
Digging Deeper into Government Information: An Introduction to Text Mining

2018 Federal Depository Library Conference

Valerie D. Glenn and Eleanor Dickson Koehl

October 22, 2018



Presenter Information

Valerie D. Glenn, Head, Map and Government Information Library, University of Georgia (previously Federal Documents Analyst, HathiTrust)

Eleanor Dickson Koehl, HTRC Digital Humanities Specialist



What is HathiTrust?

- Brief introduction to HathiTrust
- HathiTrust U.S. Federal Documents Program
 - Overall goals
 - Collection priorities
 - More information: Heather Christenson's presentation regarding the UC FedDocArc project: Tuesday, 3-4, Wilson/Harrison room
- HathiTrust Research Center (HTRC)



What is text analysis?

- Using computers to reveal information **in** and **about** text (Hearst, 2003)
 - Algorithms discern patterns in unstructured text
 - More than just search
- Everyday example
 - E-mail spam filter



How does it work?

- Break textual data into smaller pieces
- Abstract (reduce) text so that a computer can crunch it
- Count!
 - Words, phrases, parts of speech, etc.
- Develop hypotheses from relative counts



How does it impact research?

- Shift perspective → shift research questions
- Opens up:
 - Questions not provable by human reading alone
 - Larger corpora for analysis
 - Studies that cover longer time spans
- One step in the research process
 - Can be combined with close reading



Text Analysis Research Questions

- May involve:
 - Change over time
 - Pattern recognition
 - Comparative analysis



Example: “Textual Predictors of Bill Survival in Congressional Committees”

- Research question:

What makes a bill more or less likely to make it out of a Congressional committee?

Yano, et al., “Textual Predictors of Bill Survival in Congressional Committees,” *2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 793–802, Montreal, Canada, June 3-8, 2012, <https://dl.acm.org/citation.cfm?id=2382157>



Example: “Textual Predictors of Bill Survival in Congressional Committees”

- Modeled likelihood of a bill’s success
 - From 1993-2011
- Made use of metadata features
 - E.g. The sponsor, their political affiliation, and their home state; which party has majority; what month it was introduced
- Also used textual features
 - E.g. Unigrams (single words) to show which words are correlated with a bill’s success or demise in committee

<https://dl.acm.org/citation.cfm?id=2382157>



Example: “Textual Predictors of Bill Survival in Congressional Committees”

- Their results
 - Words associated with success: resources, interior
 - Words associated with failure: energy, security
- Their argument
 - “Minor legislation” more likely to survive
 - Some bills are proposed to make a political statement
 - Smaller bills get rolled into larger ones

<https://dl.acm.org/citation.cfm?id=2382157>



Building Corpora

- Identify texts through:
 - Full text search
 - Metadata (author, date, genre)
- Usually involves deduplication
 - What to keep/discard is project dependent
- Examples of deduplication criteria:
 - OCR quality
 - Earliest edition
 - Editions without forewords or afterwords



Finding Text

- Not always easy
 - copyright restrictions
 - licensing restrictions
 - format limitations
 - hard-to-navigate systems

** issues more pronounced at scale**



Sources of Text Data

Vendor Databases

Library and archives
digital collections

Social media

Digitized personal corpus

Government sources

Considerations for choosing:

Comprehensiveness

- “Official” version
- Timeliness
- Authoritative vs. Public voice
- Likelihood for clean data
- Need for unique sources
- Cost



HTRC Services

Offering	Description	Data availability	Account required
HTRC Algorithms	A set of tools for assembling collections of digitized text from the HathiTrust corpus and performing text analysis on them.	Including items in copyright for ALL USERS.	HTRC Analytics
HTRC Worksets	Worksets are sub-collections of HathiTrust volumes that can be analyzed with HTRC algorithms or used to access HTRC Extracted Features.	Including items in copyright for ALL USERS.	HTRC Analytics
Extracted Features Dataset	A dataset allowing non-consumptive analysis on specific features extracted from the full text of the HathiTrust corpus.	Including items in copyright for ALL USERS.	None
HathiTrust+Bookworm	A tool for visualizing and analyzing word usage trends in the HathiTrust corpus.	Including items in copyright for ALL USERS.	None
HTRC Data Capsule	A secure computing environment for researcher-driven text analysis on the HathiTrust corpus.	All users may access public domain items. Access to items in copyright is available ONLY to member-affiliated researchers.	HTRC Analytics ; plus additional restrictions



HTRC Analytics

HTRC Analytics | Algorithms | Data Capsules | Worksheets | Datasets | Explore | Help | About | Sign In | Sign Up

HATHI TRUST
Research Center

HathiTrust Research Center Analytics

Supports large-scale computational analysis of the works in the HathiTrust Digital Library to facilitate non-profit and educational research.

