

TRAITEMENT ET ANALYSE DES DÉBATS PARLEMENTAIRES À LA CHAMBRE DES DÉPUTÉS (1881-1899)

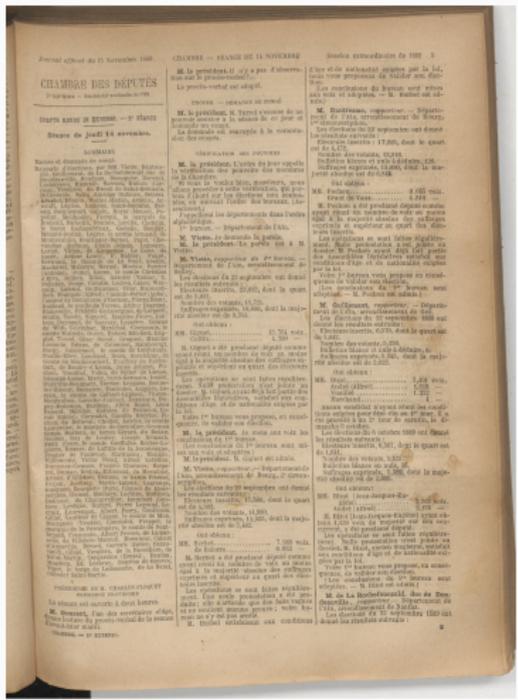
Problèmes, défis et solutions

Aurélien Pellet¹ et Marie Puren²

10 novembre 2022 - Res(t)ituer les adresses des almanachs et annuaires commerciaux parisiens du XIXe siècle. Un corpus de localisations urbaines à grande échelle. 2e Journée de l'atelier SoDUCo-BnF

¹Epitech ²LRE, EPITA

LES DÉBATS PARLEMENTAIRES DURANT LA TROISIÈME RÉPUBLIQUE



- **Projet AGODA** : Analyse sémantique et Graphes relationnels pour l'Ouverture et l'étude des Débats à l'Assemblée nationale
- Débats à la Chambre des députés (chambre basse du parlement) transcrits en détail dans le **Journal officiel de la République française. Débats parlementaires (1881-1940)**
- Disponible en ligne via **Gallica** (bibliothèque numérique de la Bibliothèque nationale de France)
- 15 législatures entre 1881 et 1940 / 10-12.000 images par législature

Figure – Séance parlementaire du 14 novembre 1889

OBJECTIFS

- Créer une plateforme de consultation
- Produire des données textuelles structurées et sémantiquement enrichies à partir de ces débats numérisés
- Contribuer à la conception d'un workflow adapté à la préparation, à la publication et à l'analyse de grands corpus de documents historiques

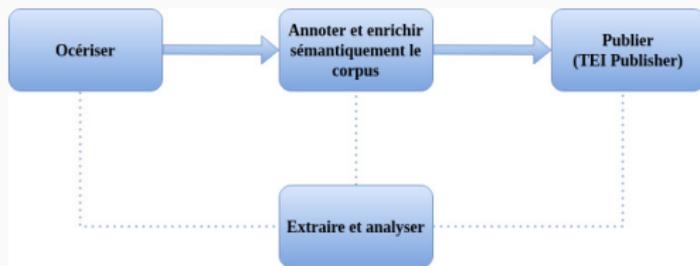


Figure – Les étapes de la chaîne de traitement

OCÉRISER LES DÉBATS

- Récupération des textes océrisés via **API Document** de Gallica => qualité inégale de l'OCR
- Erreurs dues à :
 - qualité du document : tâches et surimpression
 - la **courbure de la page** au niveau de la reliure

