

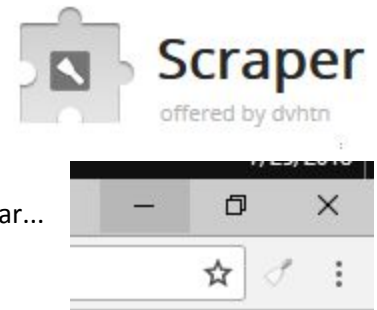
# Web Scraping with the Chrome Scraper Extension

## Set up the Scraper Extension

Install the [Scraper extension](#) in Chrome.

Test the Scraper extension

- [Highlight text and] **Right click** anywhere and choose Scrape similar...
- Highlight text and click the **icon in the toolbar** and choose Scrape similar...  
Instructions use right click, but you can always click the icon in the toolbar.



Check that the settings are correct

- Examine the window that opens, looking for **XPath** in the upper right. If it says "jQuery", click the click the Reset button at the bottom of the Scraper window and then restart Chrome.



## Using the Developer Toolbar in Chrome

- To Inspect an element, it right click on it.
- You can also press Command+Shift+C (Mac) or Ctrl+Shift+C (PC) then find the element.
- Right click on a tag and follow the menus to copy the XPath.
- In the Developer Window you can press Ctrl-F (PC) or Command+F (Mac) to open the Find toolbar and test your XPaths.

## XPath Reference

### HTML

`<p><span>this</span></p>`

`<p class="this">anything</p>`

`<p class="anything">this</p>`

`<p class="me">this</p>`

`<p class="me you">this</p>`

`<p class="you">this</p>`

`<p class="this">me</p>`

`<p id="me">this</p>`

`<div><b>Word</b>this</div>`

`<hr>Word<p>this</p>`

`<dl><dt>Word</dt><dd>this</dd></dl>`

### XPath to locate this

`p/span`

`p/@class`

`p[@class]`

`p[@class="me"]`

`p[contains(@class, "me")]`

`p[not(@class="me")]`

`p[text()='me']/@class`

`*[@id="me"]`

`div/text()`

`hr/following-sibling::p`

`dl/dt[text()='Word']/../dd`

<div><span><p><b>this</b></p></span></div> div//b

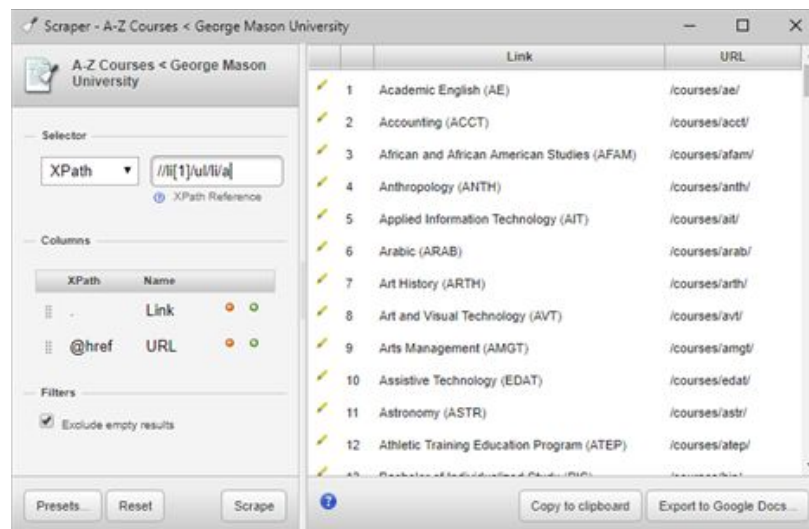
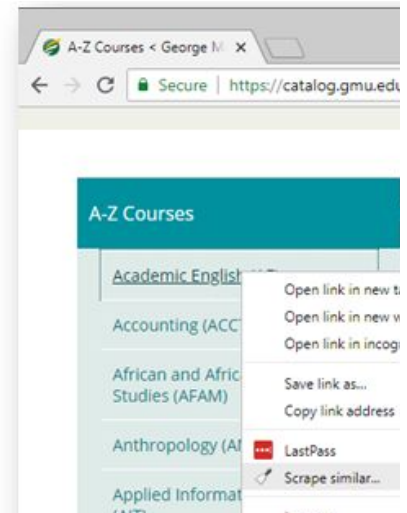
## Scrape a List

