

# Digital Trace Data

## Part 2

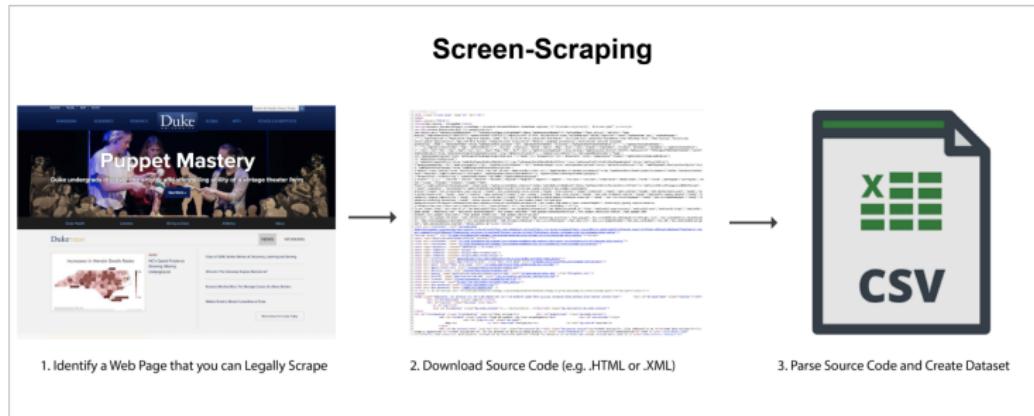
Ridhi Kashyap

University of Oxford

SICSS-Oxford  
June 18, 2019

# Web Scraping

- ▶ The process of automatically extracting data from web pages



Material on screen-scraping in these slides draw on  
[Chris Bail's tutorial on screen-scraping](#)

# Legal Issues in Web Scraping

- ▶ Website Terms of Service: are you allowed to do this?
- ▶ Websites tend to have a “robots.txt” policy that specifies rules about automated data collection on the site.
- ▶ Larger websites like Facebook, Instagram, NY Times do now allow these practices – but some tend to have Application Programming Interfaces (more later) to facilitate access to their public data.

# Robot.txt

```
# robots.txt for http://www.wikipedia.org/ and friends
#
# Please note: There are a lot of pages on this site, and there are
# some misbehaved spiders out there that go _way_ too fast. If you're
# irresponsible, your access to the site may be blocked.
#
# Observed spamming large amounts of https://en.wikipedia.org/?curid=NNNNNNN
# and ignoring 429 ratelimit responses, claims to respect robots:
# http://mj12bot.com/
User-agent: MJ12bot
Disallow: /

# advertising-related bots:
User-agent: Mediapartners-Google*
Disallow: /

# Wikipedia work bots:
User-agent: IsraBot
Disallow:
```

# Wikipedia Page

The Wikipedia logo, featuring a globe with various symbols from different scripts (Greek, Chinese, etc.) overlaid.

**WIKIPEDIA**  
The Free Encyclopedia

[Create account](#) [Log in](#)

[Article](#)

[Talk](#)

[Read](#)

[Edit](#)

[View history](#)

[Search](#)

A magnifying glass icon used for the search function.

## World Health Organization ranking of health systems in 2000

From Wikipedia, the free encyclopedia

The [World Health Organization](#) (WHO) ranked the health systems of its 191 member states in its [World Health Report](#)<sup>[1]</sup> 2000. It provided a framework and measurement approach to examine and compare aspects of health systems around the world.<sup>[2]</sup> It developed a series of performance indicators to assess the overall level and distribution of health in the populations, and the responsiveness and financing of health care services. It was the organization's first ever analysis of the world's health systems.<sup>[3]</sup>

### Contents [hide]

- [1 Ranking](#)
- [2 Methodology](#)
- [3 Criticism](#)
- [4 See also](#)
- [5 References](#)

[Main page](#)

[Contents](#)

[Featured content](#)

[Current events](#)

[Random article](#)

[Donate to Wikipedia](#)

[Wikipedia store](#)

