

Les optiques de la lecture distante : de l'occultation à la déconstruction de l'archive

Newseye — 15 mars 2021

Pierre-Carl Langlais
@Dorialexander
Alexander Doria (Wikipedia)

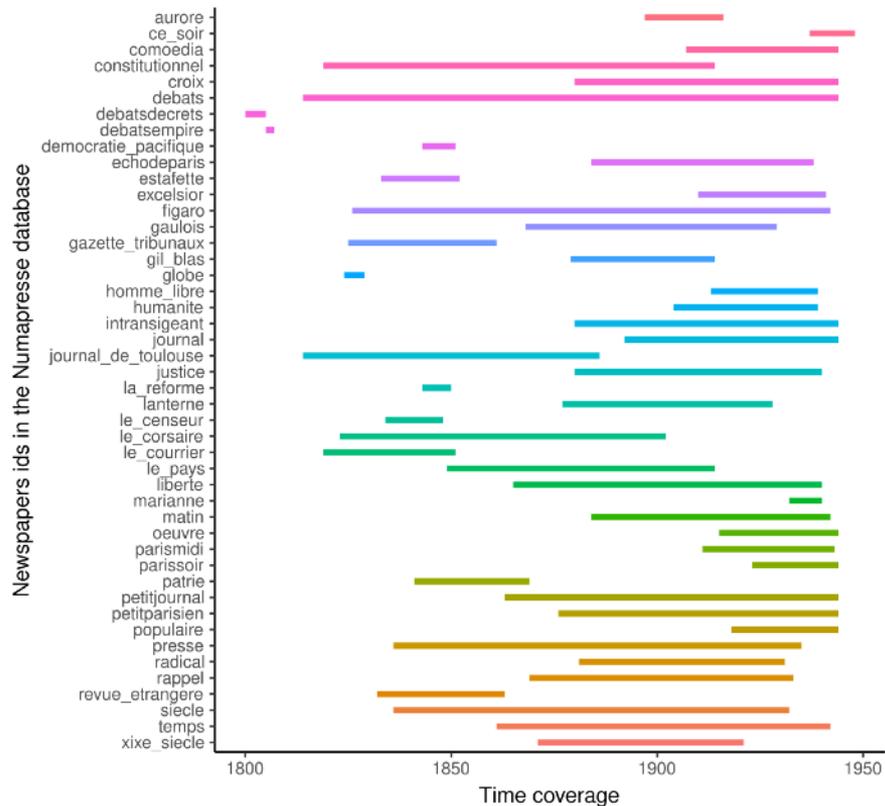




Les corpus de Numapresse

Depuis septembre 2019, Numapresse dispose des archives numériques détaillées de l'ensemble de la presse quotidienne nationale française disponibles sur Gallica.

Même si les collections numérisées restent incomplètes cette intégration à grande échelle permet d'élargir le champ de nos investigation : notre objet n'est plus seulement un journal mais un écosystème médiatique



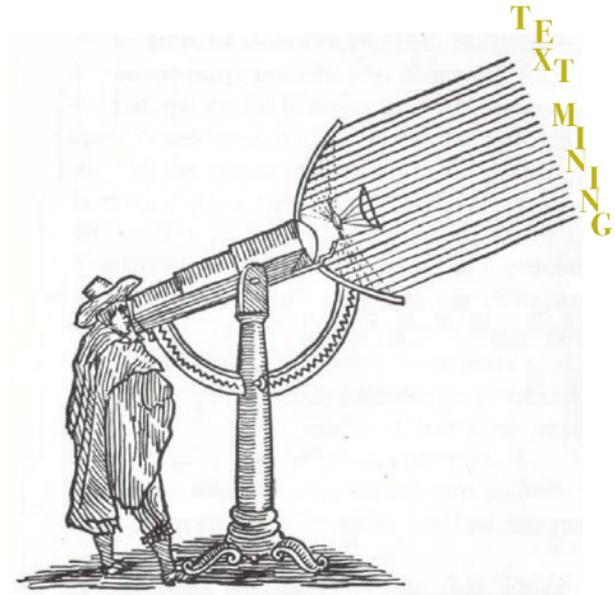


De la lecture comme un art optique

Les principaux quotidiens étudiés ont été publiés continuellement sur de très longues périodes (jusqu'à 145 ans pour le *Journal des débats*). L'évolution des formes culturelles s'intègre dans plusieurs temporalités :

- Quelques semaines pour la dissémination d'une nouvelle ou pour des collaboration de courtes durées
- Quelques mois pour les traits structurels de l'écosystème médiatique
- Quelques années pour le fonctionnement d'une rédaction.
- Quelques décennies pour l'émergence de genres journalistiques.

Pour tout prendre en compte, nous allons devoir varier nos focales...





Plan

Une présentation en trois parties :

- Lost in big data : dynamiques d'invisibilisations dans le patrimoine numérisé
- Voir l'invisible à distance : de l'archive-document à l'archive-réseau

1. Lost in big data

Dynamiques d'invisibilisations dans le patrimoine numérisé



Un patrimoine devenu “big data”...

En vingt ans, bibliothèques et acteurs privés ont développé de grands programmes de numérisation. Ces corpus sont de plus en plus disponibles pour des analyses quantitatives massives. La BNF a ainsi mis en ligne 1,5 to d'archives xml utilisées par le projet Europeana Newspaper ou 200 000 monographies utilisées par le projet OBVIL.

Documents de presse numérisés en mode « OCR » du projet Europeana Newspapers

Mots-clés: Presse, gallica, Europeana Newspapers

Ce jeu de données contient les documents numériques des collections de presse traitées durant le projet européen Europe: une reconnaissance du texte (OCR, *optical char*

Sommaire

- [Contenu du jeu de données](#)
- [Contexte de production](#)
- [Formats du jeu de données](#)
- [Exemples d'utilisation](#)
- [API et jeux de données en relation](#)

Contenu du jeu de données

Ce jeu contient la transcription réalisée par OCR d'environ 200 000 collections de presse de Gallica traitées durant le projet E

Tous les documents numérisés des titres suivants sont pr

- Le Figaro
- L'Echo de Paris
- L'Univers
- La Presse
- L'Humanité
- Le Constitutionnel
- Le Petit Journal
- Le Siècle
- L'Action Française
- L'Intransigeant

Fiche technique

Date de création ou de mise à jour : 2015

Quantité : 1 267 500 pages, 275 000 fascicules

Langue :

[TEI](#) [JSON](#) [Reconnaissance automatique des caractères \(OCR\)](#) [Textes](#)

DOCUMENTS DE GALICA PRODUITS AU FORMAT TEI PAR OBVIL

Présentation

Ce jeu de données contient le mode texte des documents de Gallica traités par l'Observatoire de la vie littéraire (Labex [OBVIL](#)). Le corpus est en français, issu majoritairement de l'édition du XIXe siècle.

N°	Auteur	Titre	Publié	Présenté	Alignement	Requ. sup.	Part.	Prés. sup.	Editeur
1	Gallica	Le Constitutionnel	1815	3245	105	2071	7	2,8 ans	L4801/0294
2	Gallica	L'Echo de Paris	1829	6591	72	2211	40	3,7 ans	L4801/0294
3	Gallica	Le Siècle	1825	1001	91	1861	28	3,2 ans	L4801/0294
4	Gallica	L'Intransigeant	1871	1001	74	2201	40	3,3 ans	L4801/0294
5	Gallica	Le Figaro	1826	4715	238	1871	16	2,7 ans	L4801/0294
6	Gallica	L'Univers	1826	14315	688	2704	12	2,5 ans	L4801/0294
7	Gallica	Le Petit Journal	1863	14611	688	2201	32	2,8 ans	L4801/0294
8	Gallica	Le Constitutionnel	1815	3245	105	2071	7	2,8 ans	L4801/0294
9	Gallica	L'Echo de Paris	1829	6591	72	2211	40	3,7 ans	L4801/0294
10	Gallica	Le Siècle	1825	1001	91	1861	28	3,2 ans	L4801/0294
11	Gallica	L'Intransigeant	1871	1001	74	2201	40	3,3 ans	L4801/0294
12	Gallica	Le Figaro	1826	4715	238	1871	16	2,7 ans	L4801/0294
13	Gallica	L'Univers	1826	14315	688	2704	12	2,5 ans	L4801/0294
14	Gallica	Le Petit Journal	1863	14611	688	2201	32	2,8 ans	L4801/0294
15	Gallica	Le Constitutionnel	1815	3245	105	2071	7	2,8 ans	L4801/0294
16	Gallica	L'Echo de Paris	1829	6591	72	2211	40	3,7 ans	L4801/0294
17	Gallica	Le Siècle	1825	1001	91	1861	28	3,2 ans	L4801/0294
18	Gallica	L'Intransigeant	1871	1001	74	2201	40	3,3 ans	L4801/0294
19	Gallica	Le Figaro	1826	4715	238	1871	16	2,7 ans	L4801/0294
20	Gallica	L'Univers	1826	14315	688	2704	12	2,5 ans	L4801/0294

TÉLÉCHARGER

Classement alphabétique (13,6 Go)

Classement Dewey (8 Go)

Classement par siècle (13.8 Go)

Métadonnées complètes (8 Mo)

FICHE TECHNIQUE

Date de mise en ligne 2018

Format

[TEI](#) [JSON](#)

Licence

[Conditions d'utilisation des contenus de Gallica](#)

Technologies

[Reconnaissance automatique des caractères \(OCR\)](#)

Sujets

[Textes](#)



La fascination du *big data*...

