l'avoir trouvée. »

Eoutez bien ceci, messieurs. l «Eh bien, messieurs, je vous dira de bonnement que si cette redevance J centimes a été imposée, c'est qu'on a, retirer aux petites bourses la faculté" une concession à l'Awa; c'est qu'on a rj;:v

Figure – OCR result (Gallica)

We took decision to re-ocerise the corpus.

1. Incorrect segmentation
2. Poor OCR performance
3. Very curved and even cut text : not very readable by a human being, even less readable by a machine
4. Automatic annotation in XML-TEI more difficult because linked to segmentation
5. Text too faulty to publish online
6. Negative influence of OCR errors on the performance of NLP analyses (NER, topic modeling, word embedding) (cf. [Assessing the Impact of OCR Quality on Downstream NLP Tasks](#))

WHAT SOLUTIONS?

- Dewarping solution : unsatisfactory as it can bend the flat parts and even cut the image
 - Consider developing a solution that removes the curvature only for the parts of the image concerned (column at the binding edge)
- Use a language model : when the text is not human readable and therefore not machine readable, use a language model to predict the missing word(s)
 - Language model for French : trained on (too much ?) contemporary French - perhaps need to re-train on late 19th century French
- Words cut because of the curvature + words ocerised but appearing in disorder (consequence of the bad segmentation due to the curvature of the page) : use an AI to restore the text, as Ithaca of Deep Mind does (cf. [Restoring and attributing ancient texts using deep neural networks](#))

IMAGE CLEANING TOOL DEVELOPED BY ANR SODICO



(a) Original image



(b) Cleaned image

Figure – Demonstration of the SODICO tool on a page

OCR TOOL DEVELOPED BY SODUCO

Tool based on the OCR engine **PERO OCR** : very efficient on historical texts

Developed within the framework of the ANR **SODUCO**

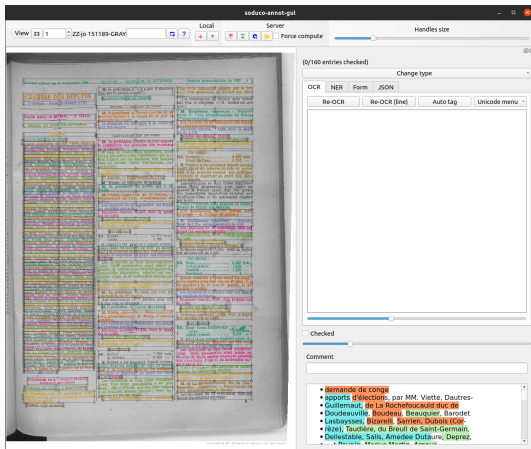
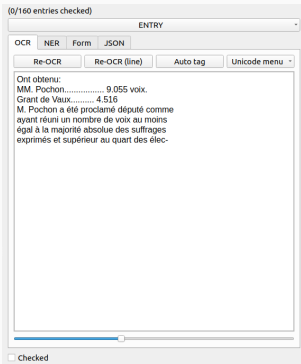


Figure – OCR tool developed by SODUCO



(a) OCR



(b) NER

Figure – OCR and NER zones

OUTPUT IN JSON

```
"box": [
  827.8018264183805,
  2374.55023568486,
  589.3638270234404,
  207.5897599646105
],
"checked": true,
"comment": "u-par",
"id": 352,
"ner_xml": "<PER>M. Francis Laur</PER>. Je n'ai pas l'intention\u2029d'apporter \u00e0 la Chambre d'autr\u2029s de",
"origin": "computer",
"parent": 274,
"persons": ["M. Francis Laur"],
"text_ocr": "M. Francis Laur. Je n'ai pas l'intention\u2029d'apporter \u00e0 la Chambre d'autres documents\u2029que ceux qui me sont fournis par la lettre\u2029de d\u00e9mission de M. Lev\u00eaque, adress\u00e9e \u00e0 M. le ministre des finances. Je demanderai\u2029la Chambre si elle conna\u00et cette lettre ou\u2029si elle veut que j'en donne lecture.",
"type": "ENTRY"
},
{
  "activities": [],
  "addresses": [],
```

Figure – Output in JSON (extract)

- Post-correction dictionary
- Use of regular expressions : manage multiple spaces, line breaks, “-”
- Endogenous corrections
- CamemBERT language model to correct faulty words (see above)
- Automatic translation ? “Translate” wrong words into the correct form

THE PRINCIPLES OF THE ENCODING IN XML-TEI

Encoding is designed according to 4 principles :