[Interaction](#)

[Help](#)

[About Wikipedia](#)

[Community portal](#)

[Recent changes](#)

[Contact page](#)

[Tools](#)

[What links here](#)

# Wikipedia Page HTML

# Wikipedia Page

```
#install package
install.packages("rvest")
library(rvest)

wikipedia_page<-read_html("https://en.wikipedia.org/wiki/World_Health_Organization_ranking_of_health_systems_in_2000")
```

# Wikipedia Page

```
{xml_document}
<html class="client-nojs" lang="en" dir="ltr">
[1] <head>\n<meta http-equiv="Content-Type" content="text/html; charset=UTF-8">\n<meta charset="UTF-8">\n<title>World Health O ...
[2] <body class="mediawiki ltr sitedir-ltr mw-hide-empty-elt ns-0 ns-subject mw-editable page-World_Health_Organization_rankin ...
```

# Wikipedia Page

The screenshot shows a Wikipedia article titled "Organization ranking of health systems in 2000" from the English Wikipedia. The browser is Google Chrome. The developer tools are open, specifically the "Developer Tools" panel, which includes the "JavaScript Console". The main content area displays the article's text, which discusses the World Health Organization's ranking of health systems across 191 member states. A table at the bottom lists the top four countries based on the ranking.

Always Show Bookmarks Bar ⌘B  
Always Show Toolbar in Full Screen ⌘F

Stop ⌘R  
Force Reload This Page ⌘ShiftR

Enter Full Screen ⌘F  
Actual Size ⌘O  
Zoom In ⌘+  
Zoom Out ⌘-

Cast...

**Developer** ►

The Free Encyclopedia

Main page  
Contents  
Featured content  
Current events  
Random article  
Donate to Wikipedia  
Wikipedia store

Interaction  
Help  
About Wikipedia  
Community portal  
Recent changes  
Contact page

Tools  
What links here  
Related changes  
Upload file  
Special pages  
Permanent link  
Page information  
Wikidata item  
Cite this page  
Print/export

View Source ⌘U  
Developer Tools ⌘I  
JavaScript Console ⌘J

## Organization ranking of health systems in 2000

The **World Health Organization (WHO)** ranked the health systems of its 191 member states in its aspects of health systems around the world.<sup>[2]</sup> It developed a series of performance indicators to assess health care services. It was the organization's first ever analysis of the world's health systems.<sup>[3]</sup>

**Contents** [hide]

- 1 Ranking
- 2 Methodology
- 3 Criticism
- 4 See also
- 5 References
- 6 External links

### Ranking [edit]

Country	Attainment of goals / Health / Level (DALE)	Attainment of goals / Health / Distribution	Attainment of goals / Overall goal at
Afghanistan	168	182	183
Albania	102	129	86
Algeria	84	110	99
Andorra	10	25	17

## Inspect

states in its **World Health Report**<sup>[1]</sup> aspects of health systems around the distribution of health in the organization's first ever analysis of the

Performance / On level of health	Performance / Overall health system performance
150	173
64	55
45	81
7	4

# Copy Xpath

The screenshot shows a browser developer tools element inspector. A context menu is open over an element with the class 'td'. The menu has two main sections: a primary section with options like 'Add attribute', 'Edit as HTML', and 'Delete element', and a secondary section with options like 'Cut element', 'Copy element', and 'Paste element'. Both sections have a 'Copy' option highlighted with a blue selection bar. The primary section also includes 'Copy outerHTML', 'Copy selector', and 'Copy XPath'.

