

<http://v.gd/webscrapping>

[Video](#)

# Web Scrapping

## Gathering Data from Websites

John Little

[Data Visualization Services](#)

Duke University Libraries



# DVS: Data & Visualization

Contact Us: [askdata@duke.edu](mailto:askdata@duke.edu)

Location: Bostock Library, the /Edge (1st floor)

Hours: Walk-in, M-F / 8-5

OR, **Anytime** the Library is open (with ID)

Walk-in, Attendant Hours:

<http://library.duke.edu/data/about/schedule.html>

Workshops:

<http://library.duke.edu/data/news>

Services:

- Data Sources
- Data Management
- Data Cleaning
- Data Analysis
- Mapping & GIS
- Data Visualization



# Web Scraping - What We'll Cover

1. **Build** a data **corpus** of congressional press releases
2. **APIs** and **gather** latitude and longitude -- using **JSON** formatted **data**
3. A brief hands-on introduction into **HTML parsing**
4. **APIs** and **Documentation** (FTP) -- OpenSecrets.org
5. Discussion of **APIs** and **Social Media** data gathering
6. A brief discussion on the **ethics** of scraping

# This is **not** a programming workshop, but...

1. We will discuss **Python** and BeautifulSoup
2. We will **not learn** or use **Python** in the workshop
3. However, **some automation** tools are used in this workshop
4. Web Scraping is about deconstructing websites. Effective scraping requires learning about technical infrastructure as well as subject content
5. Not a workshop on [Text Analysis](#) (tools that calculate or correlate your data)
6. Not a workshop on **data cleaning**

# Technical Definitions

## Deconstruction v Construction

# Definitions

- **Scraping**

*Using tools to gather data you can see on a webpage*

A wide range of web scraping techniques and tools exist. These can be as simple as copy/paste and increase in complexity to automation tools, HTML parsing, APIs and programming



Scraping propolis from the sides of the bee box

Image by [Abalg](#)-commonswiki

# Definitions

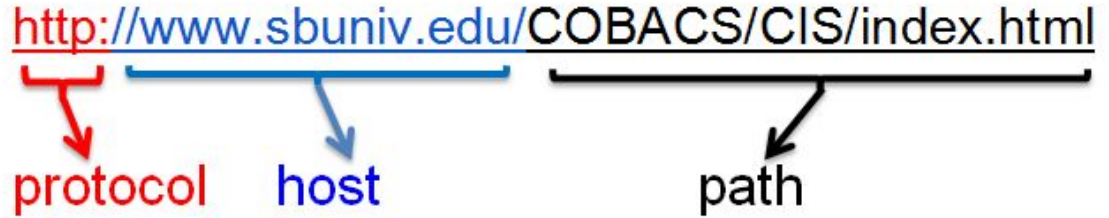
- Scraping
- **HTTP**

*HyperText Transfer Protocol*

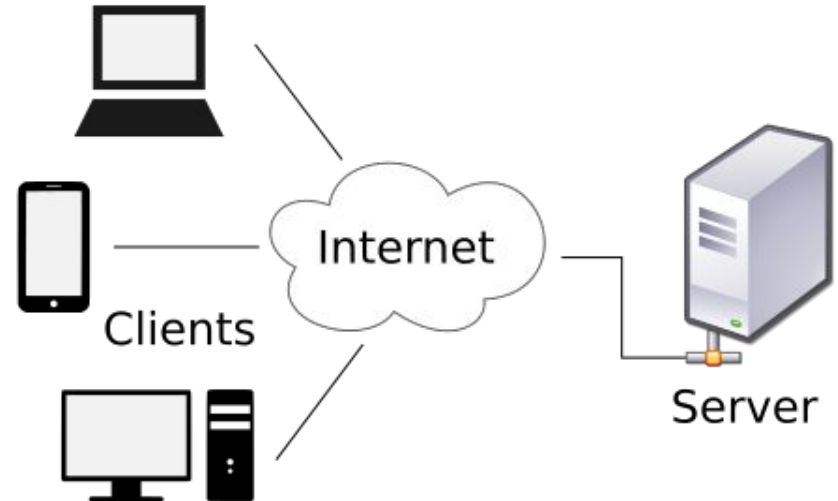
Machine interchange information transported over the Internet to enable multi-media data exchange, aka WWW. The protocol defines aspects of authentication, requests, status codes, persistent connections, client/server request/response. etc.

