

SCRAPING DATA FROM THE WEB I: FREE RESOURCES FOR STORYTELLERS

CHARLES BERRET

BROWN INSTITUTE FOR MEDIA INNOVATION

COLUMBIA UNIVERSITY

DATA STORYTELLING AT BOSTON UNIVERSITY

JUNE 6, 2017

TWO SESSIONS

**SCRAPING DATA FROM THE WEB I:
FREE RESOURCES FOR STORYTELLERS**

BREAK (10:30–10:45)

**SCRAPING DATA FROM THE WEB II:
APPLYING YOUR NEW SKILLS**

SCRAPING DATA FROM THE WEB I: FREE RESOURCES FOR STORYTELLERS

PART 1 – INTRODUCTION TO WEB SCRAPING (9-9:45)

- WHAT IS WEB SCRAPING?
- METHODS AND TOOLS FOR WEB SCRAPING
- HOW SCRAPED DATA CAN BE USED IN STORIES

PART 2 – GUIDED WALKTHROUGH (9:45-10:30)

- HOW TO SCRAPE THE WEB WITH GOOGLE SHEETS

WHAT YOU NEED FOR TODAY'S LESSON

- ▶ Modern Web Browser (Firefox, Chrome, Safari)
- ▶ Google Sheets (free with Google/GMail account)

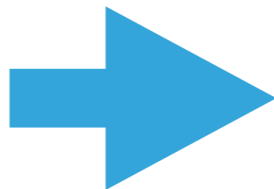
WHAT YOU WILL LEARN

- ▶ How to move information from a website to a spreadsheet
- ▶ How to make scraped data useful

WHAT IS WEB SCRAPING?

- ▶ Copying data from a website and storing it in another form so that it's useful to you.

```
<span title="ctx_ver=Z39.88-2004&amp;rft_id=info%3Aid%2Fen.wikipedia.org%3AJournalism&amp;rft.btitle=Gonzo+Journalism&amp;rft.genre=unknown&amp;rft_id=http%3A%2F%2Fwww.britannica.com%2FEBchecked%2Fto pic%2F1069436%2Fgonzo-journalism&amp;rft.pub=Encyclop%C3%A6dia+Britannica&amp;rft_val_fmt=info%3Aofi%2Ffmt%3Akev%3Amtx%3Abook" class="E3988"><span style="display:none;">#160;</span></span></li><li id="cite_note-15"><span class="mw-cite-backlink"><b><a href="#cite_ref-15"></a></b></span> <span class="reference-text"><cite class="citation journal">Robinson, Sue (2011). "<span style="padding-left:0.2em;"></span>Journalism as Process: The Organizational Implications of Participatory Online News." <i>Journalism &amp; Communication Monographs</i>. <b>13</b> (3): 137.</cite><span title="ctx_ver=Z39.88-2004&amp;rft_id=info%3Aid%2Fen.wikipedia.org%3AJournalism&amp;rft.atitle=%22Journalism+as+Process%22%3A-The+Organizational+Implications+of+Participatory+Online+News.&amp;rft.aufirst=Sue&amp;rft.aulast=Robinson&amp;rft.date=2011&amp;rft.genre=article&amp;rft.issue=3&amp;rft.jtitle=Journalism+%25-Communication+Monographs&amp;rft.pages=137&amp;rft_val_fmt=info%3Aofi%2Ffmt%3Akev%3Amtx%3AJournalism&amp;rft.volume=13" class="E3988"><span style="display:none;">#160;</span></span></li><li id="cite_note-16"><span class="mw-cite-backlink"><b><a href="#cite_ref-16"></a></b></span> <span class="reference-text"><cite class="citation web">"rst Journalism School". Columbia.: University of Missouri Press. p.&#160;1.</cite><span title="ctx_ver=Z39.88-2004&amp;rft_id=info%3Aid%2Fen.wikipedia.org%3AJournalism&amp;rft.btitle=rst+Journalism+School&amp;rft.genre=unknown&amp;rft.pages=1&amp;rft.place=Columbia.&amp;rft.pub=University+of+Missouri+Press&amp;rft_val_fmt=info%3Aofi%2Ffmt%3Akev%3Amtx%3Abook" class="E3988"><span style="display:none;">#160;</span></span> <span style="display:none;font-size:100%" class="error citation-comment">Missing or empty <code style="color:inherit;border:inherit;padding:inherit;">|url=</code> (<a href="/wiki/Help:CS1_errors#cite_web_url" title="Help:CS1 errors">help</a>)</span></li><li id="cite_note-17"><span class="mw-cite-backlink"><b><a href="#cite_ref-17"></a></b></span> <span class="reference-text"><cite class="citation journal">de Albuquerque, Afonso; Cagliardi, Juliana (2011). "THE COPY DESK AND THE DILEMMAS OF THE INSTITUTIONALIZATION OF MODERN JOURNALISM IN BRAZIL". <i>Journalism Studies</i>. <b>12</b> (1). <a href="/wiki/Digital_object_identifier" title="Digital object identifier">doi</a>:<a rel="nofollow" class="external text" href="//doi.org/10.1080%2F1461670X.2010.511956">10.1080/1461670X.2010.511956</a>.</cite><span title="ctx_ver=Z39.88-2004&amp;rft_id=info%3Aid%2Fen.wikipedia.org%3AJournalism&amp;rft.atitle=THE+COPY+DESK+AND+THE+DILEMMAS+OF+THE+INSTITUTIONALIZATION+OF+MODERN+JOURNALISM%22-IN+BRAZIL&amp;rft.aufirst=Afonso&amp;rft.au=Cagliardi%2C+Juliana&amp;rft.aulast=de+Albuquerque&amp;rft.date=2011&amp;rft.genre=article&amp;rft_id=info%3Adoi%2F10.1080%2F1461670X.2010.511956&amp;rft.issue=1&amp;rft.jtitle=Journalism+Studies&amp;rft_val_fmt=info%3Aofi%2Ffmt%3Akev%3Amtx%3AJournalism&amp;rft.volume=12" class="E3988"><span style="display:none;">#160;</span></span></li><li id="cite_note-18"><span class="mw-cite-backlink"><b><a href="#cite_ref-18"></a></b></span>
```



