

# Finding Data

*Steve Pittard*

*February 2, 2016*

## Finding Some Interesting Data

This document attempts to identify sources of information that might be of interest for those needing data to analyze as part of a larger project. In general you can search for open datasets using Google with a search string like (“finding open data sets”) so there really is no practical limit to what one can find. Note that some private data sets have crept out onto the Internet so don’t assume that all data is open or free by default. This is especially true if you encounter data found on Bit Torrent sites. For example, The Netflix Challenge involved the distribution of a large dataset to be used as part of a challenge. Unfortunately, Netflix discovered that that data could be used to reverse engineer specific user ids of actual Netflix users so they withdrew it from distribution. But that didn’t stop it from being copied a bunch of times prior to that. So, now it is not “legal” for one to use it but many still do and there are a variety of sites that still offer it for download. Proceed at your own risk.

Places like The Internet Movie Database use to have free data out there but since the rise of Big Data Analytics they have closed off access except for paying customers. They offer smaller subsets of data in unstructured format which makes it inconvenient (though not impossible) to effectively mine that kind of data. Lots of companies are moving into a direction wherein they protect their data so this isn’t peculiar to IMDB.

## Application Programming Interface ?

Some of these services offer Application Programming Interfaces (APIs) as well as built in visualization tools that can be helpful as good starting points for how one might present the data graphically. There is a web site called The Programmable Web <http://www.programmableweb.com/> that maintains a list of all data providers who offer an API. This is useful for identifying websites that allow you to programmatically access their data. They usually require some type of sign up or registration prior to use. Some of them might not be free but many are. Usually if they are free they will limit you to a specific number of access in a given time frame.

## Formats ?

Lastly, keep in mind that not all of these sites offer so called “Big Datasets” but many do. Some of these require that you navigate to a specific section of the site to obtain raw or .CSV files. It might also be that the sites offer only a specific format such as XML or JSON in which case you will have to parse that information. We’ll be looking at how to do that for some data formats later in the class. As an example of formats that are hard to parse the Million Song Database is provided in an HD5 format which can be parsed but not so easily. Amazon has offered space to host the entire 300GB dataset in a but it’s too large for most people to download conveniently. Of course Amazon wants you to use their compute resources to analyze this data which is why they offer it for free. There is a subset of the MSD at <http://labrosa.ee.columbia.edu/millionsong/pages/getting-dataset>

## Various Links

The list below is current to the best of my knowledge. That is all the links work as of this writing unless I’ve made a typo which is entirely possible. The list of interesting data grows all the time so just pick a domain of interest and Google it like “DNA Sequencing open data set” and I guarantee you that something

interesting will turn up. Also when perusing these sites if you find one of interest also search it to see if it offers an API which would simplify access. Sites are infamous for providing API access only to later revoke it usually because too many people use it. Or worse, they start charging lots of money for people to use it.

Description	Link
Medicare Data including Nursing Home information	<a href="http://data.medicare.gov">http://data.medicare.gov</a>
Amazon Public Data Sets	<a href="https://aws.amazon.com/public-data-sets/">https://aws.amazon.com/public-data-sets/</a>
Awesome Public Data Sets	<a href="https://github.com/caesar0301/awesome-public-datasets">https://github.com/caesar0301/awesome-public-datasets</a>
Stanford Large Network Data Collection	<a href="http://snap.stanford.edu/data/">http://snap.stanford.edu/data/</a>
Machine Learning Data Repository	<a href="http://mldata.org/">http://mldata.org/</a>
UC Machine Learning Repository	<a href="http://archive.ics.uci.edu/ml/">http://archive.ics.uci.edu/ml/</a>
City of Chicago	<a href="https://data.cityofchicago.org/">https://data.cityofchicago.org/</a>
IMDB Movie Review Test Data	<a href="http://ai.stanford.edu/~amaas/data/sentiment/">http://ai.stanford.edu/~amaas/data/sentiment/</a>
City of San Francisco	<a href="https://data.sfgov.org/">https://data.sfgov.org/</a>
New York City Data	<a href="https://nycplatform.socrata.com/">https://nycplatform.socrata.com/</a>
World Bank Data	<a href="http://econ.worldbank.org/datasets">http://econ.worldbank.org/datasets</a>
US Bureau of Labor Stats	<a href="http://www.bls.gov/data/">http://www.bls.gov/data/</a>
US Census Data	<a href="http://www.census.gov/">http://www.census.gov/</a>
US Energy Info Administration	<a href="http://www.eia.gov/opendata/">http://www.eia.gov/opendata/</a>
US EPA Data	<a href="http://developer.epa.gov/category/data/">http://developer.epa.gov/category/data/</a>
USDA	<a href="http://www.ers.usda.gov/data-products/.aspx">http://www.ers.usda.gov/data-products/.aspx</a>
US Open Data (Many types of Data)	<a href="http://www.data.gov/">http://www.data.gov/</a>
UK Open Data	<a href="http://data.gov.uk/">http://data.gov.uk/</a>
European Open Data Portal	<a href="http://open-data.europa.eu/en/data/">http://open-data.europa.eu/en/data/</a>
US Health Data	<a href="http://www.healthdata.gov/">http://www.healthdata.gov/</a>
Canada Open Data	<a href="http://open.canada.ca/">http://open.canada.ca/</a>
France Open Data	<a href="http://www.data.gouv.fr">http://www.data.gouv.fr</a>
Gap Minder	<a href="http://www.gapminder.org/data/">http://www.gapminder.org/data/</a>
Socrata Open Data	<a href="https://opendata.socrata.com/">https://opendata.socrata.com/</a>
World Health Organization	<a href="http://www.who.int/gho/en/">http://www.who.int/gho/en/</a>
Enron Email Data	<a href="http://www.cs.cmu.edu/~enron/">http://www.cs.cmu.edu/~enron/</a>
CDC NHANES Data	<a href="http://www.cdc.gov/nchs/nhanes/nhanes_questionnaires.htm">http://www.cdc.gov/nchs/nhanes/nhanes_questionnaires.htm</a>
Airline Challenge Data	<a href="http://stat-computing.org/dataexpo/2009/the-data.html">http://stat-computing.org/dataexpo/2009/the-data.html</a>
Kaggle	<a href="https://www.kaggle.com/">https://www.kaggle.com/</a>
UC Berkeley Data Page	<a href="http://ucdata.berkeley.edu/">http://ucdata.berkeley.edu/</a>
Links to Data on Reddit	<a href="https://www.reddit.com/r/datasets">https://www.reddit.com/r/datasets</a>
Wiki Data Downloads	<a href="https://en.wikipedia.org/wiki/Wikipedia:Database_download">https://en.wikipedia.org/wiki/Wikipedia:Database_download</a>
New York Times API	<a href="http://developer.nytimes.com/docs">http://developer.nytimes.com/docs</a>