

ANALYSIS OF THE FRENCH PARLIAMENTARY DEBATES OF THE THIRD REPUBLIC WITH TOPIC MODELLING AND WORD EMBEDDING

Methodological Challenges and First Results

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

19-20 mai 2022

Journées MASHS-19-20 mai

1. Context : the AGODA project
2. Ocerisation of digitised documents
3. Topic modelling with LDA
4. Word, Document and Topic embedding

CONTEXT : THE AGODA PROJECT

- AGODA : **A**nalyse sémantique et **G**raphes relationnels pour l'**O**uverture et l'étude des **D**ébats à l'**A**ssemblée nationale
- Project funded by the Bibliothèque nationale de France for one year
- One of the five pilot projects supported by the **DataLab** of the National Library of France
- Collaboration between Epitech (MNSHS), Inria (ALMAnaCH) and l'Université Lumière Lyon 2 (LARHRA).

FRENCH PARLIAMENTARY DEBATES OF THE THIRD REPUBLIC

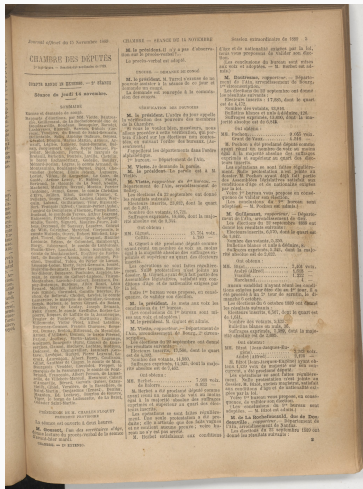


Figure – Session of November 14, 1889

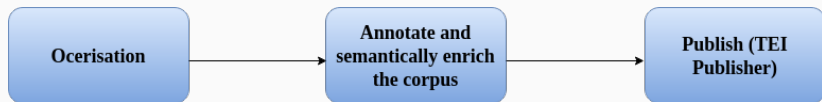
- The debates in the lower house of the French Parliament were published in the **Journal officiel de la République française. Débats parlementaires** (1881-1940)
- The issues of the Journal Officiel are available on **Gallica** (bibliothèque numérique de la Bibliothèque nationale de France)
- Difficult to work on this corpus, yet interesting for various disciplines (history, sociology, political science, linguistics)

- Facilitate access to French parliamentary debates of the Third Republic
- Facilitate research in this corpus
- Allowing the constitution of sub-corpora
- Providing new ways of viewing the documents

- Creating an online platform for consulting and exploring parliamentary debates
- Producing structured and semantically enriched textual data
- Contribute to the design of a workflow adapted to the analysis of high number of historical documents

Processing a subpart of the corpus : legislature 1889-1893 : **10418**
images to process

- Partial renewal of the political staff (boulangisme and Panama scandal)
- First signs of Catholic's rallying to the Republic
- Turnaround in customs policy (Méline laws)
- Rise of socialism and trade unionism (Fourmies)
- First anarchist terrorist attacks



OCERISATION OF DIGITISED DOCUMENTS

Mass digitization of texts : how to access/ process and analyze their content?

- OCR : Optical Character Recognition
- Processing of an image (digitized text) by an OCR engine
- Use of AI : « translation » of the image into text
 - Pages are digitised, then transformed into letters and « discreet » words » (countable : discontinuous, separate, distinct)
 - Possible to explore these pages and « differentiate » them by period, genre, language, year of publication, etc.

FRENCH PARLIAMENTARY DEBATES

- Retrieval of ocerised texts via an API of Gallica, a certain number of errors are identified
- Erros tue to the curvature of the page at the binding + stains, shadows...etc