Access a server on port 80; the declarative Document Type Definition ( HTML, XML, JSON, etc.)

<http://www.sbuniv.edu/COBACS/CIS/index.html>



protocol      host      path



# Definitions

- Scraping
- HTTP
- **HTML**

*HyperText Markup Language*

The standard markup language on the Web

As the web evolves so does the proliferation of technical wrappers surrounding the visible content of websites (text and data)

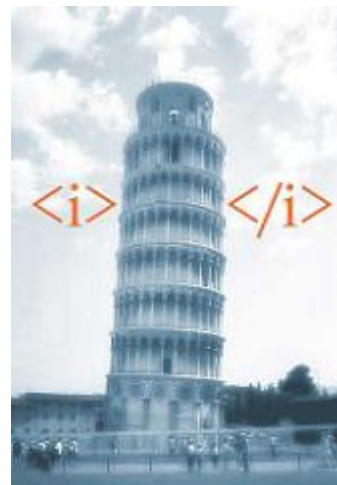


Image by [Jesper Rønn-Jensen](#)

```
<html>
  <head>
    <title>N&acute;</title>
  </head>
  <body>
    Text stránky
  </body>
</html>
```

Image by [Michaelbrabec](#)

# Definitions

- Scraping
- HTTP
- HTML
- **Parsing**

*The act of analyzing the strings and symbols to reveal only the data you need*

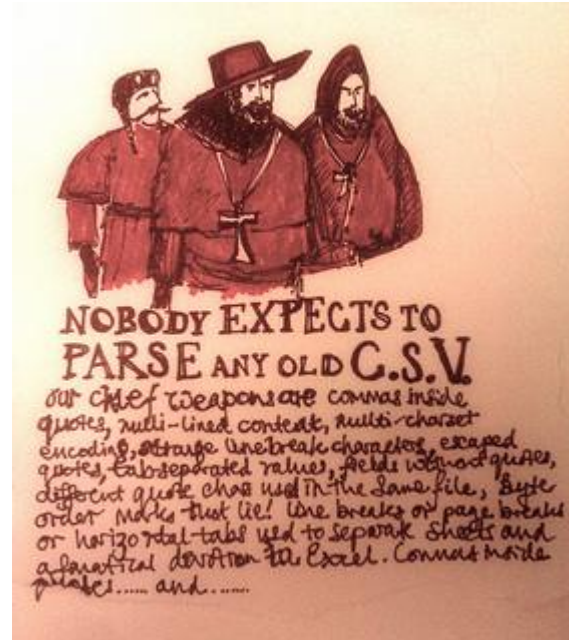


Image by [Paul Downey](#)

# Definitions

- Scraping
- HTTP
- HTML
- Parsing
- **Crawling**

*Moving across or through a website in an attempt to gather data from more than one URL or page*



Image by [Dave Gingrich](#)

# Definitions

- Scraping
- HTTP
- HTML
- Parsing
- Crawling
- **JSON**

*Javascript Open Notation*

*Readable text used to transmit data  
objects consisting of attribute-value pairs*

-- [Wikipedia](#)

```
{  
  "firstName": "John",  
  "lastName": "Smith",  
  "isAlive": true,  
  "age": 25,  
  "address": {  
    "streetAddress": "21 2nd Street",  
    "city": "New York",  
    "state": "NY",  
    "postalCode": "10021-3100"  
  },  
  "phoneNumbers": [  
    {  
      "type": "home",  
      "number": "212 555-1234"  
    },  
    {  
      "type": "office",  
      "number": "646 555-4567"  
    }  
  ],  
  "children": [],  
  "spouse": null  
}
```

# Definitions

- Scraping
- HTTP
- HTML
- Parsing
- JSON
- Crawling
- **API**

*Application Programming Interface*

A set of rules and protocols used to build a software application. In the context of Web Scraping an API is a method used to gather clean data from a website (i.e. data that is not wrapped in HTML, Javascript, bound in HTTP, etc.)

