# XPath Navigation

## WEB SCRAPING IN PYTHON

**Thomas Laetsch**
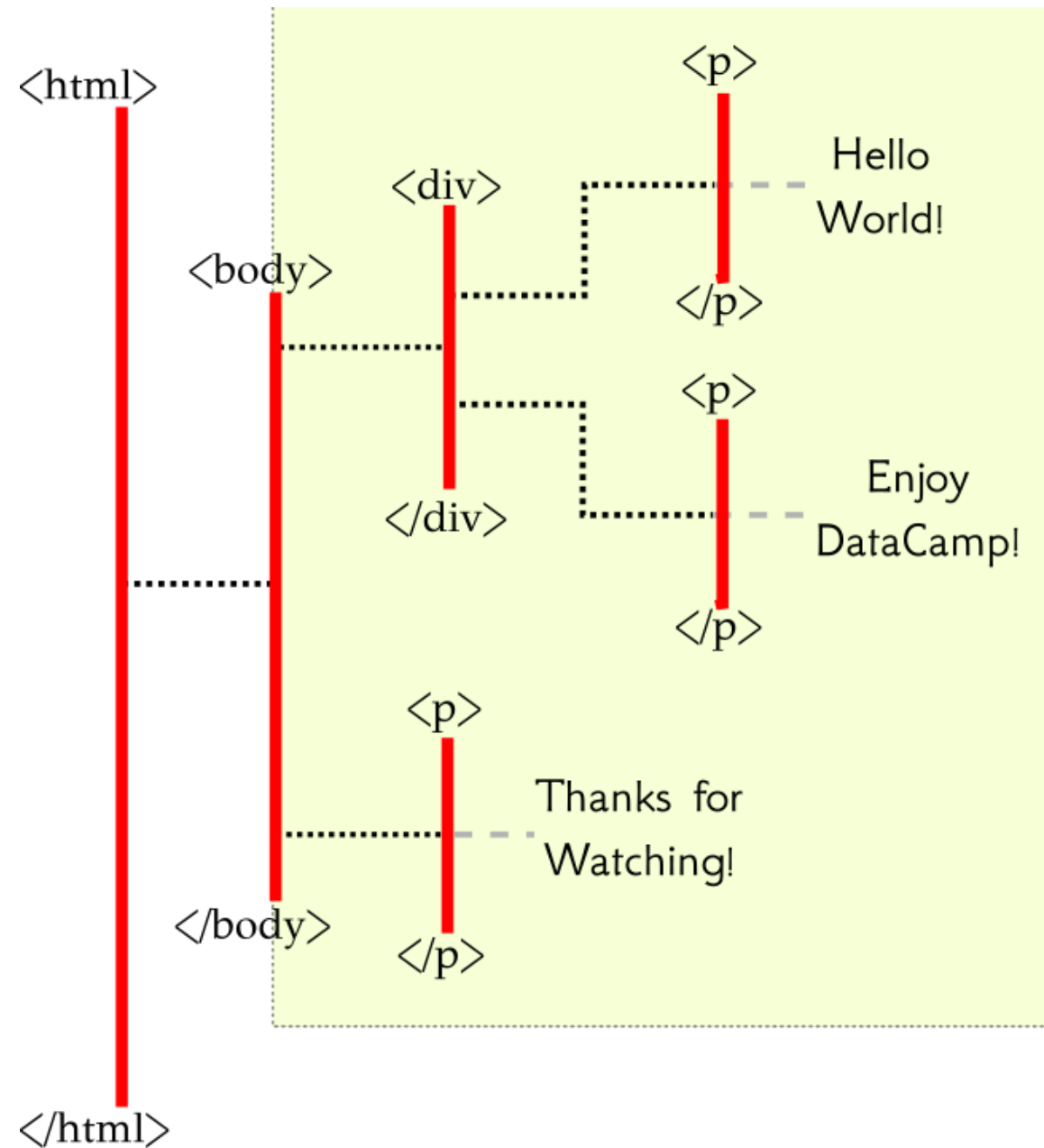Data Scientist, NYU

datacamp

# Slashes and Brackets

- Single forward slash / ==looks forward **one** generation==

- Double forward slash // looks forward **all** future generations

- Square brackets [] help ==narrow in on specific elements==

but it turns out these brackets aren't always necessary nor desired.

Note: Since "/html/body/*" selects all elements one generation below the body element without concern of the tag type, it selects all children of the body element. On the other hand, "/html/body//*" selects all elements from all future generations of the body element (that is, all descendants of the body) regardless of tag type.

