

EXTRACTING AND PROVIDING ONLINE ACCESS TO ANNOTATED AND SEMANTICALLY ENRICHED HISTORICAL DATA

The AGODA project

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

1 June 2022

DH Benelux Conference 2022, University of Luxembourg

THE AGODA PROJECT

- AGODA : Analyse sémantique et Graphes relationnels pour l’Ouverture et l’étude des Débats à l’Assemblée 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 (ALMANaCH) and the University Lumière Lyon 2 (LARHRA)

PARLIAMENTARY DEBATES DURING THE THIRD REPUBLIC

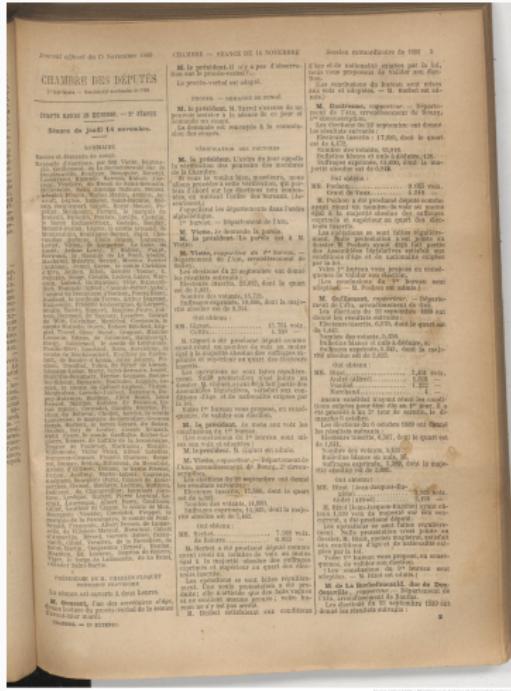


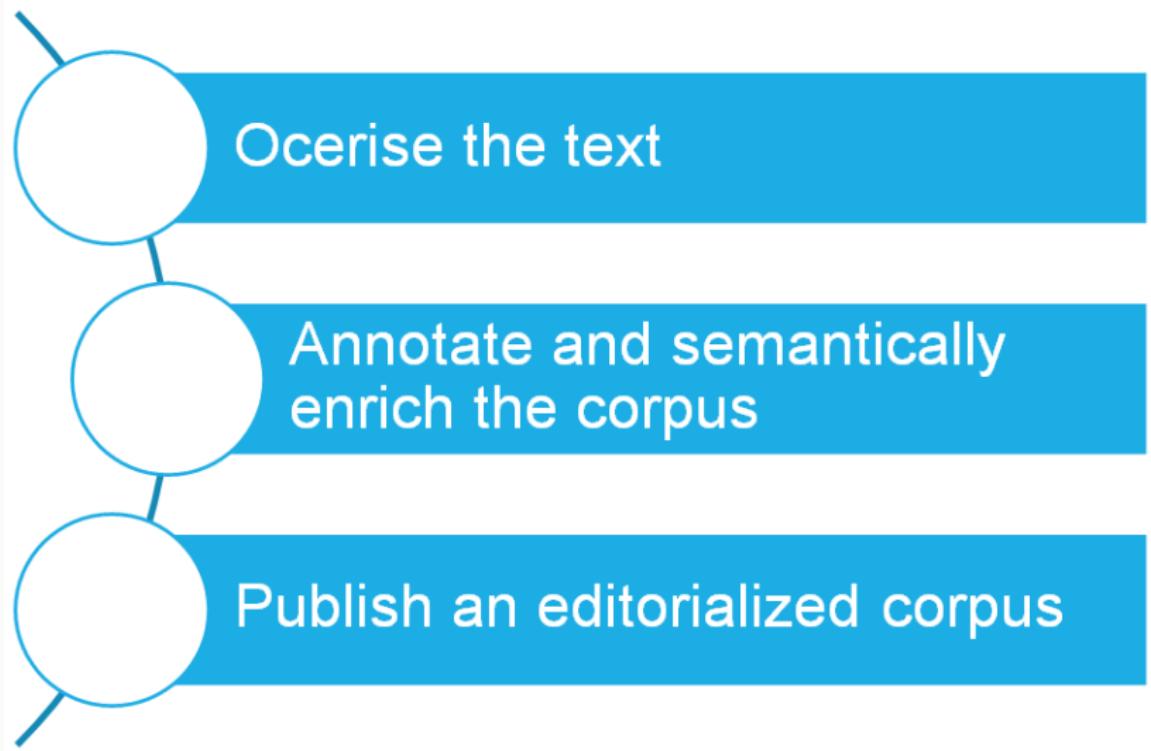
Figure – Parliamentary sitting
- 14 Nov. 1889

