# Introduction to OpenRefine: Using Open Software to Weed and Manage a Government Documents Collection

ADAPTED FROM "WORKING WITH MESSY DATA IN OPENREFINE,"
IASSIST 2018 CONFERENCE, LEANNE TRIMBLE AND KELLY SCHULTZ,
CONCORDIA UNIVERSITY, CANADA

Eimmy Solis
University of Southern California
Social Sciences Data Librarian
eimmysol@usc.edu

tinyurl.com/FDLC2019OPENREFINE

# Agenda

- Background
- What is OpenRefine?
- OpenRefine Setup
- Demonstrations and Hands-on Practice
- Additional Helpful Resources

# Learning Objectives

Participants will be able to use OpenRefine to:

- Search, sort, and filter data in a variety of ways
- Restructure and manipulate a dataset
- Perform basic data cleanup

# Background

- NEW GOV DOCS LIBRARIAN
- NO PRIOR WEEDING EXPERIENCE

# Installing OpenRefine

ALL SLIDES, HANDOUTS
AND DATASET HERE:

tinyurl.com/FDLC2019OPENREFINE

OpenRefine is installed locally on your computer, even though it uses a web browser as the user interface.

A copy of your data files are saved locally to your computer.

# What is Messy and Clean Data?

| | A | B |
|---|---|---|
| 1 | **Customer Name** | |
| 3 | John K. **Doe** Jr. | Doe, John |
| 4 | Mr. **Doe**, John | Doe, John |
| 5 | Jane A. **Smith** | Smith, Jane |
| 6 | MS. **Jane Smith** | Smith, Jane |
| 7 | **Smith, Jane** | Smith, Jane |
| 8 | Dr **Anthony** R **Von Fange** III | Von Fange, Anthony |
| 9 | **Peter Tyson** | Tyson, Peter |
| 10 | **Dan** E. **Williams** | Williams, Dan |
| 11 | **James Davis** Sr. | Davis, James |
| 12 | **James** J. **Davis** | Davis, James |
| 13 | Mr. **Donald** Edward **Miller** | Miller, Donald |
| 14 | **Miller, Donald** | Miller, Donald |
| 15 | **Rajesh Krishnan** | Krishnan, Rajesh |
| 16 | **Daniel Chen** | Chen, Daniel |

# What is OpenRefine?

Open source tool for working with messy data to clean and transform it from one format to another.