Obtain the List of Courses and URLs

### Version A

1. Go to <https://catalog.gmu.edu/courses/>
2. **Right Click** on the first link in the *left-side* list
3. Choose **Scrape similar...**
4. In the Scraper window, press **Copy to clipboard**
5. In Excel or Google Sheets, choose to **Paste** (ex Ctrl-V)

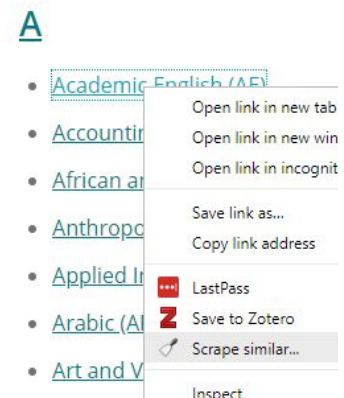
**Extra Challenge:** Change the Selector to include a tag ID instead of [1]



### Version B

1. Go to <https://catalog.gmu.edu/courses/>
2. **Right Click** on the first link in the *middle* list
3. Choose **Scrape similar...**
4. In the Scraper window, **remove the [1]** to specify that you want all letter groupings, not just A.
5. Press **Scrape**
6. Press **Copy to clipboard**
7. In Excel or Google Sheets, choose to **Paste** (ex Ctrl-V)

**Extra Challenge:** Change the Selector to include a tag ID instead of [2]



## Scrape a Table

1. Go to [https://en.wikipedia.org/wiki/List\\_of\\_countries\\_and\\_dependencies\\_by\\_population](https://en.wikipedia.org/wiki/List_of_countries_and_dependencies_by_population)
2. With your mouse, select a **whole row**, then right click and pick **Scrape Similar...**

| Rank ↕ | Country<br>(or dependent territory) ↕   | Population ↕  | Date ↕        | % of world<br>population ↕ | Source                                 |
|--------|---|---------------|---------------|----------------------------|--|
| 1      |  China <sup>[Note 2]</sup> | 1,393,110,000 | July 24, 2018 | 18.2%                      | Official population clock <sup>?</sup> |
| 2      |  India <sup>[Note 3]</sup> |               |               |                            | Official population clock <sup>?</sup> |

### Part 1

3. **Add a column** with the Green + button can click on the Name to label it Source URL
4. Click on the **XPath** area and type in the XPath you determine for the URL of the source.
  - a. It should start with \*[6] because that is the reference to the 6th column/td

### Part 2

5. Tip: Use the red minus button to **remove all columns except** for the first two (\*[1] and \*[2])
6. Use the green plus button to add columns for **Country Name**, **Note**, **Country URL**, and **Flag URL**.
7. Click on the **XPath** area and type in the XPath you determine for the each element.
  - a. The structure for Countries and Territories differ a bit, see below. Do Countries first.
  - b. A few Territories have an extra <i>: #191 Artsakh, #219 Saint-Martin, #231 Saint Helena.  
Try to include them. Hint: You may want to use a function to specify which <i> to pick.

### Country Listing:

```
<td align="left">
  <span class="flagicon" style="display:inline-block;width:25px;">
    
  </span>
  <a href="/wiki/China" title="China">China</a>
  <sup id="cite_ref-5" class="reference"><a href="#cite_note-5">[Note 2]</a></sup>
</td>
```

### Territory Listing:

```
<td align="left">
  <i>
    <span class="flagicon" style="display:inline-block;width:25px;">
      
    </span>
    <a href="/wiki/Pitcairn_Islands" title="Pitcairn Islands">Pitcairn Islands</a>
  </i>
  <i>(<a href="/wiki/United_Kingdom" title="United Kingdom">UK</a>)</i>
</td>
```

# Scrape Jobs from Indeed.com

1. Go to <https://www.indeed.com/> and search for jobs (I used “data”).

2. With your mouse, select the *whole* job listing, then right click and pick **Scrape Similar...**

- There should be a row for each job as in the screenshot. Ignore for now that only a few jobs appear.

