# DATA MINING HISTORICAL NEWSPAPERS METADATA

## Old News Teaches History

Jean-Philippe Moreux
Bibliothèque national de France,
Digitization dpt

**IFLA News Media Section,
Hamburg, April 2016**

# A True Story (@ BnF) about the Researchers' Needs

- How can we help a historian working on Stock Market quotes creation and development in French newspapers? (1800-1870)



here

# A True Story about the Researchers' Needs

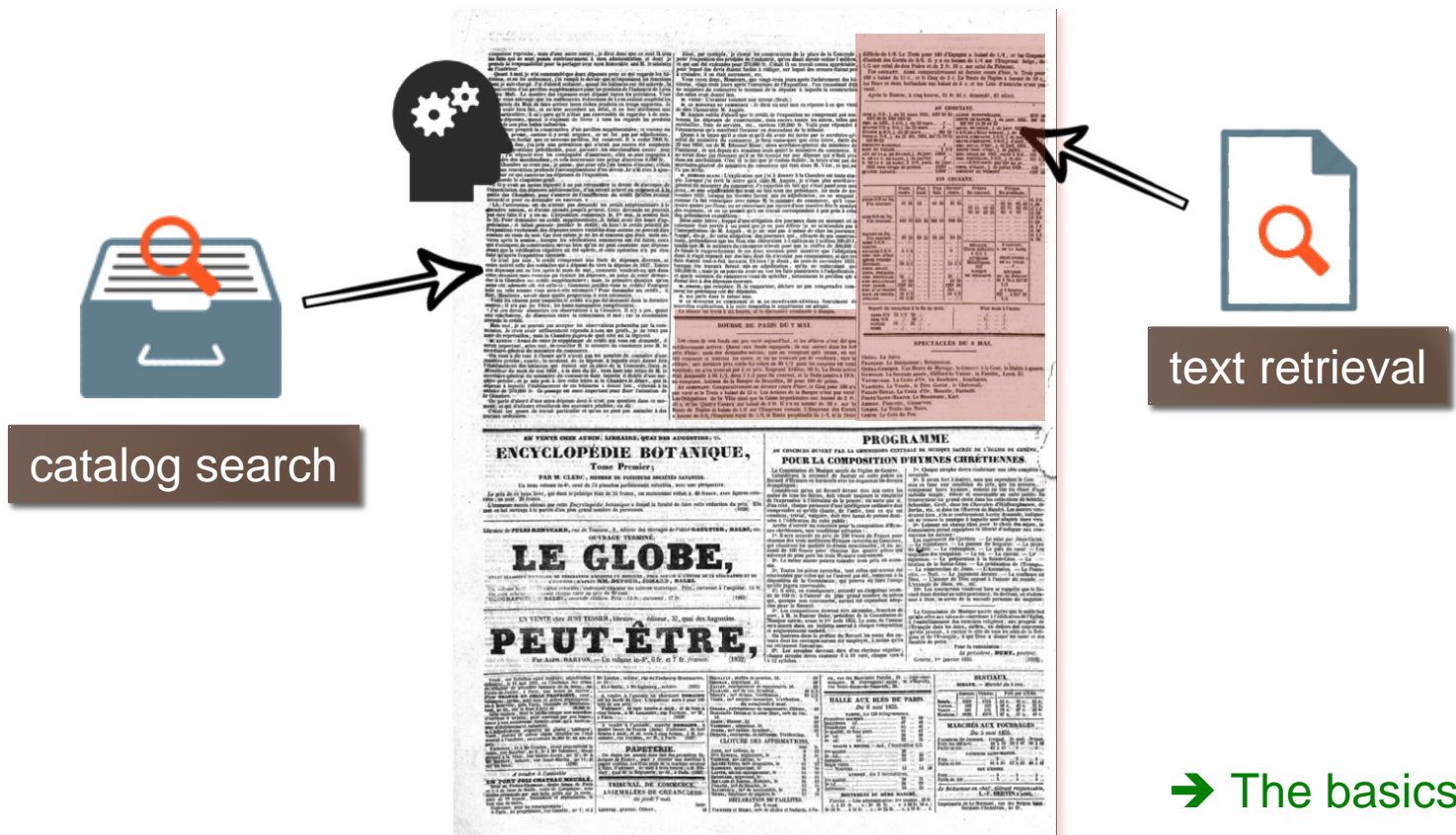- Obviously, he had to query the digital library catalog.



catalog search

# A True Story about the Researchers' Needs

- Moreover, he needed a text retrieval functionality.



catalog search

text retrieval

➔ The basics

# A True Story about the Needs of Researchers

- But is it enough? Could we do better?

text retrieval

catalog search

+ Corpora **builder**

+ Predefined qualitative
& easy-to-use corpora

+ **Advanced query**
on document
structure and layout
(to spot Stock Market
regions)

# The True Story (cont'd): unhappy Ending

"Stock Market quotes in French Newspapers (1801-1870)"
PhD in Communication and Information Science (P.-C. Langlais)

- **The creation of his corpus was very painful:**
  1. The historian had to script the DL to extract OCR and metadata from multiple newspaper titles.
  2. Then he had to refine/structure his text corpora.

