

WebScraping_Review_Lab

July 21, 2021

1 Web Scraping Lab

Estimated time needed: **30** minutes

1.1 Objectives

After completing this lab you will be able to:

Table of Contents

```
<ul>
  <li>
    <a href="BS0">Beautiful Soup Object</a>
    <ul>
      <li>Tag</li>
      <li>Children, Parents, and Siblings</li>
      <li>HTML Attributes</li>
      <li>Navigable String</li>
    </ul>
  </li>
</ul>
<ul>
  <li>
    <a href="filter">Filter</a>
    <ul>
      <li>find All</li>
      <li>find </li>
      <li>HTML Attributes</li>
      <li>Navigable String</li>
    </ul>
  </li>
</ul>
<ul>
  <li>
    <a href="DSCW">Downloading And Scraping The Contents Of A Web</a>
  </li>
</ul>
<p>
  Estimated time needed: <strong>25 min</strong>
```

</p>

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.

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

```
Collecting bs4
  Downloading https://files.pythonhosted.org/packages/10/ed/7e8b97591f6f45617413
9ec089c769f89a94a1a4025fe967691de971f314/bs4-0.0.1.tar.gz
Collecting beautifulsoup4 (from bs4)
  Downloading https://files.pythonhosted.org/packages/d1/41/e6495bd7d3781c
ee623ce23ea6ac73282a373088fcd0ddc809a047b18eae/beautifulsoup4-4.9.3-py3-none-
any.whl (115kB)
    |                                     | 122kB 22.6MB/s eta 0:00:01
Collecting soupsieve>1.2; python_version >= "3.0" (from
beautifulsoup4->bs4)
  Downloading https://files.pythonhosted.org/packages/36/69/d82d04022f02733bf9a7
2bc3b96332d360c0c5307096d76f6bb7489f7e57/soupsieve-2.2.1-py3-none-any.whl
Building wheels for collected packages: bs4
  Building wheel for bs4 (setup.py) ... done
  Stored in directory: /home/jupyterlab/.cache/pip/wheels/a0/b0/b2/4f80b94
56b87abedbc0bf2d52235414c3467d8889be38dd472
Successfully built bs4
Installing collected packages: soupsieve, beautifulsoup4, bs4
Successfully installed beautifulsoup4-4.9.3 bs4-0.0.1 soupsieve-2.2.1
```

Import the required modules and functions

```
[2]: from bs4 import BeautifulSoup # this module helps in web scrapping.
      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:

```
[3]: %%html
      <!DOCTYPE html>
      <html>
      <head>
      <title>Page Title</title>
      </head>
      <body>
```

```

<h3><b id='boldest'>Lebron James</b></h3>
<p> Salary: $ 92,000,000 </p>
<h3> Stephen Curry</h3>
<p> Salary: $85,000, 000 </p>
<h3> Kevin Durant </h3>
<p> Salary: $73,200, 000</p>
</body>
</html>

```

<IPython.core.display.HTML object>

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

```

[4]: html="<!DOCTYPE html><html><head><title>Page Title</title></head><body><h3><b id='boldest'>Lebron James</b></h3><p> Salary: $ 92,000,000 </p><h3> Stephen Curry</h3><p> Salary: $85,000, 000 </p><h3> Kevin Durant </h3><p> Salary: $73,200, 000</p></body></html>"

```

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

```

[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:

```

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

```

```

<!DOCTYPE html>
<html>
  <head>
    <title>
      Page Title
    </title>
  </head>
  <body>
    <h3>
      <b id="boldest">
        Lebron James
      </b>
    </h3>
    <p>
      Salary: $ 92,000,000
    </p>
  </body>
</html>

```

```

</p>
<h3>
    Stephen Curry
</h3>
<p>
    Salary: $85,000, 000
</p>
<h3>
    Kevin Durant
</h3>
<p>
    Salary: $73,200, 000
</p>
</body>
</html>

```

1.2 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.

```
[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

```
[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:

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

```
[9]: <h3><b id="boldest">Lebron James</b></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.

1.2.1 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:

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

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

You can access the parent with the parent

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

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

this is identical to

```
[12]: tag_object
```

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

tag_object.parent is the body element.

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