	A	B	C	D	E
1	Member	Trade mil. USD (Nom. GDP mil. U	PPP GDP mil. U:	Nom. GDP per ca
2	Argentina	142,370	594,975	924,481	13,589
3	Australia	496,700	1,343,608	1,246,480	50,962
4	Brazil	484,600	2,140,940	3,217,966	8,670
5	Canada	947,200	1,532,340	1,742,666	43,332
6	China	4,201,000	12,361,737	23,066,642	7,990
7	France	1,212,300	2,570,023	2,833,151	37,675
8	Germany	2,866,600	3,618,621	4,122,402	40,997
9	India	850,600	2,457,748	9,585,371	1,617
10	Indonesia	346,100	1,014,867	3,256,727	3,382
11	Italy	948,600	1,895,318	2,289,578	29,887
12	Japan	1,522,400	5,106,259	5,066,064	32,496
13	South Korea	1,170,900	1,521,000	2,029,861	27,195
14	Mexico	813,500	1,124,316	2,410,946	9,009
15	Russia	844,200	1,442,406	3,866,332	9,055
16	Saudi Arabia	521,600	689,004	1,803,419	20,813
17	South Africa	200,100	288,199	758,123	5,695
18	Turkey	417,000	769,474	1,756,510	9,437
19	United Kingdom	1,189,400	2,609,912	2,877,505	43,771
20	United States	3,944,000	19,377,203	19,377,203	55,805
21	European Union	4,485,000	16,970,024	20,745,303	31,918

WHAT IS WEB SCRAPING?

- ▶ Copying data from a website and storing it in another form so that it's useful to you.
 - ▶ Text
 - ▶ Numbers
 - ▶ Images
 - ▶ Videos

WHAT IS WEB SCRAPING?

- ▶ Regularities in HTML allow you to pick out just the elements you want
- ▶ Advantage: you get just the data you want, and sometimes it's data that was not previously available in a useful form
- ▶ Most scraping requires coding, or else software that handles the code behind the scenes
- ▶ No coding in today's exercise, but it will prepare you to start thinking like a programmer

WEB SCRAPING TOOLS

- ▶ **Commercial software and services** (import.io, dex.io, octoparse, parsehub, fminer)
- ▶ **Browser extensions** (OutwitHub for Firefox, Scraper for Chrome)
- ▶ **Code** (Python, Ruby, etc.)

SCRAPING WITH CODE

- ▶ Python libraries
 - ▶ **Beautiful Soup** (best for beginners)
 - ▶ **Scrapy** (crawls webpages)
 - ▶ **Selenium** (headless browser)

- ▶ Automate data gathering with **cron jobs**

LIMITS TO SCRAPING

- ▶ Poorly formatted HTML
- ▶ Authentication systems, paywalls, CAPTCHAs
- ▶ Systems that use sessions or cookies to track user activity
- ▶ Other access restrictions and usage caps
- ▶ Information spread across multiple pages

USING SCRAPED DATA

PROPUBLICA | See more at VITAL SIGNS

[f](#) [t](#) [DONATE](#)





Dollars for Docs

By Charles Ornstein, Lena Groeger, Mike Tygas, and Ryann Grochowski Jones, ProPublica. Updated December 13, 2016

Pharmaceutical and medical device companies are now required by law to release details of their payments to a variety of doctors and U.S. teaching hospitals for promotional talks, research and consulting, among other categories. Use this tool to search for general payments (excluding research and ownership interests) made from August 2013 to December 2015. | [Related Story: We've Updated Dollars for Docs. Here's What's New.](#)

Has Your Doctor Received Drug or Device Company Money?

For example: Andrew Jones, Boston, 10013

 **\$6.25B** in disclosed payments  **810,716** doctors  **1,171** teaching hospitals  **1,866** companies

Totals listed below account for all payments from August 2013 to December 2015.

Top 50 Companies

Click on a company to see how its payments break down by drug, device or doctor. Or, [see all companies](#)

Highest-Earning Doctors

NAME	PAYMENTS
ROGER JACKSON Orthopaedic Surgery of the Spine	\$54.1M

About the Dollars for Docs Data

Details behind our drug company money database.

Download the Data

The entire data set is available for purchase in the [ProPublica Data Store](#).

Source

The Centers for Medicare and Medicaid Services [Open Payments](#) data.

Archive

Search for payments made by 17 drug companies between 2009 and 2013.

Patients, Take Action

We want to know how you've used or

USING SCRAPED DATA – QUESTIONS TO ASK YOURSELF

- ▶ Is the material copyrighted?
- ▶ Is the dataset already available elsewhere?
- ▶ Have you gathered enough data?
- ▶ Are the data accurate?
- ▶ What can I learn from the data? And how can I tell a story with the data?

PART II: SCRAPING THE WEB WITH GOOGLE SHEETS

open your web browser and visit:

drive.google.com

[DRIVE.GOOGLE.COM](https://drive.google.com)

Google Drive

Search Drive



NEW

My Drive ▾

+ New folder...