## More than 100 Python scripts were needed!

☹

**Historians generally prefer to focus on research, not on writing scripts…**

# How to Satisfy Scientists' Needs?

Let's try to address this question, regarding the heritage daily corpus enriched during the Europeana Newspapers project:

- Feed the DL with <u>enriched</u> digital documents?
- Give end-users access to <u>quantitative metadata</u> describing documents structure and layout?
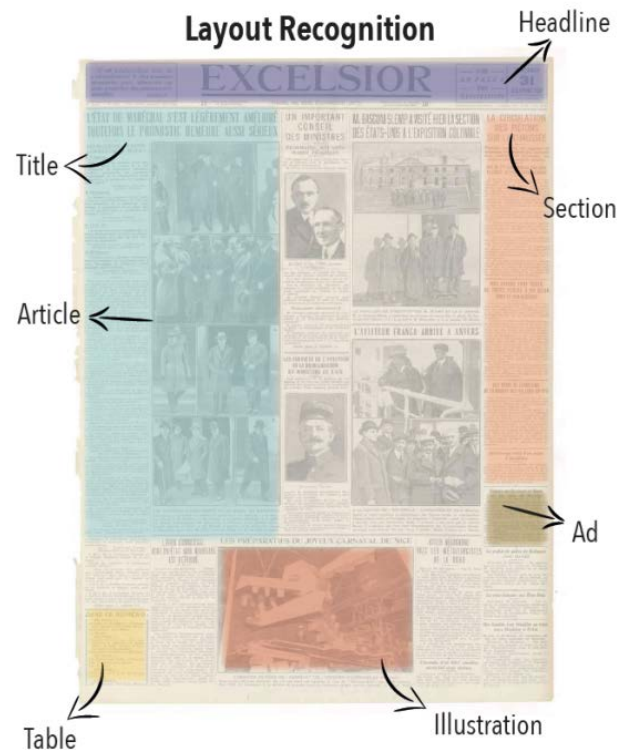- Give end-users an <u>ad hoc corpora builder</u> functionality?

**Plan**

1. The Europeana Newspapers test bed
2. Building a quantitative metadata dataset
3. Data mining and data visualization use-cases

# Enriching Digital Documents

- **Europeana Newspaper project** has enriched and aggregated millions of heritage newspapers pages with advanced refinement techniques like <u>Optical Layout Recognition</u> and <u>Named Entities Recognition</u>.

**Layout Recognition**

Headline

Title

Section

Article

Ad

Table

Illustration

**Europeana Newspapers project** (2012-2015): 11,5M OCR'ed pages, 2M OLR'ed pages from 14 European libraries

What is OLR?
- Identification of <u>structural</u> elements, including separation of <u>articles</u> and <u>sections.</u>
- Classification of <u>types of content</u> (ads, offers, obituaries…)

europeana newspapers

UIBK

CCS

# Document Analysis Technique like OLR Produce Quantitative Metadata

**The good new is OCR and OLR files are full of interesting objects tagged into the XML**:

- OCR (ALTO) is a source for quantitative metadata: number of words, illustrations & tables, paper format…

- OLR (METS) is a valuable source too for <u>high level informational objects</u>:

  - number of articles, titles, etc.

  - identification of sections (groups of articles)

  - content types classification (ads, judicial review, stock market…)

**Huge amount of valuable data for historians!**

# How to Build such Datasets?

- We have to count the number of objects in each page of the collection. Straightforward with XSLT, Java, Python, Perl, etc.
- We have to package and deliver these datasets to end-users.

**Europeana Newspapers project / BnF:** 880,000 OLR'ed pages from BnF newspapers collection, 6 titles, 1814-1944

〈BnF | Bibliothèque nationale de France



METS ALTO
documents
1 TB

→

CVS JSON XML
derived metadata
80 MB

→

Statistics

DB

**Pros**:
- Give to users <u>light</u> derived datasets, not TB of XML files!
- It's not rocket science.
- It's <u>fast</u> (2-3 h/title with an optimized NoXML parsing script)

**No Cons!**

# Who are the End-Users of the BnF Dataset?

- The EN-BnF dataset includes 5.5 M of values (150K issues, 880K p.)
- 7 metadata at issue level, 5 at page level
- XML, JSON or CSV formats

**Researchers** (Digital Humanities, History of Press, Information Science)

**Digital Curators & Mediators:** insights on the collections

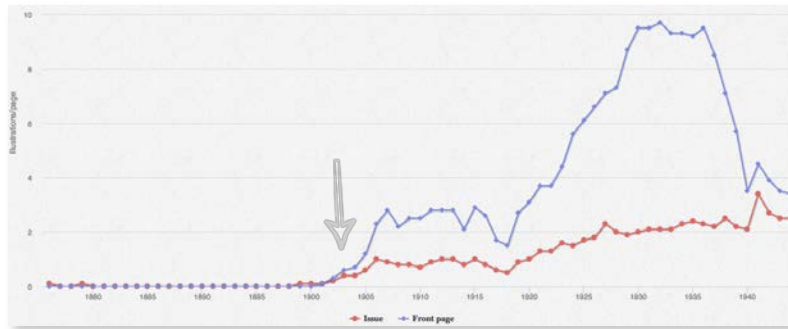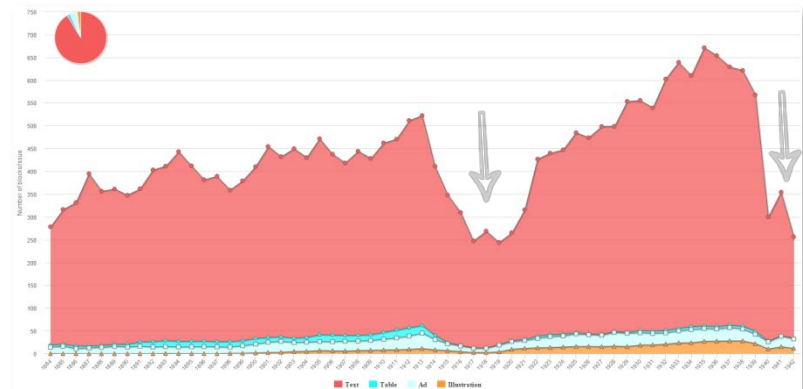**Digitization Program Managers:** statistics on digitized content

t o o l s

# Discovering Knowledge through Visualization

**Data visualization allows <u>researchers</u> to discover meaning and information hidden in large volumes of data**

tools

- **History of press/illustration:**
  Dataviz demonstrates the growing importance of illustration (blue: front page, red: inside pages).