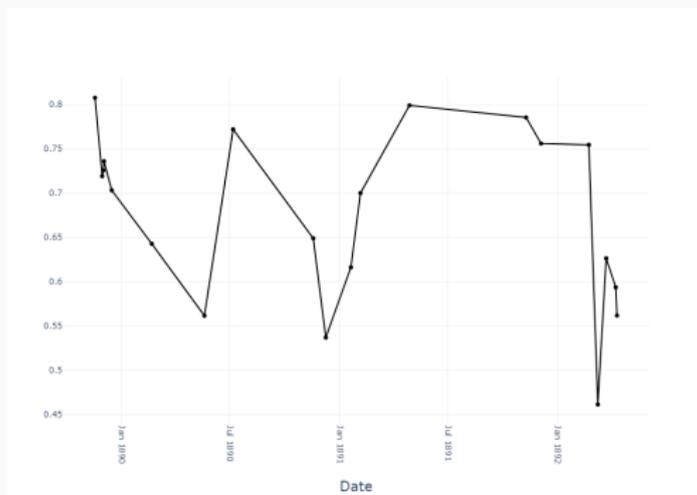


Figure – Evaluation de la qualité de l'OCR fourni par Gallica

EFFET DE LA COURBURE SUR LES RÉSULTATS DE L'OCR

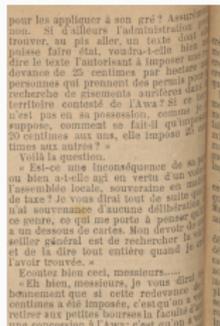


Figure – 20
octobre
1890
(p.1718)

pour les appliquer à son gré ? Assure.

non. Si d'ailleurs l'administration f, trouver, au pis aller, un texte (10111 en puisse faire état, voudra-t-elle dire le texte l'autorisant à imposer un devance de 25 centimes par hectare , personnes qui prennent des permis Pour l recherche de gisements aurifères a teSI territoire contesté de l'AwA ? Si ce tes

n'est pas en sa possession, comas le suppose, comment se fait-il qu'ilWjsj} 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 délibération ce genre, ce qui me porte à penser a un dessous de cartes. Mon devoir seiller général est de rechercher le avoir trouvés. »

Eoutez bien ceci, messieurs. l «Eh bien, messieurs, je vous dirai 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;:i}

Figure – Résultat de l'OCR

Décision de ré-océreriser le corpus

Améliorer la qualité de l'image avec une méthode de "dewarping" => résultats peu probants

- Gérer la courbure des pages avec le dewarping?
- Utiliser des outils plus avancés?



(a) Image d'origine

(b) Image "dewarpée"

Figure – Dewarping : pas adapté à nos documents

OUTIL DE NETTOYAGE - ANR SUDCO



(a) Image d'origine



Figure - Exemple d'image nettoyée

OUTIL OCR ET NER - ANR SODUCO

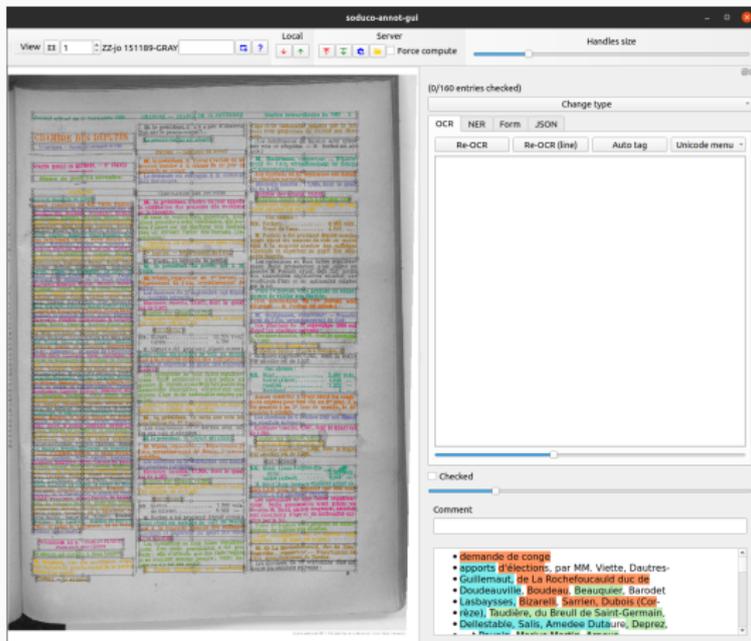
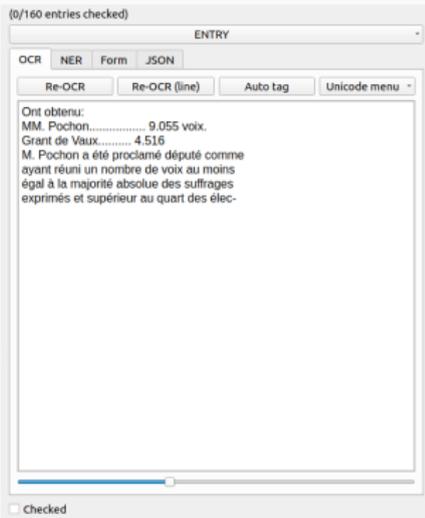