- The different uses of texts
- The particularities of the source
- Similar projects : ParlaClarín and ParlaMint
- The automatic tagging process

For more information on this aspect :

<https://github.com/mpuren/agoda/tree/ODD>

GENERAL STRUCTURE OF TEI FILES

```
<teiCorpus xmlns="http://www.tei-c.org/ns/1.0" xml:id="FR_3R_5L" xml:lang="fr">

  <!-- Métadonnées du corpus (non développées) -->
  <teiHeader>
    <fileDesc>
      <titleStmnt>
        <title></title>
      </titleStmnt>
      <publicationStmnt>
        <p></p>
      </publicationStmnt>
      <sourceDesc>
        <p></p>
      </sourceDesc>
    </fileDesc>
  </teiHeader>

  <!-- Données liées et informations contextuelles -->
  <standoff>
    <listPerson>
      <person></person>
    </listPerson>

    <listOrg>
      <org></org>
    </listOrg>

    <listPlace>
      <place></place>
    </listPlace>
  </standoff>

  <!-- Stockage du composant correspondant à la séance du 26 novembre 1889 -->
  <include xmlns:x1="http://www.w3.org/2001/XInclude" href="FR_3R_5L_1889-11-26.xml"/>

  <!-- Stockage des autres composants du corpus de façon identique -->
  <!-- ... -->
</teiCorpus>
```

(a) General structure of a corpus file

```
<TEI xmlns="http://www.tei-c.org/ns/1.0" xml:id="FR_3R_5L_1889-11-26" xml:lang="fr">

  <!-- Métadonnées du composant (non développées) -->
  <teiHeader>
    <fileDesc>
      <titleStmnt>
        <title></title>
      </titleStmnt>
      <publicationStmnt>
        <p></p>
      </publicationStmnt>
      <sourceDesc>
        <p></p>
      </sourceDesc>
    </fileDesc>
  </teiHeader>

  <text>

    <!-- Transcription du compte rendu -->
    <body>
      <div></div>
    </body>

    <!-- Annexes du compte rendu -->
    <back></back>

  </text>
</TEI>
```

(b) General structure of a component file

Figure – Encodings - General structure of corpus and component files

EXPERIMENTING WITH DIFFERENT SOLUTIONS

```
<lb/><u who="#pers_ID" xml:id="CR_1889-11-26_u5" ana="#speaker">
  <seg xml:id="CR_1889-11-26_u5.1">
    <persName ref="#pers_ID">M. Paul Déroulède</persName>. Mes amis et moi
    <lb/>avons applaudi à ce que nous avons cru, à
    <lb/>ce que nous croyons encore être la forma-
    <lb/>tion spontanée d'une majorité de con-
    <lb/>science, de justice, de tolérance et d'union.
    <lb/><incident><desc>(Applaudissements sur quelques bancs à l'ex-
    <!-- Pas de lb possible dans incident --> trémité gauche de La salle. – Exclamations
    <!-- Pas de lb possible dans incident --> au centre et à gauche.)</desc></incident>
  </seg>
</u>

<lb/><u who="pers_ID" xml:id="CR_1889-11-26_u6" ana="#speaker">
  <seg xml:id="CR_1889-11-26_u6.1">
    <persName ref="#pers_ID">Gustave Rivet</persName>. La tolérance pour
    <lb/>toutes Les ignominies qui ont été commises
    <lb/>Pendant la période électorale.
  </seg>
</u>
```

(a) Encoding model 1 : semantics and layout

```
<u who="#pers_ID" xml:id="CR_1889-11-26_u5" ana="#speaker">
  <seg xml:id="CR_1889-11-26_u5.1"><persName ref="#pers_ID">M. Paul Déroulède</persName>. Mes amis et moi
  avons applaudi à ce que nous avons cru, à ce que nous croyons encore être la formation spontanée d'une majorité de
  conscience, de justice, de tolérance et d'union. <incident><desc>(Applaudissements sur quelques bancs à l'extrémité
  gauche de la salle. –Exclamations au centre et à gauche.)</desc></incident></seg>
</u>

<u who="#pers_ID" xml:id="CR_1889-11-26_u6" ana="#speaker">
  <seg xml:id="CR_1889-11-26_u6.1"><persName ref="#pers_ID">Gustave Rivet</persName>. La tolérance pour
  toutes les ignominies qui ont été commises pendant la période électorale.</seg>
</u>
```

(b) Encoding model 2 : semantics only

Figure – Parliamentary session of 26 November 1889 (extract)