- **History of press/activity:**
  Dataviz of types of content shows the impact of the Great War on the economical activity and assesses the period of return to pre-war level activity (roughly 10 years).



©Highcharts

# Engaging new Audiences with Dataviz

**Data visualization facilitates rediscovery and reappropriation of heritage documents (by the general public)**

tools

- Data visualization of illustrations density can reveal trends or outliers, like highly illustrated issues (illustr. suppl.) or the first published illustration in a title.

Facts extracted thanks to dataviz can then enrich other digital artefacts like **timelines.**

# Engaging new Audiences with Dataviz

Interactive chart of the word density reveals breaks
due to changes in layout & paper format, outlier issues…

tools



Journal des débats politiques et littéraires, 1814-1944, 45,334 issues displayed

➜Go beyond
keyword spotting
and page flip!

➜Some users
would like to play
with those charts!

# Requesting the Dataset

**Those datasets can be requested with dedicated tools**
(statistical environments, NoSQL or XML databases...)

- **Images search solution used by Gallica Mediation Service:**
  a XQuery HTTP API identifies "graphical" pages, that is to say both
  those poor in words and including illustrations.

tools

➔ "As a digital mediator, seeking for illustrations in our 12M p. collection is a nightmare…"

http://localhost:8984/rest?run=findIllustratedPages.xq&toDate=1920-01-01&toPage=1
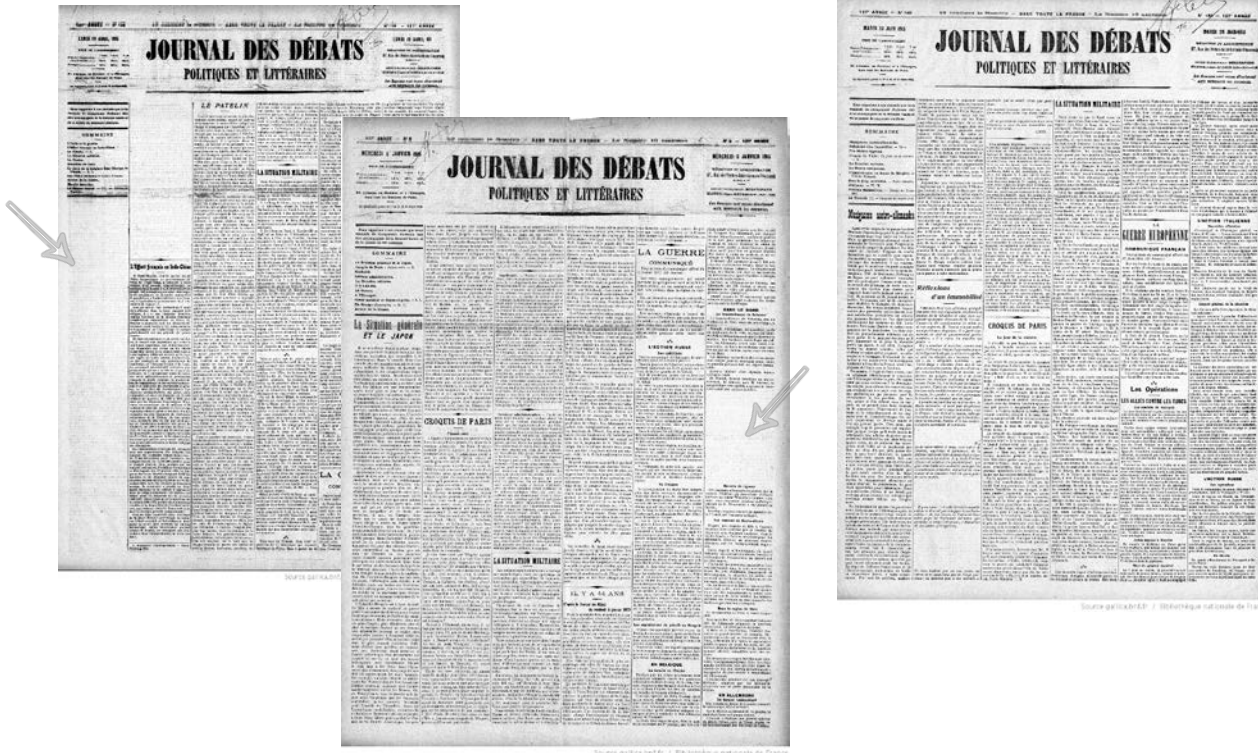
# Requesting the Dataset

- **Looking for WW1 censored front pages with BaseX:** XQueries
  can be written to dig into the data and find specific types of content, e.g.
  the front pages censored during the Great war, which have a slightly
  smaller words count than the front pages average.