Featured Services

- Extracted Features
- Text Analysis Algorithms
- Data Capsules

www.analytics.hathitrust.org

- Accounts for those from institutions of non-profit research or higher education
- Don't need to be an HT member



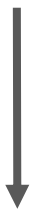
HTRC Worksets

- User-created collections of text from the HathiTrust Digital Library
 - Think of them as textual datasets
- Can be shared and cited
 - Suited for non-consumptive access



mdp.49015002221845
 mdp.49015002221837
 mdp.49015002221829
 mdp.49015002221787
 mdp.49015002221811
 mdp.49015002221761
 mdp.49015002221779
 mdp.49015002203140
 mdp.49015002203157
 mdp.49015002203033
 mdp.49015002203231
 mdp.49015002203249
 mdp.49015002203223
 mdp.49015002203405
 mdp.49015002203272
 mdp.49015002203215

HathiTrust
collection



poli_science_DDRF

Download

Description : Political science collection for DDRF workshop

Owner	Last Modified Time	Number of Volumes	Tags
rhan11	2017-10-05T18:21:35Z	16	

Filter volume by title...

Volume ID	Title	Authors	Year	Language
mdp.49015002203223	Public papers of the presidents of the United States.	United States President; Clinton, Bill 1946-; Bush, George 1924-; Reagan, Ronald; Carter, Jimmy 1924-; Ford, Gerald R. 1913-2006; Nixon, Richard M. (Richard Milhous) 1913-1994; Johnson, Lyndon B. (Lyndon Baines) 1908-1973; Kennedy, John F. (John Fitzgerald) 1917-1963; Eisenhower, Dwight D. (Dwight David) 1890-1969; Truman, Harry S. 1884-1972; Hoover, Herbert 1874-1964; United States Federal Register Division; United States Office of the Federal Register	1978	eng
mdp.49015002203272	Public papers of the presidents of the United States.	United States President; Clinton, Bill 1946-; Bush, George 1924-; Reagan, Ronald; Carter, Jimmy 1924-; Ford, Gerald R. 1913-2006; Nixon, Richard M. (Richard Milhous) 1913-1994; Johnson, Lyndon B. (Lyndon Baines) 1908-1973; Kennedy, John F. (John Fitzgerald) 1917-1963; Eisenhower, Dwight D. (Dwight David) 1890-1969; Truman, Harry S. 1884-1972;	1979	eng

https://babel.hathitrust.org/cgi/mb?a=listis;c=1848985365

Home About Collections Help Feedback

HATHI TRUST Digital Library

Search words about or within the items

Advanced full-text search | Search tips

Full view only

poli_science_DDRF

Political science collection for DDRF workshop

Search in this collection

All Items (16)

Sort by: Title A-Z | 25 per page | 1

Select all on page | Select Collection | Add Selected

Public papers of the presidents of the United States. 1971
 by United States. President.
 Published 1971

[Catalog Record](#) | [Full view](#)
[Download Extracted Features](#)



HTRC
workset



HT
volume
IDs

HTRC U.S. Federal Documents Worksets

- HathiTrust-created collections of text from the HathiTrust Digital Library
 - [U.S. Environmental Protection Agency publications](#)
 - [U.S. Bureau of Indian Affairs publications](#)
 - [U.S. Congressional Serial Set](#)
 - [Foreign Relations of the United States](#)



Building Worksets

- Ways to build:
 - Search and create in HT Collection Builder (<https://babel.hathitrust.org/cgi/mb>)
 - Compile volume IDs making use of HT metadata services (<https://www.hathitrust.org/data>)
- Import to HTRC Analytics



HTRC

Algorithms

- For analyzing HT data only
- Run against a workset
- Kind of analysis:
 - Token count and word cloud
 - Visualize themes (topic modeling)
 - List people, places, organizations, dates, etc. (named entity recognition)