Figure – Démonstration de l’outil sur une page de débat parlementaire

OUTIL OCR ET NER - ANR SODUCO



(a) OCR



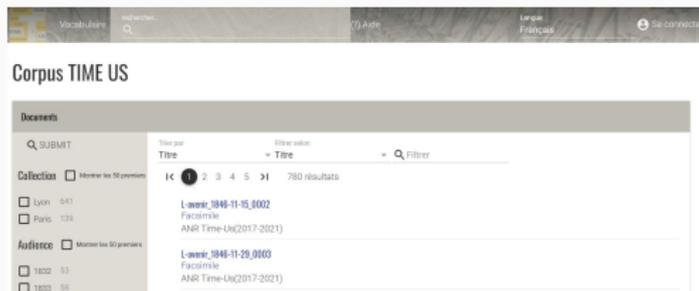
(b) NER

Figure – Zones d'OCR et de NER

- Problèmes de reconnaissance de la structure hiérarchique à cause de la courbure
- Objectif n°1 : faire disparaître au maximum cette courbure
- Développement et traitement d'images en C++

Demande de financement dans le cadre de Biblissima +

ANNOTER LES DÉBATS EN XML-TEI



(a) Interface de recherche



(b) Affichage d'un document

Figure – Corpus du projet ANR TIME-US publié avec TEI Publisher

CORPS D'UN DÉBAT ENCODÉ EN XML-TEI

```
<div type="part" corresp="#rapportelections">
  <head>SUITE DE LA VERIFICATION DES POUVOIRS</head>
  <u who="#pers_ID" xml:id="CR_1889-11-26_u23" ana="#chair">
    <seg xml:id="CR_1889-11-26_u23.1"><persName ref="#pers_ID">M. le <roleName ref="#pers_ID">président</roleName></persName>. L'ordre du jour
appelle la suite de la vérification des pouvoirs.</seg>
    <seg xml:id="CR_1889-11-26_u23.2"><persName ref="#pers_ID">M. Reybert</persName> a la parole pour donner lecture d'un rapport sur une
élection non contestée.</seg>
  </u>
  <u who="#pers_ID" xml:id="CR_1889-11-26_u24" ana="#rapporteur">
    <!-- Lecture d'un rapport -->
    <quote>
      <seg xml:id="CR_1889-11-26_u24.1"><persName ref="#pers_ID">M. Reybert, <roleName ref="#pers_ID">rapporteur</roleName></persName>.
--<placeName ref="#lieu_ID">Département de la Corrèze, arrondissement de Tulle, <num>1</num> circonscription</placeName>.</seg>
      <seg xml:id="CR_1889-11-26_u24.2">Les élections du <date when="1889-09-22">22 septembre</date> ont donné les résultats suivants :</seg>
      <seg xml:id="CR_1889-11-26_u24.3">Electeurs inscrits, <num>17,950</num>, dont le quart est de <num>4,263</num>.</seg>
      <seg xml:id="CR_1889-11-26_u24.4">Nombre des votants, <num>12,322</num>.</seg>
      <seg xml:id="CR_1889-11-26_u24.5">Bulletins blancs et nuls, à déduire, <num>133</num>.</seg>
      <seg xml:id="CR_1889-11-26_u24.6">Suffrages exprimés, <num>12,189</num>.</seg>
      <seg xml:id="CR_1889-11-26_u24.7">Ont obtenu : MM. <persName ref="#pers_ID">Borie, <roleName ref="#pers_ID">député
sortant</roleName></persName>,... <num>7,508</num> voix</seg>
      <seg xml:id="CR_1889-11-26_u24.8"><persName ref="#pers_ID">Vachal, <roleName ref="#pers_ID">ancien député</roleName></persName>...
<num>4,748</num></seg>
      <seg xml:id="CR_1889-11-26_u24.9"><persName ref="#pers_ID">M. Borie</persName> a été proclamé député comme ayant réuni un nombre de voix
au moins égal à la majorité absolue des suffrages</seg>
    </quote>
  </u>
  <floatingText><body><div><p n="176"/></div></body></floatingText>
  <!-- [...] -->
</div>
```

