

## Lab 16

# BeautifulSoup

**Lab Objective:** *Learn how to load HTML documents into BeautifulSoup and navigate the resulting BeautifulSoup object*

## HTML

HTML, or Hyper Text Markup Language is the standard markup language to create webpages. Just like XML, HTML tags describe different document content and are surrounded by angle brackets. Most tags can be combined with attributes such as `id` or `class` to help identify individual tags and make navigating the HTML tree much more simple. You can find a list of all the current HTML tags at <http://htmldog.com/reference/htmltags>. Here is an example:

```
<html>
  <body>
    <p>
      Click <a id='info' href='http://www.example.com/info'>here</a> for ↵
      more information.
    </p>
  </body>
</html>
```

The above example would output a single line

```
Click here for more information.
```

with 'here' being a clickable link to the website <http://www.example.com/info>.

This HTML code could also be written as a single line, albeit less readable, as follows:

```
<html><body><p>Click <a id='info' href='http://www.example.com/info'>here</a> for↵
more information.</p></body></html>
```

If a given tag doesn't contain any text or other tags, it could be written in a single pair of brackets as

```
<*tag_name* ... *attributes*/>
```

The HTML of a website is very easy to view. Go to your favorite website, such as <http://www.example.com/myfavoritesite>, in whatever browser you are familiar with. Once on the website, right click with the mouse pointer and look for ‘View Page Source’ or a similarly worded option. Click it and the browser will open a new browser with the HTML code for your site. Some code is easy to follow, other code not so much.

**Problem 1.** Go to the website <http://www.example.com> and open the source code. What are all the tags used? What is the value of the `type` attribute associated with the `style` tag?

## Loading HTML into BeautifulSoup

Now that we know what HTML is, we can use BeautifulSoup to create a BeautifulSoup object. BeautifulSoup is a library capable of pulling data out of HTML scripts and files, and works with a parser to provide commands to navigate and search the resulting HTML tree. Make sure the module `bs4` is installed in your Python packages. This section takes most of its material from <http://www.crummy.com/software/BeautifulSoup/bs4/doc/index.html>.

First we want a variable to store our HTML code as a string.

```
>>> doc = """
...     <html><body><p>
...     Click <a id='info' href='http://www.example.com/info'>here</a> for more ↵
...     information.
...     </p></body></html>
...     """
```

Next, import BeautifulSoup from the `bs4` module. Call `BeautifulSoup()` which takes as parameters the HTML string and an HTML parser. It returns a BeautifulSoup object, which represents the document as a nested data structure.

```
>>> from bs4 import BeautifulSoup
>>> soup = BeautifulSoup(doc, 'html.parser')
```

Including the HTML parser is optional, but you will get a warning if none is included. If that’s the case, BeautifulSoup uses the HTML parser included in Python’s standard library. Although other parsers are permitted, we have no need for them in our examples.

Once the document is stored, we can use `prettify()` to print the HTML in a readable format. This will be useful later to make sure we are getting the correct HTML from websites.

```
>>> print(soup.prettify())
<html>
```

```
<body>
  <p>
    Click <a id='info' href='http://www.example.com/info'>here</a> for more ↵
    information.
  </p>
</body>
</html>
```

**Problem 2.** Make a soup out of the following string. Then prettify it.

```
html_doc = """
<html><head><title>The Three Stooges</title></head>
<body>
<p class="title"><b>The Three Stooges</b></p>

<p class="story">Have you ever met the three stooges? Their names are
<a href="http://example.com/larry" class="stooge" id="link1">Larry</a>,
<a href="http://example.com/mo" class="stooge" id="link2">Mo</a> and
<a href="http://example.com/curly" class="stooge" id="link3">Curly</a>;
and they are really hilarious.</p>

<p class="story">...</p>
"""
```

#### NOTE

Note that the `<html>` and `<body>` tags are never actually closed. The parser used with `bs4` will automatically close these hanging tags, so don't get too stressed out by this example.

## Navigating the HTML Tree

Since `BeautifulSoup()` returns an object which acts like a nested data structure, navigating it is very intuitive. We will use the following for the rest of the section, unless otherwise specified.

```
>>> soup = BeautifulSoup(html_doc, 'html.parser')
```

where `html_doc` is the document defined in problem 2.