Output

entities.csv stdout.txt stderr.txt

[Click here to download entities.csv](#)

vol_id	page_seq	entity	type
mdp.39015037380378	00000007	December 1, 1966	DATE
mdp.39015037380378	00000007	1967	DATE
mdp.39015037380378	00000007	Regd	MISC
mdp.39015037380378	00000007	AFRICA	ORGANIZATION
mdp.39015037380378	00000008	Algeria	LOCATION
mdp.39015037380378	00000008	MPLA	ORGANIZATION
mdp.39015037380378	00000008	Sudanese	MISC
mdp.39015037380378	00000008	1966	DATE
mdp.39015037380378	00000008	Independence four years ago	DATE
mdp.39015037380378	00000008	Ghana	LOCATION
mdp.39015037380378	00000008	Central Committee	ORGANIZATION
mdp.39015037380378	00000008	Amin Shaker	PERSON
mdp.39015037380378	00000008	November 29, 1966	DATE
mdp.39015037380378	00000008	Egyptian Gazette	MISC
mdp.39015037380378	00000008	1966	DATE
mdp.39015037380378	00000008	November	DATE

Named entity extractor results



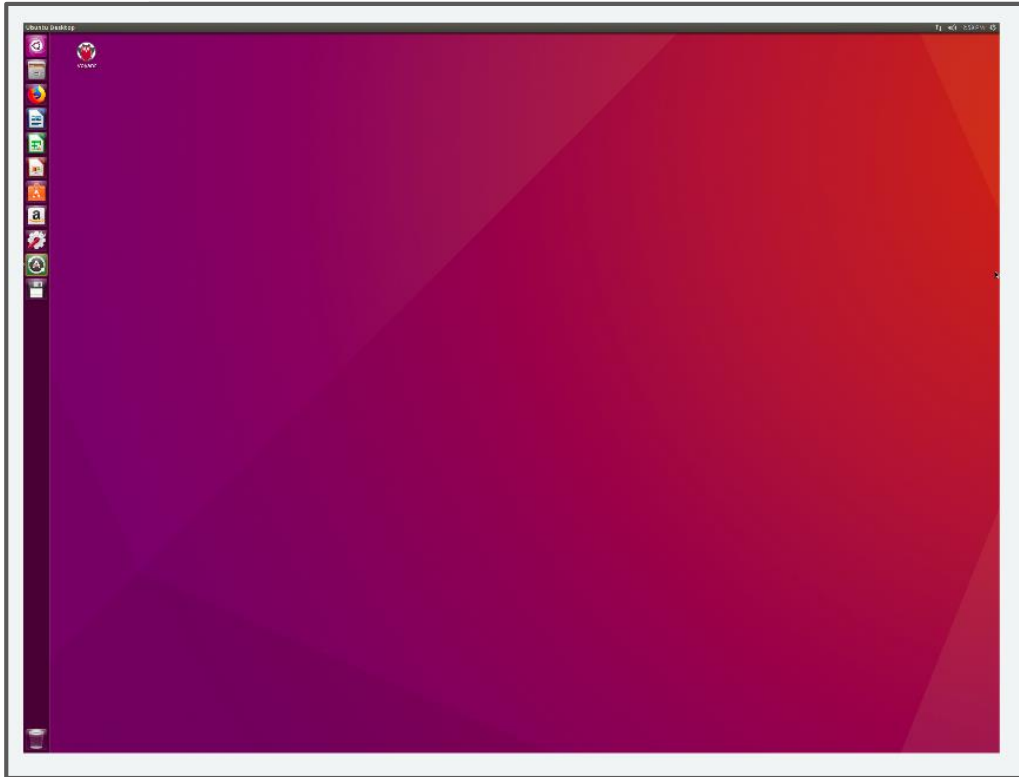
HTRC Extracted Features

- Selected data and metadata
 - Extracted from raw text
- Includes words and word counts per page

```
1  {
2    "id": "uc1.b3419888",
3    "metadata": {
4      "schemaVersion": "1.2",
5      "dateCreated": "2015-02-12T13:30",
6      "title": "Zoonomia = or The laws of organic life / by Erasmus Darwin.",
7      "pubDate": "1809",
8      "language": "eng",
9      "htBibUrl": "http://catalog.hathitrust.org/api/volumes/full/htid/uc1.b3419888.json",
10     "handleUrl": "http://hdl.handle.net/2027/uc1.b3419888",
11     "oclc": "3679915",
12     "imprint": "Thomas and Andrews, 1809."
13   },
14   "features": {
15     "schemaVersion": "2.0",
16     "dateCreated": "2015-02-20T23:58",
17     "pageCount": 616,
18     "pages": [
```