- 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!

A PROOF OF CONCEPT

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



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

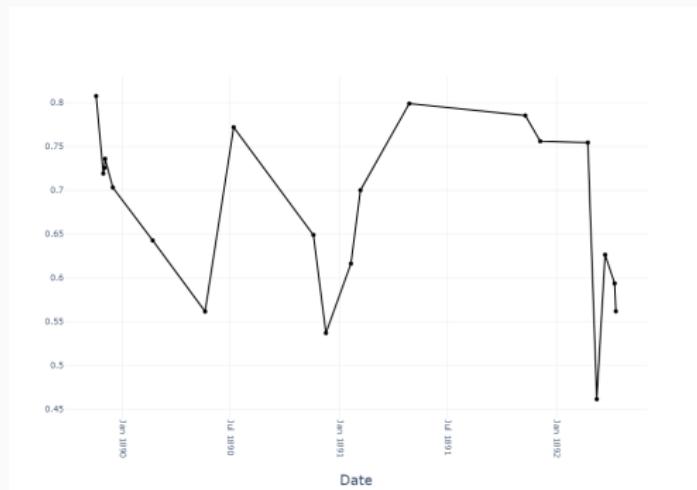


Figure – Quality assessment of the OCR provided by Gallica

EFFECT OF CURVATURE ON OCR RESULTS

Chagrin dans son cabinet et voulut à peine prononcer le discours qu'il lui tint : « Monsieur Chagny, je vous prie de faire tout ce que sera nécessaire pour que le Gouvernement et de son représentant, le ministre des Finances, soient au courant de vos vœux, mais que, si vous ne mettez pas les mesures, à mon avis, dans l'intérêt de la sécurité publique et de la survie de nos citoyens, vous n'aurez plus à compter avec moi. » Chagny, le tint pour dire à il y a deux ou trois jours, au cours d'une réunion de la commission des finances de la chambre de leurs condamnées.

M. Chagny, ministre, M. Chagny, ministre de l'Intérieur, comme on peut tous les jours ressortir à Paris, au sein de vos bons offices et des leurs, se rappellera rapidement que lorsque nous avions été nommés aux dernières élections, les deux partis avaient été unanimes que les directeurs, les administrateurs et les juges en chef devraient être élus par les électeurs, et si j'en ai eu quelques doutes, l'avertissement suffira à obéir à ces volontés, et les satisfactions que nous avons reçues.

Il est donc nécessaire d'évoquer à nouveau la question : pourquoi le principe d'une élection au suffrage universel de nos électeurs est-il aussi peu favorable pour le pays que pour le pays ? Il est tout à fait évident que le pays est menacé par les très républicains.

Il est permis de demander à nos ministres : pourquoi ne pas provoquer la voix des électeurs ? Mais il faut répondre à cette question : il faudrait faire voter à l'assemblée toutes les dissidences existantes. Ainsi, il existe de très grandes différences entre les deux partis qui se proposent ; mais quel qu'il soit, il y en a toujours un autre qui propose quelque chose qui n'est pas simplement entre les deux partis, mais qui peut également être proposé, et qui peut également être accepté. C'est pourquoi il est nécessaire que les deux partis, et des centaines de personnes qui sont dans l'opposition, doivent être libres sans défense de voter pour ces deux compagnies qui sont toutes deux très à gauche... »