### By Tag Name

To navigate the HTML parse tree, simply say the name of the tag you want. The output will be the called tag plus any nested tags or text. Below are some examples.

```
>>> soup.head
```

```
<head><title>The Three Stooges</title></head>

>>> soup.title
<title>The Three Stooges</title>
```

It is even possible to continue navigation down the tree through tags contained within tags.

```
>>> soup.body.b
<b>The Three Stooges</b>
```

Notice there are three `<a>` tags. What do you think happens when you type in `soup.a`? It only gives you the *first* tag by that name.

```
>>> soup.a
<a class="stooge" href="http://example.com/larry" id="link1">Larry</a>
```

**Problem 3.** What will the following code return?

```
>>> soup.p
```

## Tag Properties

Once a tag has been selected, you can access its properties such as its name, attributes, and strings where applicable.

A tag's name is found using `.name`.

```
>>> soup_properties = BeautifulSoup('<a href="http://example.com/larry" class="stooge" id="link1">Larry</a>')
>>> tag = soup.a
>>> tag.name
'p'
```

The attributes of a tag, if it has them, are stored in a dictionary and can be accessed as such. Accessing all the tags at once can be done through `.attrs`. Individual tags values can be reached by calling the key associated with it. If you try to access a key that is not an attribute, you get `None` in return.

```
>>> tag.attrs
{'class': ['stooge'], 'href': 'http://example.com/larry', 'id': 'link1'}

>>> tag['class']
['stooge']

>>> tag['href']
'http://example.com/larry'

>>> tag['id']
'link1'
```

```
>>> print(tag['color'])
None
```

Note that `class` returns a list. This is because the `class` attribute could have more than one value. This may show up in some HTML trees, but is not very common.

If a tag contains any text, it can be accessed with `.string`.

```
>>> tag.string
u'Larry'
```

## By Family Relations

Once we have selected a tag, we have several options available to navigate up, down, and sideways through an HTML tree.

### Going Down

Tags may contain other tags or text. These are the *children* of a tag, and the call `.contents` returns the children of the parent tag.

```
>>> head_tag = soup.head
<head><title>The Three Stooges</title></head>

>>> head_tag.contents
[<title>The Three Stooges</title>]

>>> title_tag = head_tag.contents[0]
>>> title_tag
<title>The Three Stooges</title>

>>> title_tag.contents
[u'The Three Stooges']
```

Notice in the last input, the child of `title_tag` was the string `'The Three Stooges'`. If you try to use `.contents` on this string, you would get nothing in return since strings can't contain children.

### NOTE

One thing to note is the following:

```
>>> children_doc = """
...     <html><head>The Three Little Pigs</head>
...     <body>
...     <p>The first little piggy</p>
...     <p>The second little piggy</p>
...     <p>The third little piggy</p>
...     </body>
...     </html>"""
>>> pig_soup = BeautifulSoup(children_doc, 'html.parser')
>>> pig_soup.body.contents
```

```
[u'\n',
 <p>The first little piggy</p>,
 u'\n',
 <p>The second little piggy</p>,
 u'\n',
 <p>The third little piggy</p>,
 u'\n']
```

In this example, the return carriage is counted as a child of `<body>` for each carriage return used within the `<body>` tag. This will be very important when trying to navigate between *siblings*, or children of a common tag.

Instead of creating a list, you can use `.children` to create a generator of children tags. Using the previous example, we get the following (remember the extra carriage returns):

```
>>> for pig in pig_soup.body.children:
...     print(repr(pig)) #use repr() to ignore escape sequences
u'\n'
<p>The first little piggy</p>
u'\n'
<p>The second little piggy</p>
u'\n'
<p>The third little piggy</p>
u'\n'
```

There is a `.descendants` attribute which will recursively go through a tag's children, the children's children, etc. It is left to you to look at the online documentation for this attribute.

If a tag has only one child, and that child is a string, the child is available using `.string`. If a tag has one tag, and that tag has a single string as a child, then the parent tag can use `.string` to access the string as well.

```
>>> head_tag = soup.head
>>> print(head_tag)
<head><title>The Three Stooges</head></title>
>>> title_tag = head_tag.contents[0]
>>> print(title_tag)
<title>The Three Stooges</title>
>>> head_tag.string
u'The Three Stooges'
>>> title_tag.string
u'The Three Stooges'
```

If a tag contains more than one string, `.string` return `None`. However, you can use `.strings` to return a generator that iterates through all strings contained within a tag. Check the online documentation for examples.

