



## Web Scraping Lab

Estimated time needed: 30 minutes

### Objectives

After completing this lab you will be able to:

### Table of Contents

<ul style="list-style-type: none"><li>Beautiful Soup Object<ul style="list-style-type: none"><li>Tag<ul style="list-style-type: none"><li>Children, Parents, and Siblings</li><li>HTML Attributes</li><li>Navigable String</li></ul></li><li>Filter<ul style="list-style-type: none"><li>find All</li><li>find</li><li>HTML Attributes</li><li>Navigable String</li></ul></li><li>Downloading And Scraping The Contents Of A Web</li></ul></li></ul>
Estimated time needed: 25 min

For this lab, we are going to be using Python and several Python libraries. Some of these libraries might be installed in your lab environment or in SN Labs. Others may need to be installed by you. The cells below will install these libraries when executed.

```
In [1]: !pip install bs4
# pip install requests

Collecting bs4
  Downloading bs4-0.0.1.tar.gz (1.1 kB)
  Requirement already satisfied: beautifulsoup4 in c:\programdata\anaconda3\lib\site-packages (from bs4) (4.9.3)
  Requirement already satisfied: soupsieve>1.2, python_version >= 3.0 in c:\programdata\anaconda3\lib\site-pack
  ages (from beautifulsoup4->bs4) (2.0.1)
  Building wheel for bs4: filename=bs4-0.0.1-py3-none-any.whl size=1277 sha256=01db39fbf77c8b3ba86a310c8171af571
  fec4019a9371cbbcae5637456259
  Stored in directory: c:\users\user\appdata\local\pip\cache\wheels\73\78\21\6b124549c9bdc94f822c02b9a3a3578a6
  6984f9f67776bca
  Successfully built bs4
  Installing collected packages: bs4
  Successfully installed bs4-0.0.1

Import the required modules and functions
```

```
In [2]: from bs4 import BeautifulSoup # this module helps in web scraping.
import requests # this module helps us to download a web page.
```

### Beautiful Soup Objects

Beautiful Soup is a Python library for pulling data out of HTML and XML files, we will focus on HTML files. This is accomplished by representing the HTML as a set of objects with methods used to parse the HTML. We can navigate the HTML as a tree and/or filter out what we are looking for.

Consider the following HTML:

```
In [3]: <html>
<doctype html>
<html>
<head>
<title>Page Title</title>
</head>
<body>
<div id="boldest">Lebron James</div></body>
</html>
```

#### Lebron James

Salary: \$ 92,000,000

#### Stephen Curry

Salary: \$85,000, 000

#### Kevin Durant

Salary: \$73,200, 000

We can store it as a string in the variable HTML:

```
In [4]: HTML = "<doctype html><html><head><title></head><body><div id='boldest'>Lebron James</div></html>"
```

To parse a document, pass it into the BeautifulSoup constructor, the BeautifulSoup object, which represents the document as a nested data structure.

```
In [5]: soup = BeautifulSoup(HTML, 'html5lib')
```

First, the document is converted to Unicode, (similar to ASCII), and HTML entities are converted to Unicode characters. BeautifulSoup transforms a complex HTML document into a complex tree of Python objects. The BeautifulSoup object can create other types of objects. In this lab, we will cover BeautifulSoup and Tag objects that for the purposes of this lab are identical, and NavigableString objects.

We can use the method prettify() to display the HTML in the nested structure:

```
In [6]: print(soup.prettify())

<doctype html>
<html>
<head>
<title>
</title>
</head>
<body>
<div id="boldest">
  Lebron James
</div>
</body>
</html>
```

### Tags

Let's say we want the title of the page and the name of the top paid player we can use the Tag. The Tag object corresponds to an HTML tag in the original document, for example, the tag title.

```
In [7]: tag_object = soup.title
print(tag_object, tag_object)
```

```
tag_object: <title>Page Title</title>
```

we can see the tag type bs4.element.Tag

```
In [8]: print("tag object type:", type(tag_object))
```

tag object type: <class 'bs4.element.Tag'>

If there is more than one Tag, with the same name, the first element with that Tag name is called, this corresponds to the most paid player:

```
In [9]: tag_object = soup.h3
tag_object
```

```
Out[9]: <div id="boldest">Lebron James</div></h3>
```

Enclosed in the bold attribute b, it helps to use the tree representation. We can navigate down the tree using the child attribute to get the name.

