

Converting data to friendly formats

Tools for scraping data from the web

[illegible]

Figure 1: Demographic Characteristics of Respondents by Group

1. Age

Age Group	Group 1	Group 2	Group 3	Group 4	Group 5
18-24	15%	20%	10%	25%	18%
25-34	30%	25%	35%	20%	30%
35-44	25%	30%	20%	35%	25%
45-54	20%	25%	30%	20%	25%
55-64	15%	15%	15%	15%	15%
65+	15%	15%	10%	5%	12%

2. Education

Education Level	Group 1	Group 2	Group 3	Group 4	Group 5
High School	10%	15%	5%	20%	12%
Some College	20%	25%	10%	15%	20%
Bachelor's	35%	30%	40%	35%	30%
Master's	25%	20%	30%	25%	25%
PhD	10%	10%	15%	5%	13%

3. Income

Income Level	Group 1	Group 2	Group 3	Group 4	Group 5
<\$10,000	5%	8%	3%	12%	6%
\$10,000-\$19,999	15%	18%	10%	10%	15%
\$20,000-\$29,999	25%	22%	20%	25%	22%
\$30,000-\$39,999	20%	20%	25%	20%	20%
\$40,000-\$49,999	15%	15%	15%	15%	15%
\$50,000+	20%	15%	25%	20%	22%

4. Employment

Employment Status	Group 1	Group 2	Group 3	Group 4	Group 5
Full-time	40%	35%	45%	30%	38%
Part-time	25%	30%	20%	35%	25%
Unemployed	15%	18%	10%	12%	15%
Retired	10%	10%	15%	10%	12%
Other	10%	7%	10%	13%	10%

5. Marital Status

Marital Status	Group 1	Group 2	Group 3	Group 4	Group 5
Single	30%	25%	35%	20%	28%
Married	45%	50%	40%	55%	45%
Divorced	15%	12%	18%	10%	15%
Widowed	10%	13%	7%	15%	12%

[illegible]

Tables

- Copying and pasting into Excel does not usually work

Free online tools:

- cometdocs.com
- pdftoexcelonline.com
- zamzar.com
- pdftoexcel.org

Free desktop software

- tabula.nerdpower.org

Let's try it!

- Open the File: Country Profile President's Malaria Initiative (PMI) or Maternal Index
- Make sure you understand the data: read all codes, descriptions and legends
- Convert file using one of the free online programs
- Review xls or csv file and evaluate whether the resulting tables are usable
- Clean up table, eliminating unnecessary information and note data source

Complex Formats

- Columns, data divided among various lines, mixed with text
- Automatic tools probably won't work
- The PDFs will have to be extracted and a programmer will use programming language to create a program to extract the data that you want

Scanned Images

Optical Character Recognition Software

- For documents that you can't search or select text
- Quality depends on quality of image, clarity of text
- Adobe Acrobat Professional: paid
- Google Docs: free
- [Document Cloud](#): free. Advantage is that it allows you to annotate and publish your document as part of a multimedia publication. Your media outlet must request an account

Scanned Images

Optical Character Recognition Software

- For documents that you can't search or select text
- Quality depends on quality of image, clarity of text
- Adobe Acrobat Professional: paid
- Google Docs: free
- [Document Cloud](#): free. Advantage is that it allows you to annotate and publish your document as part of a multimedia publication. Your media outlet must request an account

Let's try it!

- Open Specimen Financial Statement
- Open Google Drive
- Select “Upload”
- Choose “Convert text from PDF and image files to Google Documents”
- Review quality
- Make corrections by comparing original file to OCR output

Scraping data from the web

- Browser scrapers:
 - [Dafizilla Table2Clipboard](#) for Mozilla Firefox
 - [Scraper Extension](#) for Chrome
- Allow you to select tables on a website, including specific rows and columns
- Copy the data tables and either create a Google spreadsheet or paste into a blank Excel file

Try it!

- Google: “demographics of India”
- Open wikipedia entry
- Highlight table
- Right click
- In Chrome, select “Scrape similar” and “Export to Google Docs”
- In Firefox, select “Table2Clipboard,” “Copy whole table” and paste in Google Doc or Excel

ScraperWiki: more complicated!

Webpage Screenshot



PROFESSIONAL SERVICES HELP BLOG LOG IN

Liberate your data with ScraperWiki!

ScraperWiki helps you do data science on the web.

Get, clean, analyse, visualise and manage your data, with simple tools or custom-written code.



ScraperWiki for businesses

Sign up – Free



I usually explore data in Excel, but I want to take my data analysis to the next level.



I want a place to write and run code that does awesome stuff with data.



I'm a business and I need professional data extraction, analysis and management.

Super-charge your data analysis with tools including...

Scraping for Journalism: A Guide for Collecting Data by ProPublica

Webpage Screenshot

Don't Miss: [Fracking](#) | [IRS](#) | [Dollars for Docs](#) | [Surveillance](#) | [Patient Safety](#) | [Prescriber Checkup](#) | [Debt Inc.](#) | [990s](#) | [NSA](#)

[DONATE](#)



Journalism in the Public Interest



Receive our top stories daily

Email address

[SUBSCRIBE](#)

[Home](#)

[Our Investigations](#)

[Tools & Data](#)

[MuckReads](#)

[Get Involved](#)

[About Us](#)



Search ProPublica



The ProPublica Nerd Blog

Secrets for Data Journalists and Newsroom Developers



956



17

Google+



Scraping for Journalism: A Guide for Collecting Data

by [Dan Nguyen](#)

ProPublica, Dec. 30, 2010, 5:23 p.m.

17 Comments |

Our [Dollars for Docs news application](#) lets readers search pharmaceutical company payments to doctors. We've written a series of how-to guides explaining how we collected the data.

Most of the techniques are within the ability of the moderately experienced programmer. The most [difficult-to-scrape site](#) was actually a previous [Adobe Flash incarnation](#) of Eli Lilly's disclosure site. Lilly has since released [their data in PDF format](#).



Photo by Dan Nguyen/ProPublica

The News Apps Team



Our Tools and Style Guides

[transcribable](#)

Drop in crowdsourcing for your Rails app.
Extracted from [Free the Files](#)

Need help scraping data?

- Ask the [School of Data](#)
- Join the [data driven journalism](#) listserve
- Join the [Hacks/Hackers Nairobi Google group](#) and go to the next meeting or post your question and ask for volunteers
- Visit the Internews data journalism team for support through the whole process, from finding and scraping data to visualization and storytelling