File upload

Folder upload

Google Docs >

Google Sheets >

Google Slides >

More >

Backups

6 GB of 17 GB used

Upgrade storage

	Owner	Last modified ↓	File size
EL	nyc yak	May 3, 2017 nyc yak	—
	me	Apr 17, 2017 me	—
nMagazine	me	Mar 23, 2017 me	—
lipper	me	Feb 28, 2017 me	—
n	me	Feb 23, 2017 me	—
ssertation	me	Apr 30, 2016 me	—
Notes	me	Apr 30, 2016 me	—
Presentations	me	Apr 30, 2016 me	—
Teaching	me	Apr 30, 2016 me	—
Projects	me	Apr 30, 2016 me	—
PhD	me	Mar 18, 2016 me	—
Home	me	Mar 18, 2016 me	—
Images	me	Dec 2, 2015 me	—

COPY AND PASTE (SOMETIMES) WORKS — BUT IT'S TEDIOUS

- ▶ WIKIPEDIA.ORG —> G20
- ▶ [HTTPS://EN.WIKIPEDIA.ORG/WIKI/G20](https://en.wikipedia.org/wiki/G20)

Member country data [edit]

Member	Trade mil. USD (2014)	Nom. GDP mil. USD (2017) ^[20]	PPP GDP mil. USD (2017) ^[20]	Nom. GDP per capita USD (2015) ^[20]	PPP GDP per capita USD (2015) ^[20]	HDI (2015)	Population (2014)	Area	P5	G4	G7	BRICS	MIKTA	DAC	OECD	C'wth	N11
Argentina	142,370	594,975	904,481	13,529	22,554	0.827	42,961,000	2,790,400	✗	✗	✗	✗	✗	✗	✗	✗	✗
Australia	496,700	1,343,608	1,246,480	50,962	47,389	0.939	23,599,000	7,692,024	✗	✗	✗	✗	✓	✓	✓	✓	✗
Brazil	484,600	2,140,940	3,217,986	8,670	16,155	0.754	202,768,000	8,515,767	✗	✓	✗	✓	✗	✗	✗	✗	✗
Canada	947,200	1,532,340	1,742,656	43,332	44,967	0.920	35,467,000	9,984,670	✗	✗	✓	✗	✗	✓	✓	✓	✗
China	4,201,000	12,361,737	23,066,642	7,990	14,107	0.738	1,367,520,000	9,572,900	✓	✗	✗	✓	✗	✗	✗	✗	✗
European Union	4,485,000	16,970,024	20,745,303	31,918	37,852	0.876	505,570,700	4,422,773	✗	✗	✓	✗	✗	✓	✗	✗	✗
France	1,212,300	2,570,023	2,833,151	37,675	41,181	0.897	63,951,000	640,679	✓	✗	✓	✗	✗	✓	✓	✗	✗
Germany	2,866,600	3,618,621	4,122,402	40,997	46,893	0.926	81,940,000	357,114	✗	✓	✓	✗	✗	✓	✓	✗	✗
India	850,600	2,457,748	9,585,371	1,617	6,162	0.624	1,259,695,000	3,287,263	✗	✓	✗	✓	✗	✗	✗	✓	✗
Indonesia	346,100	1,014,867	3,256,727	3,362	11,126	0.689	251,490,000	1,904,569	✗	✗	✗	✗	✓	✗	✗	✗	✓
Italy	948,600	1,895,318	2,289,578	29,867	35,708	0.887	61,665,551	301,336	✗	✗	✓	✗	✗	✓	✓	✗	✗
Japan	1,522,400	5,106,259	5,066,064	32,486	38,054	0.903	127,061,000	377,930	✗	✓	✓	✗	✗	✓	✓	✗	✗
Mexico	813,500	1,124,316	2,410,946	9,009	17,534	0.762	119,581,789	1,964,375	✗	✗	✗	✗	✓	✗	✓	✗	✓
Russia	844,200	1,442,406	3,866,332	9,055	25,411	0.804	148,300,000	17,098,242	✓	✗	✗	✓	✗	✗	✗	✗	✗
Saudi Arabia	521,600	689,004	1,803,419	20,813	53,624	0.847	31,624,000	2,149,690	✗	✗	✗	✗	✗	✗	✗	✗	✗
South Africa	200,100	288,199	758,123	5,695	13,165	0.666	53,699,000	1,221,037	✗	✗	✗	✓	✗	✗	✗	✓	✗
South Korea	1,170,900	1,521,000	2,029,861	27,195	36,511	0.901	51,437,000	100,210	✗	✗	✗	✗	✓	✓	✓	✗	✓
Turkey	417,000	769,474	1,756,510	9,437	20,438	0.767	77,324,000	783,562	✗	✗	✗	✗	✓	✗	✓	✗	✓
United Kingdom	1,189,400	2,609,912	2,877,505	43,771	41,159	0.909	64,511,000	242,495	✓	✗	✓	✗	✗	✓	✓	✓	✗
United States	3,944,000	19,377,203	19,377,203	55,805	55,805	0.920	318,523,000	9,526,468	✓	✗	✓	✗	✗	✓	✓	✗	✗

COPY AND PASTE (SOMETIMES) WORKS — BUT IT'S TEDIOUS

- ▶ CLICK AND DRAG TO SELECT THE ENTIRE TABLE
- ▶ TYPE COMMAND+C OR CTRL+C TO COPY
- ▶ GO TO GOOGLE SHEETS AND TYPE COMMAND+V OR CTRL+V TO PASTE

Member country data [edit]

