Motivation
oooo

Get data
ooooooooooo

Extract/clean
oooooo

# Part I: Acquiring Data from the Web

Solomon Messing

Department of Communication, Statistics
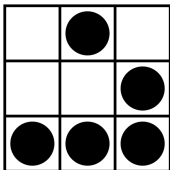Stanford Social Science Data and Software (SSDS)

September 24, 2012

**Motivation**
○●○○

Get data
○○○○○○○○○○○

Extract/clean
○○○○○○

## Overview

- How to get data from the web (cURL, APIs, JSON, XML)
- Extracting useful stuff from the results (Regex)
- Representing text as data
- Getting meaningful features (Regex, R)
- Document-term matrices
- Supervised and unsupervised approaches to analysis

**Motivation**
○●○○

Get data
○○○○○○○○○○○

Extract/clean
○○○○○○

## What to expect

- Goal: Quantitative insight from haystack of messy data
- Labs: programming, lots of R!
- Working with your neighbors
- Adapting starter code + code found on Google, StackOverflow, etc.
- You learn better, I talk less, you have code to work with.

**Motivation**
○○●○

Get data
○○○○○○○○○○○

Extract/clean
○○○○○○

## Motivation

- Massive growth in availability of text and other data
- 100K Tweets per min, 700K FB shares per min (DOMO), 200M emails—10 min = 1 LOC (Huggins)
- Proliferation of structured/semi-structured data at your fingertips: open-data, APIs, and scrape-able data sources
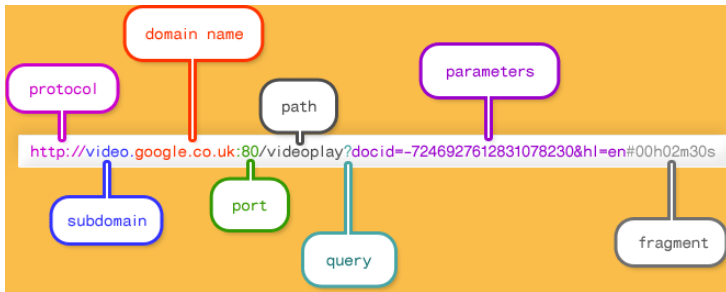- Growth of open-source tools to acquire and analyze this data

Motivation
○○○●

Get data
○○○○○○○○○○○

Extract/clean
○○○○○○

# Motivation: what's out there?

- Raw text, html tables, semi-structured HTML/XML.
- Data marketplaces, e.g.:
  http://www.infochimps.com/datasets
- APIs:
  - https://dev.twitter.com/
  - http://developer.nytimes.com/ articles, campaign finance, congress, entities, geography/population
  - http://developer.washingtonpost.com/ political speeches, campaign finance, White House visitors log
  - http://developer.yahoo.com/everything.html - search, finance, geo-coding, on-the-fly entity extraction, content analysis, term extraction.
  - https://api.facebook.com - access (a little) Facebook data.[1]

---

[1]e.g., get aggregate likes for NYTimes.com article

Motivation
oooo

Get data
●oooooooooo

Extract/clean
oooooo

## Data from the Web

1. Hit a server
2. Parse it's response
3. Clean and transform into something useful
4. Often by merging it with something else.

Motivation
oooo

Get data
o●ooooooooooo

Extract/clean
oooooo

# Hit a server



from: http://doepud.co.uk/blog/anatomy-of-a-url.php

Motivation
oooo

Get data
oo●ooooooooo

Extract/clean
oooooo

# Hit a server, in *.NIX



Unix/Linux was made for this. Windows was not.

Motivation
oooo

Get data
oooo●oooooo

Extract/clean
oooooo

## Hit a server with cURL

- Use cURL (RCurl) to send GET or POST request to server for a URL/URI
- `curl http://thecaucus.blogs.nytimes.com`
- `curl -L http://t.co/KtxsapBV`
- `curl http://search.twitter.com/search.json?q=%40obama`
- The latter is a query string, can used to return custom results from MANY websites (Twitter, nytimes.com, Library of Congress, etc.).