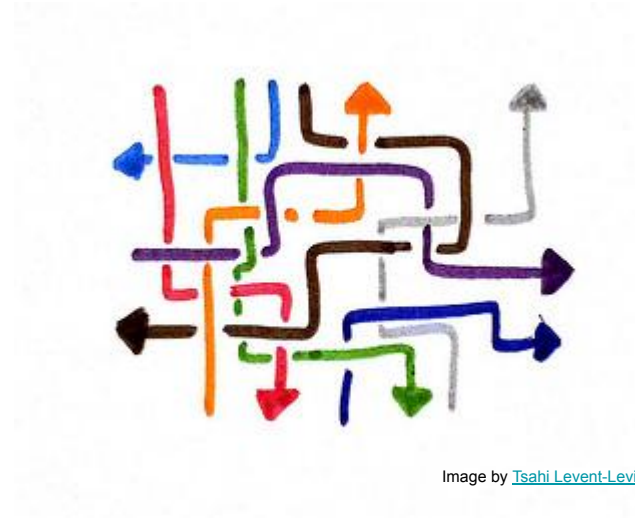


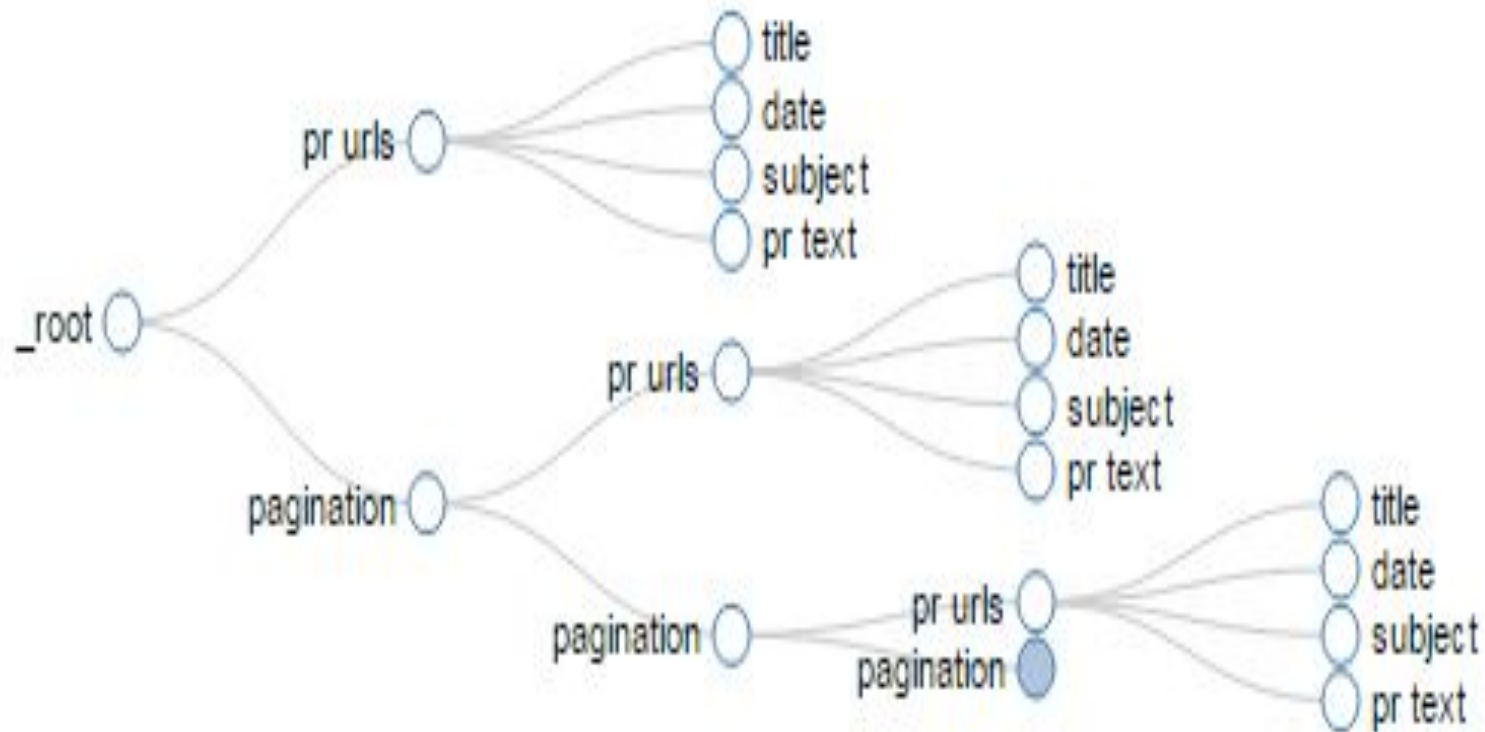
Image by [Tsahi Levent-Levi](#)

