# MESSY DATA? CLEAN IT UP WITH OpenRefine

Tricia Lampron
Metadata Services Specialist
University of Illinois Library

# GUIDES.LIBRARY.ILLINOIS.EDU/OPENREFINE

Created by Brinna Michael

# OpenRefine

A free, open source, power tool for working with messy data.

- About OpenRefine
- Hands-on:
  - *Importing your data*
  - *Manipulating and cleaning your data*
  - *Exporting your data*
- Resources
- Questions?

# OpenRefine

as a Data Exploration Tool

- OpenRefine can *normalize* and *visualize* your data!

- OpenRefine looks like a *spreadsheet*, but it acts like a *database*.

- Use OpenRefine to *explore*, *clean*, and *link* your data.
  - *Pull out greater granular details when there's a lot of information*
  - *Pinpoint very specific criteria using a combination of facets*

# Strengths:

- Free and Open Source

- Much more powerful than Excel

- Platform independent

- Great history tracking

- Can export commonly used functions for reuse

# Weaknesses:

- Can be a little unstable (some queries will run slowly)

- Some methods require light programming*

- Some tasks that are easy in Excel are more difficult/impossible in OpenRefine
  - *Adding new rows / New Data*
  - *Editing individual cells*

# Hands-On OpenRefine

- Creating a New Project

- Records vs. Rows

- Reordering Columns

- Basic Normalization

- Tracking Operation Histories

- Faceting and Clustering

- Exporting your work

# To get started:

- Download and install OpenRefine
  - http://openrefine.org/download.html

- PLEASE NOTE: OpenRefine runs in a browser, like a web app (though it does not require an internet connection). It works best with Chrome, Firefox, or Safari, and will automatically run on your default browser. OpenRefine is unstable when run on Internet Explorer and does not run at all on Edge, so it is best to choose another default browser.

# SAMPLE DATA

**goo.gl/CBqj2n**

# Creating a project

- TSV, CSV, other separated values
- Line-based; Fixed-width text files
- Spreadsheets
  - *Excel files (.xls, .xlsx)*
  - *Open document format spreadsheets (.ods)*
  - *Google Sheets*
- XML
- JSON

# Creating a project

## Project Settings
- Check your encoding!
- Choose the correct separators!
- Parse your column headers!
- Don't parse text into numbers, dates, etc.! (Unless you are VERY confident about your data)

# OpenRefine Layout

- OpenRefine displays data in a tabular format (like a spreadsheet):
  - *Each line will represent a **'record'** or **'row'** in the data*
  - *Each **column** represents a type of information*
- OpenRefine only displays a limited number of lines of data at a time.
  - *You can adjust the number to **5, 10** (default)**, 25,** or **50**.*
  - *When you manipulate your data, you are manipulating the whole set, not just the first 50!*

# Records vs. Rows

## ROWS

- Individual lines of data, not linked by relationships between the data in neighboring lines

## RECORDS

- Multi-line groupings of data, linked by relationships between data in those lines

# Reordering columns

# Basic normalization

- Trim leading and trailing whitespace
- Collapse consecutive whitespace
- Unescape HTML entities
  -  
  - &amp;
- titlecase, uppercase, lowercase
- number, date, text
- null, empty string

# Sorting Data

# Sorting Data

- Reverse sort
- Remove sort

# Splitting multiple values

| Identifier | Author |
|---|---|
| 123456 | Smith, Joe; Green, Jane |
| 123457 | Collins, Mark; Prestley, Susan |
| | |
| | |

| Identifier | Author |
|---|---|
| 123456 | Smith, Joe |
| | Green, Jane |
| 123457 | Collins, Mark |
| | Prestley, Susan |
| | |

# A note on tracking project history

**Don't be afraid to make mistakes!**

...you can always use the Undo/Redo tab to navigate to earlier stages in your project.

# Text filtering



*You can also use regular expressions in the filter!