# To Bracket or not to Bracket



```
xpath = '/html/body'
```

```
xpath = '/html[1]/body[1]'
```

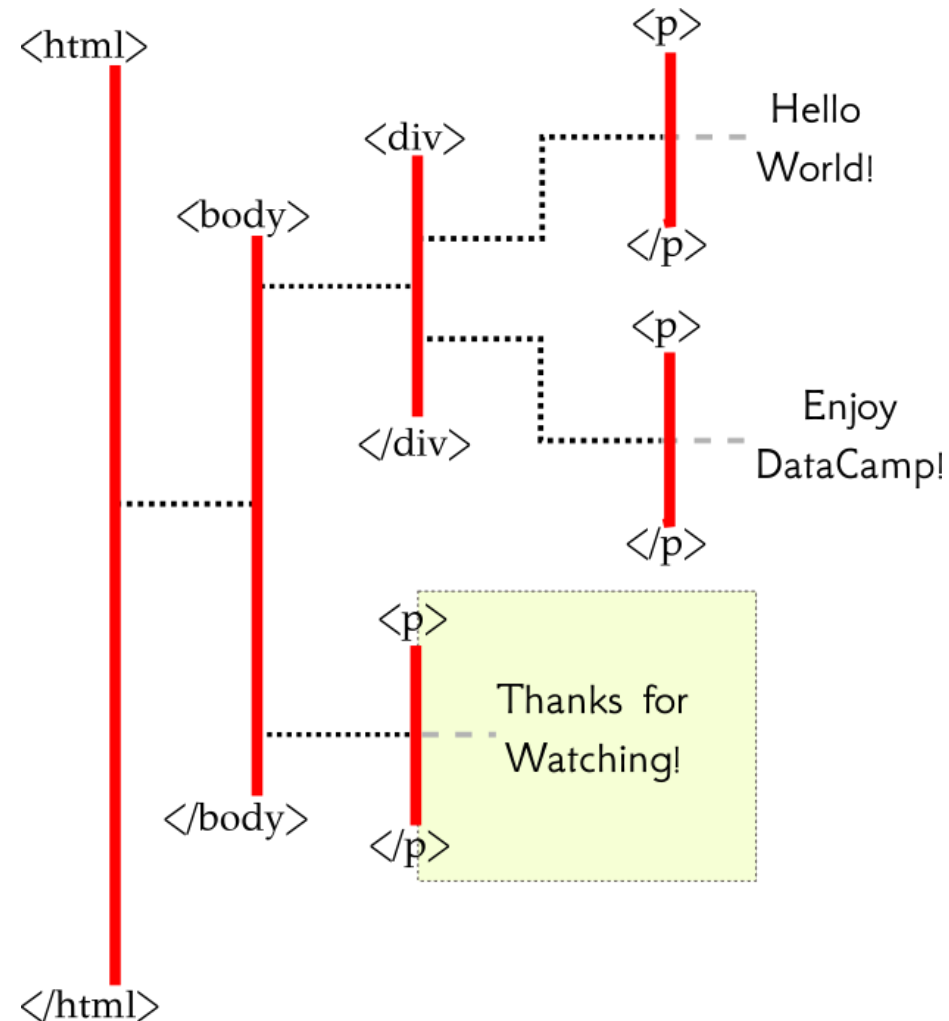- Give the same selection

3. To Bracket or not to Bracket
For example, the first XPath expression without brackets and the second with brackets lead us to the same element the body element since there is only one html element at the root level, and one body element which is a child of that html element.
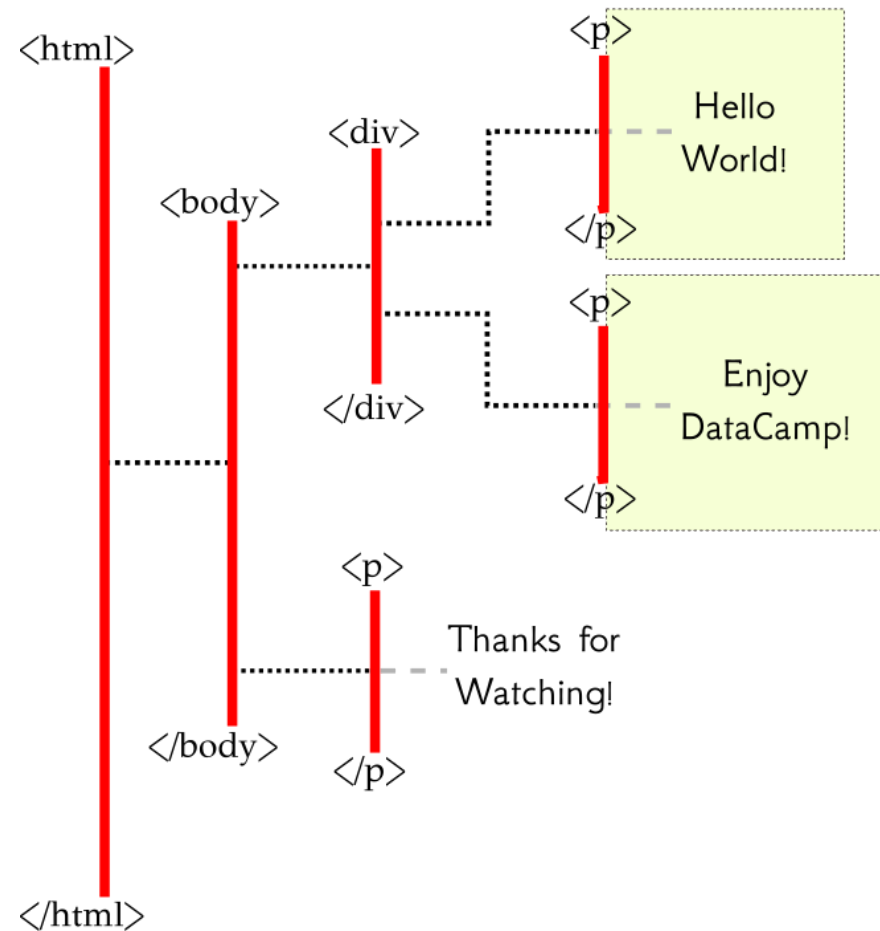
# A Body of P