# Webscraper.io

## Demonstration & Hands-on

# Demo: Scraping Congressional Press Releases

- [Representative Nancy Pelosi's Press Releases](#)
  - **CONTENT**
    - Structure of the Press Release subsection of the site
      - Pagination
      - Links to each release
      - Common elements of release
  - **TOOLS**
    - Webscraper.io tool works inside of Chrome
      - Tutorials
      - Documentation
      - Community
      - Free or, alternatively, Fee for Service



# Now You Try It

1. <http://v.gd/webcraping1111>
2. Follow the download & installation instructions: [step 9a](#)
3. Find some congressional press release sites: [step 9b](#)
4. Follow the instructions: [steps 1-8](#)

# Google Sheets

## Demonstration

ImportHTML  
ImportXML

# Google Sheets -- ImportHTML

- Example Site: <http://www.boxofficemojo.com/>
- [IMPORTHTML](#)(url, query, index)
- In Practice:

```
=IMPORTHTML ("http://www.boxofficemojo.com/", "table", 3)
```

```
http://www.boxofficemojo.com/movies/?page=intl&id=annie2014.htm
```

```
=IMPORTHTML (A6, "table", 6)
```

Help Docs: <https://support.google.com/docs/answer/3093339?hl=en>

# Google Sheets -- ImportXML

- Example Site: <http://www.nytimes.com/>
- [IMPORTXML](#)(url, xpath\_query)
- In Practice:

```
=IMPORTXML ("http://nytimes.com/", "//*/p[contains(@class, 'summary')]")
```

Help Docs: <https://support.google.com/docs/answer/3093342?hl=en>

# Resources

## YouTube Videos

- [Web Scraping with Google Sheets](#)
- [Importing Data 4 ImportXML](#)
- [Web scraping using Google Docs - Xpath](#)

## Other Resources

- [CSS Tutorial](#)
- [XPath](#)
- [XPath Language](#) defined by W3C

## Your Turn:

- <https://github.com/data-and-visualization/Rfun>

# APIs, Parsing, & JSON

# OpenRefine & JSON file format

- Demonstration (<http://v.gd/parsing3333>)
  - A step-by-step guide using OpenRefine to gather JSON data via Google Map's API; then parse the JSON for latitude & longitude

sample address data

Sequence of color and letters

- Forging Specialties
- Green Goddess
- Chang, Carolyn Esq
- Century 21 Krall Real Estate
- Thompson Steel Company Inc
- Brooks, Morris J Jr
- Dimmock, Thomas J Esq
- Shimotani, Grace T

Base map

Map showing North Carolina and surrounding areas, with several locations marked by colored pins (A through H) corresponding to the list on the left. The map includes major cities like Asheville, Charlotte, Raleigh, and Winston-Salem, and major highways like I-85, I-95, and I-40.

# Parsing HTML

# OpenRefine & ParseHtml

- **BeautifulSoup Libraries**
  - Refine uses the Jython Libraries and has Jsoup
  - Jsoup is a Java library built on BeautifulSoup -- a tool for HTML Extraction
- **[Resources](#) (OpenRefine)**
  - Step-by-step example documented in the demonstration above
- **Documentation**
  - [Refine's documentation](#) on HTML Parsing
  - [Jsoup Documentation](#)
- Now You Try it -- <http://v.gd/parsing2222>

# Case Study: OpenSecrets and documentation

# OpenSecrets API and FTP

- OpenSecrets tracks the effects of money and lobbying in elections and politics
- OpenSecrets has an API
- OpenSecrets [API Documentation](#)
- OpenSecrets [Bulk Data downloader](#)
  - Login
  - Lobby.zip

