# Scraping Data from the Web

In this lab, you will learn how to scrape data from websites using browser extensions and Google Spreadsheets.

## Getting Started



Despite business complexity, corporate data is surprisingly open in Russia. Photo: Shutterstock

In Russia, shell companies abound and political involvement in business creates significant risks for international businesses. Today, there is the added pressure of increasing US and EU sanctions. In this example, the International Consortium of Investigative Journalists scraped several Russian and international business registries to track ownership and money flows behind shell companies. The data was available online but there was no way to download the data from the site. Take a look at: http://www.icij.org/blog/2014/08/how-investigate-russias-shady-business-world-online

# Scraping Data Using Browser Extensions

Browser Extensions or Browser Add-Ons are applications that run in your internet browser that enable you to:
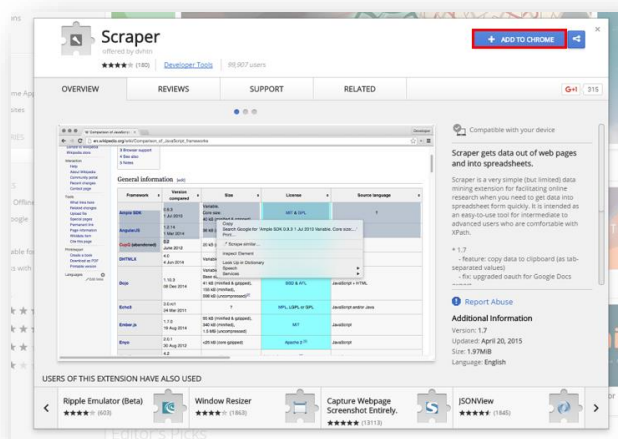
- Select data tables on a website, including specific rows and columns
- Copy data from these tables and use it in a spreadsheet application like Excel or Google Sheets

Here is an example, let's use a browser extension to scrape data about Pakistan from a United Nation's website. You can use any of these procedures – depending on the browser you are using: Google Chrome or Mozilla Firefox.
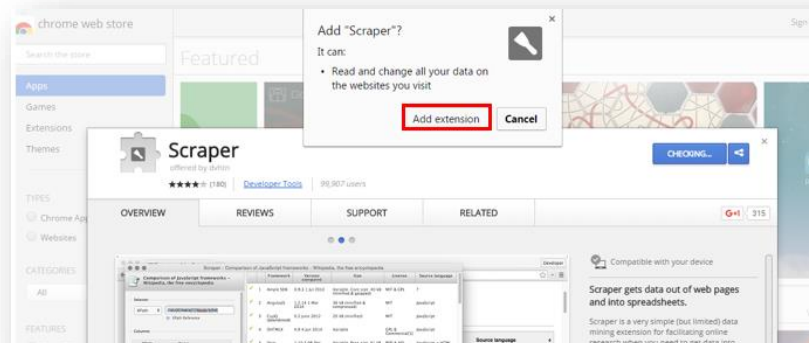
## Scraping Data using Browser Extension for Google Chrome

### Download Browser Extension to Scrape Data

1. Open your Chrome browser
2. To download **Scraper Extension**, go to
   https://chrome.google.com/webstore/detail/scraper/mbigbapnjcgaffohmbkdlecaccepngjd?hl=en
3. In the 'Scraper' window that opens, click the blue **+ ADD TO CHROME** button.