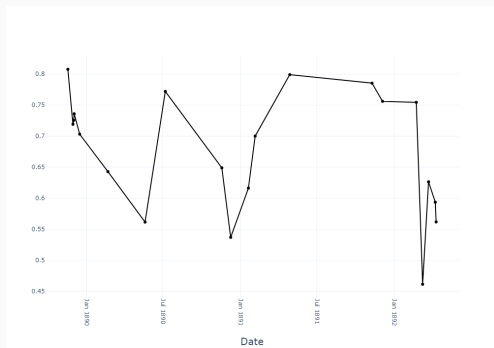
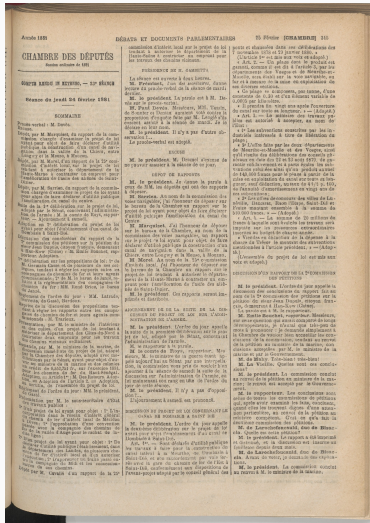


Figure – 1st evaluation of the Gallica's OCR quality

Improving OCR quality with « dewarping » methods => Inconclusive results

- Dewarping of pages?
- More advanced cleaning Tools?
 - Tesseract documentation provides a lot of tools to clean the images
 - Tools developed by other teams
- Bench marking different OCR engines

SODUCO CLEANING TOOL



(a) Original

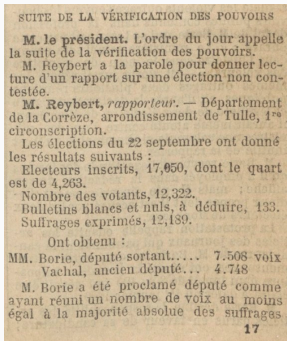


(b) (pre)Processed

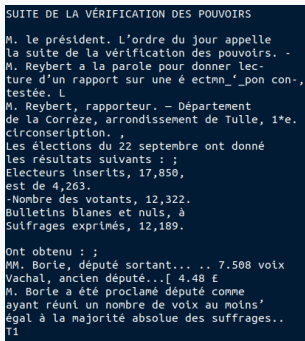
Figure 10 — SODUCO tool applied to one page of debate

Selecting and comparing multiple OCR engines :

- Tesseract (ocr-tesseract or pytesseract)
- ABBYY FineReader



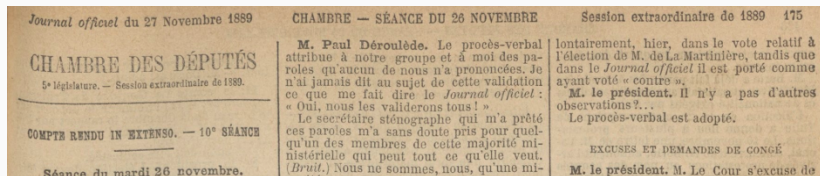
(a) original (zoom)



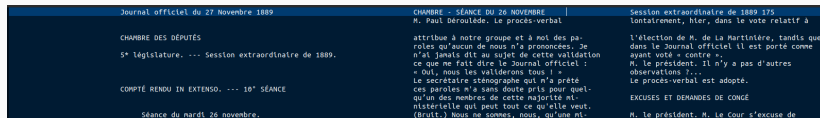
(b) OCR (zoom)

Figure – Zoom in on a text block + OCR (tesseract)

SOME SOLUTIONS



(a) original (zoom)



(b) OCR (zoom)

Figure – Zoom in on a text block + OCR (ABBYY)

TOOL DEVELOPED BY LRDE (BASED THE OCR ENGINE : PERO OCR)