Figure – Encodage - Séance parlementaire du 26 novembre 1889 (extrait)

ANNEXES D'UN DÉBAT ENCODÉES EN XML-TEI

Annexe au procès-verbal de la séance du mardi 26 novembre 1889.	
SCRUTIN	
Sur les conclusions de 1 ^{er} bureau tendant à l'annulation des opérations électorales de la 1 ^{re} circonscription de l'arrondissement de Lorient (Morbihan).	
Nombre des votants.....	506
Majorité absolue.....	254
Pour l'adoption.....	330
Contre.....	176
La Chambre des députés a adopté.	
ONT VOTÉ POUR :	
MM. Abeille, Armes (Emanuel), Armez, Arribat, Audifred, Ayraud (Edouard), Baile (Marial), Bory, Barodot, Barthon, Bartsch, Baril (Adrien), Banlard, Beauquier, Berard, Berger (Georges) (Seine), Bertrand, Bézaux, Bizarelli, Bizol, Bizouard-Bert, Bizez (Pierre), Buisy-É Anglais, Boudiguy-Silhouz, Bony-Casternes, Bourgeois, Bouchier (Volger), Boudemont, Bouderville, Bouge, Boulogne-Bermet, Boullay, Bourgeois (Jules), Bourgeois (Léon) (Marne), Bouillier de Boissière, Bortier (Pierre), Boyssac, Braud, Breton, Brézas, Brisson (Henri), Brisson (Eugène), Brégnol, Brunier, Bully, Burdeau, Buxignier.	
Ratification aux scrutins de la séance du 25 novembre 1889.	
M. Michau (Nord), porté comme s'étant abstenu dans le scrutin sur l'urgence de la proposition de M. Maxime Lecoste, déclare avoir voté pour ».	

(a) Source numérisée

```
<!-- ANNEXES -->
<back>
<head>Annexes au procès-verbal de la séance du <date when="1889-11-26">mardi 26 novembre 1889</date>.</head>
<div xml:id="vot18891126">
<!-- VOTE 1 -->
<div xml:id="vot18891126_vot1" type="voting" corresp="discussion7ebureau">
<head>
<label>SCRUTIN</label>
<note>seg-Sur les conclusions du <num>7</num> bureau tendant à l'annulation des opérations électorales de la
<placeName ref="#lieu_ID"><num>1</num> circonscription de l'arrondissement de Lorient (Morbihan)</placeName>.</seg></note>
</head>
<!-- Détail du vote -->
<desc>
<measure type="nbvotants" quantity="506">Nombre des votants <num>506</num></measure>
<measure type="maj" quantity="254">Majorité absolue <num>254</num></measure>
<measure type="yes" quantity="330">Pour l'adoption <num>330</num></measure>
<measure type="noes" quantity="176">Contre <num>176</num></measure>
</desc>
<note type="result">seg-La <orgName ref="#org_ID">Chambre des députés</orgName> a adopté.</seg></note>
<floatingText>body=>div=pb n°192</div></floatingText>
<!-- Liste des votants -->
<note type="voterslist">
<desc>Ont voté pour :</desc>
<seg>MM. <persName ref="#pers_ID">Abeille</persName>, <persName ref="#pers_ID">Armes (Emanuel)</persName>, <persName ref="#pers_ID">Armez</persName>, <persName ref="#pers_ID">Arribat</persName>, <persName ref="#pers_ID">Audifred</persName>, <persName ref="#pers_ID">Ayraud (Edouard)</persName>, <persName ref="#pers_ID">Baile (Marial)</persName>, <persName ref="#pers_ID">Bory</persName>, <persName ref="#pers_ID">Barodot</persName>, <persName ref="#pers_ID">Barthon</persName>, <persName ref="#pers_ID">Bartsch</persName>, <persName ref="#pers_ID">Baril (Adrien)</persName>, <persName ref="#pers_ID">Banlard</persName>, <persName ref="#pers_ID">Beauquier</persName>, <persName ref="#pers_ID">Berard</persName>, <persName ref="#pers_ID">Berger (Georges) (Seine)</persName>, <persName ref="#pers_ID">Bertrand</persName>, <persName ref="#pers_ID">Bézaux</persName>, <persName ref="#pers_ID">Bizarelli</persName>, <persName ref="#pers_ID">Bizol</persName>, <persName ref="#pers_ID">Bizouard-Bert</persName>, <persName ref="#pers_ID">Bizez (Pierre)</persName>, <persName ref="#pers_ID">Buisy-É Anglais</persName>, <persName ref="#pers_ID">Boudiguy-Silhouz</persName>, <persName ref="#pers_ID">Bony-Casternes</persName>, <persName ref="#pers_ID">Bourgeois</persName>, <persName ref="#pers_ID">Bouchier (Volger)</persName>, <persName ref="#pers_ID">Boudemont</persName>, <persName ref="#pers_ID">Bouderville</persName>, <persName ref="#pers_ID">Bouge</persName>, <persName ref="#pers_ID">Boulogne-Bermet</persName>, <persName ref="#pers_ID">Boullay</persName>, <persName ref="#pers_ID">Bourgeois (Jules)</persName>, <persName ref="#pers_ID">Bourgeois (Léon) (Marne)</persName>, <persName ref="#pers_ID">Bouillier de Boissière</persName>, <persName ref="#pers_ID">Bortier (Pierre)</persName>, <persName ref="#pers_ID">Boyssac</persName>, <persName ref="#pers_ID">Braud</persName>, <persName ref="#pers_ID">Breton</persName>, <persName ref="#pers_ID">Brézas</persName>, <persName ref="#pers_ID">Brisson (Henri)</persName>, <persName ref="#pers_ID">Brisson (Eugène)</persName>, <persName ref="#pers_ID">Brégnol</persName>, <persName ref="#pers_ID">Brunier</persName>, <persName ref="#pers_ID">Bully</persName>, <persName ref="#pers_ID">Burdeau</persName>, <persName ref="#pers_ID">Buxignier</persName>.
</note>
</div>
<!-- RECTIFICATIONS -->
<div corresp="vot18891125" type="rectification">
<head>Rectifications aux scrutins de la séance du <date>25 novembre 1889</date>.</head>
<note corresp="vot18891125_vot1">seg-<persName ref="#pers_ID">M. Michau</persName> <placeName ref="#lieu_ID">(Nord)</placeName>, porté comme s'étant abstenu dans le scrutin sur l'urgence de la proposition de M. Maxime Lecoste, déclare avoir voté pour ».</seg></note>
<!-- [...] -->
</div>
</back>
```