De grandes ambitions et... beaucoup de problèmes. Les promesses inabouties de Ngram Viewer et des culturonomics.

Google Books Ngram Viewer

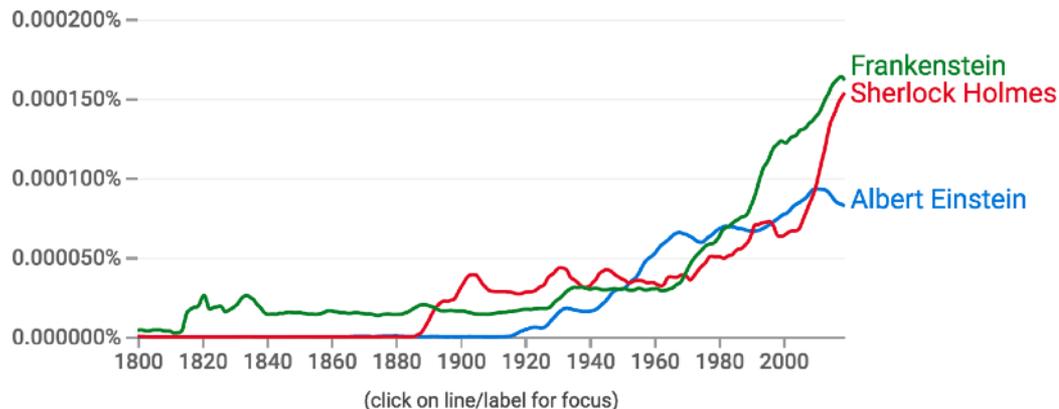
Albert Einstein, Sherlock Holmes, Frankenstein

1800 - 2019

English (2019)

Case-Insensitive

Smoothing





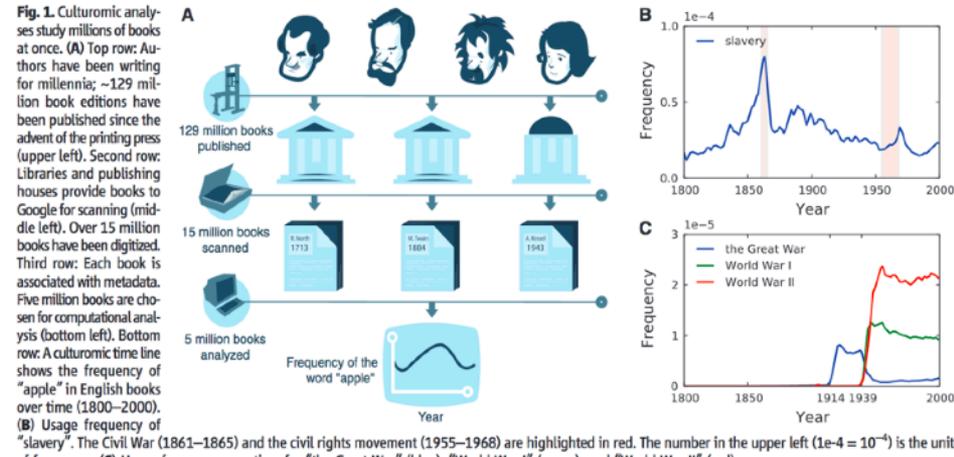
La fascination du *big data*...

« Les *culturomics* désignent l'application de la collecte et de l'analyse de données massive à l'étude de la culture humaine. Les livres sont un début, mais nous devons également intégrer les journaux, les manuscrits, les cartes, les œuvres d'art et une myriade d'autres créations humaines »

Quantitative Analysis of Culture Using Millions of Digitized Books

Jean-Baptiste Michel,^{1,2,3,4,5*}† Yuan Kui Shen,^{2,6,7} Aviva Presser Aiden,^{2,6,8} Adrian Veres,^{2,6,9} Matthew K. Gray,¹⁰ The Google Books Team,¹⁰ Joseph P. Pickett,¹¹ Dale Hoiberg,¹² Dan Clancy,¹⁰ Peter Norvig,¹⁰ Jon Orwant,¹⁰ Steven Pinker,⁵ Martin A. Nowak,^{1,13,14} Erez Lieberman Aiden^{1,2,6,14,15,16,17*}†

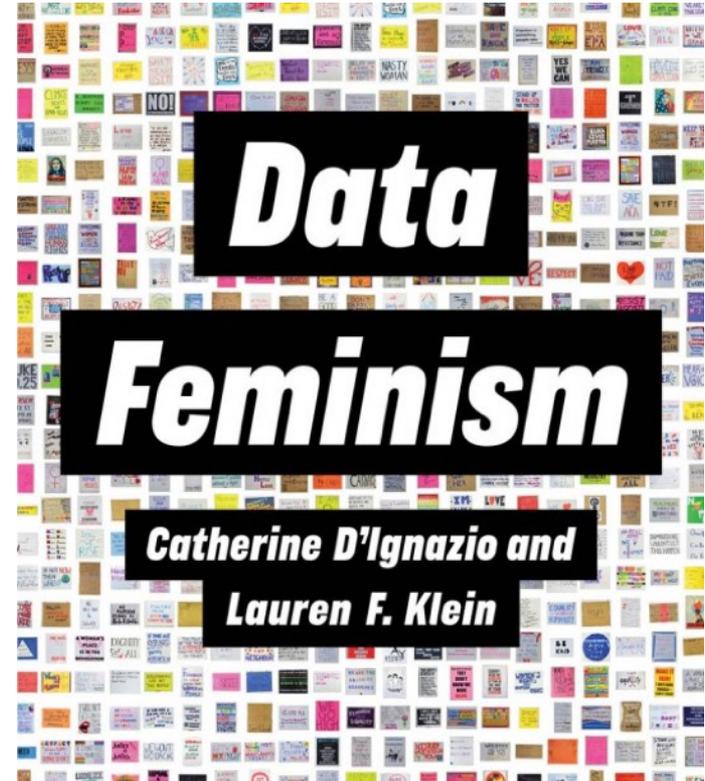
We constructed a corpus of digitized texts containing about 4% of all books ever printed. Analysis of this corpus enables us to investigate cultural trends quantitatively. We survey the vast terrain of 'culturomics,' focusing on linguistic and cultural phenomena that were reflected in the English language between





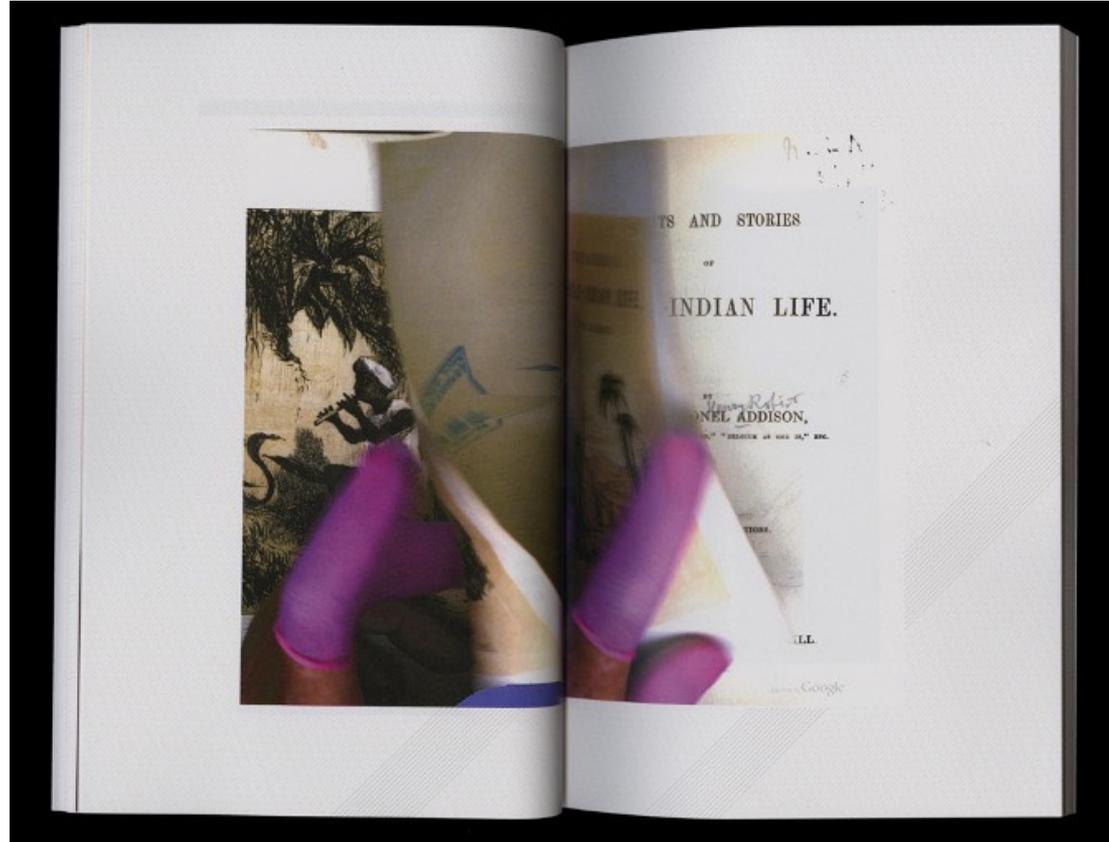
La fascination du *big data*...