4. In the popup window that appears, click the **Add extension** button.



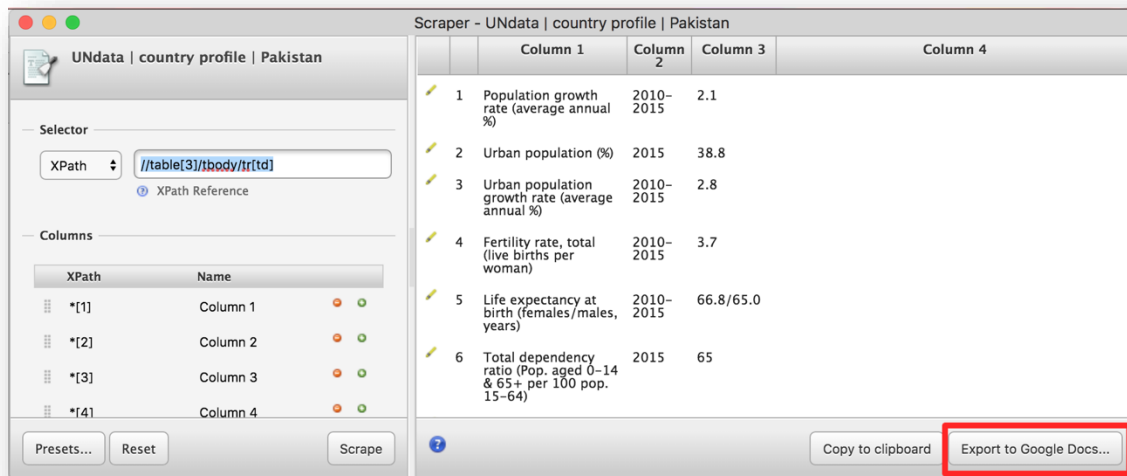5. The extension is now added to chrome.

### *Scraping Data from a Website*

**6.** Now, open a new Chrome window

**7.** Go to: http://data.un.org/CountryProfile.aspx?crName=Pakistan

**8.** Once the page opens, scroll down to reach the table with the heading "Social Indicators"

**9.** Now use the mouse pointer to highlight first three rows of the "Social Indicators" table, then right-click and select **Scrape similar**… (Note, if you select the entire table, the tool will not work. You have to select a portion of the table and Google will automatically find the rest of the table).



**10.** In the Scraper window that opens, click the **Export to Google Docs…** button.

**11.** If you are not signed in to Google, Chrome asks you to sing-in using your Gmail id. A window opens seeking permission for the Scrpaer extension. Click **Allow**.



**12.** A Google Spreadsheet now opens showing the Social Indicators data scraped from the UN website:

## Scraping Data using Browser Add-On for Mozilla Firefox

***Download Browser Add-On to Scrape Data***

**13.**   Open Firefox browser.

**14.**   To download an add-on called **Dafizilla Table2Clipboard**, go to: https://addons.mozilla.org/en-us/firefox/addon/dafizilla-table2clipboard/

**15.**   On this webpage, click **Continue to Download**

**16.**   On the next page that opens, click **+ Add to Firefox**



**17.**   In the pop-up window that appears, click **Install**.

**18.** In the next pop-up window, click **Restart Now**. Firefox closes and reopens – the add-on is now added to Firefox.

*Scraping Data from a Website*
**19.** Now, open a new Firefox window.
**20.** Go to: http://data.un.org/CountryProfile.aspx?crName=Pakistan
**21.** Once the page opens, scroll down to reach the table with the heading "Social Indicators"
**22.** Now use the mouse pointer to highlight first three rows of the "Social Indicators" table, then right-click and select T**able2Clipboard > Copy whole table**



**23.** Now, open a new Microsoft Excel file.
**24.** In the cell A1 of the spreadsheet, right-click and select **Paste Special**

**25.** In the 'Paste Special' window, check if **HTML** is selected, and then click **OK.**



**26.** The social indicators data table is now available in Excel.

# Exercise: Scrape Data Using Browser Extensions

Use the Chrome or Firefox browser to scrape data from the **Number of Health Workers** table available on this web page: https://www.quandl.com/collections/pakistan/pakistan-health-data