4. A Body of P
As another example, the XPath string here leads us to the single paragraph element which is one generation below the body element. Do you understand why there is only this one element selected? It's because there is only one paragraph element which is a child of the body element!
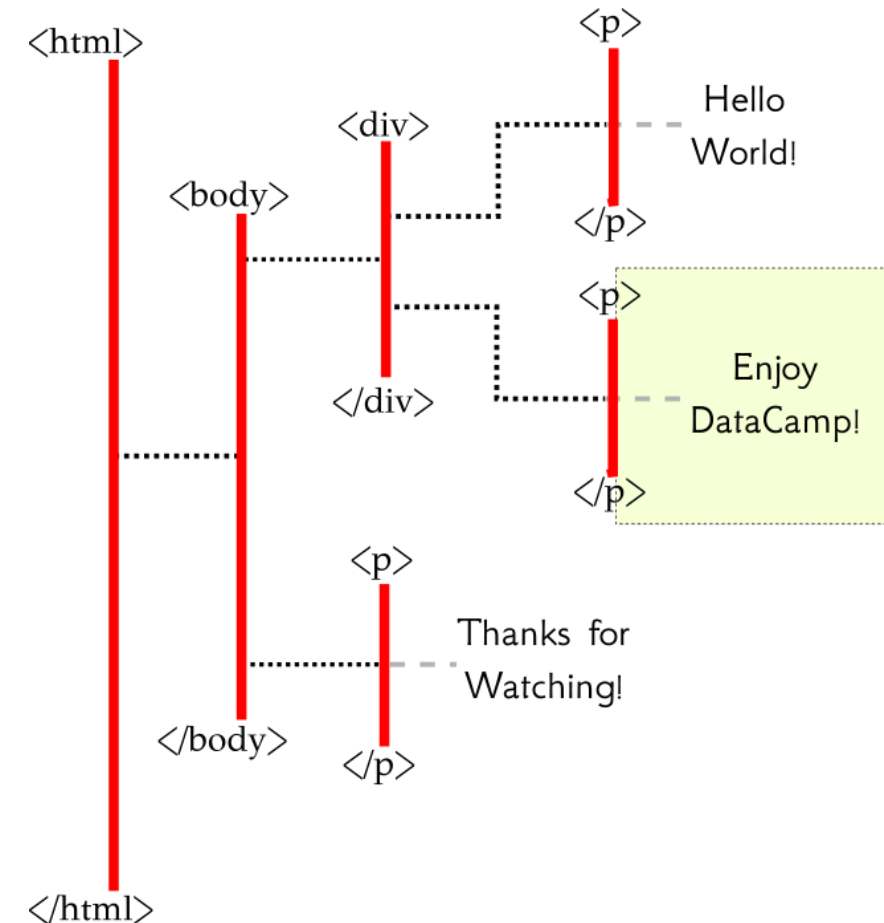
<html>

<div>

<body>

Hello
World!

</p>

<p>

Enjoy
DataCamp!

</div>

</p>

<p>

Thanks for
Watching!

</body>

</p>

</html>

# The Birds and the Ps

xpath = '/html/body/div/p'

xpath = '/html/body/div/p[2]'

# Double Slashing the Brackets

```
xpath = '//p'
```

```
xpath = '//p[1]'
```

<html>
<body>
<div>
<p>
Hello World!
</p>
<p>
Enjoy DataCamp!
</p>
</div>
<p>
Thanks for Watching!
</p>
</body>
</html>

... Honestly, I don't often mix double forward-slashes and brackets filled with numbers. We'll see in later slides there are other, more interesting ways to use brackets to select elements.

# The Wildcard

xpath = `'/html/body/*'`

- The asterisks * is the "wildcard"



7. The Wildcard
The asterisks indicates we want to ignore tag type. For example, in this expression, we are directed to both children of the body element, regardless that one is a div element and one is a paragraph element.

# Xposé

## WEB SCRAPING IN PYTHON

datacamp

# Off the Beaten XPath

## WEB SCRAPING IN PYTHON

Thomas Laetsch
Data Scientist, NYU

# (At)tribute

- @ represents "attribute"
  - @class &larr;
  - @id &larr;
  - @href &larr;

2. (At)tribute
Let's start by first pointing out that in XPath notation, the @ symbol is used to distinguish attributes. For example, if we see @class, @id, or @href, in the XPath expression, then it is referring to a class attribute, id attribute, or href attribute, respectively.

# Brackets and Attributes

3. Brackets and Attributes
We saw before that square brackets can be used in xpath syntax to hone in on a specific element or elements based on their order within a given generation. We can also include other information within square brackets to select specific elements.

<html>

<div id="uid">

<p class="class-1">

Hello
World!

</p>

<body>

<p class="class-2">

Enjoy
DataCamp!

</p>

</div>

<p class="class-1">

Thanks for
Watching!

</body>

</p>

</html>

# Brackets and Attributes

```
<html>
    <body>
        <div id="uid">
            <p class="class-1">
                Hello
                World!
            </p>
        </div>
            <p class="class-2">
                Enjoy
                DataCamp!
            </p>
        </div>
            <p class="class-1 class-2">
                Thanks for
                Watching!
            </p>
    </body>
</html>
```

```
xpath = '//p[@class="class-1"]'
```

4. Brackets and Attributesßor example, the XPath string here will direct to all paragraph elements from //p, and then reduce down to all those whose class attribute is equal to "class-1". Note that we have the class attribute in quotations. Now, my convention is to use single quotes to define the XPath string, and double quotes as needed within the XPath expression itself.
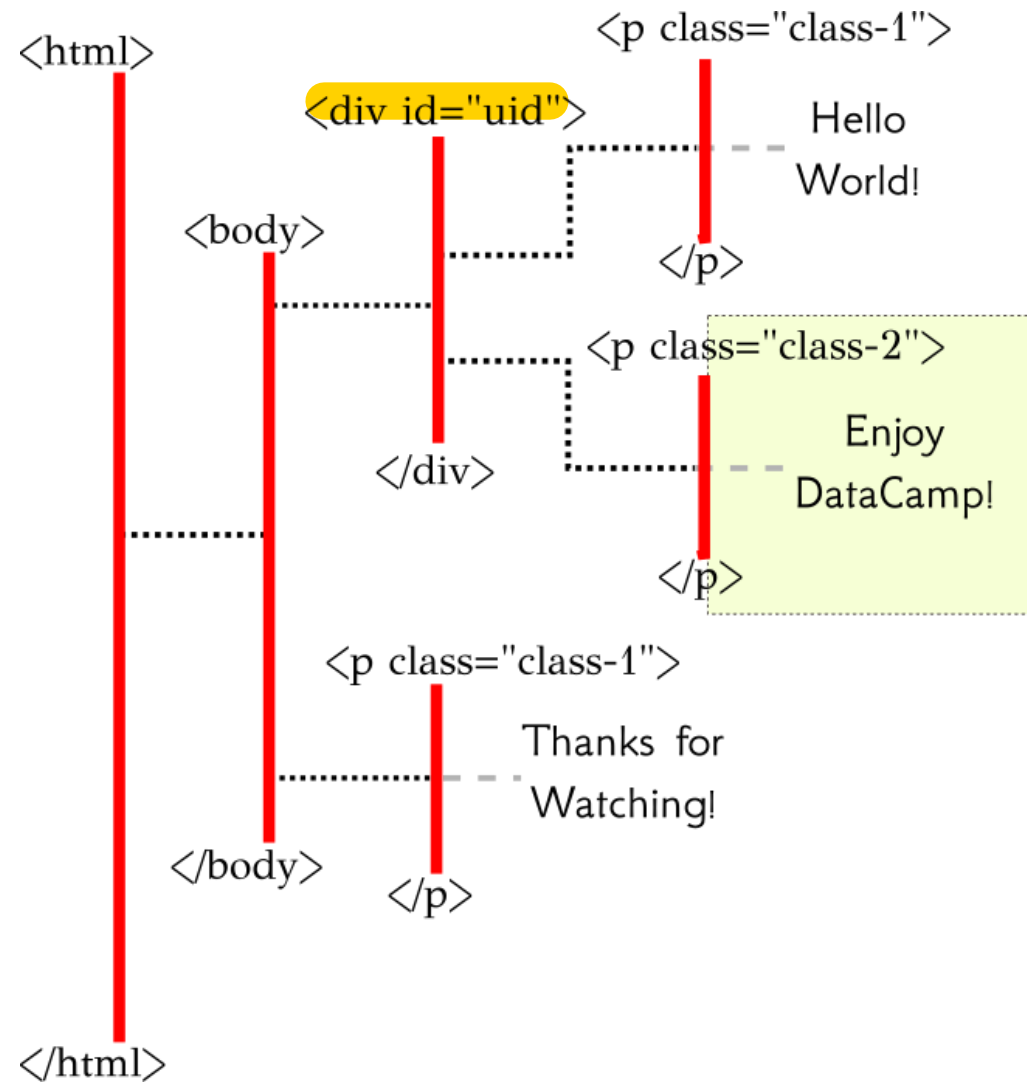
# Brackets and Attributes



```
xpath = '//*[@id="uid"]'
```