Chagot fut donc cabine et voil il pr  s le discours qu'il lui tira. « Monsieur Chagot, lui dit-il, svp avec tous les jours recou『s 『abbieau-le-Gouvernement et de son repr  s 『ent et c 『est vos affaires particuli  res. Ebhissi rappo『l 『ne de vous avertir que, si vous t『ez pas mesures, 『a mon sens i『o os blees, que vous avez pr  s contre Piter sil 『riens, vous n『avez plus 『a croper SU cette bienveillance. » Il pr  s lui. Chagot se le tint pour dit et il etelit 『a ses ouvriers de conserver et de poolpier les fonctions qu'ils tenaient de a de leurs concitoyens, « cowi Eb bien, monsieur le m『e... »

M. Chang, les compagnies de fer ont tous leurs recours à vous!om vous avons ; tous les jours, elles envoient des boîtes aux bureaux et des leuros, n'avez-vous pas seulement rappeler à MM. les îles c-contrôle qu'ils sont moins encouragés et rades que les directeurs, les conçilliers, les surveillants des ingénieurs des coH11P les surveillants des ingénieurs c'der, gnes, et si vous convainquez ce ces défauts avisez-moi suffira à obtenir mères les satisfactions que no t'mons, une Messieurs, hier, vous vous priez de venir avec moi et veiller pour but de garantir x'Uvros la liberté d'association, qui leur est reçue, puis les lois républicaines, initiatif, de Je ne demande pas à M. le II 3g, traux malice de provoquer * j'en ai loué nouvelle, je lui demande Simple gne ! faire aujou'res les lois exigeant x'Uvros et ay avoir hier des dissidences en bres de cette Chambre sur telles y'ont voté est propose. Ilia questi qu'il saurait y en avoir sur 011 que je souleve aujou'res, un t Ce n'est pas simplement en élève, C et un ouvrier que le débat s'et il se trouve entre une compagnie qui lui ren'gtzier.

vice public et des centaines de juges travaillent. La question est de auJc Jlt.

doivent être livrés sans défensee au* JieV ces des compagnies qui les ifres 01'0 et, Quant à moi, je ne le crois p *

Figure – OCR result

Figure – 14

May 1890

(p.786)

IMAGE CLEANING TOOL DEVELOPED BY ANR SODUCO



(a) Original image

(b) Cleaned image

Figure – Demonstration of the SODUCO tool on a parliamentary debate 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**

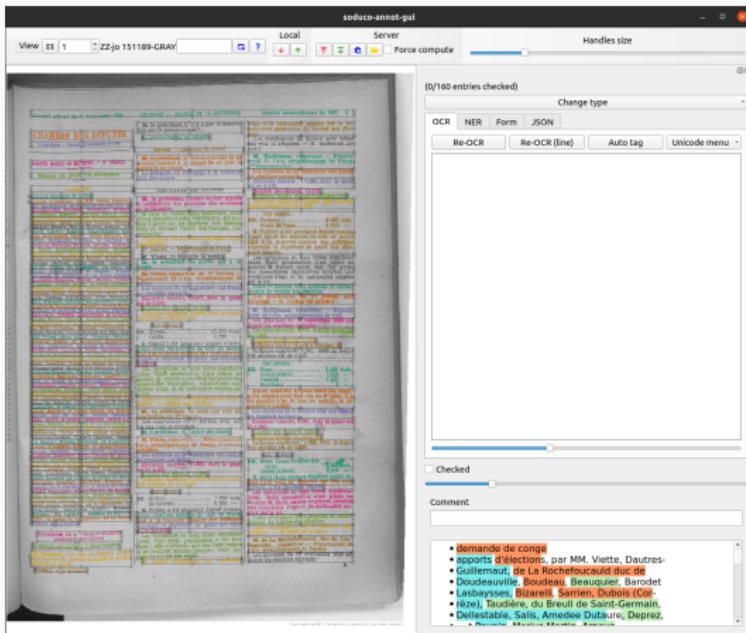


Figure – OCR tool developed by SODUCO

OCR AND NER

(0/160 entries checked)