## Number of Health Workers

| Source | Indicator | Level | Units | As Of | Api Call | vs World |
|--------|-----------|-------|-------|-------|----------|----------|
| WHO | nursing and midwifery personnel - Pakistan | 100,397.00 | unit | 2010 | JSON, CSV | |
| WHO | dentistry personnel - Pakistan | 10,508.00 | unit | 2010 | JSON, CSV | |
| WHO | pharmaceutical personnel - Pakistan | 8,102.00 | unit | 2004 | JSON, CSV | |
| WHO | laboratory health workers - Pakistan | 9,744.00 | unit | 2004 | JSON, CSV | |
| WHO | environment and public health workers - Pakistan | 106.00 | unit | 2004 | JSON, CSV | |
| WHO | community and traditional health workers - Pakistan | 11,510.00 | unit | 2010 | JSON, CSV | |
| WHO | other health workers - Pakistan | 19,082.00 | unit | 2004 | JSON, CSV | |
| WHO | health management & support workers - Pakistan | 203,337.00 | unit | 2004 | JSON, CSV | |
| WHO | dental technicians/assistants - Pakistan | 1,410.00 | unit | 2004 | JSON, CSV | |
| WHO | laboratory scientists - Pakistan | 1,751.00 | unit | 2004 | JSON, CSV | |
| WHO | laboratory technicians/assistants - Pakistan | 6,323.00 | unit | 2004 | JSON, CSV | |
| WHO | community health workers - Pakistan | 11,510.00 | unit | 2010 | JSON, CSV | |

## Density of Medical Personel

**27.** Scrape the data from Number of Health Workers using chrome Scraper Tool or Firefox Dafizilla Table2Clipboard.

**28.** Chose another table from https://www.quandl.com/collections/pakistan/pakistan-demography-data and scrape it.

## Task 3: Scraping Data Using ImportHTML

**Code Source from a webpage**

Behind of all text, table, photos, videos that a web page contain there is sourcecode. This code determines what is on the page, how it looks and what it does.

**29.** Open https://en.wikipedia.org/wiki/Pakistan_at_the_Olympics - cite_note-jang.com.pk.2Fthenews.2Fjul2-10 either in google chrome or firefox

**30.** To see the source code right click and click in Inspect (Google Chrome) or Inspect Element (Firefox)

**Medals by Games**  [ edit ]

| Games | Athletes | Athletes by sport | | | | | | | | | | | Medals | | | Total | Rank |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1948 London | 39[1] | 5 | 3 | 2 | 19 | | | | 4 | 2 | 4 | | - | - | - | 0[1] | -- |
| 1952 Helsinki | 44[2] | 16[2] | 4 | 2 | 18 | | | 1 | 2 | 1 | 1 | | - | - | - | 0[2] | -- |
| 1956 Melbourne | 62[3] | 19[3] | 6 | 4 | 18[3] | | | 2 | 3 | 3 | 6 | | - | 1[3] | - | 1[3] | 31 |
| 1960 Rome | 49[4] | 12[4] | 4 | 2 | 18 | | | 4 | | 2 | 7 | | 1[4] | - | 1[4] | 2[4] | 20 |
| 1964 Tokyo | 41 | 6 | 4[5] | 4[5] | 18 | | | 5 | | 1[5] | 6 | | - | 1[5] | - | 1[5] | 30 |
| 1968 Mexico City | 20[6] | | | | 18[6] | | | | | | 2[6] | | 1[6] | - | - | 1[6] | 29 |
| 1972 Munich | 25 | 5 | 2[7] | | 18 | | | | | 1[7] | 2[7] | | - | 1[7] | - | 1[7] | 33 |
| 1976 Montreal | 24 | 2[8] | 2 | | 16 | | | | | 2 | 2 | | - | - | 1[8] | 1[8] | 37 |
| 1980 Moscow | | | | | did not participate | | | | | | | | | | | | |
| 1984 Los Angeles | 29 | 3[9] | 4[9] | | 16 | 6[9] | | | | | 2[9] | | 1[9] | - | - | 1[9] | 25 |
| 1988 Seoul | 31 | 7 | 2[10] | | 16 | 2[10] | 1[10] | | | | 3[10] | | - | 1[10] | 1[10] | 46 | |
| 1992 Barcelona | 27 | 4 | 4 | | 16 | 2[11] | | | | | 1[11] | | - | 1[11] | 1[11] | 54 | |
| 1996 Atlanta | 24[12] | 2[12] | 4[12] | | 16[12] | | | | 1[12] | | 1[12] | | - | - | - | 0[12] | - |
| 2000 Sydney | 27[13] | 2[13] | 4[13] | | 16[13] | 3[13] | | 1[13] | 1[13] | | | | - | - | - | 0[13] | - |
| 2004 Athens | 26[14] | 2[14] | 5[14] | | 16[14] | | | 1[14] | 2[14] | | | | - | - | - | 0[14] | - |