5. Brackets and Attributes
As another example, we could use the expression here with the wildcard character to first direct to all elements within the HTML document, and reduce down to whichever element has "uid" as its id attribute.

# Brackets and Attributes

<html>

<div id="uid">

<p class="class-1">

Hello
World!

</p>

<body>

<p class="class-2">

Enjoy
DataCamp!

</div>

</p>

<p class="class-1">

Thanks for
Watching!

</body>

</p>

</html>

```
xpath = '//div[@id="uid"]/p[2]'
```

6. Brackets and Attributes
Or, we could combine what we know to first navigate to the div element with id attribute equal to "uid", and then collect the second paragraph child of that div element.

# Content with Contains

Xpath Contains Notation:

**contains( @attri-name, "string-expr" )**

7. Content with Contains
A useful tool we can include within our square-bracketed expression is the "contains" function. The format of the "contains" function is given abstractly here, with the left argument containing the attribute name (including the at symbol), and the right argument is the string expression we want to search for within the given attribute. What it does is searches the attributes of that specific attribute name and matches with those where the string expression is a sub-string of the full attribute.

# Contain This

```
xpath = '//*[contains(@class,"class-1")]'
```

✔️ `<p class="class-1"> ... </p>`

✔️ `<div class="class-1 class-2"> ... </div>`

✔️ `<p class="class-12"> ... </p>`

8. Contain This
To make this clearer, let's look at an example. The expression here will choose all elements in which the string "class-1" is contained as a substring within the full class attribute; this even includes the third paragraph belonging to class-12, because class-1 is a substring of class-12.

# Contain This

```
xpath = '//*[@class="class-1"]'
```

✓ `<p class="class-1"> ... </p>`

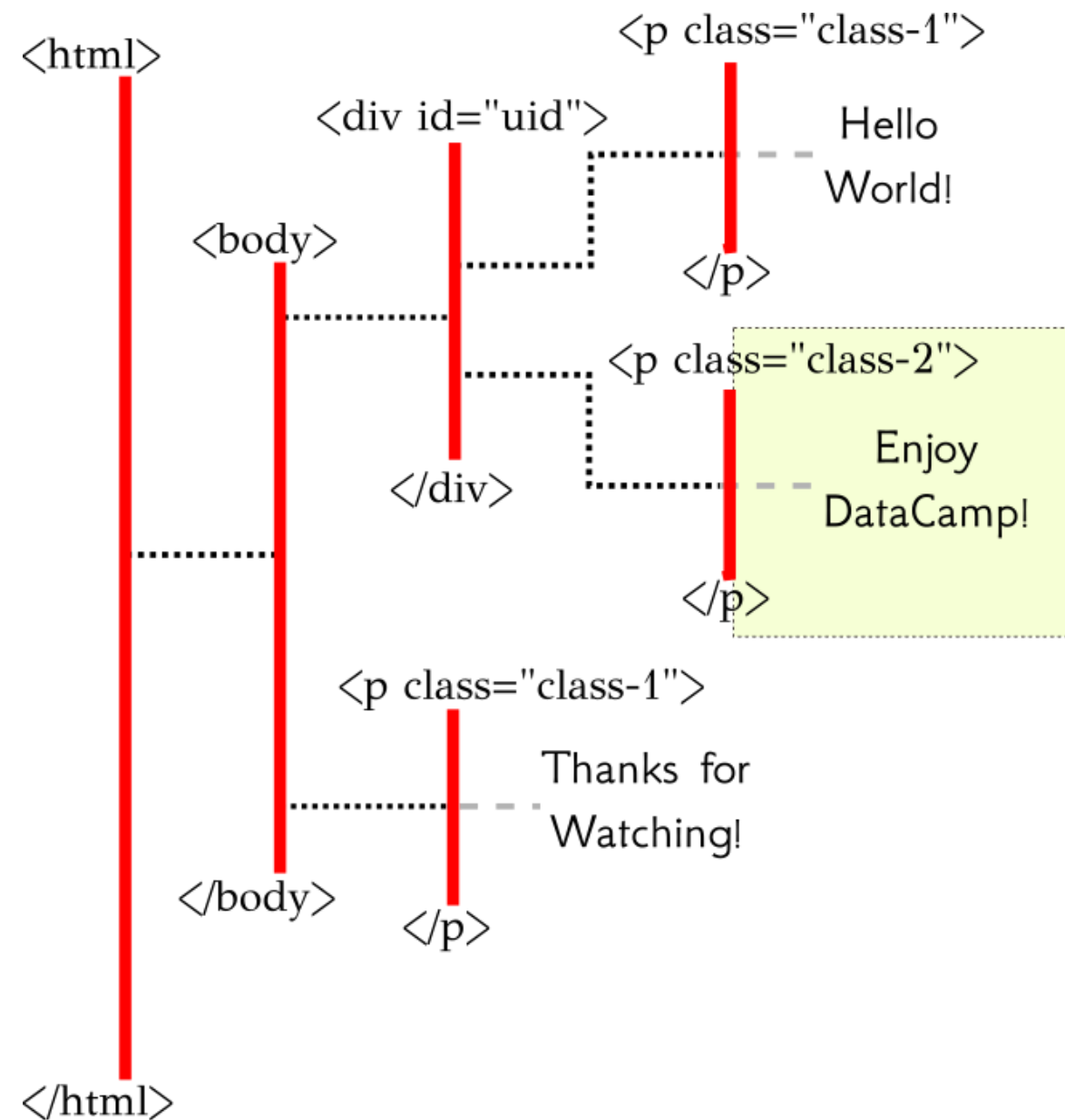✗ `<div class="class-1 class-2"> ... </div>`