# Social Media

# Social Media

1. Many ways to gather social media data
  - a. IFTTT where you compose rules to connect sites and can deposit data in a spreadsheet
  - b. APIs - often requires registered keys
  - c. Buy your data from a service such as GNIP
2. After you download it you may want to perform analysis
  - a. Sentiment Analysis, Word Frequency, Correlation, etc.
  - b. [Text Analysis tools \(from Digital Humanities LibGuide\)](#)
  - c. [Digital Studio's program on working with Texts: Comparing and choosing texts analysis tools](#)


# TAGS: a tool for collecting Twitter streams

- [TAGS](https://tags.hawksey.info/) (“New Sheets”; Version 6.0ns) - <https://tags.hawksey.info/>
  - Form driven (not command line)
  - Minimal setup
  - Data are collected in Google Sheets
  - Gather twitter stream data by type
    - screen-name stream data
    - screen-name status updates
    - twitter user favorited tweets
    - Search term for last 7 days: hashtag stream, username, boolean logic
    - Limit by date
    - Schedule to run hourly - set your interval, or run once.
    - 3 minute setup-video; easy to use - <https://youtu.be/Vm0kjAvH5HM>
    - Outputs: raw CSV structured data, plus default social graph visualizations

... make a copy and template while logged into your Google

TAGS menu click this button --> **Enable custom menu**

er run TAGS > Setup Twitter Access do so now (this should only

  <- you can use search oper  
from:BarackObama' (withou

---

off collection with TAGS > Run now! or set a trigger to collect eve  
n Tools -> Script Editor then Triggers -> Current script's triggers.

**Settings:**

filter	default	<- if search term is being must have to be included
tweets	0	
tweets	3000	<- maximum varies based
	search/tweets	<- use a search term in st