(b) Modèle d'encodage

Figure – Séance parlementaire du 26 novembre 1889 - votes, liste des votants, rectifications (extrait annexes)

STRUCTURE GÉNÉRALE DES FICHIERS XML TEI

```
<teiCorpus xmlns="http://www.tei-c.org/ns/1.0" xsl:id="FR_3R_5L" xsl:lang="fr">
  <!-- Métadonnées du corpus (non développées) -->
  <teiHeader>
    <fileDesc>
      <titleStnt>
        <title></title>
      </titleStnt>
      <publicationStnt>
        <p></p>
      </publicationStnt>
      <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 -->
  <xli:include xmlns:xli="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) Structure générale d'un fichier corpus

```
<TEI xmlns="http://www.tei-c.org/ns/1.0" xsl:id="FR_3R_5L_1889-11-26" xsl:lang="fr">
  <!-- Métadonnées du composant (non développées) -->
  <teiHeader>
    <fileDesc>
      <titleStnt>
        <title></title>
      </titleStnt>
      <publicationStnt>
        <p></p>
      </publicationStnt>
      <sourceDesc>
        <p></p>
      </sourceDesc>
    </fileDesc>
  </teiHeader>

  <text>

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

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

  </text>
</TEI>
```

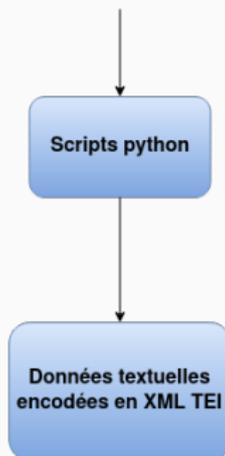
(b) Structure générale d'un fichier composant

Figure – Encodages - Structure générale des fichiers corpus et composant

APPLIQUER L'ENCODAGE : AUTOMATISATION