```
[13]: <body><h3><b id="boldest">Lebron James</b></h3><p> Salary: $ 92,000,000 </p><h3>  
Stephen Curry</h3><p> Salary: $85,000, 000 </p><h3> Kevin Durant </h3><p>  
Salary: $73,200, 000</p></body>
```

tag_object.sibling is the paragraph element

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

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

sibling_2 is the header element which is also a sibling of both sibling_1 and tag_object

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

```
[15]: <h3> Stephen Curry</h3>
```

Exercise: next_sibling

Using the object sibling_2 and the method next_sibling to find the salary of Stephen Curry:

```
[17]: sibling_2.next_sibling
```

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

Click here for the solution

sibling_2.next_sibling

1.2.2 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:

```
[18]: tag_child['id']
```

```
[18]: 'boldest'
```

You can access that dictionary directly as attrs:

```
[19]: tag_child.attrs
```

```
[19]: {'id': 'boldest'}
```

You can also work with Multi-valued attribute check out [1] for more.

We can also obtain the content of the attribute of the tag using the Python get() method.

```
[22]: tag_child.get('id')
```

```
[22]: 'boldest'
```

1.2.3 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:

```
[24]: tag_string=tag_child.string  
tag_string
```

```
[24]: 'Lebron James'
```

we can verify the type is Navigable String

```
[25]: type(tag_string)
```

```
[25]: 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:

```
[26]: unicode_string = str(tag_string)  
unicode_string
```

```
[26]: '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:

```
[27]: %%html
<table>
  <tr>
    <td id='flight' >Flight No</td>
    <td>Launch site</td>
    <td>Payload mass</td>
  </tr>
  <tr>
    <td>1</td>
    <td><a href='https://en.wikipedia.org/wiki/Florida'>Florida</a></td>
    <td>300 kg</td>
  </tr>
  <tr>
    <td>2</td>
    <td><a href='https://en.wikipedia.org/wiki/Texas'>Texas</a></td>
    <td>94 kg</td>
  </tr>
  <tr>
    <td>3</td>
    <td><a href='https://en.wikipedia.org/wiki/Florida'>Florida<a> </td>
    <td>80 kg</td>
  </tr>
</table>
```

<IPython.core.display.HTML object>

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

```
[28]: table="<table><tr><td id='flight'>Flight No</td><td>Launch site</td>_
→<td>Payload mass</td></tr><tr> <td>1</td><td><a href='https://en.wikipedia.
→org/wiki/Florida'>Florida<a></td><td>300 kg</td></tr><tr><td>2</td><td><a_
→href='https://en.wikipedia.org/wiki/Texas'>Texas</a></td><td>94 kg</td></
→tr><tr><td>3</td><td><a href='https://en.wikipedia.org/wiki/
→Florida'>Florida<a> </td><td>80 kg</td></tr></table>"
```

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

1.3 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)`

1.3.1 Name

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

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

```
[31]: [<tr><td id="flight">Flight No</td><td>Launch site</td> <td>Payload
      mass</td></tr>,
      <tr> <td>1</td><td><a
      href="https://en.wikipedia.org/wiki/Florida">Florida</a><a></a></td><td>300
      kg</td></tr>,
      <tr><td>2</td><td><a
      href="https://en.wikipedia.org/wiki/Texas">Texas</a></td><td>94 kg</td></tr>,
      <tr><td>3</td><td><a
      href="https://en.wikipedia.org/wiki/Florida">Florida</a><a> </a></td><td>80
      kg</td></tr>]
```

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

```
[33]: first_row =table_rows[0]
      first_row
```

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

The type is tag

```
[34]: print(type(first_row))

<class 'bs4.element.Tag'>

we can obtain the child
```

```
[35]: first_row.td
```

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

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

```
[36]: 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><a
href="https://en.wikipedia.org/wiki/Florida">Florida</a><a></a></td><td>300
kg</td></tr>
row 2 is <tr><td>2</td><td><a
```