tools



Is it effective?
- Recall rate: 45%
- Precision rate: 68%

(Based on a ground truth carried on the Journal des Débats front pages for 1915)
➔ Limits of a statistical approach when applied to a word based metric biased by layout singularities. Good enough for mediation:
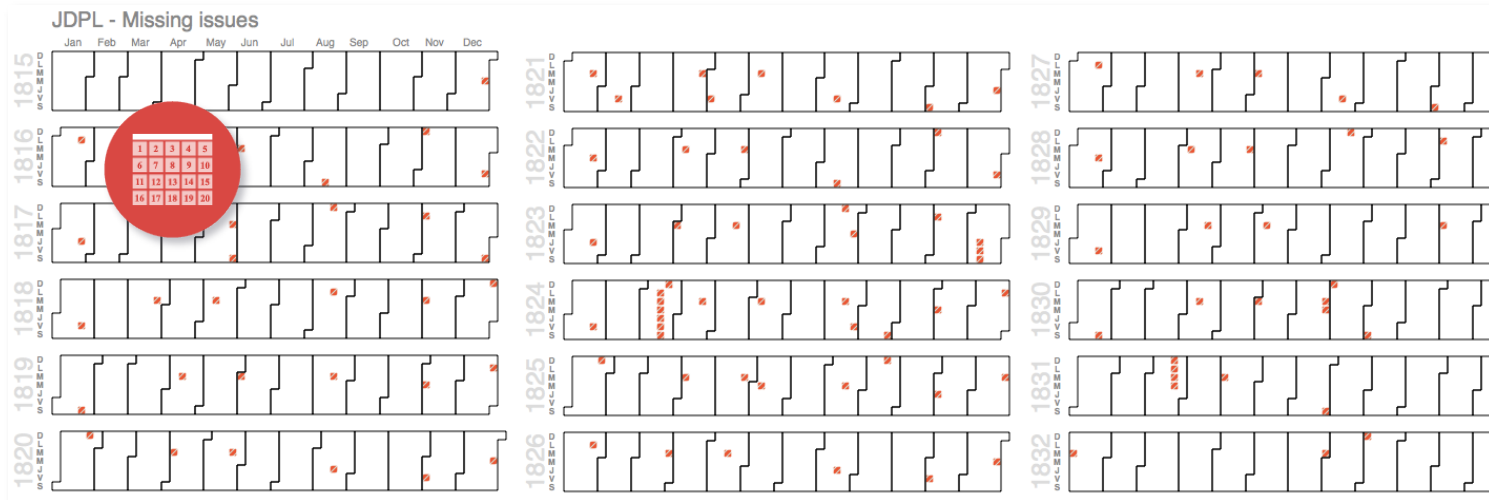Gallica blog post

# Are my Data Representative?

**The quality of datasets affects the validity of the analysis and interpretation.** Irregular data in nature or discontinuous in time may introduce bias. ➔ QA should be conducted.

<u>Data vizualisation</u> can contribute to quality control (and information of end-users)

- A <u>compact</u> calendar display of a title shows rare missing issues, which suggests that the digital collection is representative.
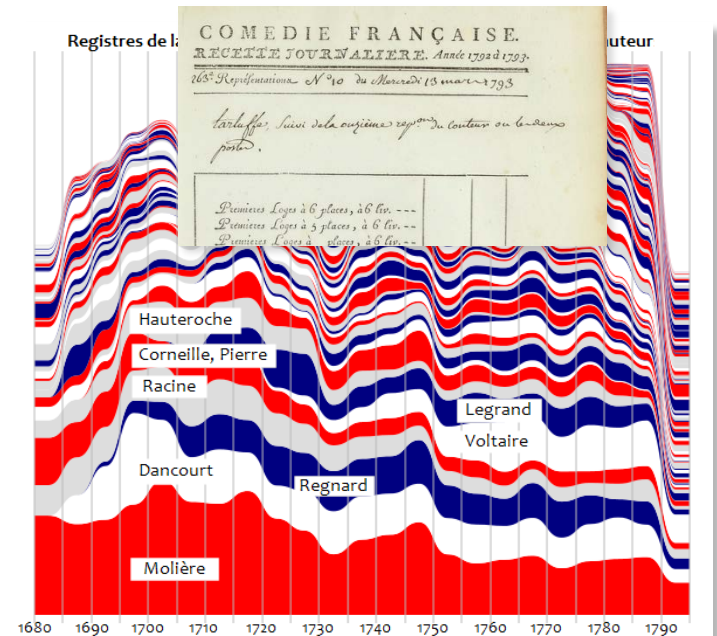


©Google Charts API

# Perspectives

- <u>Apply</u> the same data mining process to the other Europeana Newspapers OLR'ed datasets to produce more datasets. <u>Apply</u> on the on-going BnF newspapers digitization program.

- <u>Automatically build</u> the quantitative metadata datasets.

- <u>Experiment</u> on other types of materials with a temporal dimension (e.g. long life magazines or revues, early printed books).

- @BnF: Assess the opportunity of setting up a <u>data mining framework</u> targeting DH researchers ("Corpus" BnF research project, 2016-2018): Corpora builder? API? OCR dumps? Derived datasets? Remote processing?...

# Conclusion

- <u>Quantitative</u> metadata are relevant for all DLs' users: scientists, general public, institution employees

- Only <u>basic </u>data mining & dataviz methods and tools are needed.

- <u>OLR enrichment</u> provides a rich source of information for researchers.
Such data, possibly crossed with the OCRed text, usually provide a fertile ground for research hypotheses.