Copy

Copy element

Paste element

Copy outerHTML

Copy selector

Copy XPath

Add attribute

Edit as HTML

Delete element

Hide element

Force state

Break on

Expand recursively

Collapse children

Scroll into view

Focus

```
> <n>...</n>
  <table class="wikitable sortable jquery-tablesorter">
    <thead>...</thead>
    <tbody>
      <tr>
        <td>...</td>
        <td>168</td>
        <td>182</td>
        <td>183</td>
        <td>184</td>
        <td>150</td>
        <td>173</td>
      </tr>
      <tr>...</tr>
      <tr>...</tr>
```

html body #content #bodyContent #mw-content-text div table tbody tr td

Styles Event Listeners DOM Breakpoints Properties Accessibility

# HTML node

```
section_of_wikipedia <- html_node(wikipedia_page, xpath = '//*[@id="mw-content-text"]/div/table[2"]')
head(section_of_wikipedia)
```

```
$node
<pointer: 0x7feb444a56d0>
```

```
$doc
<pointer: 0x7feb4449e0a0>
```

# Extracting HTML Table

```
> health_rankings<-html_table(section_of_wikipedia)
> head(health_rankings[,c(1:2)])
      Country Attainment of goals / Health / Level (DALE)
1    Afghanistan                               168
2     Albania                                    102
3     Algeria                                    84
4     Andorra                                    10
5     Angola                                     165
6 Antigua and Barbuda                           48
```

## CSS Selector Gadget

- ▶ For more complex webpages, particularly those with CSS elements, other approaches might be needed.
- ▶ <http://selectorgadget.com/>
- ▶ Select css elements with Selector.

# Selector Gadget

The screenshot shows the Duke University homepage. At the top, there is a navigation bar with links for ADMISSIONS, ACADEMICS, RESEARCH, Duke UNIVERSITY (with a magnifying glass icon), GLOBAL, ARTS, and SCHOOLS & INSTITUTES. Below the navigation is a large banner image of three students performing puppetry. Overlaid on the banner is the title "Puppet Mastery" in large white letters, followed by a subtitle "Duke undergrads discover the artistry and storytelling ability of a vintage theater form". A yellow "See More" button is visible in the center of the banner. Below the banner is a horizontal navigation bar with links for Duke Health, Libraries, Giving to Duke, Athletics, and About. The main content area features a "Duke TODAY" header with categories NEWS and WORKING. On the left, there is a map titled "Increases in Heroin Death Rates" showing North Carolina county-level data from 1990 to 2014. To the right of the map is a news article titled "NC's Opioid Problems Growing, Moving Underground". Further down the page are other news items: "Class of 2018: Senior Stories of Discovery, Learning and Serving" (highlighted with a red border), "Who Are The Honorary Degree Recipients?", "Eduardo Bonilla-Silva: The Strange Career of a Race Scholar", and "Melton Grant to Boost Humanities at Duke". At the bottom, there is a "What's up @Duke" section and a search bar with the placeholder "Search for People, Places, Things". The search bar also includes buttons for "Clear (259)", "Toggle Position", "?", and "X".

# Selector Gadget

```
duke_page<-read_html("https://www.duke.edu")
duke_events<-html_nodes(duke_page, css="li:nth-child(1) .epsilon")
html_text(duke_events)
```

## Some Limits

- ▶ In practice, this can often yield messy data depending on web page complexity
- ▶ Web scraping can be frustrating, unfeasible or illegal. What then?
  - ▶ For complex web pages, crowd workers (Mechanical Turk) could be an option. Talk to Roberto at dinner.
  - ▶ Some websites make their content available through APIs.

## What is an API?

- ▶ Application programming interfaces or APIs are a software intermediary that allows two applications to talk to each other.
- ▶ Web APIs allow one computer (a client) to ask another computer (a server) for some resource over the internet.

## Application Programming Interface

- ▶ Modern APIs adhere to standards, that make data exchange programmatically accessible, safe and structured
- ▶ APIs generally require some kind of authentication and credentialling process that gives them more legitimacy.
- ▶ Contrast with web scraping.

# API Directory

The screenshot shows the ProgrammableWeb API Directory homepage. At the top, there's a navigation bar with links for "LEARN ABOUT APIs", "WHAT IS AN API?", "TUTORIALS", "API CHARTS & RESEARCH", "ADD API & MORE", and social media icons. A search bar says "Search over 21,780 APIs and much more". Above the search bar are links for "WRITE FOR US", "BECOME MEMBER", and "LOGIN". A banner for MuleSoft with the text "Applying DevOps to APIs" and a "Learn more" button is displayed.

