

Text classification on historical newspapers

How was the 1848 June Days uprising covered in the press?

Léo Dumont – leo.dumont@univ-paris1.fr

15 novembre 2024

Université Paris 1 – PIREH, Centre d'histoire du 19^e siècle

Historical issue

June Days were an uprising from 22 to 26 June 1848 in Paris.

- Study the event through the words used to describe it.
- Analyze the process of restoring order carried out by the triumphant republican authorities.

**In search of a narrated « law and
order fighter »**

What sources to study the discursive construction of the « law and order fighters » ?

- 1848 as an « unprecedented media explosion [...] in the history of the press (Ambroise-Rendu, 1999) »
- 450 newspaper creations

Characteristics of a daily newspaper in 1848

A newspaper issue contains a wide variety of content :

- news,
- serialized novels,
- transcription of parliamentary debates,
- letters,
- stock prices,
- advertisements,
- ...

Not always regular rules for layout.

How to find relevant content in the press ?

- Rely on digitized historical newspapers data and distant reading (Moretti, 2013) :
 - French National Library (BnF)
 - **RetroNews** : digital library dedicated to historical newspapers

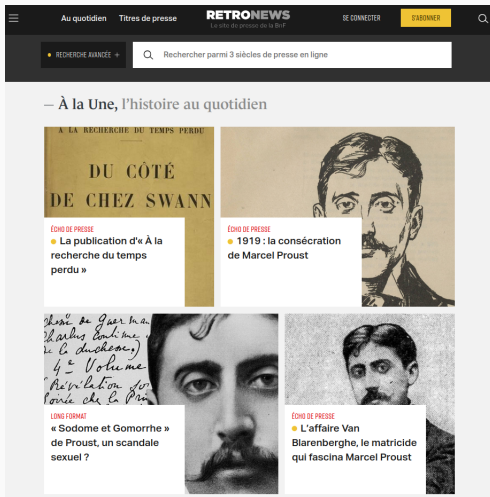


Figure 1 – RetroNews Homepage.

Digitized historical newspapers data challenges

Specific features of digitized newspapers :

OCR quality

- historical documents are still a challenge for OCR
- negative outcomes on NLP tasks

Digital materials granularity

- Two main levels of granularity (METS/ALTO XML standard) :
 - Physical Structure
 - Logical Structure

OCR quality

Le gouvernement a eu tort, dit-il, d'introduire dans le Mi-nitruv, pour le critiquer, le mot de socialisme: c'est un mot de la science qui devait être respecté. Après a été respecté, on a poursuivi de toutes parts les socialistes, la communion idéaliste; les journées de juin ont été dues à cette intolérance.



Le gouvernement a eu tort, dit-il, d'introduire dans le Mi-nitruv, pour le critiquer, le mot de socialisme: c'est un mot de la science qui devait être respecté. Après a rév -lutimon de Fév trierons a calomnié ue toues es parts les socialistes ou a poursuivi de toutes paris h s banque- , la communion s idéaliste; les journées de juin ont été dues à cette intolérance.

Figure 2 – Example of OCR quality.

A well known problem : (Mackovski, 2001), (Traub, 2015), (Mutuvi, 2018), (Hill, 2019), (Jiang, 2021), (Vitman, 2022).

Digital materials granularity

Two main levels of granularity with METS/ALTO XML standard :

- Physical Structure
→ pages
- Logical Structure
→ block layouts of text



Figure 3 – Logical Structure Detection.

Block layout is rarely a proper content.

Corpus construction

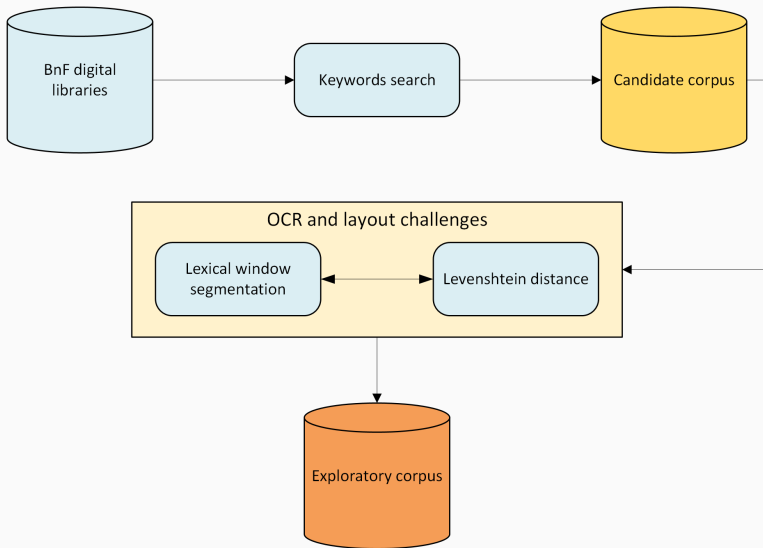


Figure 4 – Corpus construction flow chart.

Corpus construction steps :

1. Candidate corpus construction :

1.1 all occurrences of the phrase « journées de juin » in the national daily press (June 1848 – December 1852)

1.2 candidate corpus (6k results, 1.7 Gb json files)

2. Exploratory corpus construction :

2.1 Lexical window slicing of approximately 500 words around the searched keywords

2.2 1.7 million words to explore !