- Quantitative metadata is sometimes <u>enough</u>
to satisfy users. Example of a "pure"
quantitative metadata DH project  ➜

**"The Comédie-Française Registers"** project:
From 1680 until 1791, only one theater troupe
in Paris was allowed to perform the plays of
Molière, Corneille, Racine, Voltaire, Beaumarchais,
etc. This troupe played the works of these authors
over 34,000 times and kept detailed records
of their box office receipts for every single one
of those performances. (Partners: Paris-Sorbonne,
Harvard University, MIT)

©http://cfregisters.org/fr (Frédéric Glorieux)
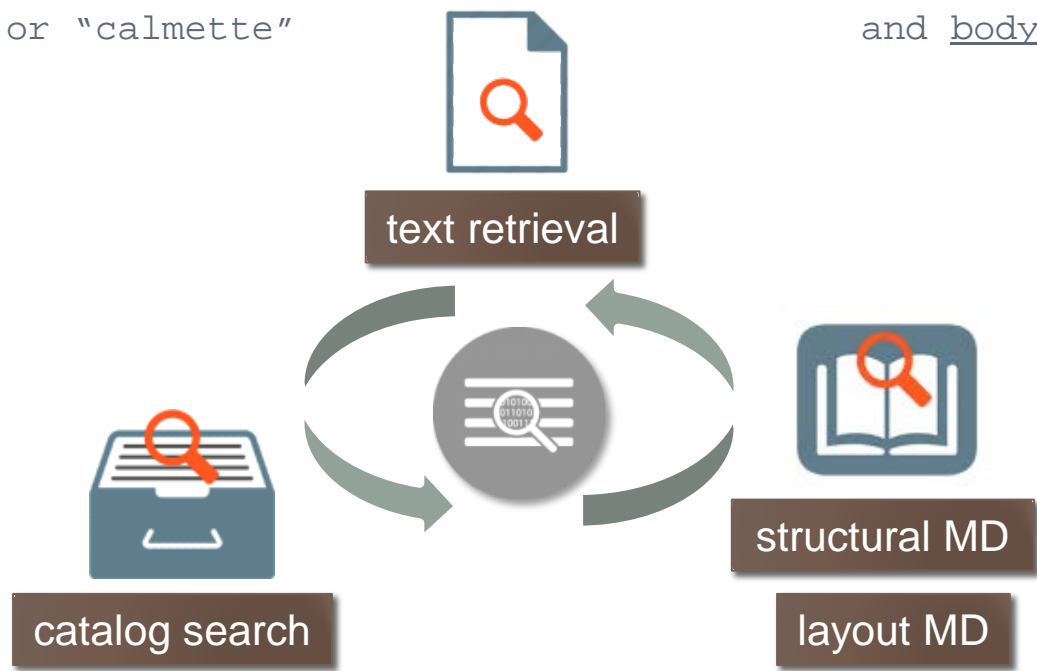
# Advanced Search in Newspapers?

- Feeding the search engine with layout and structural metadata will allow users to perform **advanced mixed queries**:

```
? illustrated articles
in Trial review section
from 1914 to 1916
where title contains
    "caillaux" or "calmette"
```

```
? articles with table
in Le Matin
where title contains
    "metal prices"
and body contains "gold"
```

text retrieval

catalog search

structural MD

layout MD

# Advanced Search in Newspapers?

- Feeding the search engine with layout and structural metadata will allow users to perform **advanced mixed queries**:

? illustrated articles ... articles with table in *Judicial review* ... *Matin* from 1914 to 1916 ... title contains where title conta ... tal prices" "caillaux" or ... y contains "gold"



→**Trove Advanced Search**

**Article Category**
Return only items in these categories

- ☑ Article
- ☐ Advertising
- ☐ Detailed Lists, Results, Guides
- ☐ Family Notices
- ☐ Literature

**Article Length**
Limit responses to articles of a particular length

- ● All
- ○ <100 Words
- ○ 100 - 1000 Words
- ○ 1000+ Words

**Illustrated Articles**
Limit to articles with or without illustrations

- ○ All
- ● Restrict to illustrated articles only
- ○ Restrict to articles without illustrations