  Query string trivia: the following string has brought down many a web server:
  `system%28%27rm+-rf+%2F%27%29`

Motivation
oooo

Get data
ooooo●oooooo

Extract/clean
oooooo

# Hit a server (or an entire site!)

- Use wget or a crawler
- `wget http://thecaucus.blogs.nytimes.com`
- See also module 3
  `http://www.stanford.edu/~seanjw/module3/#8`
- Heritrix `http://crawler.archive.org/`

Motivation
○○○○

Get data
○○○○○○●○○○○○

Extract/clean
○○○○○○

# Hit a server (on schedule!)

- Use `cron`
- `crontab -e`
- in VIM type `* * * * * /4 /path/to/R CMD myfile.R` to run every Wed
- or perhaps `* * * * * /4 /path/wget http://thecaucus.blogs.nytimes.com`
- Save and you'll see: `crontab: installing new crontab`
- Type `crontab -l` to see your crontabls
- See [http://benr75.com/pages/using_crontab_mac_os_x_unix_linux](http://benr75.com/pages/using_crontab_mac_os_x_unix_linux) for more.

Motivation
oooo

Get data
ooooooo●oooo

Extract/clean
oooooo

# Parse the server's reply: JSON

- What the **** is JSON?
- JSON: Java script object notation
- for serializing objects, for the purposes of data interchange
- semi-structured, tree-like, and pretty simple
- Nice JSON viewer for Chrome and Firefox—Use this in the R code.
- Sean's overview from module 3
  http://www.stanford.edu/~seanjw/module3/#42
- Widom's overview of JSON here:
  http://www.db-class.org/course/video/preview_list

Motivation
○○○○

Get data
○○○○○○○○●○○○

Extract/clean
○○○○○○

# Record and parse it: JSON



```
{
    "photos": {
        "page": 1,
        "pages": 94276,
        "perpage": 15,
        "total": "1414129",
        "photo": [{
            "id": "3891667770",
            "owner": "35468133120@N01",
            "secret": "4479baebf9",
            "server": "2451",
            "farm": 3,
            "title": "Mexican train dominoes with Brian and Michelle",
            "ispublic": 1,
            "isfriend": 0,
            "isfamily": 0
        },
        {
            "id": "3891661852",
            "owner": "10640301@N07",
            "secret": "79de502257",|
            "server": "2590",
            "farm": 3,
            "title": "Peaches",
            "ispublic": 1,
            "isfriend": 0,
            "isfamily": 0
        },
```
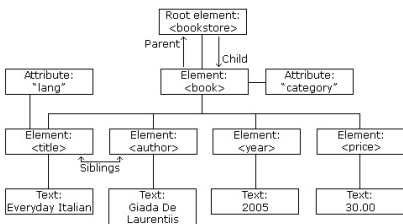
- Base values
- Objects {} - label-value pairs
- Arrays [] - list of values
- nested sets of arrays - not a table
- NO SCHEMA
- NO SQL (hard to query)

## Record and parse it: XML

- What the **** is XML?
- XML: eXtensible Markup Language
- Like HTML, but tags describe data, not formatting
- Semi-structured, tree-like, much richer, more complicated than JSON
- Nice XML viewer for Chrome.
- Sean's overview in module 3
  http://www.stanford.edu/~seanjw/module3/#21
- Widom's course:
  http://www.db-class.org/course/video/preview_list

Motivation  
oooo

**Get data**  
oooooooooooo●o

Extract/clean  
oooooo

# Record and parse it: XML



Looks like:

```
<bookstore>
  <book category="COOKING">
    <title lang="en">Everyday Italian</title>
    <author>Giada De Laurentiis</author>
    <year>2005</year>
    <price>30.00</price>
  </book>
  <book category="CHILDREN">
    <title lang="en">Harry Potter</title>
    <author>J K. Rowling</author>
    <year>2005</year>
    <price>29.99</price>
  </book>
  <book category="WEB">
    <title lang="en">Learning XML</title>
    <author>Erik T. Ray</author>
    <year>2003</year>
    <price>39.95</price>
  </book>
</bookstore>
```

Example from http://www.w3schools.com/xml/xml_tree.asp.