✗ `<p class="class-12"> ... </p>`

9. Contain This
The last example differs from what we've seen so far since the expression here without the contains function only matches elements whose entire class attribute is equal to "class-1".

# Get Classy


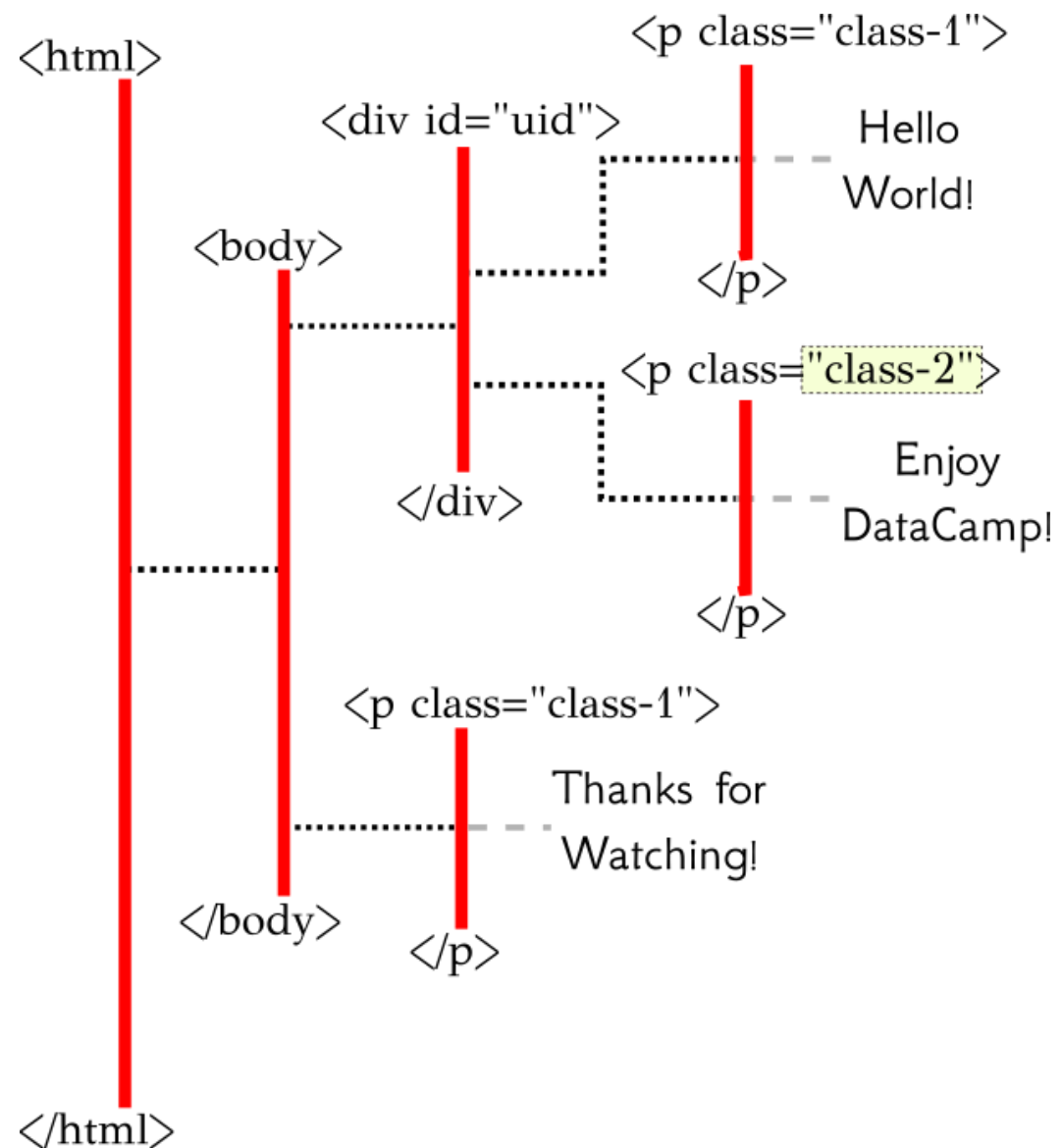
```
xpath = '/html/body/div/p[2]'
```

10. Get Classy
Now, let's consider how to direct to the attribute information itself. To do so, we first create an XPath expression to the element or elements we want to pull out some attribute information from. Say, we would like to direct to the class attribute of this highlighted paragraph element. We already know how to direct to the highlighted area.

# Get Classy

<p class="class-1">

&lt;html&gt;

&lt;div id="uid"&gt;

Hello
World!