# Why OpenRefine?

**vs**

# Demonstrations & Hands-on Practice
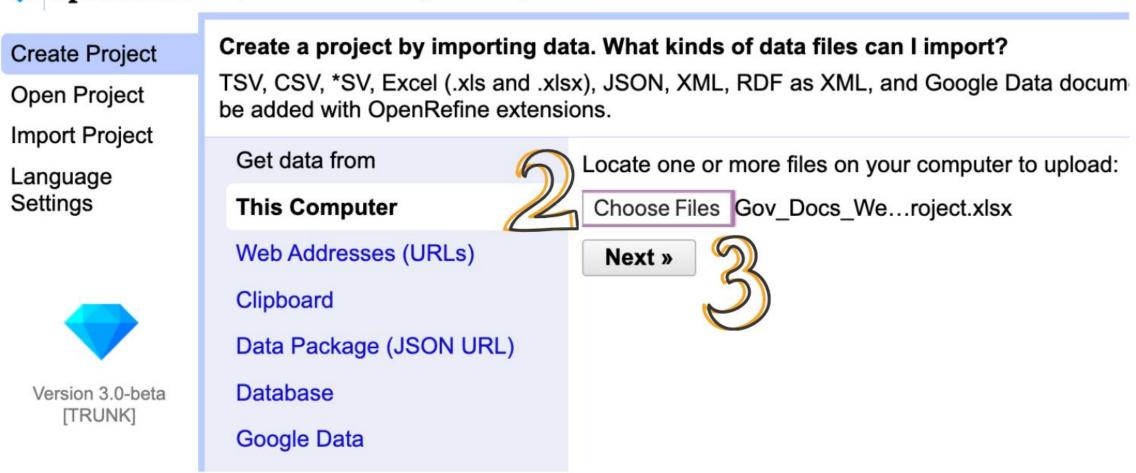
tinyurl.com/FDLC2019OPENREFINE

# IMPORTING A DATASET INTO OPENREFINE

**OpenRefine** *A power tool for working with messy data*

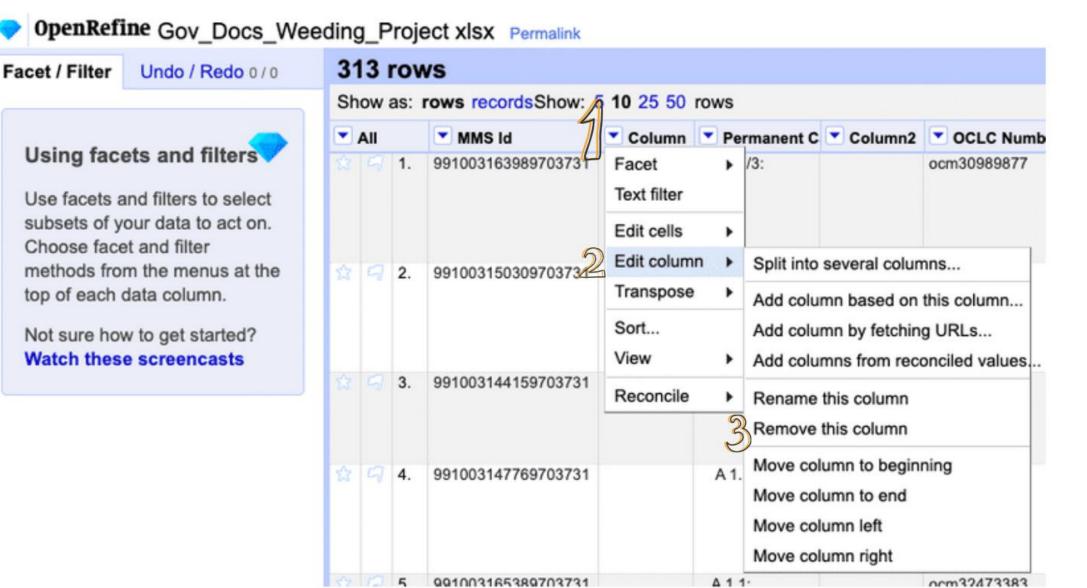New version! **Download OpenRefine v3.2 now.**

**1**

Create Project

Open Project

Import Project

Language Settings

Version 3.0-beta [TRUNK]

**Create a project by importing data. What kinds of data files can I import?**

TSV, CSV, *SV, Excel (.xls and .xlsx), JSON, XML, RDF as XML, and Google Data docum be added with OpenRefine extensions.

Get data from

**This Computer**

Web Addresses (URLs)

Clipboard

Data Package (JSON URL)