### Children, Parents, and Siblings

As stated above the Tag object is a tree of objects we can access the child of the tag or navigate down the branch as follows:

```
Out[10]: tag_child = tag_object.b
tag_child
```

```
Out[10]: <div id="boldest">Lebron James</div>
```

You can access the parent with the parent

```
In [11]: parent_tag = tag_child.parent
parent_tag
```

```
Out[11]: <div id="boldest">Lebron James</div></h3>
```

this is identical to

```
In [12]: tag_object
```

```
Out[12]: <div id="boldest">Lebron James</div></h3>
```

tag\_object.parent is the body element.

```
In [13]: tag_object.parent
```

```
Out[13]: <body><div id="boldest">Lebron James</div><div id="h3"><p>Stephen Curry</p></div><p>Salary: $ 92,000,000 </p><div id="h3"> Stephen Curry</div></body>
```

tag\_object.sibling is the paragraph element

```
In [14]: sibling_1 = tag_object.next_sibling
sibling_1
```

```
Out[14]: <p>Salary: $ 92,000,000 </p>
```

sibling\_2 is the next element which is also a sibling of both sibling\_1 and tag\_object

```
In [15]: sibling_2 = sibling_1.next_sibling
sibling_2
```

```
Out[15]: <div id="h3"> Stephen Curry</div>
```

### Exercise: next\_sibling

Using the sibling\_2 and the method next\_sibling to find the salary of Stephen Curry:

```
In [16]: sibling_2.next_sibling
```

```
Out[16]: <p>Salary: $85,000, 000 </p>
```

Click here for the solution

### HTML Attributes

If the tag has attributes, the tag id="boldest" has an attribute id whose value is boldest. You can access a tag's attributes by treating the tag like a dictionary:

```
In [17]: tag_child['id']
```

```
Out[17]: 'boldest'
```

You can access that dictionary directly as attrs:

```
In [18]: tag_child.attrs
```

```
Out[18]: {'id': 'boldest'}
```

You can also work with Multi-valued attribute check out [] for more.

```
In [19]: tag_child.get('id')
```

```
Out[19]: 'boldest'
```

### Navigable String

A string corresponds to a bit of text or content within a tag. BeautifulSoup uses the NavigableString class to contain this text. In our HTML we can obtain the name of the first player by extracting the string of the Tag object tag\_child as follows:

```
In [20]: tag_string = tag_child.string
tag_string
```

```
Out[20]: 'Lebron James'
```

we can verify the type is Navigable String

```
In [21]: type(tag_string)
```

```
Out[21]: bs4.element.NavigableString
```

A NavigableString is just like a Python string or Unicode string, to be more precise. The main difference is that it also supports some BeautifulSoup features. We can convert it to string object in Python:

```
In [22]: unicode_string = str(tag_string)
unicode_string
```

```
Out[22]: 'Lebron James'
```

### Filter

Filters allow you to find complex patterns, the simplest filter is a string. In this section we will pass a string to a different filter method and BeautifulSoup will perform a match against that exact string. Consider the following HTML of rocket launches:

```
In [23]: <html>
<tr>
<td id="flight">Flight No</td>
<td>Launch site</td>
<td>Payload mass</td>
</tr>
<tr>
<td>1</td>
<td>Florida</td>
<td>300 kg</td>
</tr>
<tr>
<td>2</td>
<td>Texas</td>
<td>94 kg</td>
</tr>
<tr>
<td>3</td>
<td>Texas</td>
<td>80 kg</td>
</tr>
</table>

Flight No Launch site Payload mass
1 Florida 300 kg
2 Texas 94 kg
3 Texas 80 kg
```

We can store it as a string in the variable table:

```
In [24]: table = "<table><tr><td id='flight'>Flight No</td><td>Launch site</td> <td>Payload mass</td></tr><tr> <td>1</td><td>Florida</td><td>300 kg</td></tr><tr> <td>2</td><td>Texas</td><td>94 kg</td></tr><tr> <td>3</td><td>Texas</td><td>80 kg</td></tr></table>"
```

```
In [25]: table_bs = BeautifulSoup(table, 'html5lib')
```

### find All

The find\_all() method looks through a tag's descendants and retrieves all descendants that match your filters.

The Method signature for find\_all(name, attrs, recursive, string, limit, \*\*kwargs)

### Name

When we set the name parameter to a tag name, the method will extract all the tags with that name and its children.

```
In [26]: table_rows = table_bs.find_all('tr')
```

```
Out[26]: [<tr><td id="flight">Flight No</td><td>Launch site</td> <td>Payload mass</td></tr>,
<tr> <td>1</td><td>Florida</td><td>300 kg</td></tr>,
<tr> <td>2</td><td>Texas</td><td>94 kg</td></tr>,
<tr> <td>3</td><td>Texas</td><td>80 kg</td></tr>]
```

The result is a Python Iterable just like a list, each element is a tag object:

```
In [27]: first_row = table_rows[0]
first_row
```

```
Out[27]: <tr><td id="flight">Flight No</td><td>Launch site</td> <td>Payload mass</td></tr>
```

The type is Tag

```
In [28]: print(type(first_row))
```

```
Out[28]: <class 'bs4.element.Tag'>
```

we can obtain the child

```
In [29]: first_row.td
```

```
Out[29]: <td id="flight">Flight No</td>
```

If we iterate through the list, each element corresponds to a row in the table:

```
In [30]: for i, row in enumerate(table_rows):
print("row", i, "is", row)
```

```
row 0 is <tr><td id="flight">Flight No</td><td>Launch site</td> <td>Payload mass</td></tr>
row 1 is <tr><td>1</td><td>Florida</td><td>300 kg</td></tr>
row 2 is <tr><td>2</td><td>Texas</td><td>94 kg</td></tr>
row 3 is <tr><td>3</td><td>Texas</td><td>80 kg</td></tr>
```

