# Academic BRASS

Published by the
BRASS Business Reference in Academic Libraries Committee

Vol 9(2), Fall 2014


Ilana Stonebraker
Business Information Specialist and Assistant Professor of Library Science
Parrish Library of Management and Economics
Purdue 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 does mean it gets any less messy. Even large, carefully curated government datasets suffer from irregularities like acronyms, open reprecord items, and missed categories. Steadfast librarians have the patience for such inaccuracies, but undergraduate students are often unprepared for the realities of the big data they crave. Teaching data cleaning and collaboration can help students better understand and use large datasets but also illustrate the importance of library-curated data, as it often has fewer of these problems than data we find on the open web. At a high level, library data and open datasets may be seem comparable, but when we give students the tools to go through the data on their own the small things start to add up.

This short article will discuss shifting the focus of a one-shot time to naway from datasets to help students do this deeper dive, using the data tools Google Fusion Tables and OpenRefine. It is my hope to make librarians more aware of the two data tools and how they can be integrated into business instruction. I first learned about Google Fusion Tables and OpenRefine from attending the Ann Arbor Data Dive in 2012. Data Dives are nonprofit weekend intensive data events. Other tools used by Data Dives are available at http://opendata-tools.org/en/. As business librarians come more to tackle issues of student data collection and reuse, approaches and tools like these more can help illustrate the importance of good sources and 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 with the libraries in the past, developed a final assignment where groups of students found and statistically

analyzed a dataset to solve a business problem. The students, all of whom were in the upper division 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 short presentation (20 minutes) to the students about library resources.

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 in class seemed better spent focusing on tools likely useful in all cases and not on the large amount of data sources the students could potentially use. The libraries also wanted to focus on Google Fusion Tables and OpenRefine.

## Google Fusion Tables

Google Fusion Tables is a web-hosted data management product (see Figure 1). Fusion 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 allow for seamless multiple-user editing of spreadsheets but also includes many basic visualization tools as well. A critical additional feature of Google Fusion Tables is the ability for "matching" two tables together if they share a common business information column such as zip-code, city, or age.

Fusion Tables was attractive to me for several reasons. First, I knew some of the students were looking to collect 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 seamless merging. Third, the interface is very simple and easy to learn. Students already have Google accounts and so do not need to sign up for any service.

**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 cleanup that operates in-browser. It is similar in many aspects to Microsoft Excel, however, it functions more like a database. For those working with large datasets, OpenRefine can help the student clean up small issues such acronym inspelling errors, etc. Several videos about how to it works are available here: http://openrefine.org/ OpenRefine allows input from Excel and CSV files as well as Google Fusion Tables

OpenRefine is a tackle in an intro to session because it helps students visualize the messy characteristics of large datasets. The state of Indiana may be listed as "IN" in some places and "INDIANA, State of" in others, but both will show up in OpenRefine's list

**Putting it all Together in a One-Shot Session**

So as prepared for the short (less/20 mins) session, T c -0.004 Teap2 (have f1 te2 (ed ) F6.6 (/.0T4bT5 (a)-6.6 ) 2 (bo)/45.9 (m)-6.6 (0.8

The short also highlighted two major biases informationists likely to play: the students: memes and collaboration. The collaborations that followed were very different from other collaborations I have had in the past. Students had questions about the things I had shown but were also more aware of the memes of the data. When I had showed them high level versions of the many datasets the library suggested, they had not seemed much different to the student compared to what they could find searching around. But once the conversation was more about memes servers access, the library datasets were a lot more attractive.

## Conclusion

This short article has covered two tools, Google Fusion Tables and OpenRefine, and their use in business undergraduate library instruction. Google Fusion Table offers collaborative and data merging features in a familiar environment students use. OpenRefine is a powerful yet simple data cleaning option. These tools highlight ways libraries can provide value to their datasets for novice big-data analysis that amplify the traditional role of access and preservation.

Works Cited:

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