Member	Trade mil. USD (2014)	Nom. GDP mil. USD (2017) ^[20]	PPP GDP mil. USD (2017) ^[20]	Nom. GDP per capita USD (2015) ^[20]	PPP GDP per capita USD (2015) ^[20]	HDI (2015)	Population (2014)	Area	P5	G4	G7	BRICS	MIKTA	DAC	OECD	C'wth	N11
Argentina	142,370	594,975	904,481	13,529	22,554	0.827	42,961,000	2,790,400	✗	✗	✗	✗	✗	✗	✗	✗	✗
Australia	496,700	1,343,608	1,246,480	50,962	47,389	0.939	23,599,000	7,692,024	✗	✗	✗	✗	✓	✓	✓	✓	✗
Brazil	484,600	2,140,940	3,217,986	8,670	16,155	0.754	202,768,000	8,515,767	✗	✓	✗	✓	✗	✗	✗	✗	✗
Canada	947,200	1,532,340	1,742,656	43,332	44,967	0.920	35,467,000	9,984,670	✗	✗	✓	✗	✗	✓	✓	✓	✗
China	4,201,000	12,361,737	23,066,642	7,990	14,107	0.738	1,367,520,000	9,572,900	✓	✗	✗	✓	✗	✗	✗	✗	✗
European Union	4,485,000	16,970,024	20,745,303	31,918	37,852	0.876	505,570,700	4,422,773	✗	✗	✓	✗	✗	✓	✗	✗	✗
France	1,212,300	2,570,023	2,833,151	37,675	41,181	0.897	63,951,000	640,679	✓	✗	✓	✗	✗	✓	✓	✗	✗
Germany	2,866,600	3,618,621	4,122,402	40,997	46,893	0.926	81,940,000	357,114	✗	✓	✓	✗	✗	✓	✓	✗	✗
India	850,600	2,457,748	9,585,371	1,617	6,162	0.624	1,259,695,000	3,287,263	✗	✓	✗	✓	✗	✗	✗	✓	✗
Indonesia	346,100	1,014,867	3,256,727	3,362	11,126	0.689	251,490,000	1,904,569	✗	✗	✗	✗	✓	✗	✗	✗	✓
Italy	948,600	1,895,318	2,289,578	29,667	35,708	0.887	61,665,551	301,336	✗	✗	✓	✗	✗	✓	✓	✗	✗
Japan	1,522,400	5,106,259	5,066,064	32,486	38,054	0.903	127,061,000	377,930	✗	✓	✓	✗	✗	✓	✓	✗	✗
Mexico	813,500	1,124,316	2,410,946	9,009	17,534	0.762	119,581,789	1,964,375	✗	✗	✗	✗	✓	✗	✓	✗	✓
Russia	844,200	1,442,406	3,866,332	9,055	25,411	0.804	148,300,000	17,098,242	✓	✗	✗	✓	✗	✗	✗	✗	✗
Saudi Arabia	521,600	689,004	1,803,419	20,813	53,624	0.847	31,624,000	2,149,690	✗	✗	✗	✗	✗	✗	✗	✗	✗
South Africa	200,100	288,199	758,123	5,695	13,165	0.666	53,699,000	1,221,037	✗	✗	✗	✓	✗	✗	✗	✓	✗
South Korea	1,170,900	1,521,000	2,029,861	27,195	36,511	0.901	51,437,000	100,210	✗	✗	✗	✗	✓	✓	✓	✗	✓
Turkey	417,000	769,474	1,756,510	9,437	20,438	0.767	77,324,000	783,562	✗	✗	✗	✗	✓	✗	✓	✗	✓
United Kingdom	1,189,400	2,609,912	2,877,505	43,771	41,159	0.909	64,511,000	242,495	✓	✗	✓	✗	✗	✓	✓	✓	✗
United States	3,944,000	19,377,203	19,377,203	55,805	55,805	0.920	318,523,000	9,526,468	✓	✗	✓	✗	✗	✓	✓	✗	✗

IT WORKED! (SORT OF) BUT IT'S BETTER TO AUTOMATE

Untitled spreadsheet ☆

File Edit View Insert Format Data Tools Add-ons Help All changes saved in Drive

Comments Share

fx Member

	A	B	C	D	E	F	G	H	I	J	K	L
1	Member	GDP mil. USD (2014)	GDP mil. USD (2015)	GDP mil. USD (2016)	GDP per capita USD (2014)	GDP per capita USD (2015)	HDI (2015)	Population (2014)	Area	P5	G4	G7
2	Argentina	142,370	591,975	921,481	13,589	22,551	0.827	42,961,000	2,780,100			
3	Australia	496,700	1,343,608	1,246,480	50,962	47,383	0.930	23,599,000	7,692,024			
4	Brazil	484,600	2,140,940	3,217,906	6,670	16,155	0.754	202,768,000	8,515,767			
5	Canada	947,200	1,532,340	1,742,856	43,332	44,967	0.92	35,467,000	9,984,870			
6	China	4,201,000	12,361,737	23,066,542	7,990	14,107	0.736	1,357,520,000	9,572,900			
7	European Union	4,485,000	16,970,021	20,745,303	31,918	37,852	0.876	505,570,700	4,422,773			
8	France	1,212,300	2,570,023	2,833,151	37,675	41,181	0.897	53,951,000	640,579			
9	Germany	2,066,600	3,618,621	4,122,402	40,997	46,893	0.926	80,940,000	357,114			
10	India	850,600	2,457,748	9,585,371	1,617	6,162	0.624	1,259,695,000	3,287,263			
11	Indonesia	346,100	1,014,867	3,256,727	3,362	11,125	0.689	251,490,000	1,904,569			
12	Italy	948,600	1,895,318	2,289,578	29,867	35,708	0.837	60,665,551	301,336			
13	Japan	1,522,400	5,106,259	5,066,064	32,486	38,054	0.903	127,061,000	377,330			
14	Mexico	813,500	1,124,316	2,410,346	9,009	17,534	0.762	119,581,789	1,964,375			
15	Russia	844,200	1,442,406	3,866,332	9,055	25,411	0.804	146,500,000	17,098,242			
16	Saudi Arabia	521,600	689,004	1,803,419	20,813	53,624	0.847	30,624,000	2,149,590			
17	South Africa	200,100	283,199	758,123	5,695	13,165	0.696	53,699,000	1,221,037			
18	South Korea	1,170,900	1,521,000	2,029,861	27,195	36,511	0.901	50,437,000	100,210			
19	Turkey	417,000	769,474	1,756,510	9,437	20,433	0.767	77,324,000	783,562			
20	United Kingdom	1,189,400	2,609,912	2,877,505	43,771	41,159	0.909	64,511,000	242,495			
21	United States	3,944,000	19,377,203	19,377,203	55,805	55,805	0.92	318,523,000	9,525,468			
22												
23												
24												