« Big Dick Data projects fetishise large size and prioritise it, along with speed, over quality, ignore context and inflate their technical capabilities. They also tend to have little consideration for inequalities or inclusion in the process. Mark Zuckerberg aiming to supersede human senses with AI might be considered one such project, along with software company Palantir's claims about massive-scale datasets. Big Dick Data projects aren't necessarily wholly invalid, but they suck up resources that could be given smaller, more inclusive projects. »



...face à l'écueil de la matérialité

La présence occasionnelle de « mains » dans Google Books souligne que la numérisation n'est pas un bouton magique mais un processus complexe trop souvent dissimulé (avec un recours fréquent au « digital labor »)



...face à l'écueil de la matérialité

Un récent tweet viral d'Internet Archive a levé un coin de voile sur les conditions de production de l'archive numérisée et du *digital textual labor* avec des divisions du travail sociales, genrées et ethniques.



At the Internet Archive, this is how we digitize a book.

We never destroy a book by cutting off its binding. Instead, we digitize it the hard way--one page at a time.

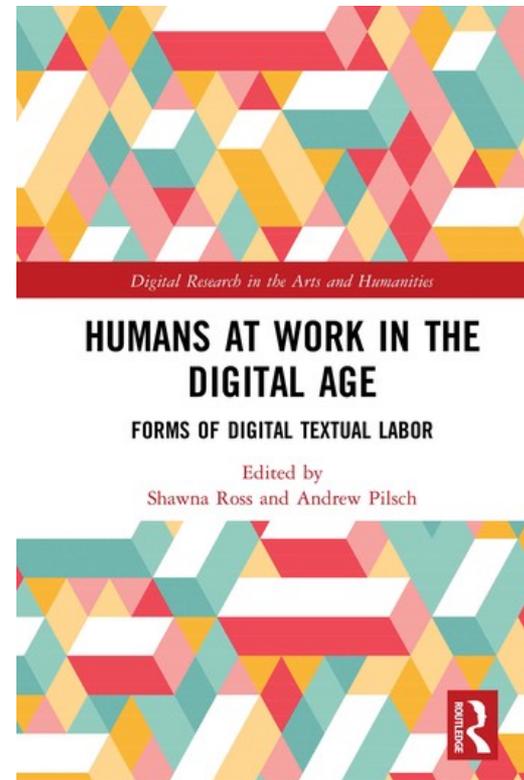
[#digitalbooks](#)

[Traduire le Tweet](#)



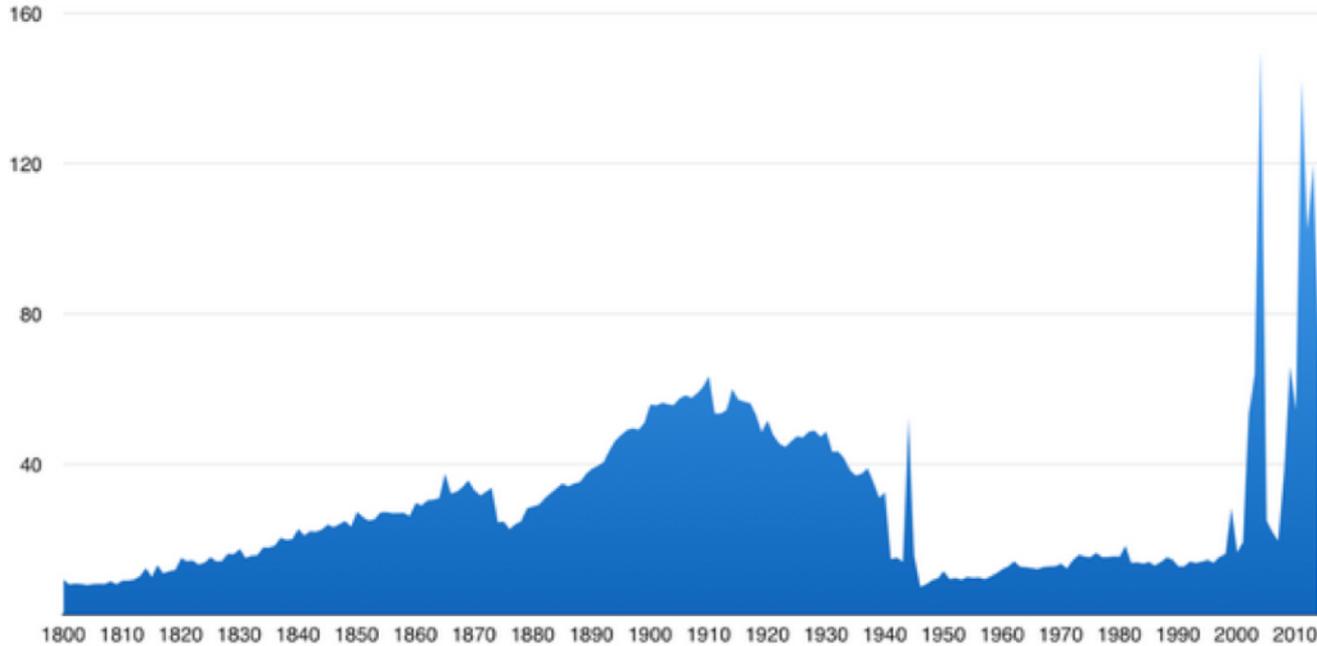
5:31 PM · 6 févr. 2021 · Twitter for iPhone