APPLYING ENCODING

```
"box": [
  827.8018264188805,
  2374.55023568466,
  589.3638270234404,
  207.5897599646105
],
"checked": true,
"comment": "u-pas",
"id": 352,
"sex_xml": "<PER>M. Francis Leur</PER>. Je n'ai pas l'intention\u0029d'apporter \u00e0 la Chambre d'autr\u0029 de",
"origin": "computer",
"parent": 274,
"persons": ["M. Francis Leur"],
"text_ocz": "M. Francis Leur. Je n'ai pas l'intention\u0029d'apporter \u00e0 la Chambre d'autres documents\u0029 que ceux qui me sont fournis par la lettre\u0029 de M. Lev\u0029, adresse\u0029 le ministre des finances. Je demanderai\u0029 la Chambre si elle conna\u0029 cette lettre ou\u0029 si elle veut que j'en donne lecture.",
"type": "ENTRY"
},
{
  "activities": [],
  "addressees": [],
```

Textual data in JSON

Python scripts

Textual data in XML-TEI

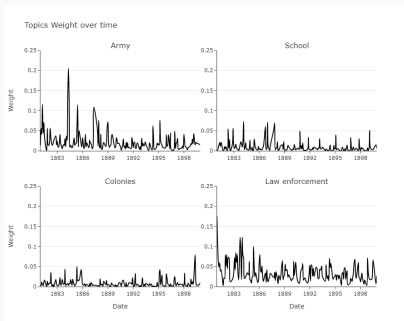
- Using **TEI Publisher** : a tool for publishing corpora in TEI
- Based on Web Components technology : allows the development of new functionalities
- Develop new modes of visualisation : for example, explore debates according to the subjects they evoke, and not only by key words

TOPIC MODELING AND WORD EMBEDDING

EXPLORING DEBATES WITH TOPIC MODELING

Topic 8	Topic 11	Topic 15
salaire	général	pari
question	commission	télégraphe
gouvernement	régiment	faire
jour	troupe	ingénieur
patron	monsieur	train
chambre	année	ligne
droit	jeune	chambre
syndicat	temps	personnel
délégué	faire	etat
monsieur	corps	administration
travail	soldat	employé
travaux	ministre	poste
ministre	homme	public
grève	loi	travaux
faire	an	service
mineur	guerre	agent
mine	service	ministre
loi	militaire	fer
compagnie	officier	chemin
ouvrier	armée	compagnie

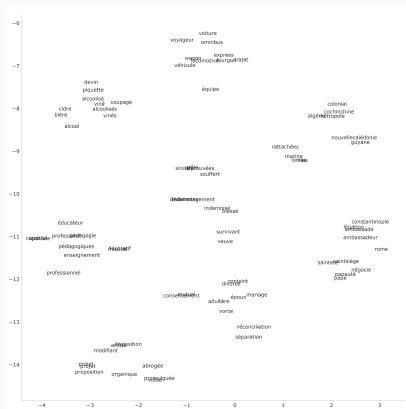
(a) 3 topics among 40 :
working class (8), army (11)
and infrastructures (15)



(b) Topic's evolution over time

Figure – Topic modeling with LDA

WORD EMBEDDING : WORD2VEC AND TOP2VEC



(a) Word vectors projected with t-SNE
(word2vec : CBOW,W=5)

Cluster 55	Cluster 68	Cluster 70
victimes	divorce	enveloppes
inondations	epoux	timbres
secourir	mariage	poste
eprouvees	conjugal	postale
orages	divorces	timbre
sinistres	adultere	recepissés
grele	conjugale	postes
secours	remarier	postaux
venir	separation	telegraphes
infortunes	indissolubilité	colis
ravages	conjoins	fixe
miseres	mutuel	recouvrements
catastrophe	separations	graphes
evenements	mari	postales
repartition	mariages	taxe
incendies	femme	decide
soulager	conjoint	soit

(b) 3 clusters among the 113 topics found by top2vec : storm (55), divorce (68) and poste (70)

Figure – Results of word2vec and top2vec



Nicolas Bourgeois : `nicolas.bourgeois@epitech.eu`

Fanny Lebreton : `fanny.lebreton@chartes.psl.eu`

Aurélien Pellet : `aurelien.pellet@epitech.eu`

Marie Puren : `marie.puren@epitech.eu`

Pierre Vernus : `pierre.vernus@msh-lse.fr`

Github repository : <https://github.com/mpuren/agoda>