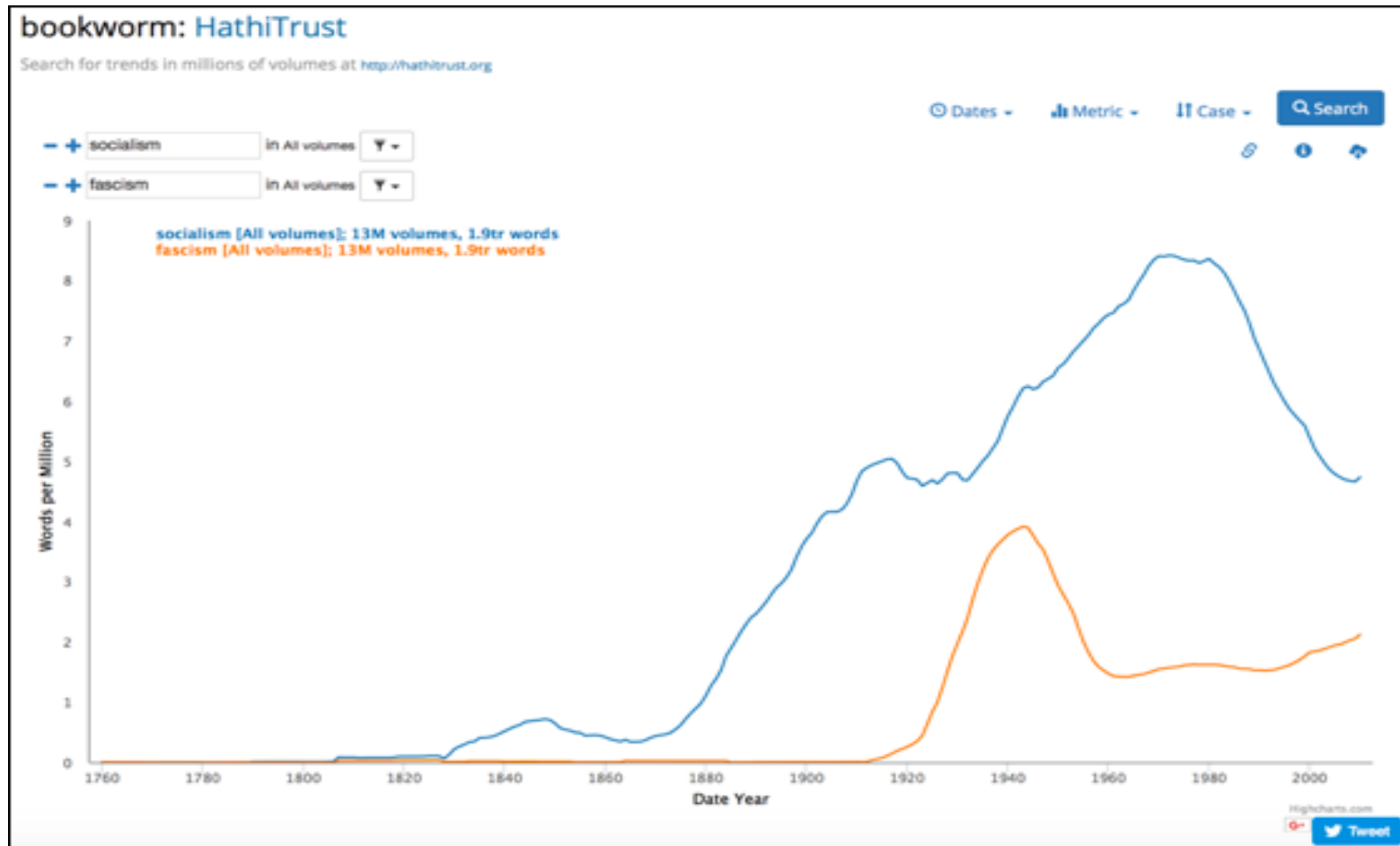
HTRC Data Capsules Environment



- Secure computing desktop environments
- Researcher-driven analysis
- In-copyright data limited to approved researchers from HathiTrust member institutions



HathiTrust + Bookworm



Getting Started

- Create an account for HTRC Analytics
 - HT-member institution affiliation NOT required
- Identify HathiTrust volumes to analyze
 - Search & create collection
 - Make use of metadata services
 - Ask for help!



Getting Started

- Download Extracted Features for the volumes and analyze locally

or

- Import the collection or volume IDs as a workset and run an algorithm

or

- Import data into a Data Capsule and analyze in-Capsule



Thank you! Questions?

Or contact htrc-help@hathitrust.org for more assistance!