23,3 k Retweets 2 816 Tweets cités 77,1 k J'aime





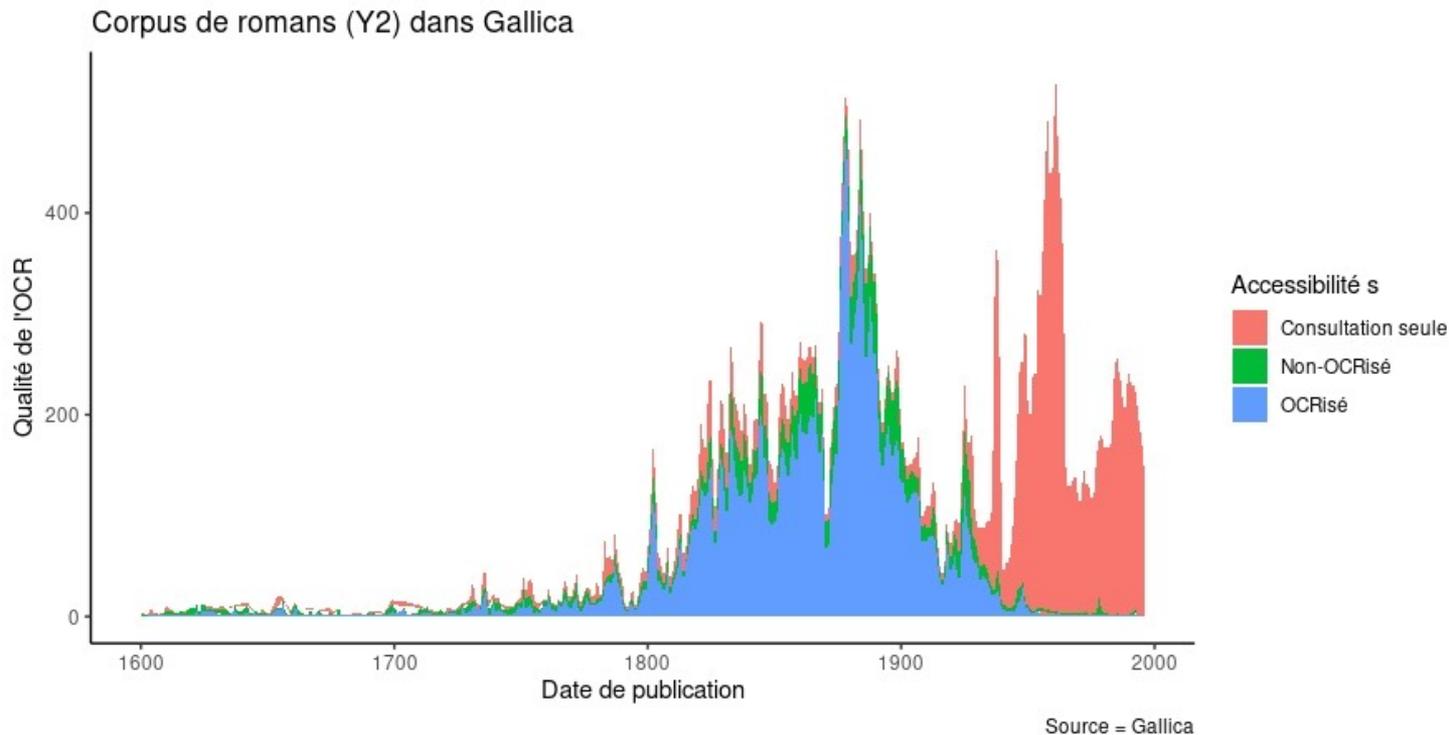
Les biais de collections



Le trou noir du web entre 1945 et 1990 dans Europeana



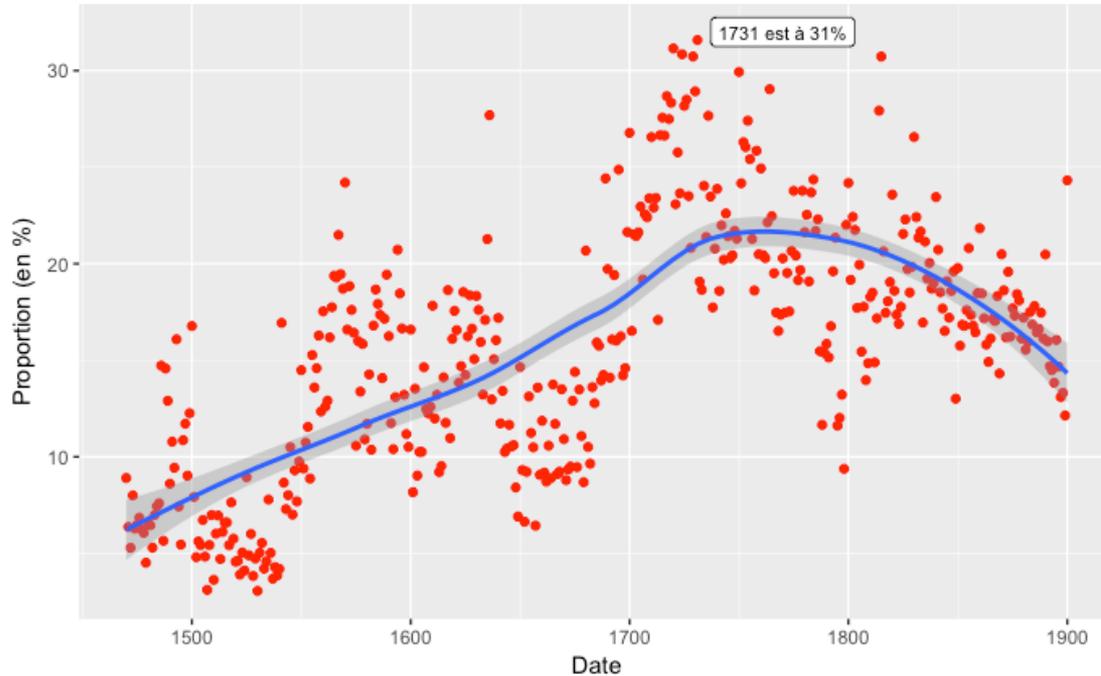
Les biais de collections



Le mille-feuille des conditions d'accès aux documents sur Gallica

Les biais de collections

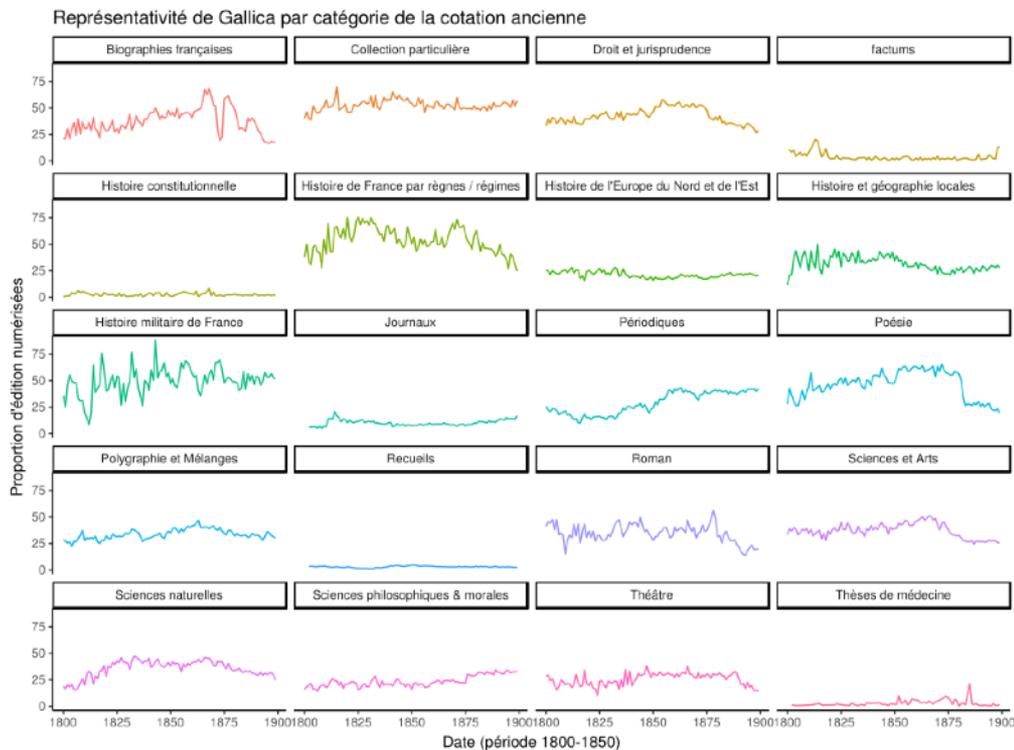
Éditions du Catalogue de la BNF numérisées dans Gallica



Source : Data BNF

Même avant 1945 la représentativité de Gallica évolue...

Les biais de collections

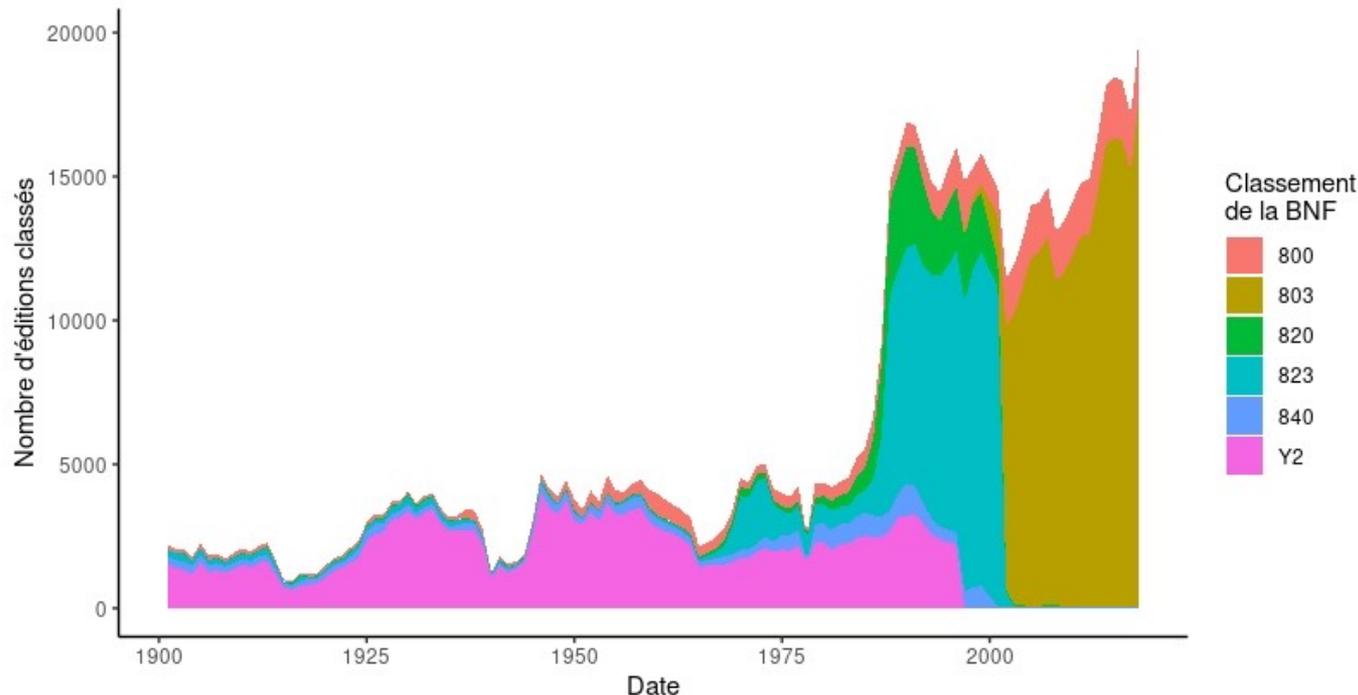


...avec d'importantes variabilités selon les thématiques.