```
▼ 19:
activities:      []
addresses:      []
▼ box:
  0:             220.12998972250773
  1:             202.1299897225077
  2:             561.8527187504491
  3:             00.4418437486525
checked:        true
comment:        "seg"
id:             276
▼ ner_xml:
  origin:        "<PER>M. Borie</PER> a «ACT>déjà fait partie des Assemblées\u2029législatives et satisfait aux conditions d«ACT>«ACT>âge\u2029et de
  parent:        261
  text_ocr:      "M. Borie a déjà fait partie des Assemblées\nlégislatives et satisfait aux conditions d'âge\net de nationalité exigées par la loi."
  type:          "ENTRY"
```

Données textuelles au format JSON



- Données textuelles dans les fichiers JSON = valeurs des clés "text_ocr"
- Clés contenues dans des objets qui correspondent à la segmentation au niveau paragraphe
- Appliquer des règles de transformation sur les valeurs des clés "text_ocr"


```
"activities": [],
"addresses": [],
"box": {
  57.6116701150507,
  1710.0,
  560.3883298849495,
  50.07766597698992
},
"checked": true,
"comment": "u seg",
"id": 305,
"key": [
  0,
  1735
],
"ner_xml": "<PER-M, Paul D roulede</PER>. Je demand la pa-rol . Tole",
"origin": "computer",
"parent": 269,
"persons": [
  "M. Paul D roulede"
],
"text_ocr": "M. Paul D roulede. Je demand la pa-rol .",
"type": "ENTRY"
},
```

(a) Fichier JSON avec une prise de parole

```
def add_utterance(data):
    """
    Ajout de l' l ment TEI "u" pour chaque bo e  tiquet e "u" ou "u-beginning" et "u-end"
    :param data: dictionnaire contenant l'ensemble des donn es issues des JSOM
    """
    for i in range(len(data)):
        if "comment" in data[i]:
            if re.search(r"^(u|u-)\b", data[i]["comment"]):
                data[i]["text_ocr"] = "".join(['<u>', data[i]["text_ocr"], '</u>'])
            elif re.search(r"u-beginning", data[i]["comment"]):
                data[i]["text_ocr"] = "".join(['<u>', data[i]["text_ocr"]])
            elif re.search(r"u-end", data[i]["comment"]):
                data[i]["text_ocr"] = "".join([data[i]["text_ocr"], '</u>'])
            else:
                pass
    return data
```

(b) Fonction add_utterance