## Going Up

Just as tags can have *children*, tags also have a *parent*. To access a tag's parent, simply use the `.parent` attribute.

```
>>> title_tag = soup.title
>>> title_tag
<title>The Three Stooges</title>
>>> title_tag.parent
<head><title>The Three Stooges</title></head>
```

The parent of a string would be the tag which contains it.

```
>>> tag = title_tag.string
>>> print(tag)
The Three Stooges

>>> tag.parent
<title>The Three Stooges</title>
```

You could use `.parents` to iterate through all parents of a given tag. Examples can be found in the online documentation.

## Going Sideways

Consider the following document, taken from the online documentation:

```
>>> sibling_soup = BeautifulSoup("<a><b>text 1</b><c>text 2</c></a>")
>>> print(sibling_soup.prettify())
<html>
  <body>
    <a>
      <b>
        text 1
      </b>
      <c>
        text 2
      </c>
    </a>
  </body>
</html>
```

Note the `<b>` and `<c>` tags are on the same level, underneath the `<a>` tag. These tags are considered *siblings*. Siblings in an HTML tree will always appear with the same indentation underneath a parent tag.

Use the attributes `.next_sibling` and `.previous_sibling` to navigate between these sibling elements. If a sibling has no next or previous sibling, these attributes are assigned `None`.

```
>>> sibling_soup.b
<b>text 1</b>

>>> sibling_soup.b.next_sibling
<c>text 2</c>

>>> sibling_soup.c.previous_sibling
<b>text 1</b>

>>> sibling_soup.c.next_sibling #<c> has no next sibling
None
```

```
>>> sibling_soup.b.string
u'text 1'

>>> print(sibling_soup.b.string.next_sibling) #text 1 and text 2 are not siblings
None
```

Recall the `pig_soup` example, how we saw extra carriage returns between the `<p>` tags.

```
>>> pig_soup.body.contents
[u'\n',
 <p>The first little piggy</p>,
 u'\n',
 <p>The second little piggy</p>,
 u'\n',
 <p>The third little piggy</p>,
 u'\n']
```

What do you expect `pig_soup.body.p.next_sibling` to return?

```
>>> pig_soup.body.p.next_sibling
u'\n'

>>> pig_soup.body.p.next_sibling.next_sibling
<p>The second little piggy</p>
```

We need to make two calls to `.next_sibling` in order to get the next `<p>` tag. Keep this in mind for future questions as you navigate across siblings.

Just as with parents and children, there are also sibling generators `.next_siblings` and `.previous_siblings` to iterate through all the siblings of a given tag. As before, check the online documentation for more information.

#### Problem 4. Answer the following questions.

1. In the Three Stooges example, what code would you use to return the following:

```
u'Mo'
```

2. How would you return the following:

```
<p class="story">...</p>
```

3. Download 'example.htm' associated with the lab into your working directory. You can go to <http://example.com> to see the site where this file originates from. The following code allows you to bring in any file into BeautifulSoup()

```
>>> example_soup = BeautifulSoup(open('example.htm'), 'html.parser')
```

This creates a `beautifulsoup` object representing the HTML source code from



the website `http://www.example.com`. Use two different methods to access the following line:

```
u'More information...'
```

## By `find()`

In actual website HTML, often there are too many tags which share a common name. It would be nice to find characteristics that might be unique for a given tag. Look back at our previous examples and think about what characteristics differentiate tags with the same name. BeautifulSoup uses `.find()` to allow you to search for a tag not only by name, but also by a specific attribute value or strings. The following examples refer back to the “Three Stooges” HTML document in problem 2.

Search by name:

```
>>> soup.find('b') #Pass in tag names, just like soup.b
<b>The Three Stooges</b>

>>> #or use the name parameter
>>> soup.find(name='a') #Still only returns the first instance
<a class="stooge" href="http://example.com/larry" id="link1">Larry</a>
```