Les biais de collections

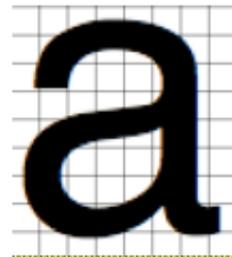
La documentation
des archives
patrimoniales
n'est jamais
simple : elle hérite
d'une longue
histoire et des
classements
successifs
adoptés par les
bibliothèques



Classifications des romans dans Gallica : de Y2 (jusqu'en 1996) jusqu'au Cadre de classement de la BNF)



Tous les textes... mais pas toutes les lettres.



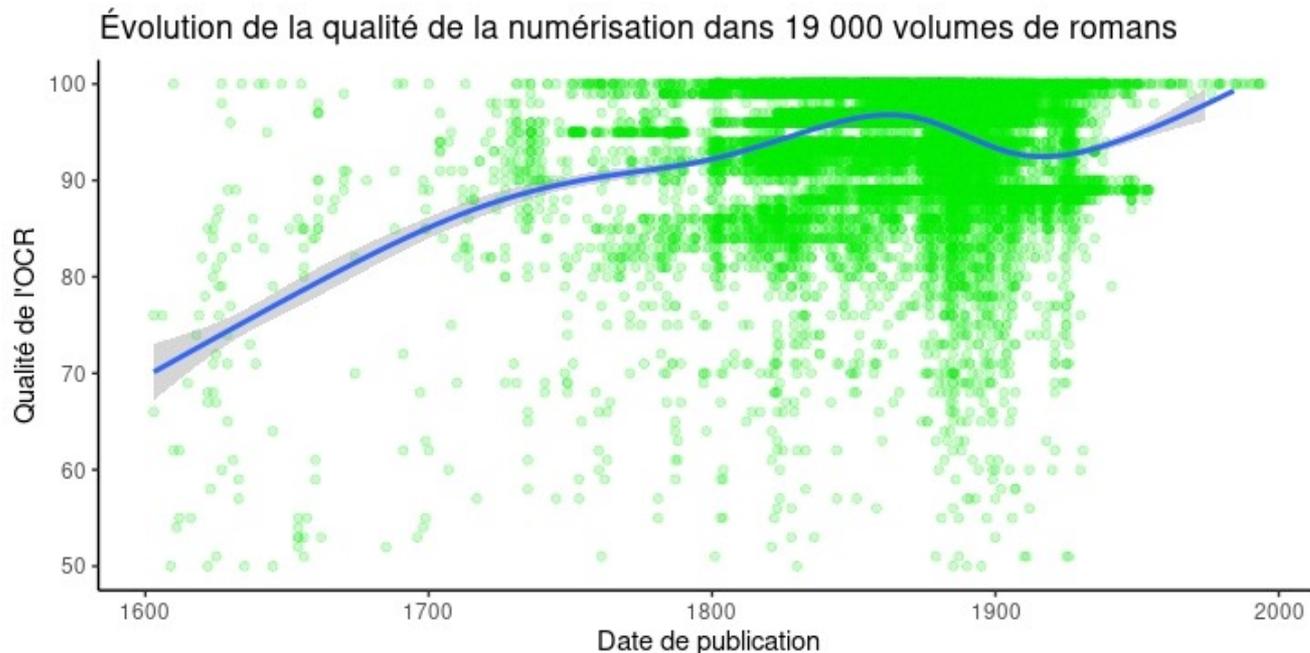
Les milliards de pages numérisés ne peuvent pas l'être à la main : on utilise des outils automatisés les OCR avec leurs propres biais (ici avec le S long de l'ancien français).





Tous les textes... mais pas toutes les lettres.

La qualité de l'océrisation est inégalement répartie selon les corpus : les logiciels ont avant tout été conçus pour des imprimés contemporains et fonctionnent de plus en plus difficilement pour les périodes anciennes.



2. Voir l'invisible à distance

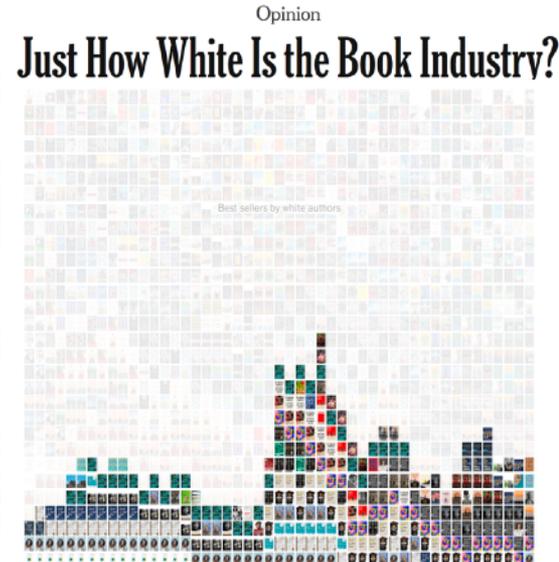
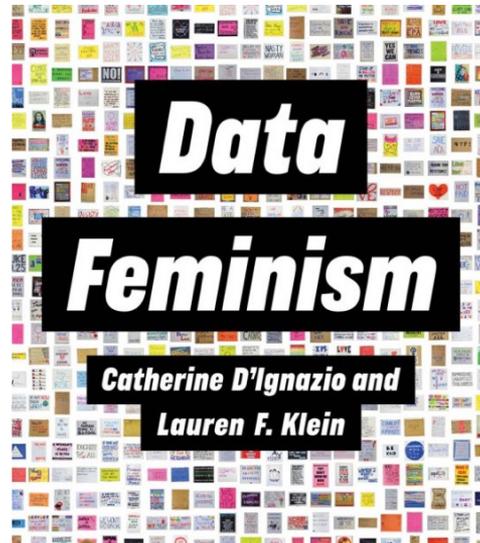
De l'archive-document à l'archive-réseau



Les archives invisibles

L'enjeu des humanités computationnelles aujourd'hui n'est plus seulement de ne pas déformer les archives existantes mais de questionner les conditions de production des archives et de l'archivage et de se donner les moyens de « produire » des archives occultées

« In calling for a conceptual reorientation from the axis of distant and close to a space defined by multiple dimensions of scale, I seek to make the case that quantitative methods can be used to probe the research questions about gender, race, and their intersection with labor that have thus far proved difficult (although certainly not impossible) to explore. » — Lauren F. Klein (« Dimensions of Scale »)



Surmonter les filtres de lectures

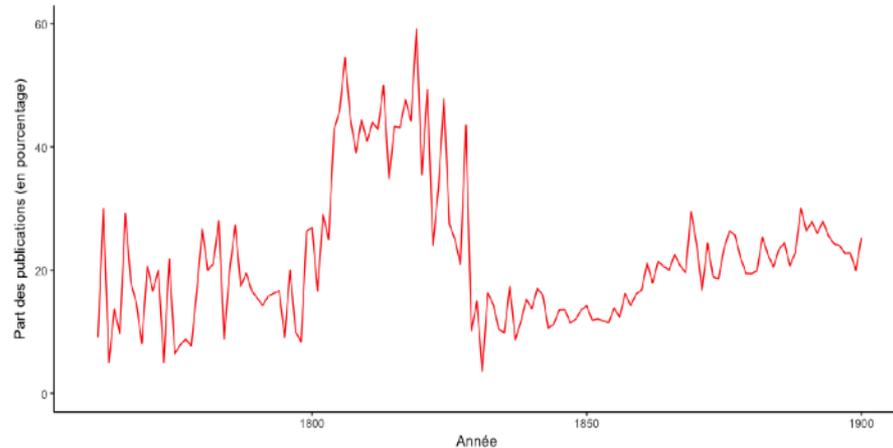
Romanrama

Classifications automatisées des romans de Gallica (1815-1850)