+ G20 Nations Sheet2

SCRAPING TABLES TO GOOGLE SHEETS

Command to scrape a table from a website:

```
=ImportHTML("http://example.com","table",1)
```

- ▶ `=ImportHTML("http://website.com","table",1)` **Command (1,2,3)**
- ▶ `=ImportHTML("http://website.com","table",1)` **"URL"**
- ▶ `=ImportHTML("http://website.com","table",1)` **"Element Type"**
- ▶ `=ImportHTML("http://website.com","table",1)` **Element Number**

SCRAPING TABLES TO GOOGLE SHEETS

To scrape a table from <https://en.wikipedia.org/wiki/G20>

Type this command directly into a spreadsheet box:

```
=ImportHTML("https://en.wikipedia.org/wiki/G20","table",1)
```

- ▶ `=ImportHTML("https://en.wikipedia.org/wiki/G20","table",1)`
- ▶ `=ImportHTML("https://en.wikipedia.org/wiki/G20","table",1)`
- ▶ `=ImportHTML("https://en.wikipedia.org/wiki/G20","table",1)`
- ▶ `=ImportHTML("https://en.wikipedia.org/wiki/G20","table",1)`

fx | =ImportHTML("https://en.wikipedia.org/wiki/G20","table",0)

	A	B	C	D	E	F	G	H	I	J	K	L
1	<i>?</i> =ImportHTML("https://en.wikipedia.org/wiki/G20","table",0)											
2												
3												
4												
5												
6												
7												
8												
9												
10												
11												
12												
13												
14												
15												
16												
17												
18												
19												
20												
21												
22												
23												
24												

fx | =ImportHTML("https://en.wikipedia.org/wiki/G20","table",0)

	A	B	C	D	E
1	<i>?</i> =ImportHTML("https://en.wikipedia.org/wiki/G20","table",0)				
2					
3					
4					

=ImportHTML("https://en.wikipedia.org/wiki/G20", "table", 0)

	A	B	C	D	E	F	G	H	I	J	K	L
1	Member	Trade mil. USD	Nom. GDP mil. U	PPP GDP mil. U	Nom. GDP per c	PPP GDP per ca	HDI (2015)	Population (2014	Area	P5	G4	G7
2	Argentina	142,370	594,975	924,481	13,589	22,554	0.827	42,961,000	2,780,400	N	N	N
3	Australia	496,700	1,343,608	1,246,480	50,962	47,389	0.939	23,599,000	7,692,024	N	N	N
4	Brazil	484,600	2,140,940	3,217,986	8,670	16,155	0.754	202,768,000	6,515,767	N	Y	N
5	Canada	947,200	1,532,340	1,742,656	43,332	44,967	0.92	35,467,000	9,984,670	N	N	Y
6	China	4,201,000	12,361,737	23,066,642	7,990	14,107	0.738	1,367,520,000	9,572,900	Y	N	N
7	France	1,212,300	2,570,023	2,833,151	37,675	41,181	0.897	63,951,000	640,679	Y	N	Y
8	Germany	2,866,600	3,618,621	4,122,402	40,997	46,893	0.926	80,940,000	357,114	N	Y	Y
9	India	850,600	2,457,748	9,585,371	1,617	6,162	0.624	1,259,695,000	3,287,263	N	Y	N
10	Indonesia	346,100	1,014,867	3,256,727	3,362	11,126	0.689	251,490,000	1,904,569	N	N	N
11	Italy	948,600	1,895,318	2,289,578	29,867	35,708	0.887	60,665,551	301,336	N	N	Y
12	Japan	1,522,400	5,106,259	5,066,064	32,486	38,054	0.903	127,061,000	377,930	N	Y	Y
13	South Korea	1,170,900	1,521,000	2,029,861	27,195	36,511	0.901	50,437,000	100,210	N	N	N
14	Mexico	813,500	1,124,316	2,410,946	9,009	17,534	0.762	119,581,789	1,964,375	N	N	N
15	Russia	844,200	1,442,406	3,866,332	9,055	25,411	0.804	146,300,000	17,098,242	Y	N	N
16	Saudi Arabia	521,600	689,004	1,803,419	20,813	53,624	0.847	30,624,000	2,149,690	N	N	N
17	South Africa	200,100	288,189	758,123	5,695	13,165	0.666	53,699,000	1,221,037	N	N	N
18	Turkey	417,000	769,474	1,756,510	9,437	20,438	0.767	77,324,000	783,562	N	N	N
19	United Kingdom	1,189,400	2,609,912	2,877,505	43,771	41,159	0.909	64,511,000	242,495	Y	N	Y
20	United States	3,944,000	19,377,203	19,377,203	55,805	55,805	0.92	318,523,000	9,526,468	Y	N	Y
21	European Union	4,485,000	16,970,024	20,745,303	31,918	37,852	0.876	505,570,700	4,422,773	N	N	Y
22												
23												
24												

finding "tables" in the HTML source code of the page

The screenshot shows a web browser displaying the G20 page. A table is visible in the main content area, listing member countries and their respective leaders and finance ministers. The browser's developer tools are open, showing the HTML source code for the table and the CSS styles applied to it.