**Sort Order**
Select how you would like your results sorted
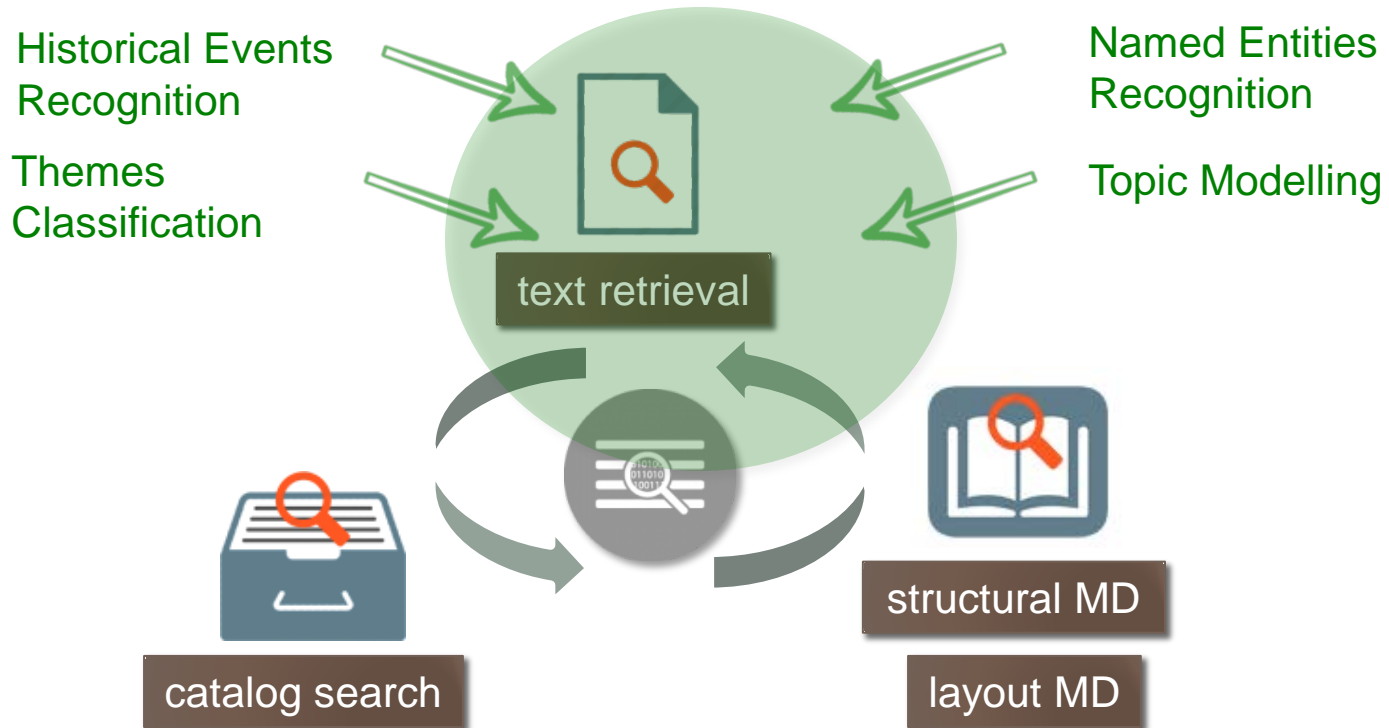
By Relevance

http://trove.nla.gov.au

catalog search

layout MD

# Advanced Search in Newspapers?

- Adding a pinch of **semantic** flavor to get closer to natural language query:

```
I'm looking for illustrated articles on front page in Trial topic
from 1914 to 1916 which contain NE.person "Henriette Caillaux"
or "Gaston Calmette"
```

Historical Events
Recognition

Themes
Classification

Named Entities
Recognition

Topic Modelling

text retrieval

catalog search

structural MD

layout MD

# Advanced Search in Newspapers?

- Adding a slice of **semantic** flavor to get closer to natural language query:

I'm ... al topic
fro ... laux"
or ...

➜ **RetroNews Advanced Search**

**484 résultats**

FILTRER VOTRE RECHERCHE ✕

Par titre de presse ▶

Par type de presse ▶

Par périodicité ▶

Par date de publication ▼

1914-1916 ✕
1914 (402)
1915 (33)
1916 (49)

Par lieu de publication ▶

**SERVICES PREMIUM**
AFFINER VOTRE RECHERCHE

Par thématique ▶

Par évènement ▶

Par sujet ▶

Par personne ▼

CAILLAUX (237)
CALMETTE (101)
POINCARÉ (88)
LABORI (70)
JAURÈS (65)
⋯ Plus de personnes

**Faceted search**: dates, NE, themes, events, topics…

**LE GAULOIS**
N° 13430 P.1
22 JUILLET 1914

prie mon excellent collaborateur M. Félix Belle de donner la physionomie de l'audience. Jusqu'au prononcé du verdict, je rougirais d'écrire une ligne qui fût de nature à impressionner le jury. Après le verdict, nous aurons tout le temps d'en dire ici notre sentiment.

Mais le spectateur n'a pas les mêmes devoirs impérieux que la neutralité impose à l'homme politique. Une chose m'a frappé à l'audience d'hier : c'est le subit apparat que le président des assises a fait déployer pour la comparution de M. Caillaux. Les postes ont été doublés ; des instructions très sévères ont été données au lieutenant de la garde républicaine ; le greffier a reçu l'ordre de lire un texte visant la répression de toute manifestation qui viendrait à se produire. Je sais qu'il fallait faire respecter la majesté du prétoire. Cela n'en créait pas moins une exception qui n'avait rien d'absolument démocratique et marquait bien aux yeux des jurés qu'un personnage d'importance dans la république allait être entendu. Ceux qui en auraient encore douté n'ont eu qu'à voir la dé-

**LE MATIN**
N° 10984 P.1
25 MARS 1914

Mme Caillaux raconte le meurtre — UN CRIME EN WAG...

**LE JOURNAL**
N° 7843 P.1
18 MARS 1914

LE JOURNAL
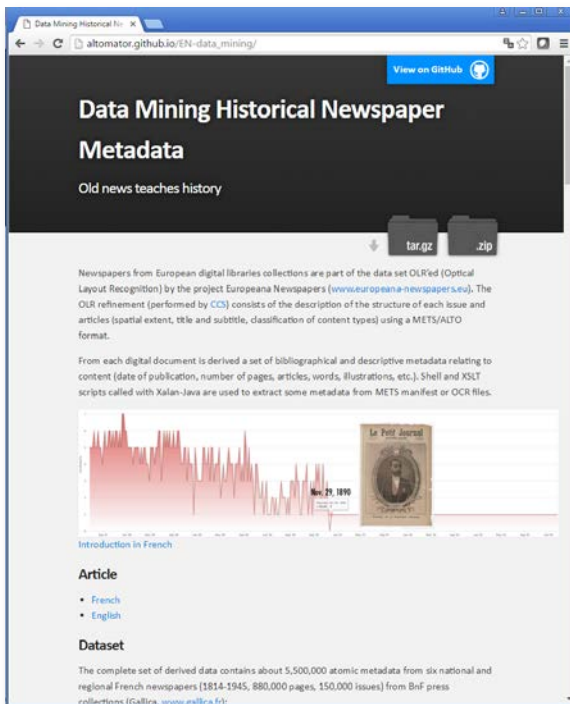Après le Meurtre de M. Calm...