---

tweets	0
	0
	30/12/1899 00:00:00
	30/12/1899 00:00:00

**Active**

hive into an interactive online resource using TAGSExplorer

Republican Debate 1pm 10/27/2016											
File Edit View Insert Format Data Tools Add-ons Help TAGS All changes saved in Drive											
Comments <span>Share</span>											
RT @GodSaveAmerica1: This morning: @CNN, @FoxNews discussing #GOPDebate ; CNBC discussing Chinese population growth. Nuff said. #Republica...											
	A	B	C	D	E	F	G	H	I	J	
1	id_str	from_user	text	created_at	time	geo_coordinates	user_lang	in_reply_to_user_id	in_reply_to_screen_n	from_user_id_str	in_reply
12802	656219	AmericanMoocher	Win a trip to the #RepublicanDebate! Cheaper than a circus ticket. <a href="https://t.co/LRKfZNoR08">https://t.co/LRKfZNoR08</a>	Mon Oct 19 21:28:3	19/10/2015 22:28:37		en			292123314	
12803	656212	jockeyje70	RT @DCadvocacy: #Republicandebate #Conservatives NEXT DEBATE - What questions do you want asked? <a href="https://t.co/oNgw04xas">https://t.co/oNgw04xas</a>	Mon Oct 19 20:58:3	19/10/2015 21:58:32		en			2381169179	
12804	656209	CUBoulder	RT @CUArtsSciences: "Issues in the Debates: Debate on the Issues" on 10/20 with #cuboulder experts - preview of #republicandebate <a href="http://t...">http://t...</a>	Mon Oct 19 20:44:0	19/10/2015 21:44:09		en			741150194	
12805	656204	EScottAdler	RT @CUArtsSciences: "Issues in the Debates: Debate on the Issues" on 10/20 with #cuboulder experts - preview of #republicandebate <a href="http://t...">http://t...</a>	Mon Oct 19 20:25:0	19/10/2015 21:25:09		en			1140438458	
12806	656202	junilio86	#republicandebate <a href="https://t.co/MiLokWkR4">https://t.co/MiLokWkR4</a>	Mon Oct 19 20:18:4	19/10/2015 21:18:41		es			2567794593	
12807	656200	PracticallyGOP	#CNBC cedes to @realDonaldTrump debate demands <a href="https://t.co/vMygQsmAH">https://t.co/vMygQsmAH</a> #trump #republicans #gop #election2016 #republicandebate #gop	Mon Oct 19 20:11:1	19/10/2015 21:11:13		en			279856723	
12808	656195	DCadvocacy	#Republicandebate #Conservatives NEXT DEBATE - What questions do you want asked? <a href="https://t.co/oNgw04xas">https://t.co/oNgw04xas</a>	Mon Oct 19 19:47:4	19/10/2015 20:47:48		en			2588074814	
12809	656184	CWABoulder	RT @CUArtsSciences: "Issues in the Debates: Debate on the Issues" on 10/20 with #cuboulder experts - preview of #republicandebate <a href="http://t...">http://t...</a>	Mon Oct 19 19:04:4	19/10/2015 20:04:47		en			83498007	
12810	656174	CUBoulderCareer	RT @CUArtsSciences: "Issues in the Debates: Debate on the Issues" on 10/20 with #cuboulder experts - preview of #republicandebate <a href="http://t...">http://t...</a>	Mon Oct 19 18:27:1	19/10/2015 19:27:16		en			51196881	
12811	656172	mccuddenm	New #BeforeTheyWereFamous is up on #BenCarson Video: <a href="https://t.co/vlOrlQDqGE">https://t.co/vlOrlQDqGE</a> #republicandebate <a href="https://t.co/S4Ryh0ihrH">https://t.co/S4Ryh0ihrH</a>	Mon Oct 19 18:18:3	19/10/2015 19:18:31		en			90972141	
12812	656168	CUBoulderAlumni	RT @CUArtsSciences: "Issues in the Debates: Debate on the Issues" on 10/20 with #cuboulder experts - preview of #republicandebate <a href="http://t...">http://t...</a>	Mon Oct 19 18:03:3	19/10/2015 19:03:35		en			25112405	
12813	656168	baenanapancakes	#election2016 #DemDebate #RepublicanDebate2015 #RepublicanDebate #DemocraticDebate <a href="http://t.co/kMViv48X2g">http://t.co/kMViv48X2g</a>	Mon Oct 19 18:02:3	19/10/2015 19:02:36		en			3296136756	
12814	656167	ruthnasrullah	Interesting: Who would have guessed that the voice frequency of #whining is 180hz? #Trump2016 #republicandebate <a href="https://t.co/E81JuzutKz">https://t.co/E81JuzutKz</a>	Mon Oct 19 17:57:4	19/10/2015 18:57:40		en			184109643	
12815	656167	cueducation	Get in the campaign spirit before the #republicandebate <a href="https://t.co/dJEHVz9m">https://t.co/dJEHVz9m</a>	Mon Oct 19 17:56:5	19/10/2015 18:56:57		en			1613099138	
12816	656156	MaryPatriotNews	#CNBC next #RepublicanDebate 10/28 PREDICTION: #MainstreamMedia will downplay their PR work 4 #ObamaAdmin #HILLARY #Bernie @PRNewswire	Mon Oct 19 17:18:2	19/10/2015 18:18:24		en			3448846420	
12817	656150	magicbeagle	RT @OpposingViews: CNBC Changes Debate Format To Accommodate #Trump And #Carson's Demands <a href="https://t.co/ZcAnVZXU5c">https://t.co/ZcAnVZXU5c</a> #RepublicanDebate #news #...	Mon Oct 19 16:52:3	19/10/2015 17:52:33		en			1113877119	
12818	656150	magicbeagle	RT @OpposingViews: CNBC Changes Debate Format To Accommodate #Trump And #Carson's Demands <a href="https://t.co/ZcAnVZXU5c">https://t.co/ZcAnVZXU5c</a> #RepublicanDebate #news #...	Mon Oct 19 16:52:3	19/10/2015 17:52:33		en			1113877119	
12819											

Add 1000 more rows at bottom.