```
href="https://en.wikipedia.org/wiki/Texas">Texas</a></td><td>94 kg</td></tr>
row 3 is <tr><td>3</td><td><a
href="https://en.wikipedia.org/wiki/Florida">Florida</a><a> </a></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.

```
[39]: 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><a
href="https://en.wikipedia.org/wiki/Florida">Florida</a><a></a></td>
column 2 cell <td>300 kg</td>
row 2
column 0 cell <td>2</td>
column 1 cell <td><a href="https://en.wikipedia.org/wiki/Texas">Texas</a></td>
column 2 cell <td>94 kg</td>
row 3
column 0 cell <td>3</td>
column 1 cell <td><a href="https://en.wikipedia.org/wiki/Florida">Florida</a><a>
</a></td>
column 2 cell <td>80 kg</td>
```

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

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

```
[40]: [<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>,
      <tr> <td>1</td><td><a
href="https://en.wikipedia.org/wiki/Florida">Florida</a><a></a></td><td>300
kg</td></tr>,
      <td>1</td>,
```

```

<td><a href="https://en.wikipedia.org/wiki/Florida">Florida</a><a></a></td>,
<td>300 kg</td>,
<tr><td>2</td><td><a
href="https://en.wikipedia.org/wiki/Texas">Texas</a></td><td>94 kg</td></tr>,
<td>2</td>,
<td><a href="https://en.wikipedia.org/wiki/Texas">Texas</a></td>,
<td>94 kg</td>,
<tr><td>3</td><td><a
href="https://en.wikipedia.org/wiki/Florida">Florida</a><a> </a></td><td>80
kg</td></tr>,
<td>3</td>,
<td><a href="https://en.wikipedia.org/wiki/Florida">Florida</a><a> </a></td>,
<td>80 kg</td>]

```

1.4 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.

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

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

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

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

```
[43]: [<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:

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

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

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



```

<td>2</td>,
<td><a href="https://en.wikipedia.org/wiki/Texas">Texas</a></td>,
<td>94 kg</td>,
<tr><td>3</td><td><a
href="https://en.wikipedia.org/wiki/Florida">Florida</a><a> </a></td><td>80
kg</td></tr>,
<td>3</td>,
<td><a href="https://en.wikipedia.org/wiki/Florida">Florida</a><a> </a></td>,
<a> </a>,
<td>80 kg</td>]

```

Click here for the solution

```
table_bs.find_all(href=False)
```

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

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

```
[48]: [<b id="boldest">Lebron James</b>]
```

Click here for the solution

```
soup.find_all(id="boldest")
```

1.4.1 string

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

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

```
[49]: ['Florida', 'Florida']
```

1.5 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:

```

[50]: %%html
<h3>Rocket Launch </h3>

<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>
<p>

<h3>Pizza Party </h3>

<table class='pizza'>
    <tr>
        <td>Pizza Place</td>
        <td>Orders</td>
        <td>Slices </td>
    </tr>
    <tr>
        <td>Domino's Pizza</td>
        <td>10</td>
        <td>100</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>

```

<IPython.core.display.HTML object>

We store the HTML as a Python string and assign two_tables:


```
[61]: for link in soup.find_all('a', href=True): # in html anchor/link is represented
      ↪ by the tag <a>