- Tagged elements
- Attributes
- Text
- Nested structure - not a table
- XML SCHEMA/DTD
- NO SQL (use XPATH/jQuery)

# Lab 1

Lab 1: Getting useful data with Curl and JSON in R.

Motivation
oooo

Get data
ooooooooooo

Extract/clean
●ooooo

# Clean things up: Regular Expressions



from: https://xkcd.com/208/

## Regex in action..

- To clean up data after scraping http://www.r-bloggers.com/scrape-web-data-using-r/.

- To extract useful information (state, latitude, longitude), when scraping a web page http://solomonmessing.wordpress.com/2011/09/18/map-of-participants/

- grep to create custom indicator variable features for analysis, agrep for approximate version of textttgrep. More on this later.

- For nice reference materials, look to http://www.regular-expressions.info/reference.html.

Motivation
oooo

Get data
ooooooooooo

Extract/clean
oooooo

## RegEx bare essentials

| RegEx | Description | Example |
|-------|-------------|---------|
| bil | Match string 'bil' | this is a 2 dollar <u>bil</u>l. |
| dollar\|bill | Match dollar or bill | this is a 2 <u>dollar</u> <u>bill</u>. |
| \d | Match any digit | this is a <u>2</u> dollar bill. |
| \w | Match any word (single letter) | <u>this</u> <u>is</u> <u>a</u> 2 <u>dollar</u> <u>bill</u>. |
| \w+ | Match at least 2 letters | <u>this</u> <u>is</u> a 2 <u>dollar</u> <u>bill</u>. |
| \s | Match any whitespace | this_is_a_2_dollar_bill. |
| \S | Match NOT whitespace | <u>this</u> <u>is</u> <u>a</u> <u>2</u> <u>dollar</u> <u>bill</u>. |
| colo?ur | Optionally match character preceding '?' | Yanks <u>color</u>, Brits <u>colour</u>. |
| col.*r | match any character between l and r 0 or more times | Yanks <u>color</u>, Brits <u>colour</u>. |

## Regex in R for Cleaning and Feature Extraction

| Command | What it does |
|---|---|
| `grep("dollar\\|bill", moneyStuff)` | return index of `moneyStuff` with item. |
| `gregexpr("\\d", "4a53f45e")` | return index of string in each match (why might this be a bad idea?) |
| `str_extract(moneyStuff, "\\d+")` | extracts groups of digits |
| `str_replace(moneyStuff, "(\\d+)", "\$ \\1")` | inserts dollar sign in front of numbers |
| `agrep("dollar\\|bill", moneyStuff, max.distance = .2)` | return index of anything with edit distance $\leq$ .2 from dollar or bill in `moneyStuff` with item. VERY useful to handle misspellings, plagiarism, etc. |

Read up on edit distance here: http://www.stanford.edu/class/cs124/lec/med.pdf

Motivation
oooo

Get data
ooooooooooo

Extract/clean
ooooo●o

## Regex == AWK

Need to quickly extract/transform a 2 TB text file?

- `awk '/pattern/'` return every line in a file matching pattern
- `cat bigfile.csv | awk '(2 > 5&&3<2) {print 1,3}'`
  `> smallerfile.csv`
  Reads in all lines from bigfile.csv where 2nd column value $> 5$
  and third column $< 2$ and prints to smallerfile.csv

Motivation
oooo

Get data
ooooooooooo

Extract/clean
oooooo●

## Lab 2

Lab 2: Scraping and Regex