# Text filtering

# Faceting

- View unique values

- Select set of data based on unique values

- Edit all cells with the same value

- Extract values and counts

# Editing Values

## In a Cell



## In a Facet

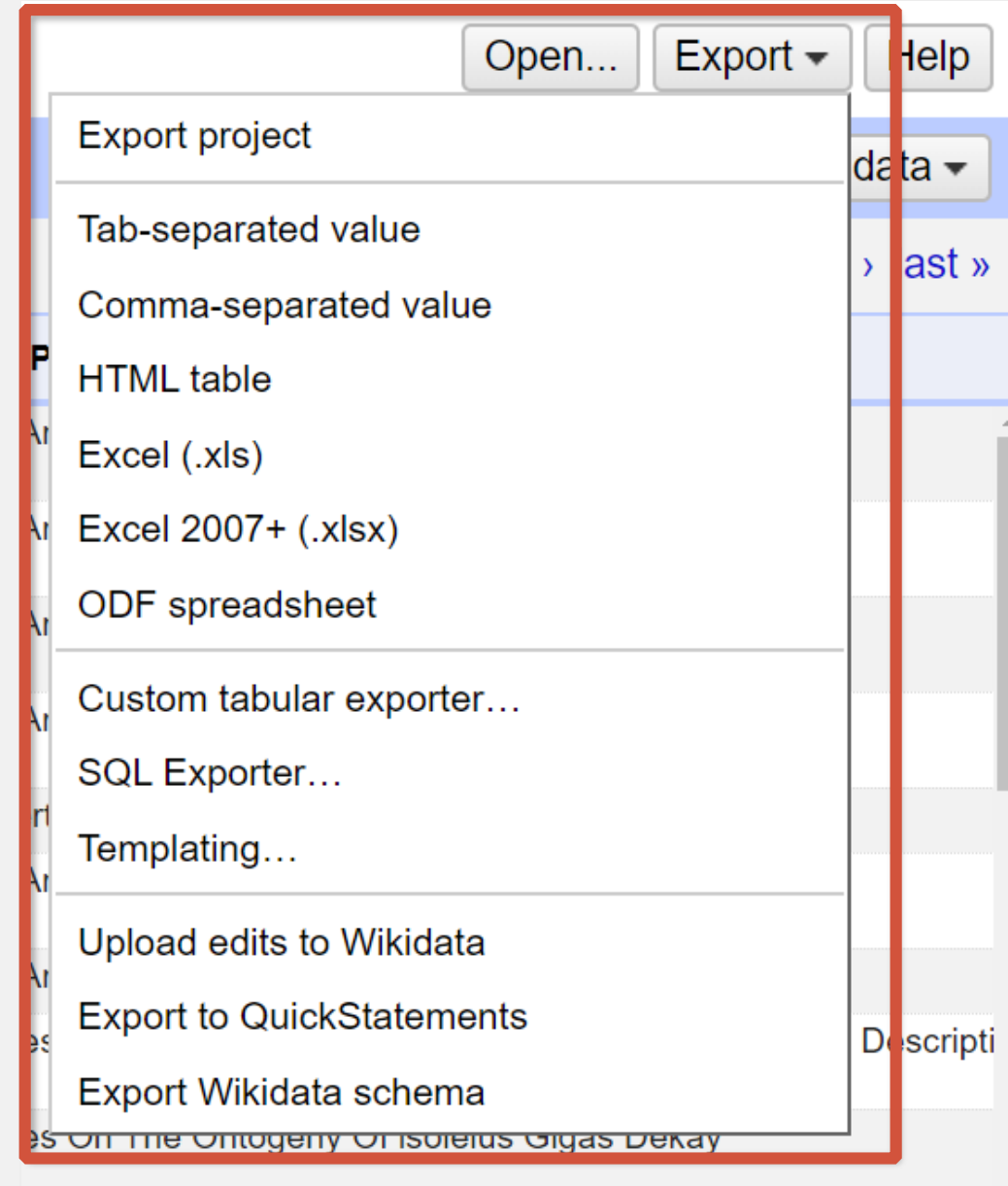# Clustering

# Clustering

## Caveat





Hotel Savoy
59th St. & 5th Ave.
New York, New York



Savoy Hotel
Strand
London WC2R 0EU
United Kingdom

# Exporting your data

# Extracting history (JSON)

# Applying previous history (JSON) to new project

# Resources     guides.library.Illinois.edu/openrefine

- *OpenRefine page*: http://openrefine.org/

- *Github repository*: https://github.com/OpenRefine/OpenRefine

- *Open Refine Documentation:*
  https://github.com/OpenRefine/OpenRefine/wiki/Documentation-For-Users

- Free Online Course:
  https://bigdatauniversity.com/courses/introduction-to-openrefine/

- *Additional Reconciliation Services:*
  https://github.com/OpenRefine/OpenRefine/wiki/Reconcilable-Data-Sources

- *More on Google Regular Expression Language:*
  https://github.com/OpenRefine/OpenRefine/wiki/GREL-String-Functions

- Understanding Expressions:
  https://github.com/OpenRefine/OpenRefine/wiki/Understanding-Expressions

- GREL Cheat Sheet:
  http://arcadiafalcone.net/GoogleRefineCheatSheets.pdf

Slides created and/or updated by: Qian Zhang, Kelly Applegate, Brinna Michael, and Tricia Lampron
Data set provided by Qian Zhang

# QUESTIONS?

Tricia Lampron

[lampron2@Illinois.edu]