PERO OCR : very efficient of historical sources

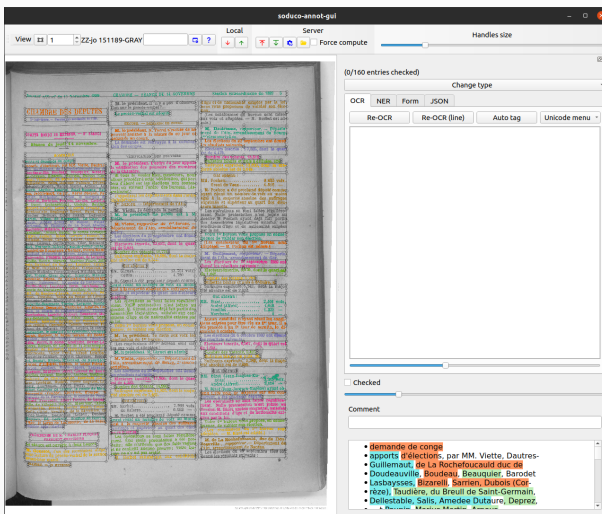
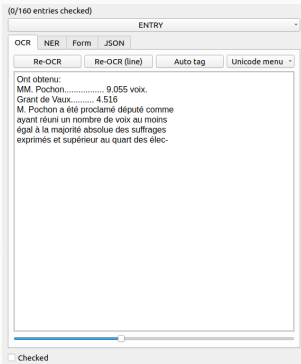
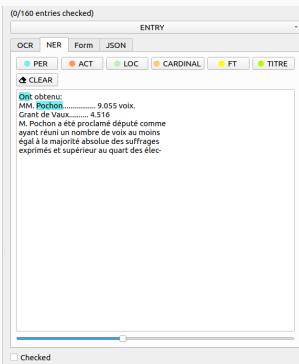


Figure – LRDE's Tool

TOOL DEVELOPED BY LRDE (BASED THE OCR ENGINE : PERO OCR)



(a) OCR



(b) NER

Figure – OCR and NER areas

2 kind of metrics :

- Unsupervised methods : Dictionary...
- Supervised methods
 - Bag Of Words : inherits from classical metrics, has limitations
 - Levenshtein's distance

Levenshtein's distance : Minimum number of insertions, deletions and substitutions on single characters to transform a text into another. Depending on the implementation, spaces can be taken into account or ignored.

OCR (string token)	Ground Truth	LEV _{1,1,1}
« des ates sont accomplis »	« Des actes sont accomplis »	3

Table – An example

$$\text{CharacterErrorRate(CER)} = \frac{\text{LEV}(\text{text}_{\text{gt}}, \text{text}_{\text{ocr}})}{\text{len}(\text{text}_{\text{gt}})}$$

$$\text{CharacterAccuracy} = \max(0, 1 - \text{CER})$$

OCR	Ground Truth	LEV	CharacterAccuracy
« des ates sont accomplis »	« Des actes sont accomplis »	3	0.88
« as aboue s0 belo »	« as above so below »	3	0.82

Table – Accuracy on 2 examples

Examples of libraries on python and in command line

- Jiwer : <https://pypi.org/project/jiwer>
- Fastwer : <https://pypi.org/project/fastwer>
- PRImA's tool : <https://www.primaresearch.org/tools/PerformanceEvaluation>

Different implementations :

1. Penalizing multiple white space
2. Pre-processing needed before the calculation of the metric
3. Possibility of not penalizing certain character
4. ...

RESULTS ON 1 PAGE OF DEBATE

OCR	Document	CharAccuracy
ocr-tesseract (eng,psm3)	26 Novembre 1889	0.93
ocr-tesseract (fr,psm3)	26 Novembre 1889	0.95
ABBY	26 Novembre 1889	0.22

- We see the main limitation of Levenshtein's distance : it's highly sensible of the detection order and layout. ABBY is overly penalized by the metric.
- Moreover, other problems aroused such as bad table detection
- We are interested in a metric more independent from the reading order

Solution : Flexible Character Accuracy (Clausner et al, 2020)

METRIC : FLEXIBLE CHARACTER ACCURACY

```
ce qui est en haut  
est comme  
ce qui est en bas
```

(a) Ground Truth

```
1 est comme  
2 ce qui est en bas  
3 ce qui est en haut
```

(b) OCR

Figure – Example A

```
Première Ligne  
row 2  
Third line  
Last row
```

(a) Ground Truth

```
Première Ligne   Third line  
row 2           Last row
```

(b) OCR

Figure – Example B

Exemple	CharAccuracy	FlexCharAccuracy
A	0.45	1
B	0.39	0.95

Table – Character Accuracy vs Flexible Character Accuracy

FLEXIBLE CHARACTER ACCURACY (CLAUSNER, 2020)

A metric, less sensible to reading order, over and under segmentation :