**L'ÉCHO DE PARIS**
N° 10935 P.1
22 JUILLET 1914

galère ?
— M. Albanel croyais pourtant bien avoir conjuré les mauvais sorts, les manifestations et tous les tumultes par une innova-

...ntities
...ion
...delling

http://www.retronews.fr

# Thank you for your attention!

- Dataset (CSV, XML, JSON) and charts are publicly available. Just play with it! (no language barrier: not a single word of French inside)



Thanks to all
the EN partners!

http://altomator.github.io/EN-data_mining

# The True Story (cont'd)

- Could we have helped thim?

tools

**OLR facilitates the corpus creation task** ☺

→ Content Types classification, Section identification
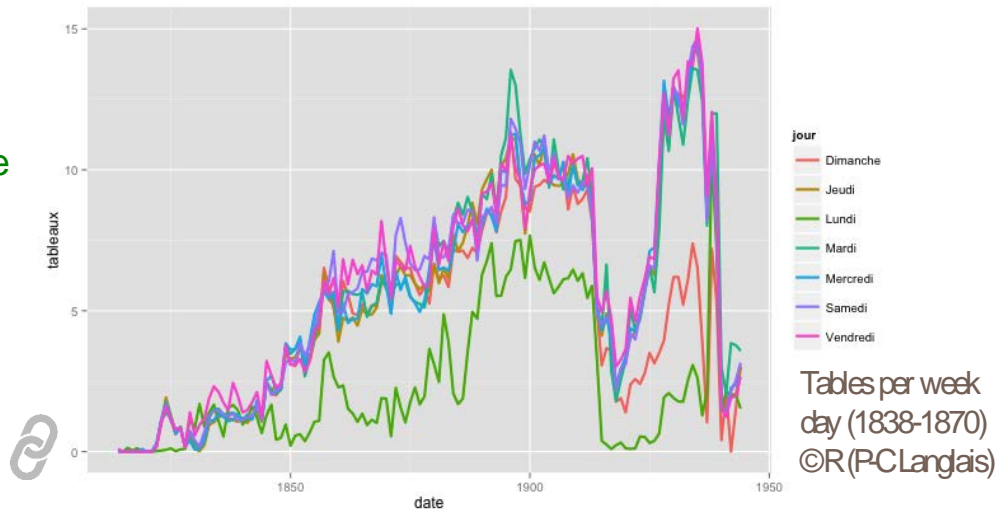
Types of content are tagged

**The quantitative dataset is of a great help** ☺

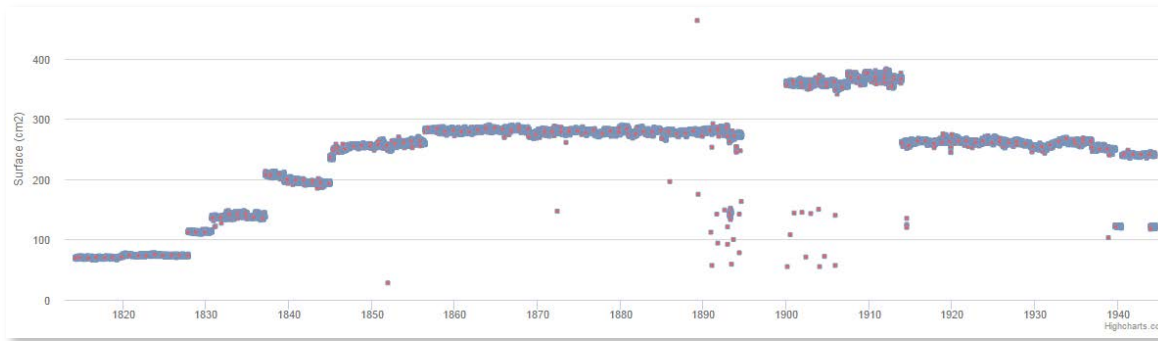"Tables" in newspapers are predominantly used in Stock Market quotes → instant use of this metadata!

Tables per week day (1838-1870) ©R (P-C Langlais)

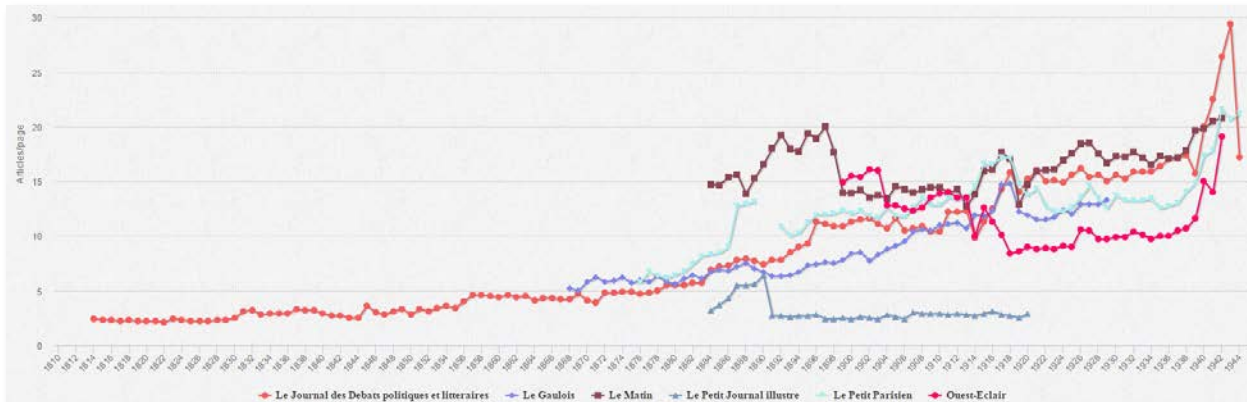# Discovering Knowledge through Visualization