ENTRY

OCR NER Form JSON

Re-OCR Re-OCR (line) Auto tag Unicode menu ▾

Ont obtenu:
MM. Pochon..... 9.055 voix.
Grant de Vaux..... 4.516
M. Pochon a été proclamé député comme
ayant réuni un nombre de voix au moins
égal à la majorité absolue des suffrages
exprimés et supérieur au quart des élec-

Checked



(a) OCR

(0/160 entries checked)

ENTRY

OCR NER Form JSON

PER ACT LOC CARDINAL FT TITRE

CLEAR

Ont obtenu:
MM. Pochon..... 9.055 voix.
Grant de Vaux..... 4.516
M. Pochon a été proclamé député comme
ayant réuni un nombre de voix au moins
égal à la majorité absolue des suffrages
exprimés et supérieur au quart des élec-

Checked



(b) NER

Figure – OCR and NER zones

OUTPUT IN JSON

```
      50
],
"id": 302,
"ner_xml": "<PER>Suirrages</PER> exprim\u00e9s, 9,90<CARDINAL>8</CARDINAL>, dont la majo-\u2029rit\u00e9 absolue est de <CARDINAL>4,455</CARDINAL>.",
"origin": "computer",
"parent": 263,
"persons": ["Suirrages"],
"text_ocr": "Suirrages exprim\u00e9s, 9,908, dont la majo-\nrit\u00e9 absolue est de 4,455.",
"type": "ENTRY",
"activities": [],
"comment": "",
"checked": false
}
```

Figure – Output in JSON (extract)

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 : **ParlaClarin** and **ParlaMint**
- The automatic tagging process

For more information on this aspect :

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

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
  avions 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'extrême
  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>
```

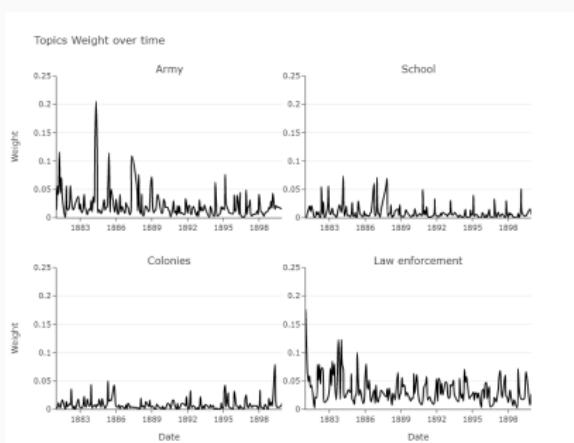
(b) Encoding model 2 : semantics only

Figure – Parliamentary session of 26 November 1889 (extract)

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	époux	timbres
secourir	mariage	poste
éprouvées	conjugal	postale
orages	divorces	timbre
sinistres	adultere	recepisses
grele	conjugale	postes
secours	remarier	postaux
venir	separation	telegraphes
infortunes	indissolubilité	colis
ravages	conjoints	fixe
misères	mutuel	recouvrements
catastrophe	separations	graphes
événements	mari	postales
répartition	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

FOR MORE INFORMATION

- Nicolas Bourgeois, Aurélien Pellet, Marie Puren. "Using Topic Generation Model to explore the French Parliamentary Debates during the early Third Republic (1881-1899)". (hal-03526254v2)
- Marie Puren, Nicolas Bourgeois, Aurélien Pellet, Pierre Vernus, Fanny Lebreton. "Between History and Natural Language Processing : Study, Enrichment and Online Publication of French Parliamentary Debates of the Early Third Republic (1881-1899)". ParlaCLARIN III at LREC2022 - Workshop on Creating, Enriching and Using Parliamentary Corpora, Jun 2022, Marseille, France. (hal-03623351)

THANK YOU FOR YOUR ATTENTION !



Marie Puren : marie.puren@epitech.eu

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

Aurélien Pellet : aurelien.pellet@epitech.eu

Nicolas Bourgeois : nicolas.bourgeois@epitech.eu

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