1. Split the two input text into line
2. Sort the ground truth text lines by length (in descending order)
3. For the first ground truth line, find the best matching OCR result line segment (by minimising a penalty that is partly based on string edit distance)
4. If full match :
 - Mark as done and remove line from list
 - Else subdivide and add to respective list of text lines; resort
5. If any more lines available repeat step 3
6. Count non-matched lines / strings as insertions or deletions (depending on origin : ground truth or result)
7. Sum up all partial edit distances and calculate overall character accuracy

ANOTHER EXAMPLE

```
1 " Stagiaires ..... 1.125
2 " 5e classe ..... 1.500
3 " 4e classe ..... 1.750
4 " 3e classe ..... 1.875
5 " 2e classe ..... 2.250
6 " 1re classe ..... 2.500
```

(a) Ground Truth

```
1 " Stagiaires .....
2 " 5e classe .....
3 " 4e classe .....
4 " 3e classe .....
5 " 2e classe .....
6 " 1re classe .....
7 1.125
8 1.600
9 1.750
10 1.875
11 2.250
12 2.500
```

(b) Example A

```
1 " Stagiaires .....
2 1.125
3 " 5e classe .....
4 1.600
5 " 4e classe .....
6 1.750
7 " 3e classe .....
8 1.875
9 " 2e classe .....
10 2.250
11 " 1re classe .....
12 2.500
```

(c) Example B

```
1 " Stagiaires .....
2 " 5e classe .....
3 " 4e classe .....
4 " 3e
5 classe .....
6 " 2e classe .....
7 " 1re classe .....
8 .....
9 1.125
10 1.500
11 1.750
12 1.875
13 2.250
14 2.500
```

(d) Example C

Figure – Our ground truth with 3 different OCR

Example	CharAccuracy	FlexCharAccuracy
A	0.59	0.96
B	0.84	0.96
C	0.58	0.96-

FLEXIBLE CHARACTER ACCURACY FOR 1 PAGE OF DEBATE

OCR	Document	CharAccuracy	FlexCharAccuracy
ocr-tesseract	1889/11/26	0.95	0.95
ABBY	1889/11/26	0.22	0.97

- .
- Results are similar to character accuracy when reading order is good
- Measurement of correct character detection somehow independent from the reading order
- Time consuming : according to the original paper it took around 17 s for one pages of 900 characters. Our debates have a average of 10000 characters per pages

- Post correction dictionary
- Using Regex
- Endogenous corrections

Results are to be improved

- OCR quality is highly impacted by the digitization process
- Post correction needed
- Defining a well suited metric

TOPIC MODELLING WITH LDA

LDA : Latent Dirichlet Allocation

- A corpus is written based on a finite number of topics. The number is a parameter of the model
- A Topic can be seen as a semantic field, a list of words linked by their meaning
- We want to find Topic as distributions over words and documents as a mixture of topics, conditional on the observed words
- Each document is then produced by selecting words from a subset of topics given a certain probability distribution that we aim to find

Bag of Words method, order of the word doesn't impact; OCR quality does.

PREPARATION STAGE

Pre-processing

1. Removing stop word
2. 50 Topics
3. Whole corpus is divided in similar size chunks, we hope to find a better way

	topic_0	topic_1	topic_2	topic_3	topic_4	topic_5	t
0	and	and	the	you	of	no	
1	was	the	of	to	had	gutenberg	
2	the	of	and	it	she	project	
3	miss	they	her	not	the	of	
4	in	to	to	be	his	the	
5	all	in	was	is	her	or	
6	of	as	his	but	to	is	
7	from	his	in	have	been	you	
8	elizabeth	their	she	will	in	indeed	
9	he	he	elizabeth	for	that	tm	
10	but	was	their	that	and	this	
11	to	were	by	if	was	work	
12	his	it	with	do	he	and	
13	cried	mr	were	the	as	works	
14	at	she	for	of	which	any	
15	an	had	on	we	not	terms	
16	am	all	which	me	could	foundation	
17	bingley	could	be	can	with	are	
18	it	that	from	so	mr	electronic	
19	bennet	not	as	must	have	all	

