

Cleaning and enriching data with OpenRefine

ETHAN YOO, DATA SCIENCE GRADUATE SPECIALIST

FEBRUARY 7, 2023

ETHAN.YOO@RUTGERS.EDU

Outline

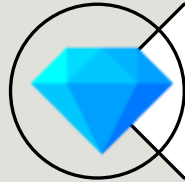
- OpenRefine and its core functions
- Cleaning data with OpenRefine
- Enriching data with OpenRefine
- Demonstration/walkthrough
- Additional resources

OpenRefine

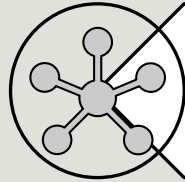
- “a Java-based power tool that allows you to load data, understand it, clean it up, reconcile it, and augment it with data coming from the web” ([OpenRefine](#))
- Free and open source software



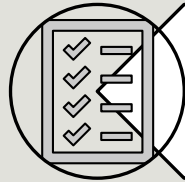
OpenRefine's Functionality



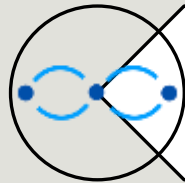
Faceting



Clustering and editing



Reconciling



History

Cleaning data



Enriching data (data augmentation)



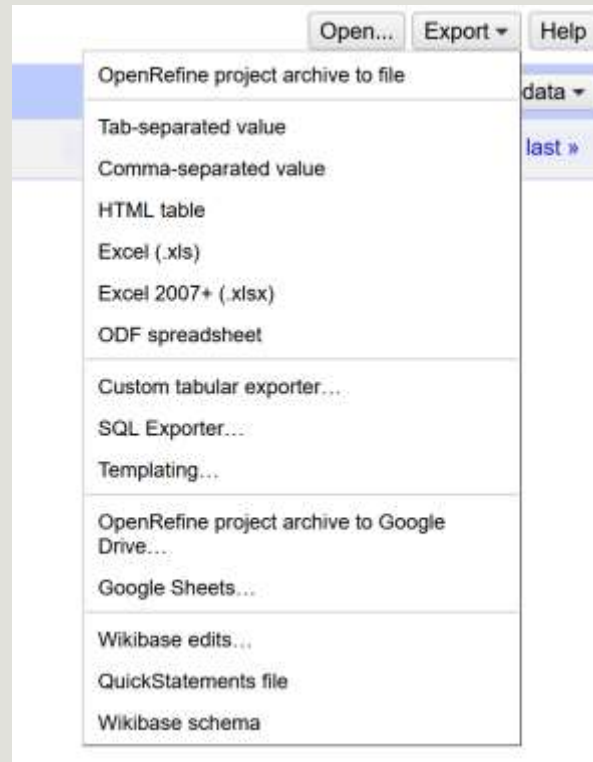
The screenshot shows a web application for data enrichment. On the left, there is a sidebar with a search bar and a list of categories. The main area displays a table with multiple columns, including 'Original Data', 'Enriched Data', and 'Status'. The table contains several rows of data, with some cells highlighted in yellow. The interface is clean and modern, with a blue header and a white body.

Original Data	Enriched Data	Status
100	100	Success
200	200	Success
300	300	Success
400	400	Success
500	500	Success
600	600	Success
700	700	Success
800	800	Success
900	900	Success
1000	1000	Success

Running OpenRefine for the first time

- Requires a browser but an Internet connection is not required for most actions
 - Open any web browser (e.g., Mozilla Firefox or Google Chrome)
 - Navigate to <http://127.0.0.1:3333> (or <http://localhost:3333>)
 - Create a project by importing data
 - **Supported:** TSV, CSV, and other delimited files, Excel files, JSON, XML
 - **Sources:** Local computer, web address, clipboard, relational database (SQLite, PostgreSQL, MySQL/MariaDB), Google Sheets (public spreadsheet or authenticated account)

Exporting or backing up work



Select “Export” in the upper-right corner and choose an option

Additional Resources

- **Data Carpentry**
 - [OpenRefine for Social Science Data](#)
 - [Data Cleaning with Open Refine for Ecologists](#)
- [**OpenRefine user manual**](#)

Microsoft Excel

- **Data**
 - Get & Transform Data
 - Get Data From Other Sources



Microsoft Excel

- **Data**
 - Get & Transform Data
 - Get Data From Other Sources