Explorer tous les romans Trouver un roman Genres et arcs narratifs Les coulisses de la classification Féminin/masculin

Show 25 entries

Titre	Date	Probabilité	Genre
Le Solitaire de Colonna par le Cte Janus Binski	1835	99.96502	Roman chrétien/conte
Wat-Tyler ou Dix jours de révolte roman historique par A-J-B Defauconpret	1825	99.91555	Roman historique
Entrée dans le monde par Miss Jane Porter traduit de l'anglais par Madame***	1828	99.90106	Roman sentimental
Célon ou Entretiens d'un vieillard avec son fils prêt à entrer dans le monde Traduction du Théophrone allemand de M Campe	1820	99.88244	Roman sentimental
Le Masque de veilleurs par Jules Lacroix	1844	99.86958	Drame
Les Fêtes des enfants ou Recueil de petits contes moraux par M Ducrest-Duminiil	1822	99.85908	Roman de mœurs
Le parc de Mansfield ou Les trois cousines par l'auteur de Raison et sensibilité ou Les deux manières d'aimer traduit de l'anglais par M Henri V*****N (Villemain)	1816	99.83142	Roman sentimental
Chinki ou Les maîtrises en Cochinchine Histoire cochinchinoise	1824	99.86429	Roman de mœurs



En donnant accès un grand nombre de publications oubliées la numérisation s'accorde avec une vision « décentralisée » de l'histoire culturelle, potentiellement en rupture avec les représentations canoniques.

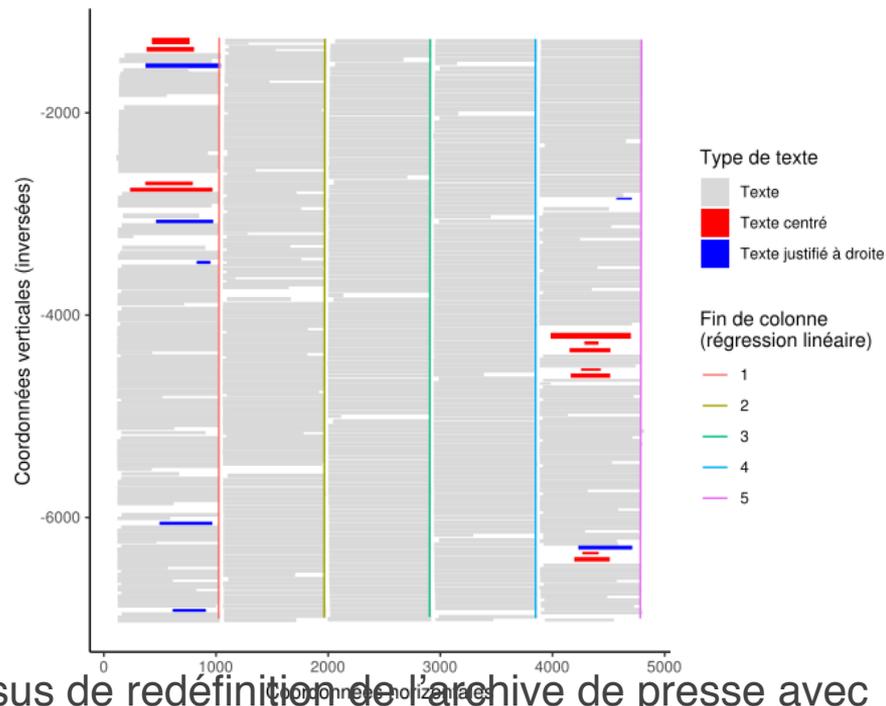
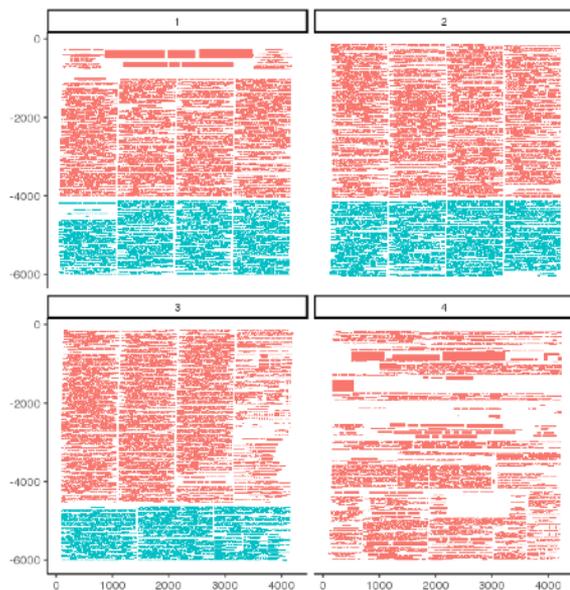
Inventer des archives : le cas de la presse

Dans le cas de la presse, l'indexation des archives n'a jamais été satisfaisante... jusqu'à aujourd'hui. Les essais n'ont pas manqué de Deschiens (1828) aux Tables du temps, avec des résultats au mieux limités.

The screenshot shows a digital archive interface. At the top, it says 'COLLECTION DE MATÉRIAUX POUR L'HISTOIRE DE LA RÉVOLUTION DE FRANCE, DEPUIS 1787, JUSQU'A CE JOUR.' Below this, there is a search bar with the number '5' and a dropdown menu. A list of search results is visible, including 'LES ÉLEGANCES', 'À l'écoute...', 'Il y a cent ans', and 'CE QUE L'ON DIT...'. A small box at the bottom left indicates 'SpecialPages' with a list of items: '[7] Article (1)', '[8] Article (1)', '[9] Article (1)', and '[10] Article (1)'.

The screenshot shows a historical newspaper page from 'L'ECHO DE PARIS'. The page is dated '20 Juin 1937' and is the 'Édition de 5 heures'. The main headline is 'L'Allemagne et l'Italie se retirent définitivement du système de contrôle'. Other headlines include 'LES CHANCES DU CABINET CHAUTEMPS', 'LES ÉLEGANCES', 'LA GRANDE COURSE DE HAIES GAGNÉE PAR "SYOHU II"', and 'ON GARDE LES MÊMES'. The page features several columns of text, a large illustration of a man in a suit, and a smaller illustration of a man in a hat. The page is highlighted in yellow.

Inventer des archives : le cas de la presse



Le projet Numapresse est investi dans un processus de redéfinition de l'archive de presse avec l'introduction de nouveaux découpages (par article, par paragraphes...) et la production de nouveaux corpus qui n'ont jamais été édités isolément (comme les romans-feuilletons)



Vers une « bibliothèque » de modèles

Genres journalistiques "1840-1860"

Un exemple de fiche : le modèle des genres journalistiques de 1840-1860 [http://](http://www.numapresse.org/generotheque/items/show/4)

www.numapresse.org/generotheque/items/show/4

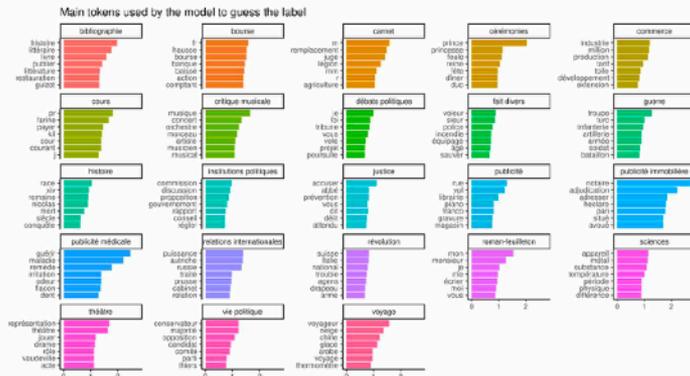
Corpus d'entraînement

1 444 blocs de textes de plus de 100 mots collectés à partir d'une sélection aléatoire de 500 exemplaires de presse numérisés par Gallica entre 1840 et 1860. Le corpus comprend quatre titres actifs sur l'ensemble de la période (*le Journal des débats*, *la Presse*, *le Siècle* et *le Constitutionnel*), deux titres partiellement publiés (*le Pays*, créé en 1849 et *le Courrier français*, disparu en 1851), ainsi que titres avec des numérisations parcellaires sur cette période (*le Figaro*, *la Démocratie pacifique*, *le Corsaire* et *l'Estafette*)

Du fait de la sélection aléatoire les titres les plus tardifs ou les titres avec un volume de publication plus faibles sont moins représentés (par exemple 5 textes pour *l'Estafette* avec des archives présentes pour la seule année 1852).

Le corpus a été annoté manuellement par blocs de 250 blocs de textes dans un tableur. En raison des imperfections du processus de numérisation, certains blocs de textes ont été subdivisés (s'ils comportaient plusieurs nouvelles différentes). Les textes difficiles à classer ont été écartés.