```
<u>seg>M. Paul D roulede. Je demande la parole.</seg></u>
```

(c) R sultat

Figure – Exemple d'application de la fonction add_utterance

1. Récupération des données JSON stockées dans une variable "data" sous forme d'un dictionnaire.
2. Gestion de l'encodage des changements de page.
3. Intégration des scripts contenant les règles de transformation dans la chaîne de traitement. Permettent d'encoder les données contenues dans la variable "data".
4. Création de l'élément <teiHeader> contenant les métadonnées :
 - Rédaction d'une partie du <teiHeader> à la main contenant les métadonnées "fixes"
 - Construction de règles permettant de rechercher les métadonnées propres à chaque compte rendu, présentes dans la variable "data", par la suite intégrées au <teiHeader> préalablement établi
5. Création des fichiers XML valides.
6. Nettoyage des fichiers XML (gestion des césures).

1. Génération du fichier corpus.
2. Intégration des métadonnées.
3. Gestion des `<xi :include>`.

1. Annotation manuelle très consommatrice de temps et non exempte d'erreurs
 - Ajout de manière semi-automatique des étiquettes
 - Permettre aux utilisateurs de réutiliser les étiquettes et de créer leurs propres étiquettes (vocabulaire contrôlé?)
2. Utiliser ces scripts pour proposer un export XML-TEI (outil SODUCO)
 - Ajouter les entités nommées (<persName>, <orgName>, etc.)

Demande de financement dans le cadre de Biblissima +

TOPIC MODELING ET WORD EMBEDDING

DÉCOMPOSITION EN VALEURS SINGULIÈRES

- On observe deux blocs de documents :
 - Mots positivement corrélés à la première composante : martin, paul, jules, léon, vicomte, seine, baron, comte...
 - Mots positivement corrélés à la seconde composante : monsieur, commission, ministre, projet, chambre, président, loi...

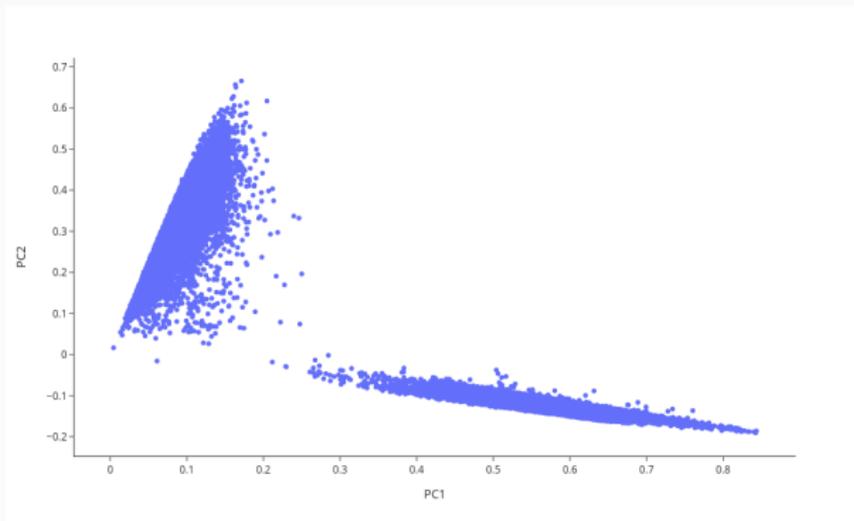


Figure – Décomposition en valeur singulière de la matrice de fréquence (+ tfidf)

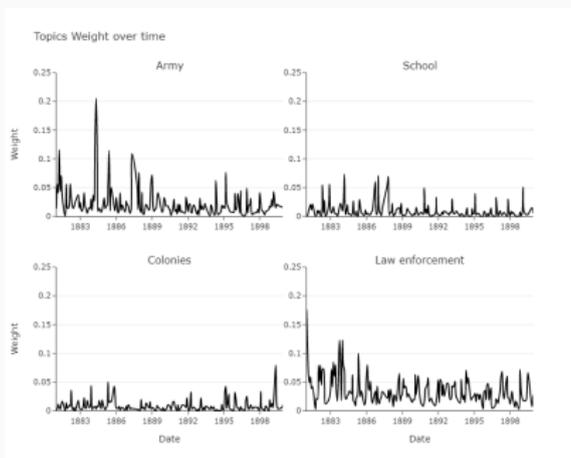
ONT VOTÉ POUR :

MM. Abbatucci. Achard. Agniel. Allain-
 Targé. Allègre. Allemand. Amat. André
 (Jules). Andrieux. Anisson-Duperron. An-
 thoard. Arenberg (prince d'). Ariste (d'). Ar-
 mez. Arnoult. Arrazat. Audiffred. Aulan
 (marquis d'). Azémar.
 Baduel d'Oustrac. Baïhaut. Ballue. Bam-
 berger. Bansard des Bois. Barascud. Barbe-
 dette. Bardoux. Barodet. Barthe (Marcel).
 Bastid (Adrien). Baudry-d'Asson (àe). Baurry.
 Beauchamp (de). Beauquier. Beaussire. Bel
 (François). Belle. Bénazet. Benoist. Berger.
 Bergerot. Berlet. Bernard. Bernier. Bert
 (Paul). Bertholon. Bianchi. Bienvenu. Bi-
 nachon. Bizarelli. Bizot de Fonteny. Blanc
 (Louis) (Seine). Blanc (Pierre) (Savoie). Blan-
 din. Blin de Bourdon (vicomte). Bonnaud.
 Bonnet-Duverdier. Borrighione. Bosc. Bou-
 chet. Boudeville. Boulard (Cher). Boulart
 (Landes). Bouquet. Bourgeois. Bousquet.
 Bouteille. Bouthier de Rochefort. Boyer
 (Ferdinand). Brame (Georges). Bravet. Bre-
 lay. Bresson. Brice (René). Brierre. Bris-
 son (Henri). Brossard. Bruneau. Buyat.
 Cadot (Louis). Caduc. Cantagrel. Carnot (Sadi).
 Casabianca (vicomte de). Casimir-Perier (Aube).
 Casimir-Perier (Paul) (Seine-Inférieure). Casse
