# Digging Deeper into Text and Data Mining

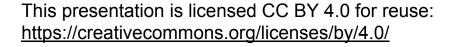
I MIXXXXIII

By Inga Haugen, Edward F. Lener, Virginia Pannabecker, & Philip Young
Virginia Tech, University Libraries

Presentation available in VTechWorks institutional repository <a href="http://hdl.handle.net/10919/79483">http://hdl.handle.net/10919/79483</a>









### Roadmap / Outline



#### **Introductions - Presenters**



**Inga Haugen** Agriculture, Life Sciences, and Scholarly Communication Librarian; Interim liaison to the College of Natural Resources and the Environment (CNRE)

**Edward Lener** Associate Director for Collection Management and College Librarian for the Sciences

**Ginny Pannabecker** Associate Director for Research Collaboration and Engagement; Liaison Librarian for life sciences and biomedical programs

**Philip Young** Institutional Repository Manager



# What is Text and Data Mining (TDM)?

Text and data mining (TDM) uses methods of automated extraction, combination, and analysis of data to create new information by revealing trends, patterns, and relationships. The mining of text and of data usually require different considerations. Text mining, sometimes called text analytics, can be viewed as a subset of data mining.

### Examples of TDM in research / scholarship

#### From VT

<u>TDM Forum</u> - variety of topics and examples - <u>business</u>, <u>cyber bullying</u>, <u>statistics</u>

<u>History - 1918 influenza pandemic project</u>

#### Other examples

Health

<u>Discovering associations between adverse events from electronic health record</u> <u>data</u>

**Digital Humanities** 

<u>Librarian collaborations for research and training in text encoding</u>

**Opioid Crisis** 

A Text Mining Analysis of Public Reactions to Opioid Crisis

## **TDM: Academic library support opportunities**



# **Identifying & Sharing TDM Sources and Tools**

- Tool and Methods Guides
  - UC Berkeley Library guide
  - MIT Libraries APIs guide
  - University of Melbourne Text Mining Tools list
  - Carnegie Mellon University Libraries guide
  - VT Libraries
- Conversations with researchers / Community of practice
- Open educational training options
- Conducting literature reviews for updates to to stay current in best practices and tools



#### **Expanding library licensing permissions**

- Each license represents an opportunity
- Vendors are at widely different places on this issue
- <u>Liblicense model agreement</u> and others like it can be a very helpful starting point
- Need to decide how important this issue is to you and what you are prepared to accept



### **Clarifying legal aspects\***



- Subject to U.S. law
- No specific exemption for TDM in U.S. copyright law\*\*
- License agreement overrides Fair Use provisions of copyright law
- Researchers may need data from multiple sources with varying rights

#### **ARL Issue Brief**

Text and Data Mining and Fair Use in the United States

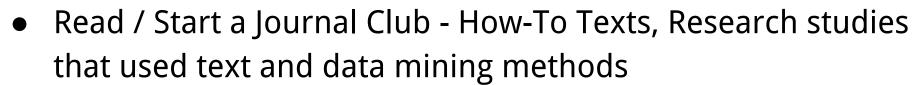
- \* Note We are not attorneys!
- \*\*You may hear about a specific tdm exemption in the UK more info here.

Sample license language from one major library vendor

"TDM Output" means the result of any TDM activity by Researcher, such as the creation of an index, abstract, relative or absolute description or representation of the Content; any algorithm, metrics, method, standard or taxonomy describing or based on the Content, ... whether in a the form of a direct extraction or a representation in any form which is based on any portion of the Content. Any quotes from the Content shall be limited to fifty (50) words or less.

#### **Developing expertise**

- Join Communities
  - Text and Data Mining Research Support List (JISC)\*



- Introduction to Text Analysis
- Tutorials / Training / Webinars
  - Coursera Data Science course via Johns Hopkins
  - Hathi Trust Research Center | Training
  - Programming Historian
  - R and Data Mining Courses Directory of free options
  - Software Carpentry / Data Carpentry / Library Carpentry



<sup>\*</sup>JISC (formerly Joint Information Systems Committee) is a UK non-profit supporting higher education

#### **Outreach and Training**

Host events / Workshops / Discussions

- Open Data Day
  - Invite <u>a local Code for America brigade</u> to co-host or facilitate a hack-a-thon or workshop

#### ContentMine

 Identify researchers or applications producers you'd like to work with and work with them to provide an event.