**31.** When you move the cursor inside the code the parts in the webpage will be selected. Inside the source code go to body class - content – mw-contect-text – wiketable

**32.** You will see the Table selected. As you can see every part in the table is selected as you move the cursor in the code, links, pictures, and data.

table.wikitable | 780 × 986

| Games | Athletes | Athletes by sport | | | | | | | | | | Medals | | | Total | Rank |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1948 London | 39[1] | 5 | 3 | 2 | 19 | | | 4 | 2 | 4 | | - | - | - | 0[1] | – |
| 1952 Helsinki | 44[2] | 16[2] | 4 | 2 | 18 | | 1 | 2 | 1 | 1 | | - | - | - | 0[2] | – |
| 1956 Melbourne | 62[3] | 19[3] | 6 | 4 | 18[3] | | 2 | 3 | 3 | 6 | | - | 1[3] | - | 1[3] | 31 |
| 1960 Rome | 49[4] | 12[4] | 4 | 2 | 18 | | 4 | | 2 | 7 | 1[4] | - | 1[4] | 2[4] | 20 |
| 1964 Tokyo | 41 | 6 | 4[5] | 4[5] | 18 | | 5 | 1[5] | 6 | | - | 1[5] | - | 1[5] | 30 |
| 1968 Mexico City | 20[6] | | | | 18[6] | | | 2[6] | | 1[6] | - | - | 1[6] | 29 |
| 1972 Munich | 25 | 5 | 2[7] | 18 | | 1[7] | 2[7] | | - | 1[7] | - | 1[7] | 33 |
| 1976 Montreal | 24 | 2[8] | 2 | 16 | | 2 | 2 | | - | - | 1[8] | 1[8] | 37 |
| 1980 Moscow | *did not participate* | | | | | | | | | | | | | | |
| 1984 Los Angeles | 29 | 3[9] | 4[9] | 16 | 6[9] | | 2[9] | 1[9] | - | - | 1[9] | 25 |
| 1988 Seoul | 31 | 7 | 2[10] | 16 | 2[10] | 1[10] | 3[10] | - | - | 1[10] | 1[10] | 46 |
| 1992 Barcelona | 27 | 4 | 4 | 16 | 2[11] | | 1[11] | - | - | 1[11] | 1[11] | 54 |
| 1996 Atlanta | 24[12] | 2[12] | 4[12] | 16[12] | | 1[12] | 1[12] | - | - | - | 0[12] | - |
| 2000 Sydney | 27[13] | 2[13] | 4[13] | 16[13] | 3[13] | 1[13] | 1[13] | - | - | - | 0[13] | - |
| 2004 Athens | 26[14] | 2[14] | 5[14] | 16[14] | | 1[14] | 2[14] | - | - | - | 0[14] | - |
| 2008 Beijing | 21 | 2 | 16 | | 1 | 2 | - | - | - | 0 | - |

**Using ImportHTML**

The Google Spreadsheet function **=importHTML("","table",N)** helps scrape a table from an HTML web page into a Google spreadsheet. Within this function:

- The URL of the target web page, and the target table element both need to be in double quotes.
- The number **N** identifies the table in the page (counting starts at 0) as the target table for data scraping.

Here is an example. Let's use Google Spreadsheets to download data about the medals won by Pakistan in the Olympics:

**33.** Open https://en.wikipedia.org/wiki/Pakistan_at_the_Olympics#cite_note-jang.com.pk.2Fthenews.2Fjul2-10 Scroll down the page to find "Medal by Games" table.