3. Identify an XPath to obtain the Job Title

- To get you started, here is sample HTML for a title with some elements that you can use highlighted, see if you can use each one (e.g., make 4 different XPaths that work). Remember, classes are not separated, so use the contains() function to reference one.

```
<a target="_blank" id="sja1" data-tn-element="jobTitle"
class="jobtitle turnstileLink"
href="https://www.indeed.com/viewjob?jk=c74e528069aba
&from=tp-serp&tk=1c3gl72080j0f119&tk=1c3gl0j0f119&jsa=6104" title="Digital Network Analyst" rel="noopener nofollow"
onmousedown="sjomd('sja1'); clk('sja1');"
onclick="setRefineByCookie([]);sjoc('sja1',0);
convCtr('SJ')">Digital Network Analyst</a>
```

4. **Add more columns** with the Green + button and identify XPaths to locate the other job elements, such as URL, Company, Location, and Summary

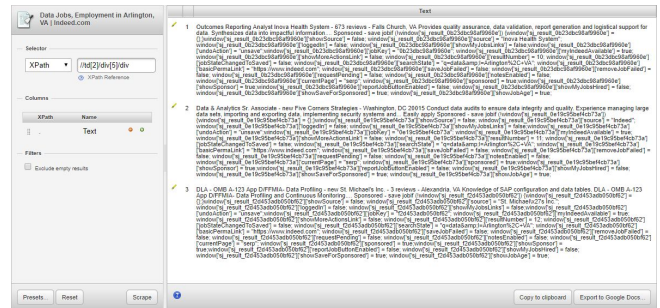
5. **Fix the XPath selector** to include all the jobs on the page by Inspecting the job listing as a whole and finding a more specific identifier for each individual job listing.

- The issue is that the Sponsored jobs are all in a special, extra div.
- Note:** You may also have to fix some XPaths that previously worked because the structure of sponsored and unsponsored jobs also differs (uses div instead of span, etc)

6. **Display more job listings** by opening the Advanced Job Search (next to the Find button). Scroll to the bottom and have it display 50 results, find jobs and re-scrape.

7. **Remove extra spaces** in the Company and Summary columns (to see why, feel free to copy the data to Excel first)

- The XPath function to do this is called normalize-space(). Put the path inside the ().
- You could also use the TRIM function in Excel or Find-and-Replace



# Get Country Flags from the State Dept using Google Sheets

In **Chrome**:

1. Go to <https://www.state.gov/misc/list/>
2. Right click one of the country names and click “**Scrape similar...**”
3. *If you want to practice XPath*, modify the Selector to
  - a. Show all letters groups, not just A’s
  - b. Use another attribute to refer to the div instead of [2]
  - c. Avoid including other links like the individual letters

*Note:* you can use not() as in a[not(@id)] to find a link without an id attribute
4. Re-scrape if you made changes, then press **Copy to clipboard**
5. **Paste** into a Google Sheet so that the data is in columns A and B
6. Think ahead and create an XPATH to get the contents of the Flag image’s **src** attribute
  - a. Go to an individual country’s page and **Inspect** the first image, which is of their flag. For more ideas, use the Scraper Extension and see what is put in the Selector box.
  - b. Test it using the Developer Windows Find (Ctrl-F) or the Scraper Extension

In **Google Sheets**:

7. In column C, use the **IMPORTXML** function to fetch the image url for the flag.
  - a. Refer to the URL of the country information
  - b. Be sure to put the XPath you created in step 6 in quotes
8. In Column D, use **CONCAT** or **&** to add “<http://www.state.gov>” to the value in column C
  - a. For a few countries, the URL already has the domain. If you wish, try to fix this.  
Hint: You can check if first/left character is a slash and thus needs the domain using LEFT
9. In Column E, use **IMAGE** to get the image from the URL in column D
  - a. Tip: To make the images bigger, make the rows taller

## “Answers” - Scrape a Table

7. Here are some ways to get the elements of the Country column:

|                                      |                                    |
|--------------------------------------|------------------------------------|
| <code>*[2]//@title</code>            | Country or Territory               |
| <code>*[2]//a[@title]</code>         | Country or Territory (for sorting) |
| <code>*[2]/i[last()]/a</code>        | County (if Dependent)              |
| <code>*[2]/i[last()]/a/@title</code> | County (if Dependent) - Full       |
| <code>*[2]//a[@title]/@href</code>   | Country URL                        |
| <code>*[2]//sup</code>               | Note                               |
| <code>*[2]//img/@src</code>          | Flag URL                           |

|   |                              |
|---|------------------------------|
|  | Georgia <sup>[Note 15]</sup> |
|  | Moldova <sup>[Note 16]</sup> |
|  | Bosnia and Herzegovina       |
|  | Uruguay                      |
|  | Puerto Rico (U.S.)           |