Readme/Settings Archive

File -> make a copy and template while logged into your Google

TAGS menu click this button --> **Enable custom menu**

er run TAGS > Setup Twitter Access do so now (this should only

<- you can use search operator from:BarackObama' (without

off collection with TAGS > Run now! or set a trigger to collect every  
Tools -> Script Editor then Triggers -> Current scripts triggers..

tings:

default

t filter

0

<- if search term is being must have to be included

sets

3000

<- maximum value is based

search/tweets

<- use a search term in s

weets

0

;

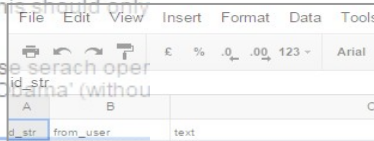
0

30/12/1899 00:00:00

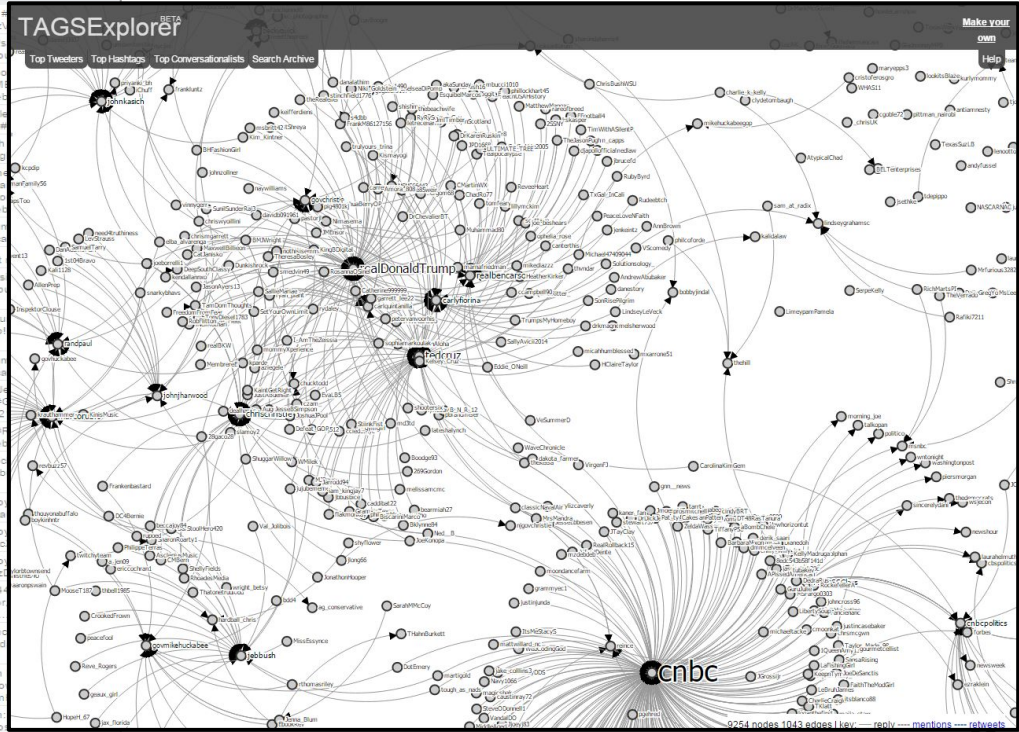
30/12/1899 00:00:00

tive

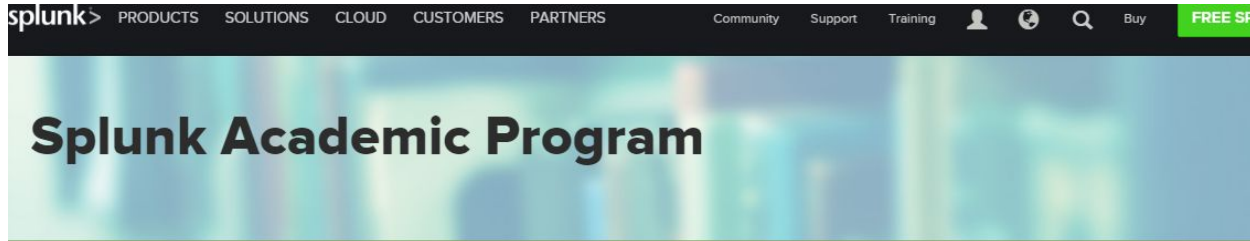
tive into an interactive online resource using TAGSExplorer