&lt;/p&gt;

&lt;body&gt;

<p class="class-2">

Enjoy
DataCamp!

&lt;/div&gt;

&lt;/p&gt;

<p class="class-1">

Thanks for
Watching!

&lt;/body&gt;

&lt;/p&gt;

&lt;/html&gt;

```
xpath = '/html/body/div/p[2]/@class'
```

**11. Get Classy**

To direct to the attribute itself, we take the XPath, follow it by a forward slash, and follow that by the @ symbol connected to the attribute name of interest, in this case, class. As a quick note. If we were instead to use a **double forward slash before the @ symbol** with the attribute name, we would not only direct to the attribute of the elements selected in the XPath, but also all of those attributes in their future generations too.

datacamp

# End of the Path

WEB SCRAPING IN PYTHON

# Introduction to the scrapy Selector

## WEB SCRAPING IN PYTHON

**Thomas Laetsch**
Data Scientist, NYU

1. Introduction to the scrapy Selector
In this lesson we will begin to familiarize ourselves with scrapy's Selector object, so named because it is the scrapy object used to select portions of the HTML using XPath (or a so-called CSS Locator -- something we'll learn about later). Some of what we see in this lesson may seem a bit gnarly. But keep in mind that once we master this step, we will have already learned how to use the main parsing tool scrapy offers, letting us actually read in an HTML document and access the inner elements we want.

datacamp

# Setting up a Selector

```python
from scrapy import Selector
```

```python
html = '''
<html>
  <body>
    <div class="hello datacamp">
      <p>Hello World!</p>
    </div>
    <p>Enjoy DataCamp!</p>
  </body>
</html>
'''
```

2. Setting up a Selector
Through this lesson, we will be using the Selector we set up in this slide as our running example. We will import Selector from scrapy. We have made a string of HTML, which we pass to the Selector as text, creating a selector object "sel", which is the object we'll be learning to use. It will become clearer as we move along, but we can think of the Selector "sel" as having "selected" the entire HTML document. Before moving on, let's note that the html has two paragraph elements, the first saying "Hello World!"; the second saying "Enjoy DataCamp!".

```python
sel = Selector( text = html )
```

- Created a scrapy Selector object using a string with the html code

- The selector sel has selected the **entire** html document

# Selecting Selectors

- We can use the `xpath` call within a `Selector` to create new `Selector` s of specific pieces of the html code

- The return is a `SelectorList` of `Selector` objects

```
sel.xpath("//p")
# outputs the SelectorList:
[<Selector xpath='//p' data='<p>Hello World!</p>'>,
 <Selector xpath='//p' data='<p>Enjoy DataCamp!</p>'>]
```

3. Selecting Selectors
To put to work all the XPath notation we've built-up, we can call the xpath Selector method to create new Selector objects, selecting the pieces of HTML we are interested in. When doing so, the return value is a SelectorList (a list with some "scrapy extras") containing new Selector objects. For example, if we use an xpath to select all paragraph objects from our running example, we will have a SelectorList of two Selector objects, one for each of the paragraphs.

# Extracting Data from a SelectorList

- Use the `extract()` method

```
>>> sel.xpath("//p")
out: [<Selector xpath='//p' data='<p>Hello World!</p>'>,
      <Selector xpath='//p' data='<p>Enjoy DataCamp!</p>'>]
```

```
>>> sel.xpath("//p").extract()
out: [ '<p>Hello World!</p>',
       '<p>Enjoy DataCamp!</p>' ]
```

- We can use `extract_first()` to get the first element of the list

```
>>> sel.xpath("//p").extract_first()
out: '<p>Hello World!</p>'
```

4. Extracting Data from a SelectorList Selectors and SelectorLists are nice, but at the end of the day, we really want to access the data inside the Selector or SelectorList. We can do this easily by using the "extract" method. For a SelectorList, by calling the extract method, we are left with a list of strings, where each of the strings is the data from the Selectors which were originally in the SelectorList. If we only want the data (as a string) from the first element of a SelectorList, we can also call the quick extract_first method.

XPath Chaining
Selector and SelectorList objects allow for chaining when using the xpath method. What this means is that you can apply the xpath method over once you've already applied it. For example, if sel is the name of our Selector, then

# Extracting Data from a Selector

```
ps = sel.xpath('//p')

second_p = ps[1]
```

```
second_p.extract()

out: '<p>Enjoy DataCamp!</p>'
```

5. Extracting Data from a Selector
Although extract_first is convenient if you want the first piece of data in a SelectorList, we could grab the data from any other Selector within our SelectorList. To give an example of this, lets say we create a SelectorList named ps, then take the second Selector in the list (remembering that python indexes lists starting at 0, so we use the index 1 rather than 2). Then we apply the extract function to this Selector. Note that a Selector only has one piece of data, so extract leaves us with the string of this data (rather than a list of strings as was the case with SelectorLists).

