

Academic BRASS

Published by the
BRASS Board of Reference in Academic Libraries Committee

Volume 19(2), Fall 2014

Editorial Board
Information Specialist and Assistant Professor of Library Science
Paris Library of Management and Economics
Paris University

Good Library Data Made Better With Technology! Using OpenRefine and Google Fusion Tables in Academic Business Libraries Instruction

Introduction

Big data just seems to get bigger all the time, but that doesn't mean it gets any less messy. Enlarge, carefully curated, and go to the heart of the matter. Data is no longer just a collection of items and metadata, but a complex web of relationships. Librarians have the patience for rich, accurate, and detailed data. They are often prepared for the realities of the big data heyday. Teaching data cleaning and collaboration can help them better understand and manage large datasets. It also illustrates the importance of library-curated data, as it often has fewer of the problems associated with online open data. At a high level, library data and open data may seem comparable, but they are not. The big data is not going through the data on their own, but the small things that add up.

This article will discuss the focus of a recent book titled *From Data to Knowledge*, which helps to do this deeper dive into the data tools Google Fusion Tables and OpenRefine. It is my hope to make librarians more aware of the tools and how they can be integrated into instruction. I first learned about Google Fusion Tables and OpenRefine from attending the Annual Arbor Data Day in 2012. Data Day is a not-for-profit weekend intensive data event. Other tools used by Data Day are available at <http://opendata-toolbox.org/en/>. As librarians continue to tackle issues of data collection and reuse, approaches and tools like this can help illustrate the importance of good research methods.

Setting

During the spring 2015 semester, the library was approached by a professor of an Electronic Commerce and Information Strategies course. The professor, who had worked in the libraries in the past, developed a final assignment where groups of students would actually

analyzed a dataset to solve a business problem. The data, all of which were in the period of the management program, were given a large amount of freedom to choose any business problem that interested them. The professor invited the libraries to give a 30-minute presentation (20 minutes) to the students at the library.

To prepare for the course, my colleague and I made a library guide of potential sources of data. Given the breadth of the project and the collaborative nature of the course, the time included seemed better spent for creating a number of likely influential cases and not on the large amount of data sources that would be available. The libraries also provided information on Google Fusion Tables and OpenRefine.

Google Fusion Tables

Google Fusion Tables is a web-based data management tool (see Figure 1). Fusion Tables was first introduced in a scientific paper in 2010 (Gonzalez et al., 2010). Google Fusion Tables can be accessed through Google Drive, Google Sheets, or from <https://www.google.com/fusiontables>. Like Google Docs, Google Fusion Tables allows for samples multiple users editing of spreadsheets in real time. It includes many basic visualization tools. A critical additional feature of Google Fusion Tables is the ability to filter data by tables together if they share a common dimension such as zip-code, city, or age.

Fusion Tables was a topic for me for several reasons. First, I knew of the tool from a colleague who collected their own data, and I knew they would probably use Google Docs for that. Second, I knew many datasets had common attributes that would allow for interesting comparisons. Third, the interface is simple and easy to learn. Students already have Google accounts and so do not need to sign up for anything.

Figure 1 Google Fusion Tables example data. I pulled the data from Wikipedia and LexisNexis and limited by temperature and country.

OpenRefine

OpenRefine is an open source desktop application for data cleaning and analysis. It is similar to many aspects of Microsoft Excel, however, it acts more like a database. For those working with large datasets, OpenRefine can help them clean up small issues across a large dataset. Several video tutorials are available here: <http://openrefine.org/>. OpenRefine allows you to import Excel and CSV files as well as Google Fusion Tables.

OpenRefine is a great tool to use in a scenario because it helps identify the messy characteristics of large datasets. The state of Indiana may be listed as "IN" in one place and "INDIANA, State of" in another. In OpenRefine, it is

Putting it all Together in a One-Shot Session

To prepare for the hour (under 20 minutes) session, I pulled together a Google Fusion Tables dataset of 4,728 citations in Alabama (see Figure 1). I first imported the data into Google Fusion Tables, then I cleaned up the data, and finally I

and finally I

The results highlighted by major bibliometric studies likely page the
these : research collaboration The countries that follow are very different from
other countries have had in the past. Studies have shown that the high
so where also more aware of the measures of the data. When I had to read them high
level of the many data the library get, they had to be more different
to be able to compare to what they could find searching around. But what the
collaboration is more about the research access the library database a lot more re
active.

Conclusion

This article has covered by tools, Google Fusion Tables and OpenRefine, and their
in the undergraduate library into a Google Fusion Table offers collaboration and data
merging features a familiar environment. OpenRefine is a powerful tool
data cleaning. These tools highlight libraries can provide to their data for
to use big-data analysis that amplify the traditional role of access and research.

Works Cited:

Gonzalez H., Haley A., Jensen C. S., Larsen A., Madhavan J., Shapley R., & Shen W.
(2010). Google fusion tables data management integrated collaboration in the cloud.
In *Proceedings of the 1st ACM symposium on Cloud computing* (pp 175– 180). ACM.