Catégorisation



Description

Le modèle 1840-1860 appartient à la série des modèles "générationnels" de Numapresse couvrant les genres journalistiques de la presse quotidienne nationale française du début du 19e siècle à la Seconde Guerre Mondiale.

Modèle

[Télécharger le modèle au format R](#)

Corpus

[Télécharger le corpus d'entraînement](#)

Format original

Modèle SVM enregistré avec R et Tidysupervise (format .rda)

Auteur

Pierre-Carl Langlais



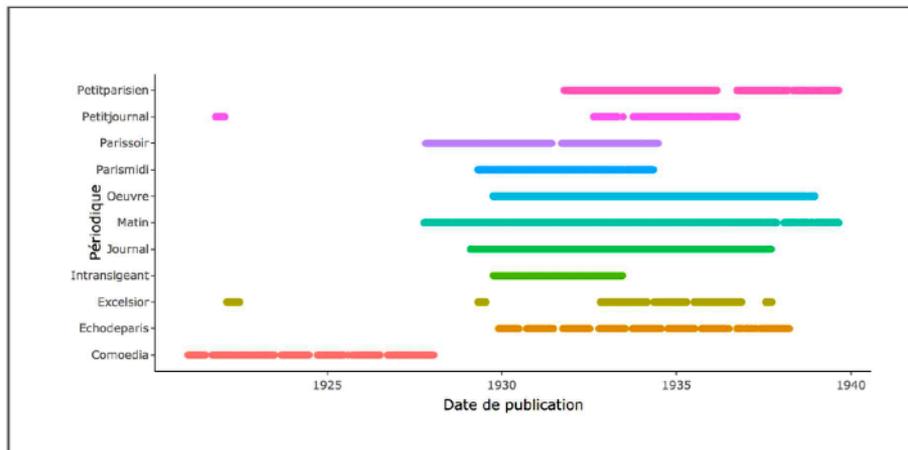
Créer des corpus contextualisé par IA



Présentation Requetes sur les corpus Requetes sur les illustrations

La Page du cinéma est un [projet Numapresse](http://www.numapresse.org) recensant 4000 suppléments hebdomadaires dédiés au cinéma publiés dans la presse généraliste nationale française pendant la première moitié du XXe siècle.

À partir du début des années 1920, plusieurs quotidiens commencent à consacrer des pages entières à l'actualité cinématographique. Ces initiatives sont d'abord temporaires (*L'Excelsior*, *Le Petit Journal*) à l'exception du quotidien culturel *Comoedia*. Elles connaissent un essor considérable lors du passage du muet au parlant, alors que le cinéma se trouve à l'avant-garde de la modernité médiatique. Tous les quotidiens à grands tirages (*Paris-Soir*, *Paris-Midi*, *Le Matin*, *Le Petit Journal*, *Le Petit Parisien*, *L'Intransigeant*, *Le Journal*) et d'autres titres importants (*L'Œuvre*, *L'Écho de Paris*, *Le Figaro*, *La Liberté*) se dotent alors de leur page cinéma. Certaines rédactions adoptent même une page cinéma quotidienne telles que *Comoedia* ou, beaucoup plus brièvement, *Paris-Midi*. La pratique régresse en partie à la fin des années 1930, en partie en raison de l'affaiblissement du rythme hebdomadaire : l'actualité cinématographique devient quotidienne et est diluée dans des rubriques moins spécialisée (Spectacles, Culture...)



La classification de masse rend également possible l'exploration de corpus quasi-exhaustifs à l'image du projet « La Page de cinéma » qui contient désormais 4000 suppléments cinéma hebdomadaires publiés dans 10 quotidiens différents : http://www.numapresse.org/exploration/cinema_pages/presentation.php

Créer des corpus contextualisé par IA

Pour les illustration,
l'impact de la classification
est encore plus profond :
ces documents ne sont
tout simplement pas
indexés par défaut [http://
www.numapresse.org/
exploration/cinema_pages/
query_illustration.php](http://www.numapresse.org/exploration/cinema_pages/query_illustration.php)



Présentation Requetes sur les corpus Requetes sur les illustrations

Ce formulaire permet de faire des requêtes sur un corpus Numapresse.

Date de début

1920-09-01

Date

1938

Requête libre sur la légende

Votre requête

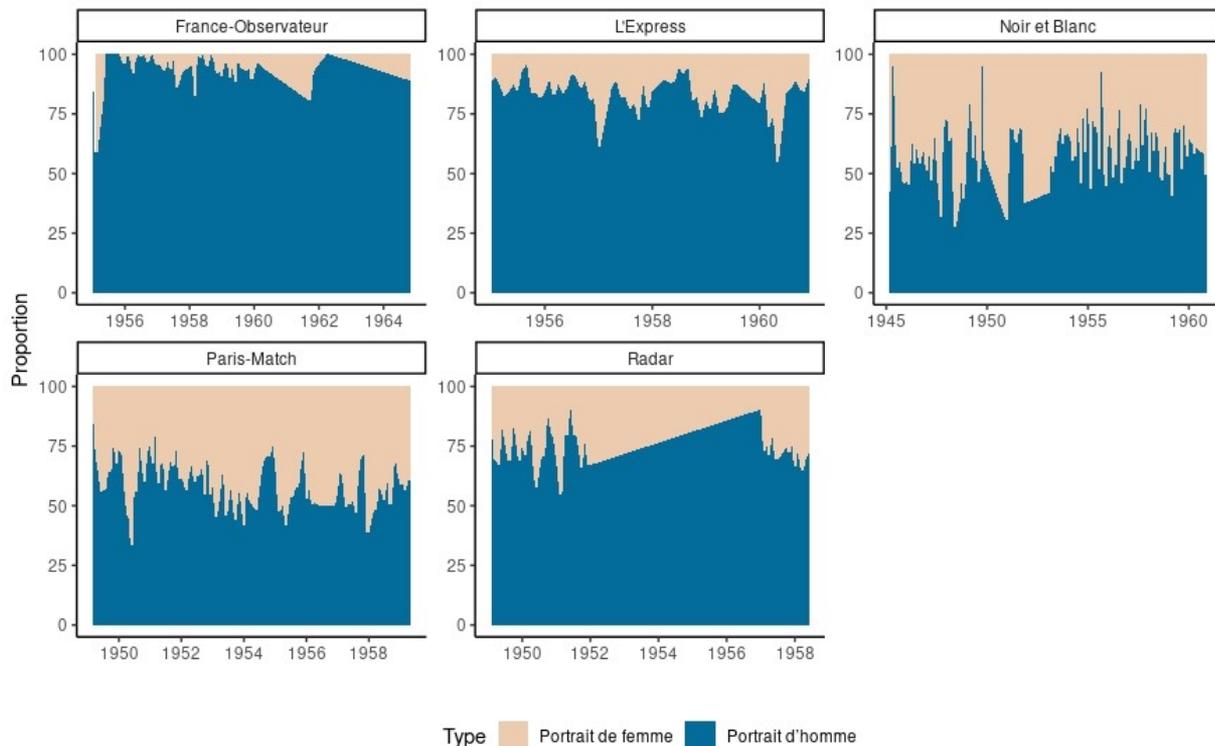
Soumettre la requête





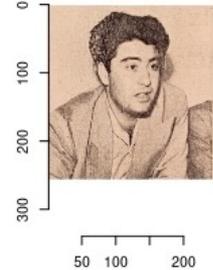
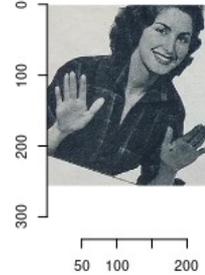
Trouble dans les catégories

Les biais des modèles peuvent être détournés pour analyser des biais sociaux. Ces graphes montrent les taux de « féminisation » dans cinq hebdomadaires des années 1950. varient très fortement selon les corpus ce qui reflète assez nettement le rapport entre les politiques éditoriales et les normes de genre