As row is a cell object, we can apply the method find\_all to it and extract table cells in the object cells using the tag td, this is all the children with the name td. The result is a list, each element corresponds to a cell and is a Tag object, we can iterate through this list as well. We can extract the content using the string attribute.

```
In [31]: for i, row in enumerate(table_rows):
print("row", i)
cells = row.find_all('td')
for j, cell in enumerate(cells):
print("column", j, "cell", cell)
```

```
row 0
column 0 cell <td id="flight">Flight No</td>
column 1 cell <td>Launch site</td>
column 2 cell <td>Payload mass</td>
row 1
column 0 cell <td>1</td>
column 1 cell <td>Florida</td>
column 2 cell <td>300 kg</td>
row 2
column 0 cell <td>2</td>
column 1 cell <td>Texas</td>
column 2 cell <td>94 kg</td>
row 3
column 0 cell <td>3</td>
column 1 cell <td>Texas</td>
column 2 cell <td>80 kg</td>
```

If we use a list we can match against any item in that list.

```
In [32]: list_input = table_bs.find_all(name=["tr", "td"])
list_input
```

```
Out[32]: [<tr><td id="flight">Flight No</td><td>Launch site</td> <td>Payload mass</td></tr>,
<td id="flight">Flight No</td>,
<td>Launch site</td>,
<td>Payload mass</td>,
<td id="1">1</td>,
<td>Florida</td>,
<td>300 kg</td>,
<td id="2">2</td>,
<td>Texas</td>,
<td>94 kg</td>,
<td id="3">3</td>,
<td>Texas</td>,
<td>80 kg</td>]
```

### Attributes

If the argument is not recognized it will be turned into a filter on the tag's attributes. For example the id argument, BeautifulSoup will filter against each tag's id attribute. For example, the first td elements have a value of id of flight, therefore we can filter based on that id value.

```
In [33]: table_bs.find_all(id = "flight")
```

```
Out[33]: [<td id="flight">Flight No</td>]
```

We can find all the elements that have links to the Florida Wikipedia page:

```
In [34]: list_input = table_bs.find_all(href = "https://en.wikipedia.org/wiki/Florida")
list_input
```

```
Out[34]: [<a href="https://en.wikipedia.org/wiki/Florida">Florida</a>,
<a href="https://en.wikipedia.org/wiki/Florida">Florida</a>]
```

If we set the href attribute to True, regardless of what the value is, the code finds all tags with href value:

```
In [35]: table_bs.find_all(href = True)
```

```
Out[35]: [<a href="https://en.wikipedia.org/wiki/Florida">Florida</a>,
<a href="https://en.wikipedia.org/wiki/Florida">Florida</a>,
<a href="https://en.wikipedia.org/wiki/Texas">Texas</a>]
```

There are other methods for dealing with attributes and other related methods; Check out the following link

### Exercise: find\_all

Using the logic above, find all the elements without href value

```
In [36]: table_bs.find_all(href=False)
```

```
Out[36]: [<td id="flight">Flight No</td>,
<td>Launch site</td>,
<td>Payload mass</td>,
<td>1</td>,
<td>Florida</td>,
<td>300 kg</td>,
<td>2</td>,
<td>Texas</td>,
<td>94 kg</td>,
<td>3</td>,
<td>Texas</td>,
<td>80 kg</td>]
```

Using the soup object soup, find the element with the id attribute content set to "boldest".

```
In [37]: soup.find_all(id="boldest")
```

```
Out[37]: [<div id="boldest">Lebron James</div>]
```

### string

With string you can search for strings instead of tags, where we find all the elements with Florida:

```
In [38]: table_bs.find_all(string="Florida")
```

```
Out[38]: ['Florida', 'Florida']
```

### find

The find\_all() method scans the entire document looking for results, it's if you are looking for one element you can use the find() method to find the first element in the document. Consider the following two table:

```
In [39]: <html>
<div>Rocket Launch </div>
<p>
<table class="rocket">
<tr>
<td>Flight No</td>
<td>Launch site</td>
<td>Payload mass</td>
</tr>
<tr>
<td>1</td>
<td>Florida</td>
<td>300 kg</td>
</tr>
<tr>
<td>2</td>
<td>Texas</td>
<td>94 kg</td>
</tr>
<tr>
<td>3</td>
<td>Florida</td>
<td>80 kg</td>
</tr>
</table>
</p>
<div>Pizza Party </div>

<table class="pizza">
<tr>
<td>Pizza Place</td>
<td>Orders</td>
<td>Slices</td>
</tr>
<tr>
<td>Domino's Pizza</td>
<td>100</td>
<td>10</td>
</tr>
<tr>
<td>Little Caesars</td>
<td>12</td>
<td>144</td>
</tr>
<tr>
<td>Papa John's</td>
<td>15</td>
<td>165</td>
</tr>
```