## Search the Largest API Directory on the Web

Search Over 21,780 APIs SEARCH APIs

Filter APIs

By Category Include Deprecated APIs

API Name	Description	Category	Submitted
Google Maps	[This API is no longer available. Google Maps' services have been split into multiple APIs, including the Static Maps API.]	Mapping	12.05.2005
Twitter	[This API is no longer available. It has been split into multiple APIs, including the Twitter Ads API, Twitter Search Tweets.]	Social	12.08.2006
YouTube	The Data API allows users to integrate their program with YouTube and also have a summary of the operations available on the website. It provides the capability to search for videos, retrieve...	Video	02.08.2006
Flickr	The Flickr API can be used to retrieve photos from the Flickr photo sharing service using a variety of feeds - public photos and videos, favorites, friends, group pools, discussions, and more. The...	Photos	09.04.2005
Facebook	[This API is no longer available. Its functions have been split among the following APIs: Facebook Ads,...	Social	08.16.2006

**API UNIVERSITY**

FEATURED LATEST

**FOR API PROVIDERS**

What Are APIs and How Do They Work?  
8 Real World API Strategies and the Keys to Their Success  
Microservices 101: Understanding and Leveraging Microservices  
[More for API Providers >](#)

**FOR DEVELOPERS**

How to Get Started With Google Actions  
How to Build a Monitoring Application With the Google Cloud Vision API  
How to Access Any RESTful API Using the R Language  
[More for Developers >](#)

**Today in APIs**  
Latest news about the API economy and newest APIs, delivered daily:  
 SUBSCRIBE

Figure: <https://www.programmableweb.com/apis/directory>

## Extracting Data from APIs

- ▶ Web APIs tend to use HTTP methods.
- ▶ HTTP is the network protocol used to deliver virtually all files and other data on the World Wide Web, such as HTML files, image files, query results, etc
- ▶ These are verbs: GET(), POST(), DELETE(), and so on.
- ▶ GET() and POST()
  - ▶ GET() is used by your browser when requesting a page
  - ▶ POST() is (usually) used when submitting a form to a server.

## httr in R

- ▶ The `httr` package in R is useful to work with APIs
- ▶ You can make a request (to a url) and get a response.
- ▶ Response contains a status, header and body.

# httr in R

```
> library(httr)
>
> trial <- GET("https://http.cat/100")
> status_code(trial)
[1] 200
>
> trial_content <- content(trial, as = 'raw')
> head(trial_content)
[1] ff d8 ff e0 00 10
>
> write_file(trial_content, "100.jpg")
```

httr in R



100

Continue

# Facebook Marketing API

- ▶ You need:
  - ▶ Facebook account
  - ▶ Marketing app with token and an ad account number ("act")
  - ▶ See this helpful tutorial by Sofia Gil-Clavel from MPIDR on how to obtain these credentials –  
[https://github.com/SofiaG11/Using\\_Facebook\\_API/  
blob/master/First\\_Step.pdf](https://github.com/SofiaG11/Using_Facebook_API/blob/master/First_Step.pdf)

# Facebook Marketing API

[Create New Audience](#)   [Use Saved Audience](#)

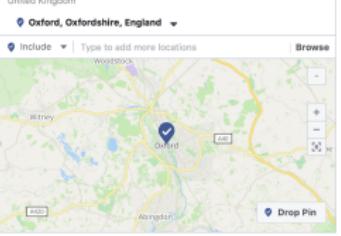
[Custom Audiences](#)  Add a previously created Custom or Lookalike Audience

[Exclude](#)   [Create New](#)

[Locations](#)  Everyone in this location

United Kingdom  
Oxford, Oxfordshire, England [Browse](#)

Include  Type to add more locations



Add Locations in Bulk

Age: 18 - 65+ [Edit](#)

Gender: All Men Women [Edit](#)

Languages:  Enter a language... [Edit](#)

Detailed Targeting [Edit](#) Include people who match

Interests > Additional Interests  
Cricket  
Cricket World Cup

Add demographics, interests or behaviors [Suggestions](#) [Browse](#)

[Exclude People](#) or [Narrow Audience](#)

Connections [Edit](#) Add a connection type

**Stories**

Because Facebook Stories is a new placement being released gradually, audience and reach estimates aren't currently available. These estimates are based on the other placements you've selected.