Figure – LDA over Pride and Prejudice, without removing stop words

4 STRAIGHTFORWARD TOPICS

Topic 8	Topic 11	Topic 13	Topic 15
salaire	général	adoption	pari
question	commission	absolue	télégraphe
gouvernement	régiment	ouvert	faire
jour	troupe	votant	ingénieur
patron	monsieur	majorité	train
chambre	année	nombre	ligne
droit	jeune	secrétaire	chambre
syndicat	temps	député	personnel
travail	soldat	article	employé
travaux	ministre	adopté	poste
grève	loi	vote	travaux
mineur	guerre	demande	agent
mine	service	chambre	ministre
loi	militaire	voix	fer
compagnie	officier	scrutin	chemin
ouvrier	armée	président	compagnie

HAND LABELLING OF TOPICS

Label	Topics	Examples of words	Weight in the corpus
Names of MPs	0,5,10,14,18, 23,35,37,39	Duval, Sigismond, Jules, Martin	0.101
government/parliament	1,6,9,13,17, 19,22,36,38, 41,45,46,49	tribune, projet, adoption, majorité	0.284
economy	2,4,16	agriculture, commerce, patente, betterave	0.069
working class	7,8,31,34	travailler, salaire, usine, mutuelle	0.070
army	11,48	général, régiment, contrôle, militaire	0.041
department	12	Calais, Alpes, Saône, Charente	0.006
trains/communications	15,44	télégraphe, ingénieur, train, travaux	0.069
local politics	20, 33	ville, arrondissement, local, département	0.030
law enforcement	21,40	police, préfet, tribunal, délit	0.055
school	24	lycée, faculté, classe, enfant	0.023
alcohol	25	bouilleur, degré, raisin, octroi	0.019
budget	26,29,30,43	chiffre, budget, dépense, exercice	0.097
colonies	28	métropole, juif, algérien, tonkin	0.018
navy	32	marin, flotte, mer, bâtiment	0.021
building works	27, 42	construction, théâtre, hectare, terrain	0.024
foreign affairs	47	puissance, Madagascar, Angleterre, traité	0.034

Figure – topic-classes

Topics a quite consistent with the major themes that marked political life during the early third republic

TOPIC EVOLUTION OVER TIME (BY MONTH)

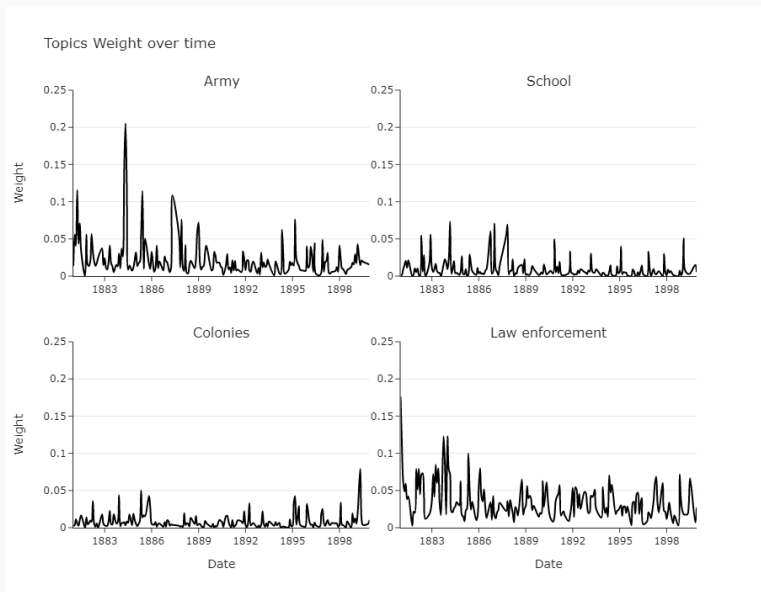
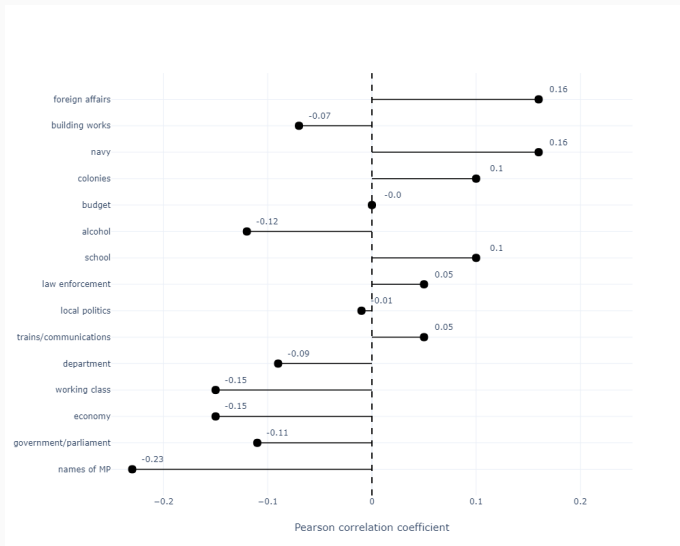