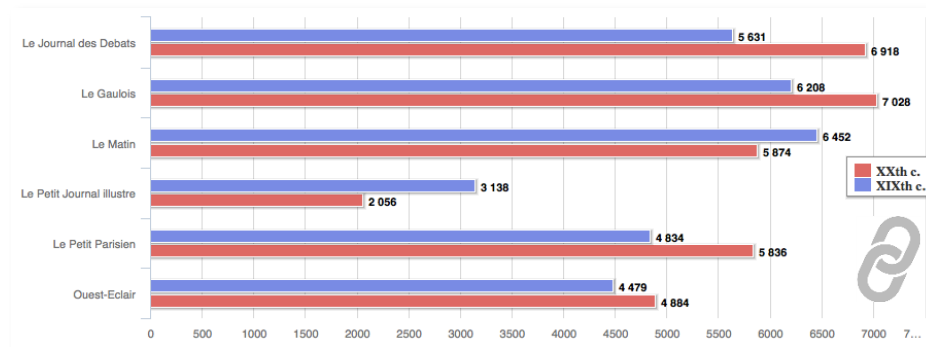- **History of press/page format:** Digital archeology of papermaking and printing.



tools

- **History of press/layout:** Visualization of the articles density per page reveals the shift from XVIIth "gazettes" to modern daily.
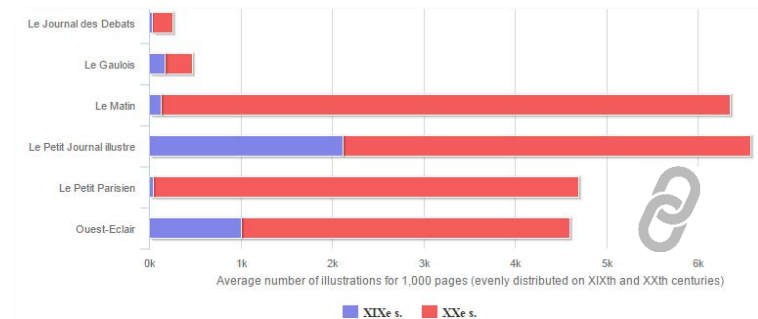


©Highcharts

# Other Users might be Interested by those Metadata: Digitization dpt

**Statistical information on digitized content for <u>project managers</u>.**
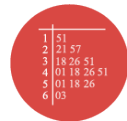
tools

- **OCR Crowdsourcing project:** What is the average density in words of these documents? What text correction efforts will be required?



- **Image bank**: What titles contain illustrations? What is the total number of images one can expect?
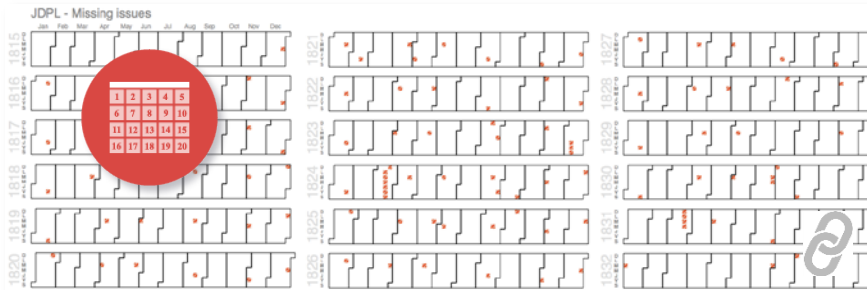


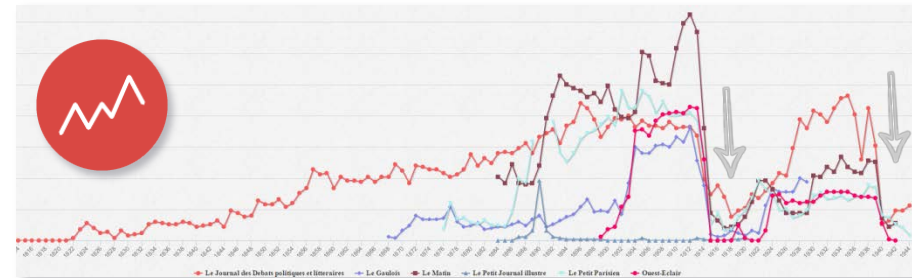©Highcharts

# Are my Data Representative?

**The quality of datasets affects the validity of the analysis and interpretation.** Irregular data in nature or discontinuous in time may introduce bias. ➔A qualitative assessment should be conducted.

<u>Data vizualisation</u> can contribute to quality control (and information of end-users)

- A calendar display of a title data shows rare missing issues, which suggests that the digital collection is representative.



©Google Charts API

- Stock Market quotes study based on the content tagged "table": one can empirically validate this hypothesis by the sudden inflections recorded in 1914 and 1939 for all titles, being known and established the historical fact of the virtual halt of trading during the two World Wars



©Highcharts