**Audience Size**

Your audience is defined.



Potential Reach: 38,000 people

**Estimated Daily Results**

Reach: 3.2K - 7.6K

The accuracy of estimates is based on factors like ad placement, budget, and market data. Numbers are provided to give you an idea of performance for your budget, but are only estimates and don't guarantee results.

Were these estimates helpful?

## Facebook Marketing API

- ▶ We can programmatically make a query to the API to retrieve these ad audience estimates.
- ▶ For more information on how targeting specifications work – see <https://developers.facebook.com/docs/marketing-api/buying-api/targeting/>
- ▶ To search for available targeting options – see <https://developers.facebook.com/docs/marketing-api/targeting-search>

# Facebook Marketing API

```
[credentials_2 <- paste0('https://graph.facebook.com/',version,'/act_',act,'/delivery_estimate?access_token=', token)
```

- ▶ version refers to the version of the API (v 3.3. is current)
- ▶ act is the ad account number
- ▶ token is the access token
- ▶ Important to remember rate limiting when working with APIs

**New campaign**

None of your ads are running - Your campaigns and ad groups are paused or removed. Enable them to begin showing your ads.

[LEARN MORE](#)

1 Create your campaign 2 Confirmation

**Edit targeted demographics**

Gender	Age	Parental status	Household income
<input checked="" type="checkbox"/> Female	<input checked="" type="checkbox"/> 18 - 24	<input checked="" type="checkbox"/> Not a parent	<input checked="" type="checkbox"/> Top 10%
<input checked="" type="checkbox"/> Male	<input checked="" type="checkbox"/> 25 - 34	<input checked="" type="checkbox"/> Parent	<input checked="" type="checkbox"/> 11 - 20%
<input checked="" type="checkbox"/> Unknown <small>(?)</small>	<input checked="" type="checkbox"/> 35 - 44	<input checked="" type="checkbox"/> Unknown <small>(?)</small>	<input checked="" type="checkbox"/> 21 - 30%
	<input checked="" type="checkbox"/> 45 - 54		<input checked="" type="checkbox"/> 31 - 40%
	<input checked="" type="checkbox"/> 55 - 64		<input checked="" type="checkbox"/> 41 - 50%
	<input checked="" type="checkbox"/> 65+		<input checked="" type="checkbox"/> Lower 50%
	<input checked="" type="checkbox"/> Unknown <small>(?)</small>		<input checked="" type="checkbox"/> Unknown <small>(?)</small>

**DONE**

**Your targeting's reach (?)**

Impressions  
**10B+**

What's defining your reach (?) ▼

**Your weekly estimates (?)**

Enter a bid and budget to see your estimated performance

⚠ Note: Household income targeting is only available in select countries. [Learn more](#)

## Exercise

- ▶ Familiarise yourself (and obtain credentials) to work with one of the following three data sources:
  1. Facebook marketing API
  2. Twitter – use Chris Bail's helpful tutorial using the `rtweet` package.
  3. ...Or another API of your choosing.

## Group Exercise

- ▶ Divide yourselves into groups of four by counting off in order around the room.
- ▶ Work together to identify a research question that you believe could be answered using one of the following three sources of digital trace data:
  1. Facebook marketing API
  2. Twitter
  3. Google Trends
- ▶ Think of a research design using digital trace data either on its own, or in combination with existing observational or survey data.
- ▶ Collect preliminary data to test the feasibility of your approach.
- ▶ What are the strengths and limitations of this approach?
- ▶ How could it be improved with a hybrid design?