## “Answers” - Getting Jobs from Indeed.com

3. Any of these will work:

|   |   |
|---|---|
| <code>./a</code>                              | Reference an element                            |
| <code>./@title</code>                         | Reference an attribute                          |
| <code>./*[@data-tn-element='jobTitle']</code> | Filters to tags with a specific attribute value |
| <code>./*[contains(@class,"jobtitle")]</code> | Filters to tags with part of an attribute value |

4. Here are some examples

```

./@href
./*[@class="company"]
./*[@class="location"]
./*[@class="summary"]

```

| XPath                               | Name     |   |   |
|-------------------------------------|----------|---|---|
| <code>./@title</code>               | Title    |  |  |
| <code>./@href</code>                | URL      |  |  |
| <code>./*[@class="company"]</code>  | Company  |  |  |
| <code>./*[@class="location"]</code> | Location |  |  |
| <code>./*[@class="summary"]</code>  | Summary  |  |  |

5. Because the class for each job listing is “row result clickcard”, you must use the contains function:

```
//div[contains(@class,'clickcard')]
```

7a. Remove extra spaces with

```
normalize-space(./*[@class="summary"])
```

7b. In Excel

```
=TRIM(C2)
```

## “Answers” -- Get Country Flags

3a. `//div[2]/div/div/ul/li/a`

3b. `//div[@class="l-wrap"]/ul/li/a`

3c. To clean up the list of countries, here are some options for the Selector:

```
//div[@class="l-wrap"]/ul[@class="no-bullet"]/li/a
//div/ul[@class="no-bullet"]/li/a[not(@target)]
//div[@class="l-wrap"]/ul/li/a[contains(@href,"state.gov")]
```

6. To fetch the flag image's URL, here are some options:

```
//p[1]/img/@src
//img[contains(@title,"Flag")]/@src
```

7. To Fetch the image src for the flag image in google sheets, here is an example:

```
=IMPORTXML(B2,"//p[1]/img/@src")
```

8. To do a straight concatenate, use one of the following in cell D2:

```
"http://www.state.gov"&C2
=CONCAT("http://www.state.gov", C2)
```

8a. To concatenate only if the first character is a slash, put this in cell D2













```
=IF(LEFT(C2, 1)="/", "http://www.state.gov"&C2, C2)
```

9. To Fetch the image at that URL put this in cell E2

```
=IMAGE(D2)
```

To see it in action, go to:

<https://goo.gl/nBJJU3>

| fx =IMPORTXML(B2, "//div[@class='l-wrap']/p[1]/img/@src") |                     |   |   |   |   |
|---|---------------------|---|---|---|---|
|   | A                   | B   | C   | D   | E   |
| 1   | Link                | URL   | Scraped Flag URL  | Full Flag URL   | Image   |
| 2   | Afghanistan         | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/14/57790/afgha">/img/14/57790/afgha</a>     | <a href="http://www.state.gov">http://www.state.gov</a> |  |
| 3   | Albania             | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/14/57819/allgfla">/img/14/57819/allgfla</a> | <a href="http://www.state.gov">http://www.state.gov</a> |  |
| 4   | Algeria             | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/14/58003/algeri">/img/14/58003/algeri</a>   | <a href="http://www.state.gov">http://www.state.gov</a> |  |
| 5   | Andorra             | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/14/57820/anlgfla">/img/14/57820/anlgfla</a> | <a href="http://www.state.gov">http://www.state.gov</a> |  |
| 6   | Angola              | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/12/48337/Angol">/img/12/48337/Angol</a>     | <a href="http://www.state.gov">http://www.state.gov</a> |  |
| 7   | Antigua and Barbuda | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/14/57692/antigu">/img/14/57692/antigu</a>   | <a href="http://www.state.gov">http://www.state.gov</a> |  |
| 8   | Argentina           | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/14/57696/argen">/img/14/57696/argen</a>     | <a href="http://www.state.gov">http://www.state.gov</a> |  |
| 9   | Armenia             | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/14/55802/armen">/img/14/55802/armen</a>     | <a href="http://www.state.gov">http://www.state.gov</a> |  |
| 10  | Aruba               | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/14/57974/nllgfla">/img/14/57974/nllgfla</a> | <a href="http://www.state.gov">http://www.state.gov</a> |  |
| 11  | Australia           | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/13/56062/aslgfla">/img/13/56062/aslgfla</a> | <a href="http://www.state.gov">http://www.state.gov</a> |  |
| 12  | Austria             | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/14/57821/aulgfla">/img/14/57821/aulgfla</a> | <a href="http://www.state.gov">http://www.state.gov</a> |  |
| 13  | Azerbaijan          | <a href="http://www.state.gov">http://www.state.gov</a> | <a href="http://img/14/55803/azerb">/img/14/55803/azerb</a>     | <a href="http://www.state.gov">http://www.state.gov</a> |  |