Figure – 4 topic and the evolution of their weight in the corpus

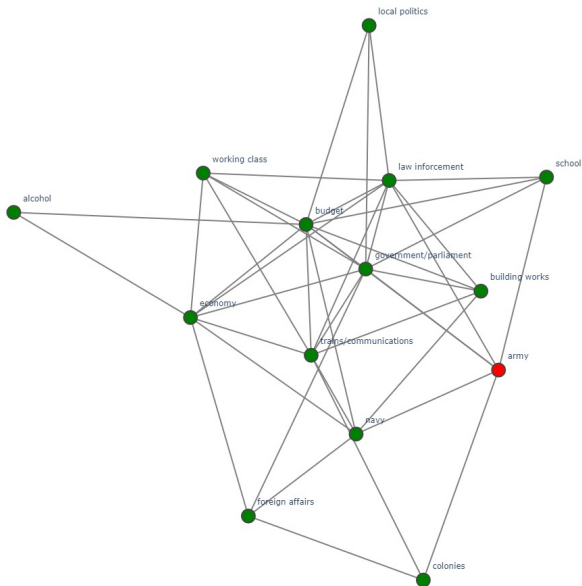
Popularity of the topic army at 4 four different times :

- In 1881 began the conquest of Tunisia
- The year 1884 was marked by the Tonkin campaign and discussion on reducing military service to 3 years
- In 1887 and 1888 : new discussions on the reform of military service.
- In 1895 the difficult conquest of Madagascar gave rise to new debates

CORRELATION WITH THE TOPIC ARMY



STRONGER CORRELATION



WORD EMBEDDING

A non bag of words method

LDA doesn't specify a notion of similarity between words and also between topics

Stop words have a huge impact, even if a Tf-idf is possible, we look for other possibilities

We look for a word encoding that capture a certain form of semantic similarity

- One-Hot Encoding of vectors (too simple)
- Singular Value Decomposition (too complex)

- Neural Network architecture with an hidden layer
- 2 methods : Continuous Bag of Words and Skip-Gram
- CBOW : predict the target word based on the context words
- Skip-Gram : predict the context words based on a target word

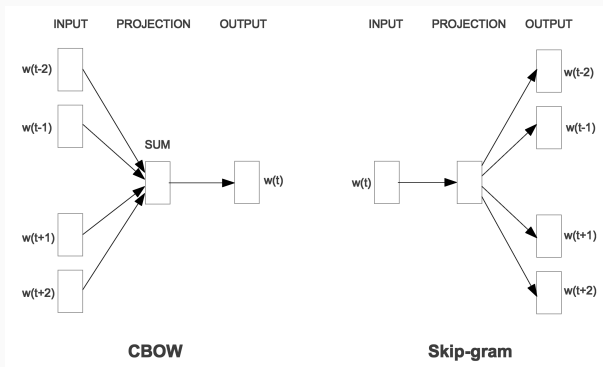


Figure – Word2Vec : CBOW et Skip-gram

FIRST RESULTS

église	train	alcool	colonie
presbytère	express	vin	colonial
cimetière	wagon	vinaigre	métropole
temple	locomotive	cidre	nies
cathédrale	omnibus	alcoolisés	marine
geneviève	véhicule	viné	lonies
basilique	voyageur	bière	cochinchine
chantre	fourgon	hectolitre	guyane
chapelain	trajet	poirés	algérie
sacristie	équipe	hydromel	rattachées
rabbin	voiture	hecto	nouvellecalédonie

Table – 4 words, with the top10 most similar words (cosine similarity)