memberships or other international organisations, such as the [G7](#) and [BRICS](#). Total GDP figures are given in millions of US dollars.

`table.wikitable.sortable.jquery-tablesorter` 1040 x 683

Member	Leader position	Head of government	Finance portfolio	Finance minister	Central bank governor
Argentina	President	Mauricio Macri	Minister of Public Finances Minister of the Treasury	Luis Caputo Nicolás Dujovne	Federico Sturzenegger
Australia	Prime Minister	Malcolm Turnbull	Treasurer	Scott Morrison	Philip Lowe
Brazil	President	Michel Temer	Minister of Finance	Henrique Meirelles	Ilan Goldfajn
Canada	Prime Minister	Justin Trudeau	Minister of Finance	Bill Morneau	Stephen Poloz
China	President ^(note 1)	Xi Jinping ^(note 1)	Minister of Finance	Xiao Jie	Zhou Xiaochuan
France	President	Emmanuel Macron	Minister of the Economy	Bruno Le Maire	François Villeroy de Galhau
Germany	Chancellor	Angela Merkel	Minister of Finance	Wolfgang Schäuble	Jens Weidmann
India	Prime Minister	Narendra Modi	Minister of Finance	Arun Jaitley	Urjit Patel

```
<table class="wikitable sortable jquery-tablesorter" style="width:100%;"> == $0
  <thead>
    <tr>
      <th class="headerSort" tabindex="0" role="columnheader" button title="Sort ascending">Member</th>
      <th class="headerSort" tabindex="0" role="columnheader" button title="Sort ascending">Leader position</th>
      <th class="headerSort" tabindex="0" role="columnheader" button title="Sort ascending">Head of government</th>
      <th class="headerSort" tabindex="0" role="columnheader" button title="Sort ascending">Finance portfolio</th>
      <th class="headerSort" tabindex="0" role="columnheader" button title="Sort ascending">Finance minister</th>
      <th class="headerSort" tabindex="0" role="columnheader" button title="Sort ascending">Central bank governor</th>
    </tr>
  </thead>
  <tbody>_</tbody>
  <tfoot></tfoot>
</table>
<h3>_</h3>
<table class="wikitable sortable jquery-tablesorter" style="font-size:90%; width:100%;">_</table>
<p>_</p>
<p>_</p>
```

Styles Computed Event Listeners >>

Filter :show .cls +

```
element.style {
  width: 100%;
}
table.wikitable load.php?debug=...tor.desktop...:1
{
  margin: 1em 0;
  background-color: #f8f9fa;
  border: 1px solid #e2e9e1;
  border-collapse: collapse;
  color: #000;
}
table {
  font-size: 100%;
}
table {
  user agent stylesheet
```

html body #content #bodyContent #mw-content-text div.mw-parser-output table.wikitable.sortable.jquery-tablesorter thead tr th.headerSort

SCRAPING LISTS TO GOOGLE SHEETS

To scrape a list from <https://en.wikipedia.org/wiki/G20>

Enter the last command, changing "table" to "list" and 1 to 2:

```
=ImportHTML("https://en.wikipedia.org/wiki/G20","list",2)
```

- ▶ `=ImportHTML("https://en.wikipedia.org/wiki/G20","list",2)`
- ▶ `=ImportHTML("https://en.wikipedia.org/wiki/G20","list",2)`
- ▶ `=ImportHTML("https://en.wikipedia.org/wiki/G20","list",2)`
- ▶ `=ImportHTML("https://en.wikipedia.org/wiki/G20","list",2)`

finding "lists" in the HTML source code of the page

The screenshot shows a web browser displaying the Wikipedia page for G20. The table of contents is visible, listing sections such as History, Summits, Organisation, List of members, Invitees, and Criticisms. The developer tools are open, showing the HTML source code for the table of contents. The code includes a list of links for each section, with the 'List of summits' link selected. The style pane on the right shows the default styles for the selected list element.

Page information
Wikidata item
Cite this page
Print/export
Create a book
Download as PDF
Printable version
In other projects
Wikimedia Commons
Languages
Afrikaans
Alemannisch
العربية
Azərbaycanca
বাংলা
Bân-lâm-gú

Contents [hide]

- History
- Summits
 - 2.1 List of summits
 - 2.2 Chair rotation
- Organisation
 - 3.1 Proposed permanent secretariat
- List of members
 - 4.1 Leaders
 - 4.2 Member country data
 - 4.3 Role of Asian countries
- Invitees
 - 5.1 Permanent guest invitees
- Criticisms
 - 6.1 Exclusivity of membership
 - 6.1.1 Norwegian perspective

Membership 20 [show]
Chairperson Angela Merkel (2017)
 Mauricio Macri (2018)
Staff None^[2]
Website g20.org