Search by attribute:

```
>>> soup.find(id='link3') #Search by unique id attribute
<a class="stooge" href="http://example.com/curly" id="link3">Curly</a>

>>> #class is a Python keyword. Use 'class_' for the attribute key.
>>> soup.find(class_='title')
<p class="title"><b>The Three Stooges</b></p>

>>> #use the attrs parameter
>>> soup.find(attrs={'id':'link3'})
<a class="stooge" href="http://example.com/curly" id="link3">Curly</a>

>>> #combine attributes
>>> soup.find(attrs={'class':'stooge', 'href':'http://example.com/curly'})
<a class="stooge" href="http://example.com/curly" id="link3">Curly</a>

>>> soup.find(class_='stooge', href='http://example.com/curly')
<a class="stooge" href="http://example.com/curly" id="link3">Curly</a>

>>> #use True to find the first tag containing a given attribute
>>> soup.find(href=True)
<a class="stooge" href="http://example.com/larry" id="link1">Larry</a>
```

Search by string:

```
>>> soup.find(string='Mo') #Recall strings act as individual units
u'Mo'
```

```
>>> soup.find(string='Mo').parent #access the tag through the parent
<a class="stooge" href="http://example.com/mo" id="link2">Mo</a>
```

Search by combining parameters:

```
>>> soup.find('a', attrs={'id':'link2','class':'stooge'})
<a class="stooge" href="http://example.com/mo" id="link2">Mo</a>
```

**Problem 5.** Refer to the `example.htm` file. Load it using BeautifulSoup. What is the website associated with the “More information...” link? Find two methods to access the website string.

**Problem 6.** Download ‘SanDiegoWeather.htm’ and load it into BeautifulSoup. You can find the corresponding website at `http://www.wunderground.com/history/airport/KSAN/2015/1/1/DailyHistory.html?req_city=San+Diego&req_state=CA&req_statename=California&reqdb.zip=92101&reqdb.magic=1&reqdb.wmo=99999&MR=1`.

1. What is the tag which contains the date, Thursday, January 1, 2015?
2. What are the tags which contain the links ‘Previous Day’ and ‘Next Day’?
3. What is the tag which contains the number associated with the Actual Max Temperature.

[Hint: You can do a `ctrl+f` to find where the text is in the HTML, then study the tags around it.]

## By `find_all()`

Recall that when a tag appeared multiple times, calling that tag name would return the first tag of that name.

```
>>> soup.a
<a class="stooge" href="http://example.com/larry" id="link1">Larry</a>
```

If you need to get all instances of a certain tag, use the `find_all()` command.

```
>>> soup.find_all('a')
[<a class="stooge" href="http://example.com/larry" id="link1">Larry</a>,
  <a class="stooge" href="http://example.com/mo" id="link2">Mo</a>,
  <a class="stooge" href="http://example.com/curly" id="link3">Curly</a>]
```

This works with all the same arguments as the `.find()` function. You may refer to the online documentation for explicit examples.

## Advanced Techniques

The following examples are techniques that can aid you in your search for specific tags. Consider the “Three Stooges” example from before. Suppose you want to find the tag that includes the url `http://example.com/curly`. You could search for it using the following:

```
>>> soup.find(href='http://example.com/curly')
<a class="stooge" href="http://example.com/curly" id="link3">Curly</a>
```

This could be annoying to type out the whole website. Instead you could use regular expressions as follows:

```
>>> import re
>>> soup.find(href=re.compile('curly')) #find href containing 'curly' in it.
<a class="stooge" href="http://example.com/curly" id="link3">Curly</a>
```

This method can be used for tag names, attributes, and strings as well.

```
>>> soup.find(string=re.compile('^Cu')) #find string that starts with 'Cu'.
<a class="stooge" href="http://example.com/curly" id="link3">Curly</a>
```

If you want to find a tag that has an attribute with a value, but don’t care what the value is, you can use `True` and `False` in place of actual values. The following returns all tags that even have the `href` attribute:

```
>>> soup.find_all(href=True)
[<a class="stooge" href="http://example.com/larry" id="link1">Larry</a>,
  <a class="stooge" href="http://example.com/mo" id="link2">Mo</a>,
  <a class="stooge" href="http://example.com/curly" id="link3">Curly</a>]
```

**Problem 7.** Use BeautifulSoup to load the ‘Big Data dates’ file. This page is found at [www.federalreserve.gov/releases/lbr](http://www.federalreserve.gov/releases/lbr), although the actual website may include more dates. Notice all the release dates of bank data, ranging from 2003 to 2015 in the file you downloaded. Using `find_all()` and `re`, find all the links to bank data from September 30, 2003 to December 31, 2014. (Don’t include the 2015 links) There should be 46 links total that you list.