id	text
59751	My take on the https://t.co/Uj4z
59751	Obfuscation - it's https://t.co/mHo
59751	RT @kurbraun: PLEASE LET M #RepublicanDel
59750	JamesHowdenIllil Job Bush handi Johnny Cueto. #
59750	illieK Stoked to watch to miss it last nig
59750	Attention all Am Just because I suggest you c
59750	CalmeDowns RT @estheram way. #Republica
59750	ArvelMauldin The debate last Fuming that bas watch the #Rep
59750	StarDonovan #realDonaldTrump night Mr. Trump #Trump2016
59750	B_N_R_12 RT @estheram way. #Republica
59750	Political_Bill #JohnKasich #J #Republicans # https://t.co/LPc2
59749	Giliganista Last night the # #RepublicanDel
59749	RevMattPurkey @GovMikeHud @goldencorall b trots.
59749	cristinapalumbo Ted Cruz destro https://t.co/k4Za
59749	Jacque_Kubin Ted Cruz destro https://t.co/Ku3e
59749	WTCPolitics1 Ted Cruz destro https://t.co/f192at
59749	CommDigNews RT @MWM444: helping the poor more tax cuts."
59749	daronich Biggest differenc laughs at candid response.
59749	l1nydanca The Republican https://t.co/dp to https://t.co/UXm
59749	uSeenBitter RT @AsBillJosh They Live so clo



# Splunk Academic Program



## Evaluate. Deploy. Develop.

Faculty and students at more than [100 universities around the world use Splunk's software platform](#) to understand how real-time Operational Intelligence drives innovation in the 21st century. Learning applied big data analytics, gaining practical cybersecurity skills and gathering research insights from sources as diverse and specialized as medical records or weather data are all examples of how these cutting-edge institutions leverage Splunk software.

With an active [community](#), free [documentation](#) and plenty of [training and education courses](#) available, Splunk offers resources to help you make the most of your Splunk license.

The following resources are available to assist your institution in its mission of

### Get Started

-  [Press Release  
Universities Worldwide Take  
Splunk to the Classroom](#)
-  [Resource  
Splunk Community Wiki](#)

[Free Download](#)

# Ethics

# Ethics & Technical Challenges of Web Scraping

Like copyright, there are rules, protocols, and best practices

1. Do you identify yourself?
2. Do you document your downloading steps?
3. Do you document your analysis steps?
4. Is your research reproducible?
5. Are your methods and processes confidential or proprietary?
6. How do you protect and document your sources?
7. How do you give credit to your sources?
8. Articles:
  - a. [To Scrape or Not to Scrape](#): Technical and Ethical Challenges of Collecting Data off the Web
  - b. On the [ethics of web scraping from a data journalism](#) perspective

# Resources

- Catalog of [Web Scraping tools](#)  
[https://docs.google.com/spreadsheets/d/1A\\_9wBEmc8VP6HUm2R-7jdPU9FWY8dSfNsfIjQ6cz638/edit#gid=0](https://docs.google.com/spreadsheets/d/1A_9wBEmc8VP6HUm2R-7jdPU9FWY8dSfNsfIjQ6cz638/edit#gid=0)

## Tools used or demonstrated in this Workshop

- Webscraper.io - <http://webscraper.io/> (script and crawl a website)
- OpenRefine - <http://openrefine.org> (scripting, cleaning, parsing, utility data tool)
  - Regular Expressions
  - GREL
  - Jsoup
- TAGS - <https://tags.hawksey.info/> (Twitter Stream collector)
- Splunk - [Academic Program](#) (Twitter Stream collector / Search and Discover )

## Ethics

- [To Scrape or Not to Scrape](#): Technical and Ethical Challenges of Collecting Data off the Web
- <http://gijn.org/2015/08/12/on-the-ethics-of-web-scraping-and-data-journalism/>

# Thank You!

Please complete feedback forms

# John Little

<http://v.gd/webscraping>

[Data & Visualization Services](#)  
[Duke University Libraries](#)

This [presentation](#) contains links to directions and sample data. Each section of the workshop can be completed at your own pace.

**Data, presentation, and handouts are shareable under CC BY-NC license:**

<https://creativecommons.org/licenses/by-nc/4.0/>