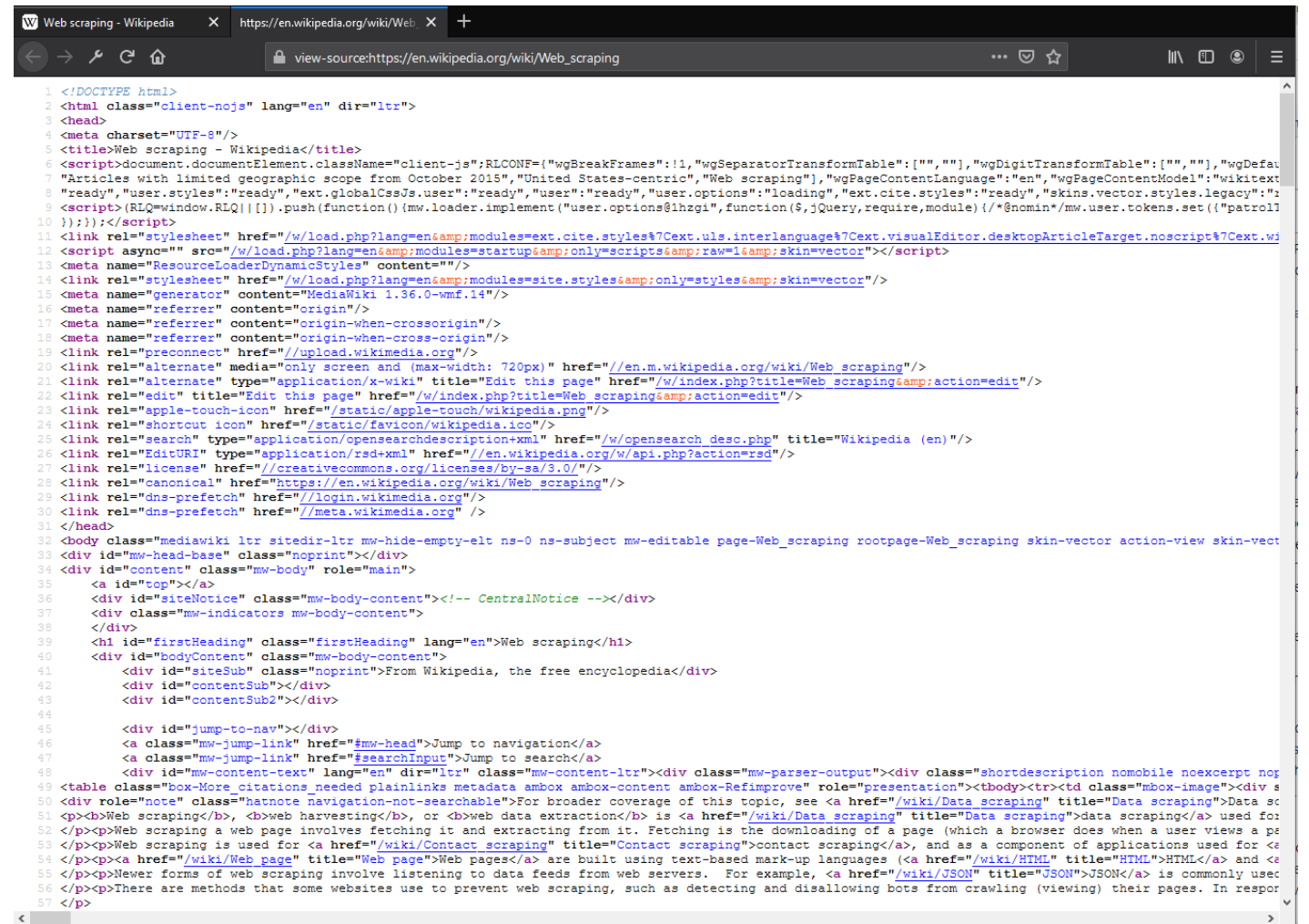# Select This Course!