Database

Google Data

**2** Locate one or more files on your computer to upload:

Choose Files   Gov_Docs_We…roject.xlsx

**Next »** **3**

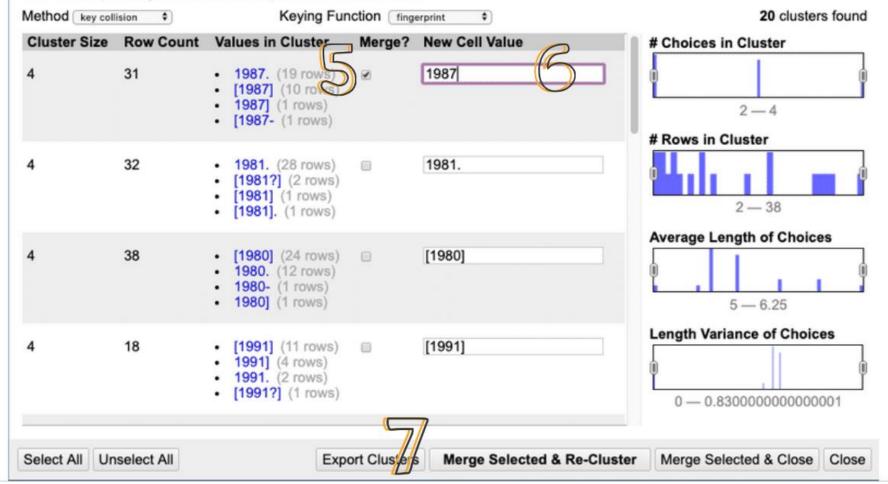# IMPORTING A DATASET INTO OPENREFINE

# REMOVING A COLUMN

# CLUSTERING

# CLUSTERING

**Publication Date** change

61 choices  Sort by: **name** count  Cluster

[1976] 1
[1979] 1
[1980] 24
[1981?] 2
[1981] 1
[1981]. 1
[1983] 3
[1984] 18
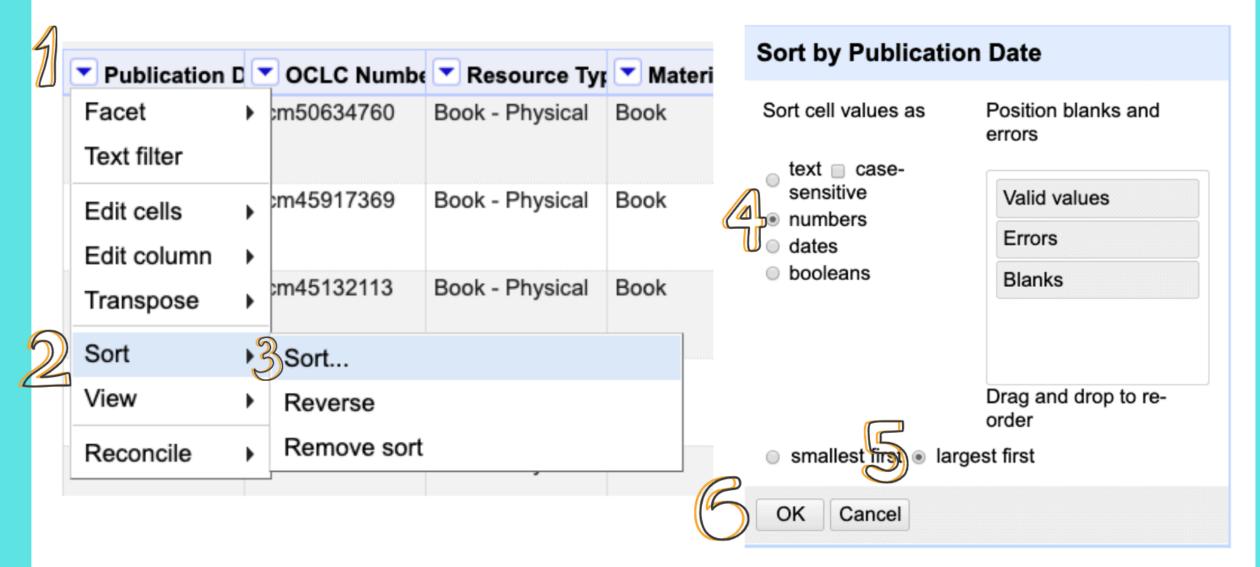[1985] 5
[1986] 1
[1987- 1
[1987] 10

# CLUSTERING

# Sort

# FILTER

# FACET

# RE-ORDER / REMOVE COLUMNS

**313 rows**

Show as: **rows** records Show: 5 10 25 **50** rows

*1*

▼ All | ▼ MMS Id | ▼ Permanent Call

| Transform | |
| Facet ▶ | 13233719703731 | A 1.2:ST 8/3 |
| Edit rows ▶ | |
| Edit columns ▶ *3* Re-order / remove columns... | |
| View ▶ | 08612489703731 | A 1.2:L 75/3 |

*2*

*4* Drag columns to re-order          Drop columns here to remove

| Title |
| Permanent Call Number |
| Publication Date |
| OCLC Number |
| Resource Type - Bibliographic Details |
| Material Type - Bibliographic Details |
| Material Type - Physical Item Details |
| Receiving Date |
| Library Name |
| Column3 |
| Column4 |
| Column5 |

| Location Name |
| Network Number |
| MMS Id |

*5* OK   Cancel

# Closing OpenRefine

- Click on OpenRefine icon and type Command- Q.

- Wait until there's a message that says the shutdown is complete.

# Helpful Resoures

- OpenRefine documentation wiki: https://github.com/OpenRefine/OpenRefine/wiki/Documentation-For-Users

- OpenRefine Tutorial from John Little (Duke University): https://libjohn.github.io/openrefine/index.html

- Software Carpentry OpenRefine Workshop: https://data-lessons.github.io/library-openrefine/

- Cleaning Data with OpenRefine from the Programming Historian: https://programminghistorian.org/lessons/cleaning-data-with-openrefine

- Fetching and Parsing Data from the Web with OpenRefine from the Programming Historian: https://programminghistorian.org/lessons/fetch-and-parse-data-with-openrefine

- Regex Cheat Sheet: http://www.rexegg.com/regex-quickstart.html

# Questions?

**Eimmy Solis**
**eimmysol@usc.edu**
**tinyurl.com/FDLC2019OPENREFINE**