# OpenRefine Workshop

Meredith L. Hale, Metadata Librarian
mhale16@utk.edu
Open Sandbox Series event on February 16, 2020
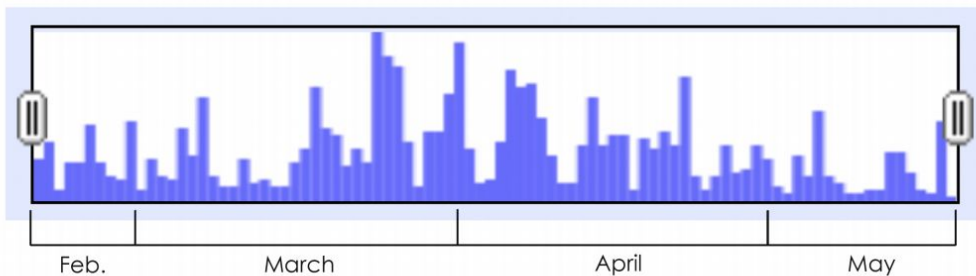
# What is OpenRefine Good For?

- Analyzing Data
- Cleaning Data
- Enhancing Data (with URIs, coordinates, etc.)
- Changing the Format of Data
- Automating Routine Data Processes

# OpenRefine for Analysis

- My first foray into OpenRefine focused on analyzing log data from the Ackland Art Museum



Timeline of the Ackland's Search Log
(Feb. 19 to May 19, 2015)

Feb.    March    April    May

# Finding Your Data

- Pull from an endpoint using the Open Archives Initiative Protocol for Metadata Harvesting (OAI-PMH)
  - https://github.com/vphill/pyoaiharvester
  - UTK Libraries endpoint - http://dloai.lib.utk.edu/cgi-bin/XMLFile/dlmodsoai/oai.pl
- Application Programming Interfaces (APIs)
  - Digital Public Library of America (DPLA)
- Other Open Data Sites
  - https://www.data.gov/
  - https://www.pewresearch.org/internet/datasets/
  - https://data.worldbank.org/
- Create your own data set

# Consider OpenRefine's Limits

- Question whether your data will be effectively manipulated using OpenRefine

- Memory allocation can be increased from the standard 1024MB to enhance performance speeds
  - https://docs.openrefine.org/manual/installing/#increasing-memory-allocation
- Limits to OpenRefine
  - No more than one million total cells
  - Input file size maximum of 50 megabytes (MB)
- My personal experience has shown that spreadsheets with more than 15,000 rows or 256 columns cannot be easily accommodated

# Not all Data are Diamonds . . .

Katie Rawson and Trevor Muñoz, "Against Cleaning", http://curatingmenus.org/articles/against-cleaning/,



VS

# Cleaning your Data

## Built-in functions

- Facet (by text, number, timeline, word)
- Change data format (text, date, number)
- Normalize spacing
- Normalize case
- Split values into different columns
- Clustering to normalize values

## Documentation

- [OpenRefine User Manual](#)
- [OpenRefine Lessons for Digital Humanities](#)

# Enhancing Your Data - Reconciliation

- Reconciliation means standardizing value variations in accordance with an established controlled form of the value

Waterhouse, John
Waterhouse, John W.
Waterhouse, J. W.

Waterhouse, John William
http://vocab.getty.edu/ulan/5
00027032



John William Waterhouse, Destiny, 1900.

# Enhancing Your Data - Reconciliation

Library of Congress Authorities

- [Christina Harlow's service](#)
- [Michael Phillip's service](#)
- [Jeff Chiu's service for LoC Name Authority Files](#)
- [University of Tennessee, Knoxville, Heroku-hosted service](#)

Getty Vocabularies

- [Getty Reconciliation Instructions](#)

Virtual International Authority File (VIAF)

- [Jeff Chiu's service](#)

Geonames

- [Christina Harlow's service](#)
- [University of Tennessee, Knoxville, Heroku-hosted service](#)

# Automating Processes

# OpenRefine in My Work Today



- GitHub - https://github.com/UTKcataloging

- Example export template - https://github.com/UTKcataloging/archivision/blob/master/cleaned_data/remediation_files/open_refine_template.md

# Resources

- "Free Your Metadata",  https://freeyourmetadata.org/
- Library Carpentry's "Introduction to Working with Data - Regular Expressions", https://librarycarpentry.org/lc-data-intro/01-regular-expressions/index.html
- "Open Refine User Manual", https://docs.openrefine.org/
- Verborgh, Ruben and  Max De Wilde. *Using OpenRefine.* Olton: Packt Publishing, Limited, 2013. (Library Ebook)
- Workshop GitHub Repository, https://github.com/mlhale7/OpenRefineWorkshop

Demo!