Previous research : (Wilson Black , 2022), (Ehrmann et al., 2022), (Beach & Walker Hanlon, 2023), (Owens & Padilla, 2021), (Sternfeld, 2011).

Corpus Exploration with Distant Reading

Which Distant Reading methodology ?

Two approaches have been chosen :

1. Supervised classification from a training dataset
2. Unsupervised classification from the labeled dataset

Each utilizing **BERT** based language models :

« fine-tuning pre-trained BERT on OCR'd texts significantly improves its resilience to OCR noise in classification tasks » (Jiang et al., 2021)

Corpus Exploration with Distant Reading

Supervised classification

Generating A Training dataset

Random Sampling of a **training dataset** (5% of the candidate corpus) :

- 373 documents
- 90k words

Iterative labeling heuristic trials :

- From 5 to 3 labels

Final tradeoff :

- Political debates (n = 185)
- Law and order fighters (n = 72)
- Insurgents (n = 116)

SetFit framework (Tunstall et al. 2022) :

- based on *Sentence Transformers*
- fast to fine-tune with few-shot
- prompt-free
- well-documented

Based-model used for sentence embeddings : **DistilCamemBERT**
(Delestre & Amar, 2022)

- distillation version of **CamemBERT** (Martin et al., 2020)

Model evaluation

- training dataset : 20 labeled examples per class (60 documents)
- test dataset : 313 documents

Class	Precision	Recall	F1-Score	Support
Political debates	0.92	0.80	0.86	96
Law and order fighters	0.86	0.96	0.91	165
Insurgents	0.85	0.75	0.80	52
Accuracy			0.88	313
Macro Avg	0.88	0.84	0.85	313
Weighted Avg	0.88	0.88	0.87	313

Table 1 – Classification Report

Corpus Exploration with Distant Reading

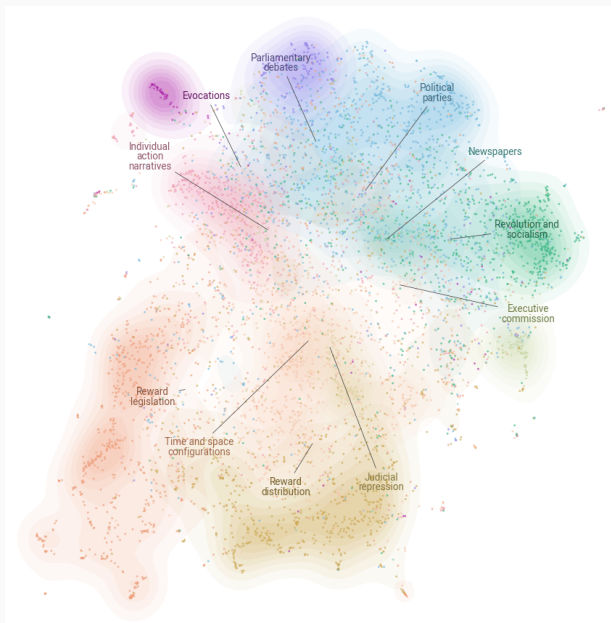
Unsupervised classification

BERTopic (Grootendorst, 2022)

Topic modeling approach based on 5 main steps :

1. Sentence embeddings
2. Dimensionality Reduction
3. Clustering
4. Tokenizer
5. Weighting scheme
6. Fine-tune topics representations

Documents and topics map



Topics and labels

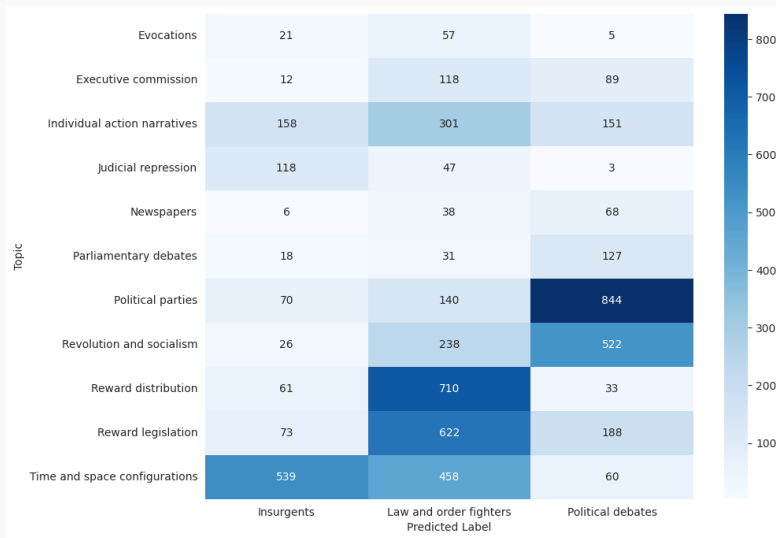


Figure 5 – Topics and predicted labels heatmap
(Pearson's Chi-squared test :X-squared = 3491.2, df = 20, p-value < 2.2e-16).

Conclusions

- Usable results with sentence-transformers based models and noisy textual data.
- A combination of supervised and unsupervised classification is a way to get a better sense of the discursive phenomena at work.
- Main future research directions :
 - expand the corpus,
 - OCR post-correction based on LLMs before the candidate corpus stage.