Elements Console Sources Network Performance Memory Application Security Audits HTTPS Everywhere

```
>>><sup id="cite_ref-g20members_2-2" class="reference"></sup>  
</p>  
<p></p>  
>>><div id="toc" class="toc">  
>>><div id="toctitle" class="toctitle"></div>  
...>>><ul> == $8  
>>><li class="toclevel-1 toctection-1"></li>  
>>><li class="toclevel-1 toctection-2">  
>>><a href="#Summits"></a>  
>>><ul>  
>>><li class="toclevel-2 toctection-3">  
>>><a href="#List_of_summits">  
>>><span class="tocnumber">2.1</span>  
>>><span class="toctext">List of summits</span>  
>>></a>  
>>></li>  
>>><li class="toclevel-2 toctection-4"></li>  
>>></ul>  
>>></li>  
>>></ul>
```

html body div#content.mw-body div#bodyContent.mw-body-content div#mw-content-text.mw-content-ltr div.mw-parser-output div#toc.toc ul

list 9 of 108 [Cancel]

Style Computed Event Listeners >>

Filter :show .cls +

```
element.style {  
}  
.mw-content-ltr load.php?debug=...tor.desktop...:1  
.toc ul, .mw-  
content-ltr #toc ul, .mw-content-rtl .mw-  
content-ltr .toc ul, .mw-content-rtl .mw-  
content-ltr #toc ul {  
  text-align: left;  
}  
#toc ul, .toc load.php?debug=...tor.desktop...:1  
ul {  
  list-style-type: none;  
  list-style-image: none;  
  margin-left: 0;  
  padding: 0 0;  
  text-align: left;
```

MORE POWERFUL SCRAPING COMMANDS WITH GSHEETS

type command directly into a spreadsheet box:

```
=IMPORTFEED("https://example.com/whatever")
```

- ▶ `=IMPORTFEED("https://example.com")` **Command (1)**
- ▶ `=IMPORTFEED("https://example.com")` **"URL"**

Note: the `IMPORTFEED` command takes only **ONE** argument

SCRAPING FROM AN RSS FEED WITH GSHEETS

type command directly into a spreadsheet box:

```
=IMPORTFEED("https://example.com/whatever")
```

Try one of these RSS feeds from the Library of Congress:

▶ <https://www.congress.gov/rss>

RSS and Email Alerts

Keeping up with Congress is easy with RSS and Email Alerts from Congress.gov. This page allows you to subscribe to a variety of RSS and Email alerts related to Congressional activity and legislation.

Jump to: [How to Get Email Alerts](#)

Subscribe to RSS

RSS (Really Simple Syndication) is a technology that delivers news to a computer or mobile device. Congress.gov offers several RSS feeds for use in an RSS reader or RSS-enabled Web browser. For details about RSS, see the [RSS help page on loc.gov](#).

Most-Viewed Bills

A weekly top-ten list of the bills measured by page views on Congress.gov.



Search Tips

Tips on how to perform effective, advanced searches on Congress.gov.



Appropriations Tables

Congress.gov includes appropriation tables to simplify the process of tracking legislative appropriation bills. Each time the table is updated, this feed will let you know what has been added to the table.



Bills Presented to the President

When a piece of legislation that requires the President's signature to become law passes Congress, this feed will alert you that the legislation has been submitted to the President for his or her signature.



On the House Floor Today

Legislation brought to the floor of the United States House of Representatives.



On the Senate Floor Today

Legislation brought to the floor of the United States Senate.



In Custodia Legis: Law Librarians of Congress

The latest posts from the award winning Law Library of Congress blog, *in Custodia Legis*, featuring news about Congress.gov, topical articles concerning foreign and domestic law, and an insider's view of the happenings in and around the Library of Congress and Capitol Hill.



Help

[Saved Search Help](#)

[Alerts Help](#)

 [Creating and Using Email Alerts](#)

fx =IMPORTFEED("https://www.congress.gov/rss/house-floor-today.xml")

	A	B	C	D	E	F	G	H	I
1	S.1083	https://www.congress.gov/bill/115	A bill to amend section 1214 of title 5, United States Code, to provide for stays during a period that the Merit Systems Protection Board lacks a quorum. (05/25/2017 legislative day)						
2	H.R.1973	https://www.congress.gov/bill/115	Protecting Young Victims from Sexual Abuse Act of 2017 (05/26/2017 legislative day)						
3	H.R.1761	https://www.congress.gov/bill/115	Protecting Against Child Exploitation Act of 2017 (05/25/2017 legislative day)						
4	S.Con.Res.14	https://www.congress.gov/bill/115	A concurrent resolution authorizing the use of Emancipation Hall in the Capitol Visitor Center for an event to celebrate the birthday of King Kamehameha I. (05/24/2017 legislative day)						
5	H.Res.350	https://www.congress.gov/bill/115	Permitting official photographs of the House of Representatives to be taken while the House is in actual session on a date designated by the Speaker. (05/24/2017 legislative day)						
6	H.R.953	https://www.congress.gov/bill/115	Reducing Regulatory Burdens Act of 2017 (05/24/2017 legislative day)						
7	H.R.624	https://www.congress.gov/bill/115	To restrict the inclusion of social security account numbers on Federal documents sent by mail, and for other purposes. (05/24/2017 legislative day)						
8	H.R.1293	https://www.congress.gov/bill/115	To amend title 5, United States Code, to require that the Office of Personnel Management submit an annual report to Congress relating to the use of official time by Federal employees. (05/24/2017 legislative day)						
9	H.Res.352	https://www.congress.gov/bill/115	Providing for consideration of the bill (H.R. 1973) to prevent the sexual abuse of minors and amateur athletes by requiring the prompt reporting of sexual abuse to law enforcement authorities, and for other purposes; providing for consideration of the bill (H.R. 1761) to amend title 18, United States Code, to criminalize the knowing consent of the visual depiction, or live transmission, of a minor engaged in sexually explicit conduct, and for other purposes; and providing for proceedings during the period from May 26, 2017, through June 5, 2017. (05/24/2017 legislative day)						
10	H.R.2473	https://www.congress.gov/bill/115	Enforcing Justice for Victims of Trafficking Act of 2017 (05/23/2017 legislative day)						