# Useful Links

See also <http://infoguides.gmu.edu/data-work/scraping>

## Selectors

XPath -- [https://www.w3schools.com/xml/xpath\\_syntax.asp](https://www.w3schools.com/xml/xpath_syntax.asp) and <https://www.w3.org/TR/xpath/all/>

CSS - [https://www.w3schools.com/cssref/css\\_selectors.asp](https://www.w3schools.com/cssref/css_selectors.asp) and <http://learnlayout.com/display.html>

JQuery - [https://www.w3schools.com/jquery/jquery\\_ref\\_selectors.asp](https://www.w3schools.com/jquery/jquery_ref_selectors.asp)

Regular Expressions - <http://infoguides.gmu.edu/data-work/regex>

JQuery vs XPath: <https://www.ibm.com/developerworks/library/x-xpathjquery/index.html>

<https://genius.com/Mat-brown-xpath-is-actually-pretty-useful-once-it-stops-being-confusing-annotated>

CSS vs XPath: [https://en.wikibooks.org/wiki/XPath/CSS\\_Equivalents](https://en.wikibooks.org/wiki/XPath/CSS_Equivalents)

## XPath Hierarchies

<http://dh.obdurodon.org/introduction-xpath.xhtml>

<http://dh.newtfire.org/explainXPath.html>

<https://dpastov.blogspot.com/2015/10/preceding-sibling-and-following-signling-xpath.html>

## Reference

Tester: <https://extendsclass.com/xpath-tester.html> or <http://xpather.com/>

## Node / Path

|                  |                                 |
|------------------|---------------------------------|
| <b>element</b>   | div p a span                    |
| <b>attribute</b> | @class @href                    |
| <b>text</b>      | text() <div> <i>this</i> </div> |

|                                      |        |
|--------------------------------------|--------|
| any <i>element</i>                   | *      |
| any <i>attribute</i>                 | @*     |
| either <i>element</i> or <i>text</i> | node() |

## Axis / Direction

|                     |                                 |
|---------------------|---------------------------------|
| self::              | .                               |
| parent::            | ..                              |
| child::             | / (also for attributes or text) |
| descendant::        | //                              |
| ancestor::          |                                 |
| attribute::         | @                               |
| preceding-sibling:: |                                 |
| following-sibling:: |                                 |
| preceding::         |                                 |
| following::         |                                 |

## Predicates / Filter

|   |
|---|
| div[1]                                  |
| div[b]                                  |
| div[@id = "me"]                         |
| div[@id <b>eq</b> "me"]                 |
| div[@class != "me"]                     |
| div[@class <b>ne</b> "me"]              |
| div[@class = ("you", "me")]             |
| div[@class="you" <b>or</b> @class="me"] |
| div[@class="you"]   div[@class="me"]    |
| div[@class="you" <b>and</b> @id="me"]   |

## Functions

|                                  |
|----------------------------------|
| contains(@class, "me")           |
| starts-with(@class, "me")        |
| ends-with(@class, "me")          |
| matches(@id, "[0-9][a-f]+", "i") |
| td[last()]                       |
| not(@id)                         |
| concat(" ", @class, " ")         |
| normalize-space(/div/p)          |