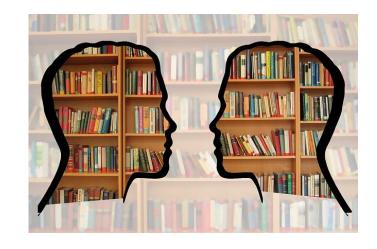
#### TDM Forum

 Invite researchers in your community/institution to share their experiences and expertise



# **Employ TDM to improve library services**

'Bibliomining'



#### **Examples**

"Gaining strategic advantage through bibliomining: Data mining for management decisions in corporate, special, digital, and traditional libraries."

"Use and understand: the inclusion of services against texts in library catalogs and "discovery systems"

"Towards linking libraries and Wikipedia: automatic subject indexing of library records with Wikipedia concepts"

"Data-mining the Library"

# Activity Consulting on a TDM Project: Questions to Ask



If you knew someone were coming to you for a 'TDM Consultation,' what questions might you ask?:

- Ahead of time (if you have the opportunity)
- At the start of the consultation
- During the consultation (what are key info points you'd want to be sure to touch on during the conversation)

# Activity Consulting on a TDM Project: Questions to Ask - Examples



- Does the content source provider support TDM projects, provide policies, or offer contacts for requests?
- How will the file type or extraction method affect the TDM project goals?
- How large is the expected data? How will it be transferred? Where will it be stored?
- If the data is digitized text, what is the quality of the Optical Character Recognition (OCR)? Is it possible to get a sample of the OCR to check for quality?

### → What would you add?

#### **Additional References**

**Report:** Young, P., Brittle, C., Haugen, I., Lener, E., Pannabecker, V. (2017). **Library support for text and data mining: A report for the University Libraries at Virginia Tech**. Retrieved from <a href="http://hdl.handle.net/10919/78466">http://hdl.handle.net/10919/78466</a>

This Presentation: <a href="http://hdl.handle.net/10919/79483">http://hdl.handle.net/10919/79483</a>

Libraries, licensing, and TDM

Text & Data Mining Clauses in Academic Library Licenses: A Case Study

#### **Organization to Follow**

Future TDM

#### Selected recent items of interest since our report was shared

- <u>Text mining of 15 million full-text scientific articles</u>
- <u>JSTOR Labs Text Analyzer</u> for Topics and Recommended Readings
- <u>Legal Analytics Lab</u> at Georgia State University
- <u>Common Crawl</u> web archiving organization | <u>Related article</u>

#### **Images**

#### Sourced from **Pixabay** - CC0 - public domain license

https://pixabay.com/en/binary-binary-system-data-dataset-2728121/ https://pixabay.com/en/road-asphalt-space-sky-clouds-220058/ https://pixabay.com/en/balloon-discussion-comment-2223048/ https://pixabay.com/en/hello-bonjour-hi-greeting-foreign-1502369/ https://pixabay.com/en/chat-multiple-icon-symbol-message-2389223/ https://pixabay.com/en/contract-consultation-pen-signature-1332817/ https://pixabay.com/en/weight-scale-equal-arm-balance-scale-2402966/ https://pixabay.com/en/home-office-workstation-office-336377/ https://pixabay.com/en/people-girls-women-students-2557396/ https://pixabay.com/en/silhouette-head-bookshelf-know-1632912/ https://pixabay.com/en/banner-header-question-mark-1090830/ https://pixabay.com/en/icon-feedback-message-cloud-data-1968237/

# **Questions?**





#### **Evaluation Link**



# Please let us (and VLA) know your thoughts on this presentation!

tinyurl.com/th2017vla



Thank You word cloud by Ashashyou [CC BY-SA 4.0], via Wikimedia Commons