MORE POWERFUL SCRAPING COMMANDS FOR GSHEETS

Pull data from an API (but don't try this just yet):

```
=IMPORTDATA("https://data.cityofboston.gov/resource/427a-3cn5.json")
```

- ▶ `=IMPORTXML("https://example.com", "//a/@href")` **Command (1)**
- ▶ `=IMPORTXML("https://data.cityofboston.gov/resource/427a-3cn5.json")` **"URL"**

=IMPORTDATA("https://data.cityofboston.gov/resource/427a-3cn5.json")

	A	B	C	D	E	F	G	H	I	J
1	{":@computed_region_aywg_kpfl	address:"1000 Washington ST"	businessname:"1000 Wash	city:"Roxbury"	comments:"Prep	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalowner:"TER	licenset
2	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washingto	city:"Roxbury"	comments:"Clear	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
3	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"clear	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
4	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"Provi	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
5	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalowner:"TER	licenset	
6	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"clear	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
7	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"Provi	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
8	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"Provi	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
9	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"Keep	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
10	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"Keep	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
11	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalowner:"TER	licenset	
12	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"Provi	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
13	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"Cove	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
14	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"Keep	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
15	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"Clear	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
16	{":@computed_region_aywg_kpfl	address:"1000 Washingto	businessname:"1000 Washington Ci	city:"Roxbury"	comments:"Cove	descript:"Eating &	expdtm:"2011-12	issdtm:"2011-11-	legalow	
17	{":@computed_region_aywg_kpfl	address:"635 Hyde Park A	businessname:"100 Percent Delicia	city:"Roslindale"	comments:"Low	descript:"Eating &	expdtm:"2017-12	issdtm:"2013-04	legalow	
18	{":@computed_region_aywg_kpfl	address:"635 Hyde Park A	businessname:"100 Percent Delicia	city:"Roslindale"	comments:"PIC	descript:"Eating &	expdtm:"2017-12	issdtm:"2013-04	legalow	
19	{":@computed_region_aywg_kpfl	address:"635 Hyde Park A	businessname:"100 Percent Delicia	city:"Roslindale"	comments:"Repa	descript:"Eating &	expdtm:"2017-12	issdtm:"2013-04	legalow	
20	{":@computed_region_aywg_kpfl	address:"635 Hyde Park A	businessname:"100 Percent Delicia	city:"Roslindale"	comments:"Pain	descript:"Eating &	expdtm:"2017-12	issdtm:"2013-04	legalow	
21	{":@computed_region_aywg_kpfl	address:"635 Hyde Park A	businessname:"100 Percent Delicia	city:"Roslindale"	comments:"Certi	descript:"Eating &	expdtm:"2017-12	issdtm:"2013-04	legalow	
22	{":@computed_region_aywg_kpfl	address:"635 Hyde Park A	businessname:"100 Percent Delicia	city:"Roslindale"	comments:"Clear	descript:"Eating &	expdtm:"2017-12	issdtm:"2013-04	legalow	
23	{":@computed_region_aywg_kpfl	address:"635 Hyde Park A	businessname:"100 Percent Delicia	city:"Roslindale"	comments:"All fo	descript:"Eating &	expdtm:"2017-12	issdtm:"2013-04	legalow	
24	{":@computed_region_aywg_kpfl	address:"635 Hyde Park A	businessname:"100 Percent Delicia	city:"Roslindale"	comments:"Base	descript:"Eating &	expdtm:"2017-12	issdtm:"2013-04	legalow	

MORE POWERFUL SCRAPING COMMANDS FOR GSHEETS

Scrape all links on a webpage:

```
=IMPORTXML("https://example.com","//a/@href")
```

- ▶ `=IMPORTXML("https://example.com","//a/@href")` **Command** (1,2)
- ▶ `=IMPORTXML("https://example.com","//a/@href")` **"URL"**
- ▶ `=IMPORTXML("https://example.com","//a/@href")` **"Query"** (for links)

FIFTEEN

MINUTE

BREAK

10:30-45

SCRAPING DATA FROM THE WEB II: APPLYING YOUR NEW SKILLS

CHARLES BERRET

BROWN INSTITUTE FOR MEDIA INNOVATION

COLUMBIA UNIVERSITY

DATA STORYTELLING AT BOSTON UNIVERSITY

JUNE 6, 2017

FOUR WAYS TO SCRAPE WITH GOOGLE SHEETS

- ▶ **Tables** (=IMPORTHTML)
- ▶ **Lists** (=IMPORTHTML)
- ▶ **RSS Feeds and APIs** (=IMPORTFEED)
- ▶ **Advanced** (=IMPORTXML)

EXERCISE: SCRAPE DATA, COLLECT SPREADSHEETS, FIND STORIES

- ▶ Choose a topic or question that interests you
- ▶ Find a website with information on that subject
- ▶ Use Google Sheets commands to scrape data into separate tabs of the spreadsheet
- ▶ Gather a personal collection of scraped data that you can study to find patterns, regularities, outliers, or other data insights.