      print(link.get('href'))
```

```
#main-content
http://www.ibm.com
https://www.ibm.com/cloud/paks?lnk=ushpv18l1
https://developer.ibm.com/blogs/ibm-announces-first-machine-learning-end-to-end-
pipeline-starter-kit/?lnk=ushpv18f1
https://www.ibm.com/blogs/systems/boost-cyber-resilience-and-more-with-ibm-
storage/?lnk=ushpv18f2
https://www.ibm.com/services/data-analytics?lnk=ushpv18f3
https://www.ibm.com/blogs/internet-of-things/geospatial-data-the-really-big-
picture/?lnk=ushpv18f4
https://www.ibm.com/products/offers-and-
discounts?link=ushpv18t5&lnk2=trial_mktpl_MPDISC
https://www.ibm.com/cloud/watson-
assistant?lnk=ushpv18t1&lnk2=trial_WatAssist&psrc=none&pexp=def
https://www.ibm.com/products/hosted-security-
intelligence?lnk=ushpv18t2&lnk2=trial_QRadarCloud&psrc=none&pexp=def
https://www.ibm.com/products/cloud-pak-for-
data?lnk=ushpv18t3&lnk2=trial_CloudPakData&psrc=none&pexp=def
https://www.ibm.com/cloud/cloud-pak-for-
integration?lnk=ushpv18t4&lnk2=trial_CloudPakInt&psrc=none&pexp=def
https://www.ibm.com/search?lnk=ushpv18srch&locale=en-us&q=
https://www.ibm.com/products?lnk=ushpv18p1&lnk2=trial_mktpl&psrc=none&pexp=def
https://developer.ibm.com/depmoels/cloud/?lnk=ushpv18ct16
https://developer.ibm.com/technologies/artificial-intelligence?lnk=ushpv18ct19
https://developer.ibm.com/videos/?lnk=ushpv18ct12
https://developer.ibm.com/?lnk=ushpv18ct9
https://www.ibm.com/docs/en?lnk=ushpv18ct14
https://www.redbooks.ibm.com/?lnk=ushpv18ct10
https://www.ibm.com/mysupport/s/?language=en_US&lnk=ushpv18ct11
https://www.ibm.com/training/?lnk=ushpv18ct15
https://www.ibm.com/cloud/hybrid?lnk=ushpv18ct20
https://www.ibm.com/cloud/learn/public-cloud?lnk=ushpv18ct17
https://www.ibm.com/cloud/redhat?lnk=ushpv18ct13
https://www.ibm.com/artificial-intelligence?lnk=ushpv18ct3
https://www.ibm.com/quantum-computing?lnk=ushpv18ct18
https://www.ibm.com/cloud/learn/kubernetes?lnk=ushpv18ct8
https://www.ibm.com/products/spss-statistics?lnk=ushpv18ct7
https://www.ibm.com/blockchain?lnk=ushpv18ct1
https://www-03.ibm.com/employment/technicaltalent/developer/?lnk=ushpv18ct2
https://www.ibm.com/search?lnk=ushpv18srch&locale=en-us&q=
https://www.ibm.com/products?lnk=ushpv18p1&lnk2=trial_mktpl&psrc=none&pexp=def
https://www.ibm.com/cloud/hybrid?lnk=ushpv18pt14&bv=true
```

```

https://www.ibm.com/watson?lnk=ushpv18pt17&bv=true
https://www.ibm.com/it-infrastructure?lnk=ushpv18pt19&bv=true
https://www.ibm.com/us-en/products/categories?technologyTopics%5B0%5D%5B0%5D=cat
.topic:Blockchain&isIBMOffering%5B0%5D=true&lnk=ushpv18pt4&bv=true
https://www.ibm.com/us-
en/products/category/technology/security?lnk=ushpv18pt9&bv=true
https://www.ibm.com/us-
en/products/category/technology/analytics?lnk=ushpv18pt1&bv=true
https://www.ibm.com/cloud/automation?lnk=ushpv18ct21
https://www.ibm.com/quantum-computing?lnk=ushpv18pt16&bv=true
https://www.ibm.com/support/home/?lnk=ushpv18ct11
https://www.ibm.com/training/?lnk=ushpv18ct15
https://www.ibm.com/demos/?lnk=ushpv18ct12
https://developer.ibm.com/?lnk=ushpv18ct9
https://www.ibm.com/garage?lnk=ushpv18pt18
https://www.ibm.com/docs/en?lnk=ushpv18ct14
https://www.redbooks.ibm.com/?lnk=ushpv18ct10
https://www-03.ibm.com/employment/technicaltalent/developer/?lnk=ushpv18ct2
https://www.ibm.com/

```

1.6 Scrape all images Tags

```

[62]: for link in soup.find_all('img'):# in html image is represented by the tag <img>
      print(link)
      print(link.get('src'))

```

```




https://1.dam.s81c.com/public/content/dam/worldwide-
content/homepage/ul/g/ce/a9/20210517-ls-cloud-paks-25904-720x360.jpg


