# Websites, Twitter & Facebook - oh my!
## How to start gathering data from the web

Michael Beckstrand, Ph.D.
*Mixed Methods Research Associate*

Alicia Hofelich Mohr, Ph.D.
*Research Data Manager*

# Lots of data "out there" to harvest

How can you access it?

How can efficiently download and handle data once accessed?

How can you get downloaded data into a usable format?

# Lots of tools "out there" to help you

What to ask yourself to choose:

1. What do you want to do?

2. How much time/effort are you willing to put into learning a tool?

3. How extensible do you need the tool to be?

# What do you want to do?

Three major categories:

- Save what you see on a page

- Capture what you see on the screen in a structured format

- Extract specific sets of information from a website

# How easy will it be to learn?

Continuum from click-based to fully coded

# How extensible?

Can you rely on the listed, built-in feature set or do you need customizability to very specific scenarios

| Tool | | Use | Ease to Learn | Extensibility |
|---|---|---|---|---|
|  | Get Them All | Save what you See | Easy | Limited |
|  | HTTrack | | Easy | Limited |
|  | If this then that | | Medium | Some |
|  | Scraper | Structure what you see | Medium | Some |
|  | Twitter Archiving Google Sheet | Query APIs | Easy | Limited |
|  | FacePager | | Medium | Some |
|  | Python | Structure what you see & query APIs | Difficult | High |
|  | R | | Difficult | High |

# Saving what you see

Download a bunch of stuff - 'batch' the operations to:

- Get all files linked on a page

- Mirror whole (or sections) of sites

- Create a workflow of tools to store/organize articles

# Example Tools for Batching the Web



**GetThemAll**

Chrome extension for downloading multiple links on a single page all at once ('Chrono Download Manager' also an excellent choice)



**HTTrack**

Stand-alone crawler that can mirror a whole (or part of) website onto your local computer



**IFTTT**

'Recipes' for automating actions between web services, like cloud storage, social media, article managers/RSS feeds

File Types

# Turning what you see (and don't see) into structured data

'Web scraping' relies on underlying HTML

This hierarchical/marked-up data can be converted to spreadsheets and tables

*Big caveat*: 'scraping' only grabs what is on screen or has been loaded by the browser

# Example Tools for Scraping



**Scraper**
Chrome extension with point & click interface & regular expression-style queries



**Python (among other coding solutions)**
Programming language with libraries for extracting data from the web (Beautiful Soup)



**R**
Packages for scraping and harvesting webpages including RSelenium, Rvest, scrapeit

# For more advanced selection

Select target area, right click, choose 'Inspect'



Area of interest is Highlighted

Inspect to bring up Developer Console

# Use the details to specify a path

# Forget what you see - ask the computer to get your data

Using an Application Programming Interface (API)

- Structure that defines how programs can communicate with each other
  - Sharing information
  - Taking in information

- Structured requests; structured responses

# Application Programming Interface (API)

# Defined Queries & Limitations

- Set actions with defined parameters
  - Twitter: https://dev.twitter.com/rest/public
  - Facebook: https://developers.facebook.com/docs/graph-api/
- Set Limitations
  - Privacy: can't get everything you may be able to see
  - Rate limitations: can't make unlimited calls
  - Return limitations: can't get all the data at once

# Example Tools for Querying APIs

**FacePager**
Point & click interface for querying Twitter, Facebook, and general APIs

**Twitter Archiving Google Spreadsheet**
Very easy-to-use script for Google sheets that accesses Twitter's streaming API to collect and archive tweets

**R**
Packages for accessing APIs for Twitter, Facebook, Elsevier, and general sites (TwitteR, Rfacebook, rscoipus, httr)

# TAGS

**Face Pager**

Nodes View

Data View

Query Set Up

Status Log

Column Set Up

Open Database | New Database | Export Data | Add Nodes | Delete Nodes | Presets | Help

Expand nodes | Collapse nodes | Copy Node(s) to Clipboard

Add Column | Unpack List | Copy JSON to Clipboard

Object ID | Object Type | Query Status | Query Time | Query Type

Key | Value

Facebook | Twitter | Generic | Files | Twitter Streaming

Resource: <checkin>/comments

Parameters: <checkin> | <Object ID>

Maximum pages: 100

Access Token: •••••••••••••••••••••••••••••• | Login to Facebook

For all selected nodes and their children of level | 1 | Fetch Data

Custom Table Columns (one key per line)

name
screen_name
type
metadata.type
location
ids

Apply Column Setup

/Users/ahofelich/Desktop/Example_Data/MacalesterWorkshop/blank.db

Timer stopped | 0 node(s) selected

| Tool | | Use | Ease to Learn | Extensibility |
|---|---|---|---|---|
|  | Get Them All | Save what you See | Easy | Limited |
|  | HTTrack | | Easy | Limited |
|  | If this then that | | Medium | Some |
|  | Scraper | Structure what you see | Medium | Some |
|  | Twitter Archiving Google Sheet | Query APIs | Easy | Limited |
|  | FacePager | | Medium | Some |
|  | Python | Structure what you see & query APIs | Difficult | High |
|  | R | | Difficult | High |

# Thank you! Questions?

Michael Beckstrand

[mjbeckst@umn.edu](mailto:mjbeckst@umn.edu)


Alicia Hofelich Mohr

[hofelich@umn.edu](mailto:hofelich@umn.edu)