**34.** Now, open a blank Google Spreadsheet from: https://docs.google.com/spreadsheets/ (if you are not signed-in, Google will ask you to sign in).

**35.** This is the format for formulas in Google sheets: =function("url", "object" , number)

**36.** In the new Google Spreadsheet that opens, type the following function: =**importHTML("https://en.wikipedia.org/wiki/Pakistan_at_the_Olympics#cite_note-jang.com.pk.2Fthenews.2Fjul2-10","table",3)**

=importHTML("https://en.wikipedia.org/wiki/Pakistan_at_the_Olympics","table",3)

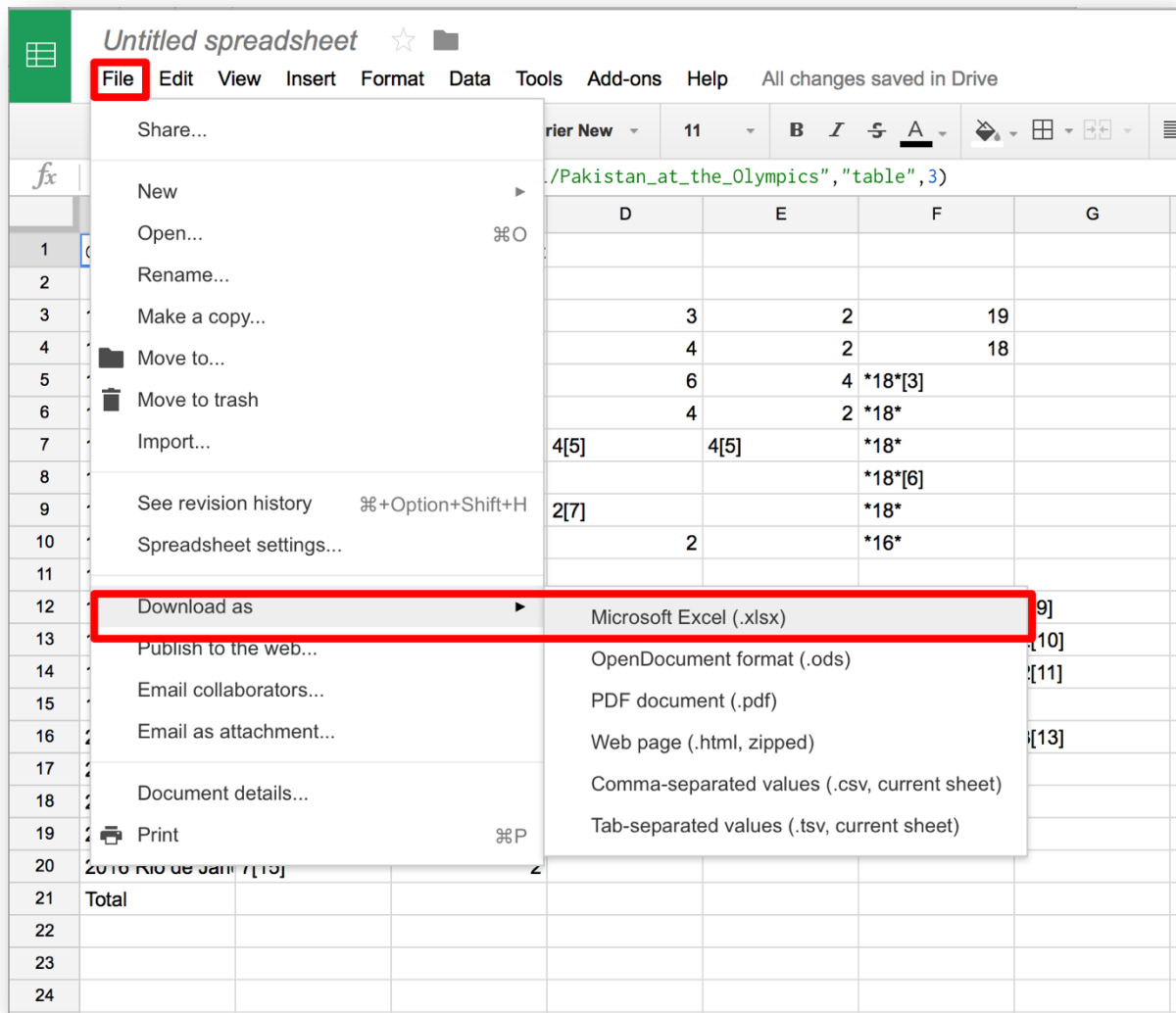| | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Games | Athletes | | Athletes by sport | | | | | | | | |
| 2 | | | | | | | | | | | | |
| 3 | 1948 London | 39[1] | 5 | 3 | 2 | 19 | | | | 4 | 2 | 4 |
| 4 | 1952 Helsinki | 44[2] | 16[2] | 4 | 2 | 18 | | | 1 | 2 | 1 | 1 |
| 5 | 1956 Melbourne | 62[3] | 19[3] | 6 | 4 | *18*[3] | | | 2 | 3 | 3 | 6 |
| 6 | 1960 Rome | 49[4] | 12[4] | 4 | 2 | *18* | | | 4 | | 2 | *7* |
| 7 | 1964 Tokyo | 41 | 6 | 4[5] | 4[5] | *18* | | | 5 | 1[5] | | 6 |
| 8 | 1968 Mexico City | 20[6] | | | | *18*[6] | | | | | 2[6] | |
| 9 | 1972 Munich | 25 | 5 | 2[7] | | *18* | | | | 1[7] | 2[7] | |
| 10 | 1976 Montreal | 24 | 2[8] | | 2 | *16* | | | | | 2 | 2 |
| 11 | 1980 Moscow | *did not participate* | | | | | | | | | | |
| 12 | 1984 Los Angeles | 29 | 3[9] | 4[9] | | *16* | 6[9] | | | | 2[9] | |
| 13 | 1988 Seoul | 31 | 7 | *2*[10] | | 16 | 2[10] | 1[10] | | | 3[10] | |