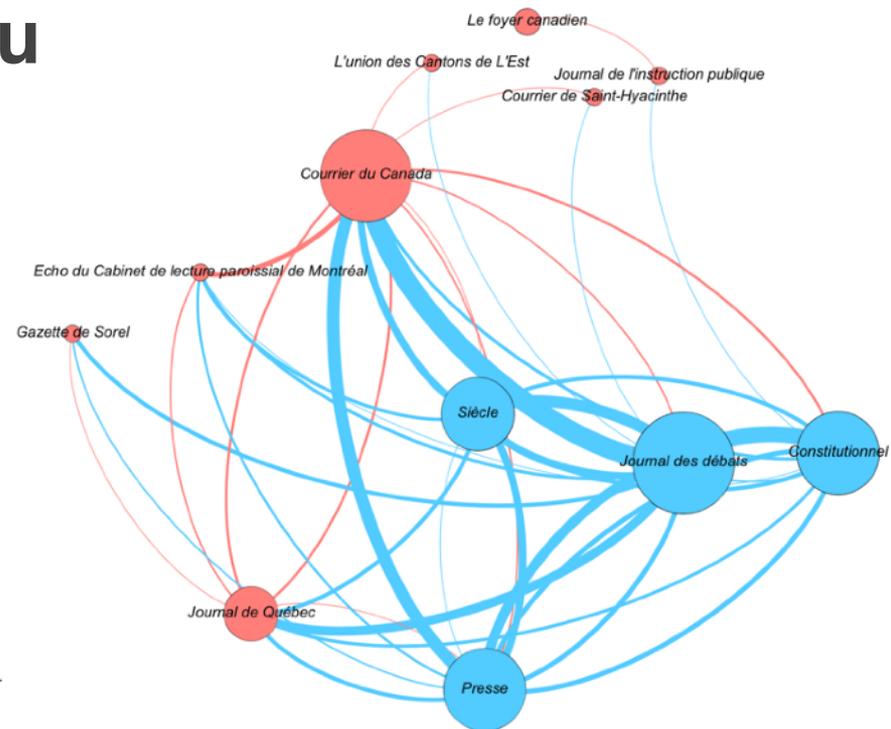
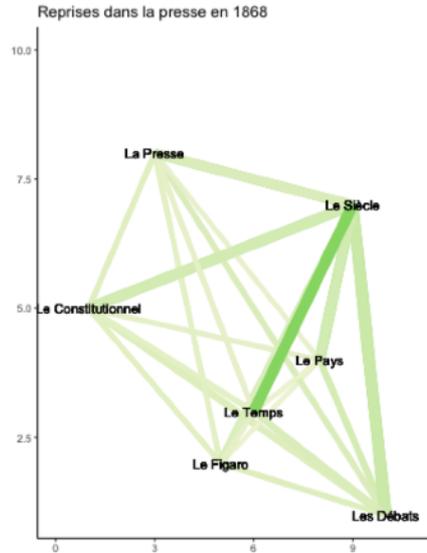
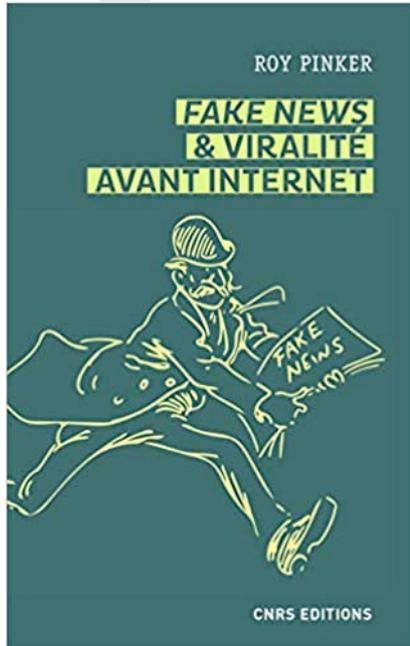


Trouble dans les catégories

En cherchant des images associant des traits féminins et masculins (> 30% dans les deux cas) nous parvenons à isoler des cas limites : représentations historicisées, ambiguës voire *queer*



Vers l'archive-réseau

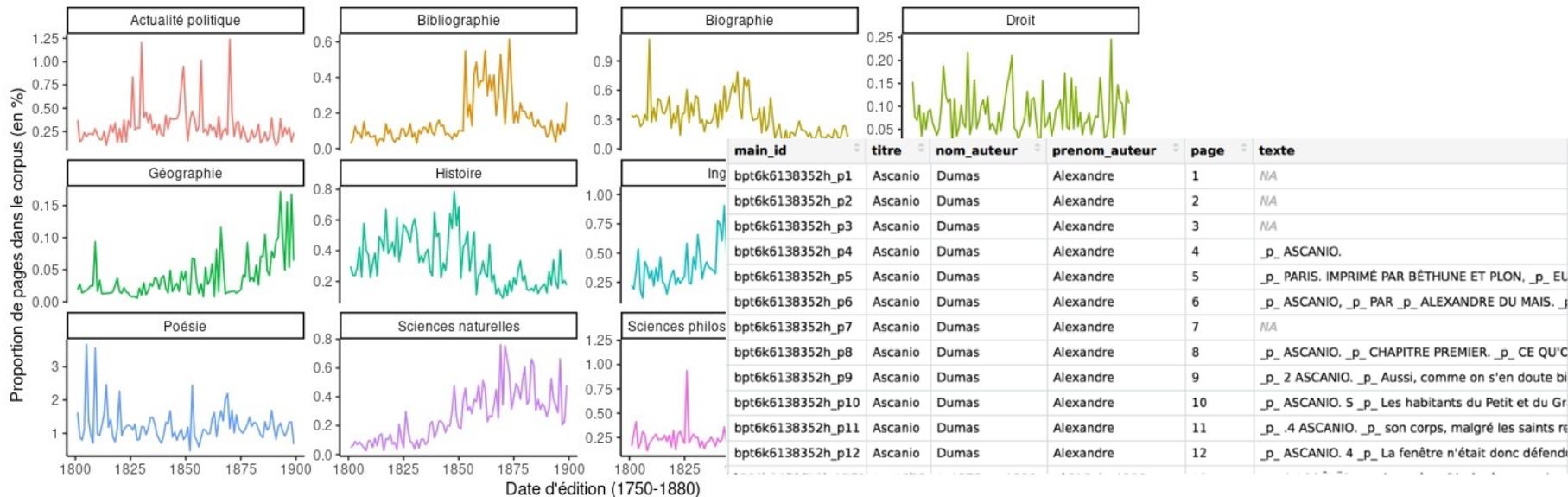


L'éclatement de l'unité documentaire du journal permet d'opérer d'autres rapprochements. D'autres projets de Numapresse portent ainsi sur l'identification automatique des nouvelles virales, ce qui permet de reconstituer des chaînes de transmission successive.

Vers l'archive-réseau

Classifications intertextuelles du corpus de roman

>50% de probabilité par page pour un modèle entraîné sur d'autres cotes historiques de Gallica

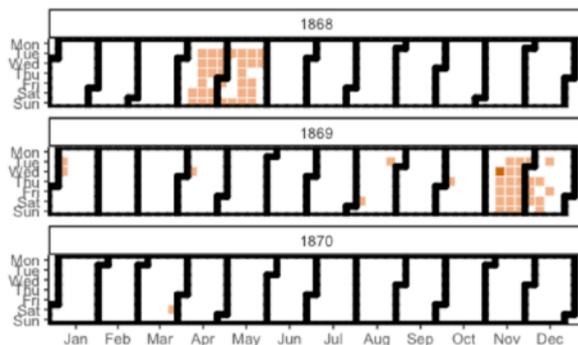


L'éclatement des cadres usuels de l'archivage affecte aussi les monographies. Les projets de cultural analytics comme le notre privilégient l'indexation au niveau de la page voire de sous-éléments internes à la page (comme les illustrations).



Vers l'archive-réseau

L'extraction des signatures rend possible l'identification de profils sociologiques à des périodes où la documentation externe fait défaut (avant 1880). Dans la *Liberté* de 1865 à 1870, les femmes occupent des positions beaucoup plus précaires que les hommes.



Judith Gauthier



Wilfried de Fonvielle

Conclusion

Du journalisme à la littérature et au-delà





Ce que l'archive change à la classification

Un possible effet retour des études en histoire culturelle avec l'intelligence artificielle : mieux informer les biais sous-jacents aux outils contemporains.



Im 
@willie_agnew



Buried in the recent trillion parameter language model paper is how the dataset to train it was created. Any page that contained one of these words was excluded: [github.com/LDNOOBW/List-o...](https://github.com/LDNOOBW/List-of-Dirty-Naughty-Obscene-and-Otherwise-Bad-Words) Two sample banned words: "twink" and "sex"

[Traduire le Tweet](#)



[LDNOOBW/List-of-Dirty-Naughty-Obscene-and-Otherwise...](https://github.com/LDNOOBW/List-of-Dirty-Naughty-Obscene-and-Otherwise-Bad-Words)
List of Dirty, Naughty, Obscene, and Otherwise Bad Words -
[LDNOOBW/List-of-Dirty-Naughty-Obscene-and-Otherwise...](https://github.com/LDNOOBW/List-of-Dirty-Naughty-Obscene-and-Otherwise-Bad-Words)
github.com



Au-delà des images réelles : vers l'exploration des images possibles.

De nouveaux projets d'intelligence artificielle sont capables de créer des images qui n'ont jamais existé mais qui font sens dans l'univers culturel d'où elles sont extraites. Ici, il y a une représentation très rétrofuturiste des Iphone par décennie.