SOME COMMENTS - 1

- A word is represented with a unique vector (size 300 typically)
- Vector coordinates do not have an intuitive interpretation
- A word (vector) is defined (learned) by the context
- The size of the context, i.e. the **Window**, is a parameter of the model

Boire (Window=1)	Boire (Window=15)
chercher	indécis
rire	vodka
trouver	allumer
changer	fuma
garder	commanda
entendre	pastis
marquer	coktail
choisir	bière
souffir	comptoire
accepter	champagne

Table – Closest words to « Boire » from a 1 and 15 window size model, French author corpus

- Vectors can encode more subtle semantic relationships : small is to big what smallest is to biggest (Mikilov et al, 2013)
- It's not systematic, both large corpus and dimensions are needed
- Analogy reported may be biased (Nissim et al, 2019) :
 - **Man** is to **King** what **Woman** is to **Queen** ? In fact King is the first answer returned, Queen comes in second place
- It does not work as well on our corpus
- Both dimension and corpus size need to be improved to keep increasing the quality of our vectors
- Skip-Gram has been shown to give better result with respect to the semantic analogies

We want to assign vectors to document

- Paragraph Vector : Distributed Memory model
- Paragraph Vector : Distributed Bag of Words

Creating a jointly embedded topic, documents and word vectors

The main assumption is that similar documents indicate a topic. A Topic can be seen as an area of close documents in the projected words/documents space

1. We start by finding vectors for words and documents (ex : Doc2Vec)
2. Similar documents will be close to each other and close to similar words
3. Find dense areas of documents, i.e. a topic
4. Topic vector is then computed as the centroid of the documents, the closest words form the vocabulary of the topic

We obtained a large number of Topics : 113. They are quite coherent and precise. We can study them in relation to each other, but also in relation with words and documents

SOME TOPICS OBTAINED

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

Table – Three clusters among 113 : storms (55), divorce (68) and the post office (70).

Topic 80	Topic 32	Topic 64	Topic 81
pedagogique	lycées	polytechnique	séminaristes
enseignement	enseignement	ecole	ecclesiastiques
aptitude	collèges	élèves	vocations
licencie	secondaire	cyr	prêtes
secondaire	proviseur	examens	clergé
université	internants	jeunes	paroissial
diplome	universités	destinent	caserne
bachelier	boursiers	caserne	sacerdoce
baccalauréat	aggregation	forestieres	dispense
infortunes	indissolubilité	études	religion
ravages	conjointes	bachelier	vouer

Table – 5 closest Topics to the word « Education »

Cluster 57	Cluster 97	Cluster 54
paquebots	trains	strategiques
postal	freins	loulan
escale	wagons	ligne
messageries	mecaniciens	chemins
antilles	signaux	timbre
transatlantiques	train	tronçons
west	voyageurs	kilomètres

Table – A cluster made a three topics

LIMITATIONS

By definition, each document is assigned to at most 1 Topic. Therefore we lose a property that LDA had to get a mixture of Topics for each document. Yet with the possibility to look for similarity between topics/words, topics/topics, document/topics we think we can retrieve some of the temporal information that we had with LDA

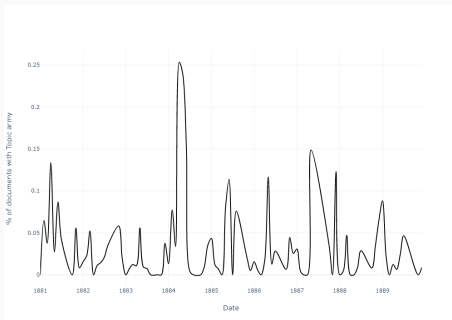
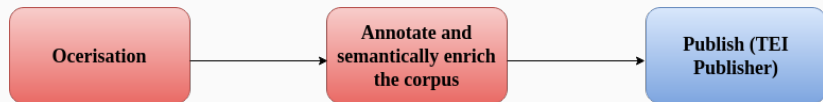


Figure – Evolution of the « Topic Army » We can find similarities with the previous graph

- First results obtained with LDA on our « raw » corpus
- Word, document and topic embedding can lead to more interpretations
- The analysis is limited by the quality of the OCR : we must improve that part





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

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

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

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