 (Germain). Castagnède. Caurant. Caze. Ca-
 zeaux.

EXPLORER LES DÉBATS AVEC LA MODÉLISATION DE SUJETS

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

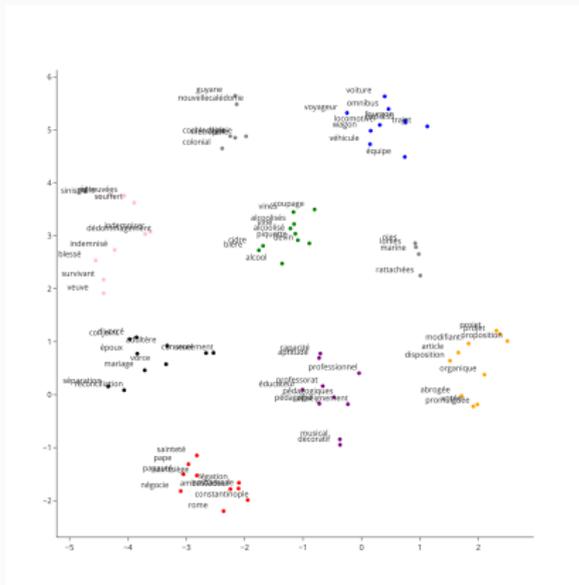
(a) 3 sujets parmi 40 : classe ouvrière (8), armée (11) et infrastructures (15)



(b) Evolution de quatre sujets au cours du temps

Figure – Résultats de la modélisation de sujets

LES PLONGEMENTS DE MOTS : WORD2VEC ET TOP2VEC



(a) Projection t-SNE des centroïdes des vecteurs (word2vec)

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

(b) 3 clusters parmi les 113 obtenus avec top2vec : tempêtes (55), divorce (68) et poste (70)

Figure – Résultats obtenus avec word2vec et top2vec

Possibilité de calculer la cosinus similarity entre un topic et des documents

QUESTION A M. LE MINISTRE DE L'INTÉRIEUR

M. le président. La parole est à M. Margaine pour adresser une question à M. le ministre de l'intérieur, qui l'a acceptée.

M. Margaine. Messieurs, d'accord avec M. le ministre de l'intérieur, je viens lui demander s'il n'entend pas appliquer le secours de 500,000 fr. qu'il a demandé aux quelques départements qui, en dehors des inondations, ont été victimes d'orages, d'incendies ou de sinistres quelconques. (Mouvements divers.)

M. de Châtensy. Il y a un crédit spécial voté pour cela !

M. le président. La parole est à M. le ministre de l'intérieur.

M. Sarrien, ministre de l'intérieur. Messieurs, le crédit de 500,000 fr. que je demande à la Chambre a uniquement pour objet de secourir les misères résultant des inondations. (Très bien ! très bien !)

En ce qui concerne les départements qui ont été éprouvés par des orages ou par la grêle, — et je crois que c'est le cas du département de la Marne, que représente l'honorable M. Margaine. . .

Sur divers bancs. Et des Ardennes ! — Et de l'Aisne ! — Et de la Meuse !

(a) $\text{cos sim} = 0.73$

M. le président. L'ordre du jour appelle la première délibération sur le projet de loi tendant à créer des timbres spéciaux pour la constatation des versements sur les livrets de la caisse d'épargne postale.

La parole est à M. le rapporteur.

M. Balhaut, rapporteur. Messieurs, au nom de la commission du budget, d'accord avec M. le ministre des postes et télégraphes, j'ai l'honneur de demander à la Chambre de vouloir bien prononcer l'urgence sur le projet de loi relatif à la création de timbres spéciaux pour les livrets de la caisse d'épargne postale.

M. le président. Je consulte la Chambre sur la demande de la déclaration de l'urgence. (L'urgence, mise aux voix, est déclarée.)

M. le président. Personne ne demandant la parole pour la discussion générale, je consulte la Chambre sur la question de savoir si elle entend passer à la lecture des articles.

(b) $\text{cos sim} = 0.81$

Figure – Résultats obtenus avec word2vec et top2vec

- 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](#))



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

Marie Puren : `marie.puren@epita.fr`

Répertoire Github du projet : <https://github.com/mpuren/agoda>

Mémoire de Fanny Lebreton :

<https://github.com/FannyLbr/Memoire-AGODA-TNAH2022>