| 14 | 1992 Barcelona | 27 | 4 | 4 | | *16* | 2[11] | | | | 1[11] | |
| 15 | 1996 Atlanta | 24[12] | 2[12] | 4[12] | | 16[12] | | | | 1[12] | 1[12] | |
| 16 | 2000 Sydney | 27[13] | 2[13] | 4[13] | | 16[13] | 3[13] | | 1[13] | 1[13] | | |
| 17 | 2004 Athens | 26[14] | 2[14] | 5[14] | | 16[14] | | | 1[14] | 2[14] | | |
| 18 | 2008 Beijing | 21 | 2 | | | 16 | | | 1 | 2 | | |
| 19 | 2012 London | 21 | 2 | | | 16 | | | 1 | 2 | | |
| 20 | 2016 Rio de Jan | 7[15] | 2 | | | | | | 2 | 2 | | |
| 21 | Total | | | | | | | | | | | |

**37.** If you prefer to use Excel, click on File- Download as-Microsoft Excel.

## Exercise

Scrape http://privatisation.gov.pk/?page_id=125

# Complex Scraping: References

If you would want to learn to about more complex scraping applications, refer to these resources:

- **Import.io** is a browser based web scraping tool: https://www.import.io/ By following their easy step-by-step method you select the data you want to scrape and the tool does the rest.
- **Outwit Hub** is a software package that you can use on your PC or laptop: http://www.outwit.com/products/hub/ Note that the free version can only scrape 100 rows of data.
- **Regular Expressions** is a coding language that looks for a particular pattern of words, digits, or characters that recurs and extracts data using code: http://www.regular-expressions.info/
- **Morph.io** offers pre-written, editable scripts that use code not only to scrape but also to schedule and reformat scraping with reusable libraries: https://morph.io/