## WEB SCRAPING IN PYTHON
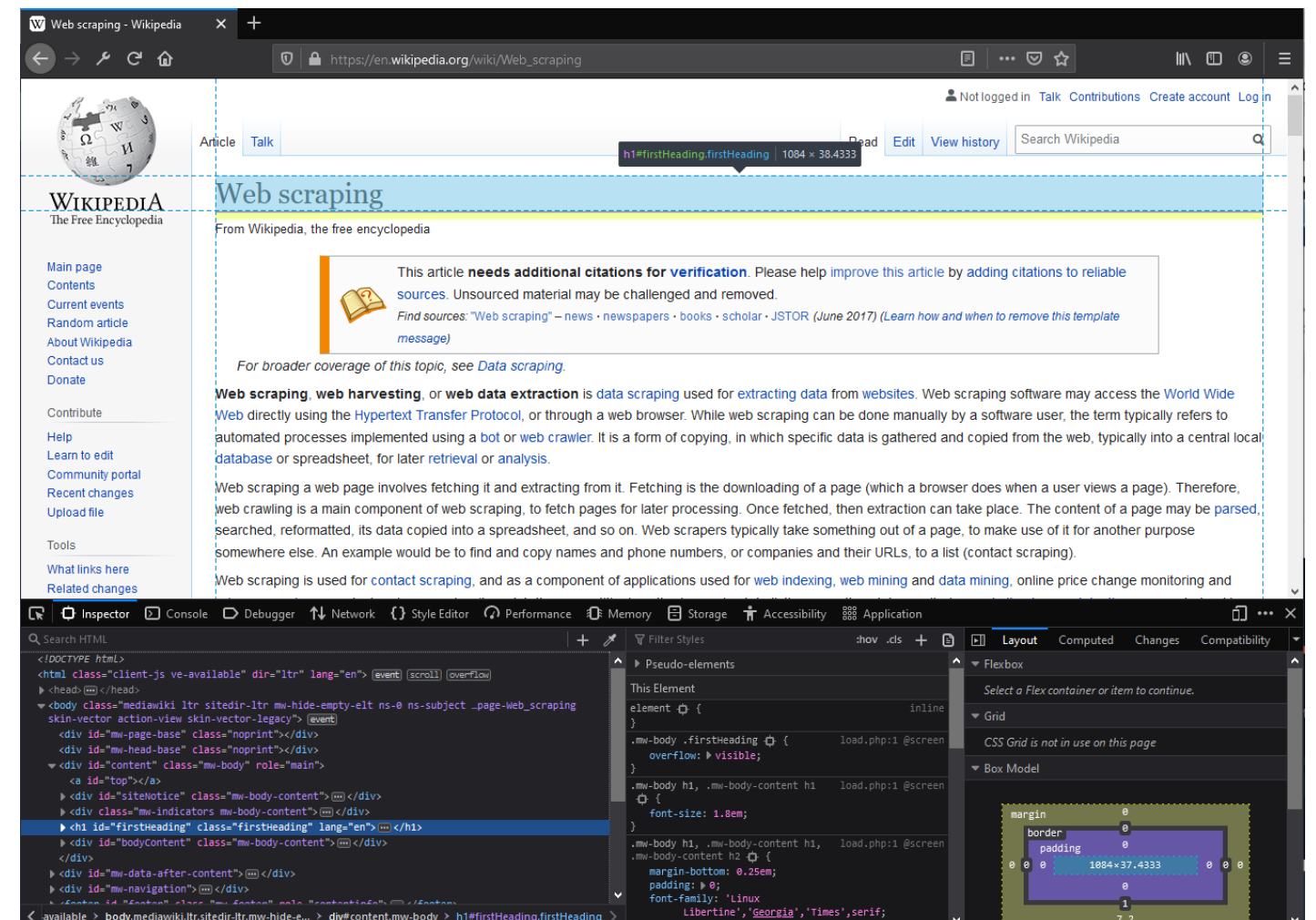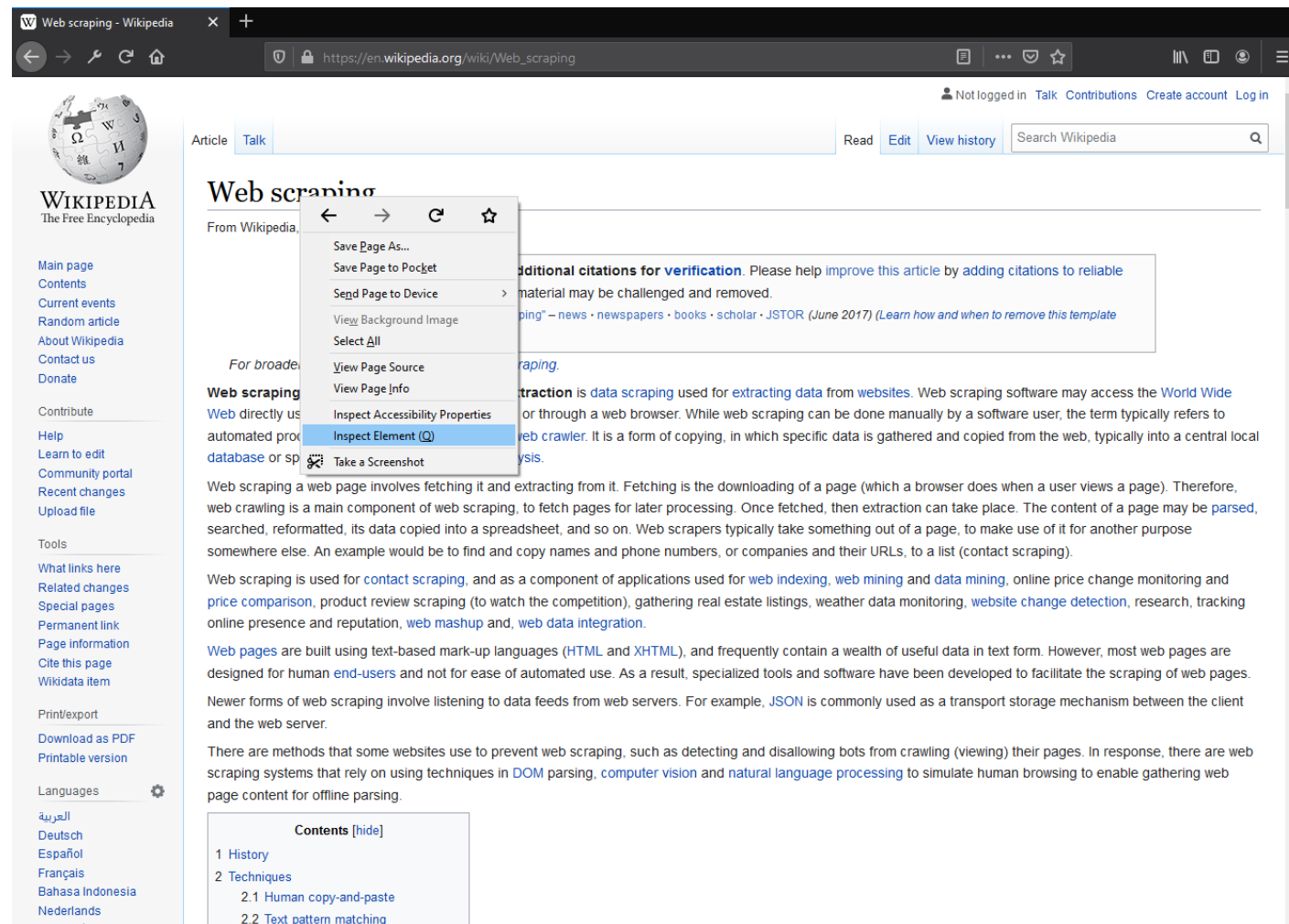
# "Source" = HTML Code

# Inspecting Elements



## 3. Inspecting Elements
Another useful tool which Firefox provides (as does Chrome, and I'm sure many other browsers), is the ability to "inspect an element". What this means is that you can select an element on the website and ask to be directed to the actual HTML Code for that specific element. In Firefox, you let your mouse hover over the element you're interested in, then right-click (or Control+Click on a Mac), and select "Inspect Element". This opens a second pane in the browser and highlights the HTML for that specific element. This is extremely useful when trying to figure out characteristics of an element which you might want to use in your scraper.

# HTML text to Selector

```python
from scrapy import Selector
```

```python
import requests

url = 'https://en.wikipedia.org/wiki/Web_scraping'
html = requests.get( url ).content
```

```python
sel = Selector( text = html )
```

4. HTML text to Selector
The final piece of this lesson is to give you a quick understanding of how we get the HTML text into a Selector object. The way we do this is using the python requests library. We won't delve into requests more than this short introduction here. In fact, in the last chapter, you will learn how to do everything within scrapy itself, without using requests. But this is still a nice piece of information to carry with you. In this example we will get the HTML from the DataCamp course directory. The requests library makes it easy to get the HTML as a string by first passing to requests.get the url (as a string) of the site you want the HTML from, then looking at the content, as we have done here. You'll notice that we assigned the content to the variable html to pass to the Selector.

# You Know Our Secrets

## WEB SCRAPING IN PYTHON