<img alt="New machine learning starter&nbsp;kit" class="ibm-resize ibm-ab-
image featured-image" decoding="async"

```


src="https://1.dam.s81c.com/public/content/dam/worldwide-content/homepage/ul/g/fd/b4/20210719-f-machine-learning-starter-kit-26012.png" style="position:absolute;top:0;left:0;bottom:0;right:0;box-sizing:border-box;padding:0;border:none;margin:auto;display:block;width:0;height:0;min-width:100%;max-width:100%;min-height:100%;max-height:100%"/>
https://1.dam.s81c.com/public/content/dam/worldwide-content/homepage/ul/g/fd/b4/20210719-f-machine-learning-starter-kit-26012.png



https://1.dam.s81c.com/public/content/dam/worldwide-content/homepage/ul/g/23/78/20210719-f-safeguarded-copy-ann-25952.jpg



https://1.dam.s81c.com/public/content/dam/worldwide-content/homepage/ul/g/b3/43/20210719-f-data-and-analytics-consulting-25989.jpg



https://1.dam.s81c.com/public/content/dam/worldwide-content/homepage/ul/g/78/1d/20210712-f-geospatial-analytics-25957.jpg















```
A6Ly93d3cudzMub3JnLzIwMDAv3ZnIiB2ZXJzaW9uPSIxLjEiLz4=


```

1.7 Scrape data from HTML tables

```
[68]: #The below url contains an html table with data about colors and color codes.
url = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/
↳IBM-DA0321EN-SkillsNetwork/labs/datasets/HTMLColorCodes.html"
```

Before proceeding to scrape a web site, you need to examine the contents, and the way data is organized on the website. Open the above url in your browser and check how many rows and columns are there in the color table.

```
[69]: # get the contents of the webpage in text format and store in a variable called
↳data
data = requests.get(url).text
```

```
[70]: soup = BeautifulSoup(data,"html5lib")
```

```
[71]: #find a html table in the web page
table = soup.find('table') # in html table is represented by the tag <table>
```

```
[72]: #Get all rows from the table
for row in table.find_all('tr'): # in html table row is represented by the tag
↳<tr>
    # Get all columns in each row.
    cols = row.find_all('td') # in html a column is represented by the tag <td>
    color_name = cols[2].string # store the value in column 3 as color_name
    color_code = cols[3].string # store the value in column 4 as color_code
    print("{}--->{}".format(color_name,color_code))
```

```
Color Name--->None
lightsalmon--->#FFA07A
salmon--->#FA8072
darksalmon--->#E9967A
lightcoral--->#F08080
coral--->#FF7F50
tomato--->#FF6347
orangered--->#FF4500
gold--->#FFD700
orange--->#FFA500
darkorange--->#FF8C00
```

```

lightyellow--->#FFFFE0
lemonchiffon--->#FFFACD
papayawhip--->#FFEFD5
moccasin--->#FFE4B5
peachpuff--->#FFDAB9
palegoldenrod--->#EEE8AA
khaki--->#F0E68C
darkkhaki--->#BDB76B
yellow--->#FFFF00
lawngreen--->#7CFC00
chartreuse--->#7FFF00
limegreen--->#32CD32
lime--->#00FF00
forestgreen--->#228B22
green--->#008000
powderblue--->#B0E0E6
lightblue--->#ADD8E6
lightskyblue--->#87CEFA
skyblue--->#87CEEB
deepskyblue--->#00BFFF
lightsteelblue--->#B0C4DE
dodgerblue--->#1E90FF

```

1.8 Scrape data from HTML tables into a DataFrame using BeautifulSoup and Pandas

```
[73]: import pandas as pd
```

```
[74]: #The below url contains html tables with data about world population.
url = "https://en.wikipedia.org/wiki/World_population"
```

Before proceeding to scrape a web site, you need to examine the contents, and the way data is organized on the website. Open the above url in your browser and check the tables on the webpage.

```
[75]: # get the contents of the webpage in text format and store in a variable called
      ↪ data
data = requests.get(url).text
```

```
[78]: soup = BeautifulSoup(data,"html5lib")
```

```
[80]: #find all html tables in the web page
tables = soup.find_all('table') # in html table is represented by the tag
      ↪ <table>
```

```
[81]: # we can see how many tables were found by checking the length of the tables
      ↪ list
```

```
len(tables)
```

[81]: 26

Assume that we are looking for the 10 most densely populated countries table, we can look through the tables list and find the right one we are look for based on the data in each table or we can search for the table name if it is in the table but this option might not always work.

```
[82]: for index,table in enumerate(tables):  
        if ("10 most densely populated countries" in str(table)):  
            table_index = index  
print(table_index)
```

5

See if you can locate the table name of the table, 10 most densely populated countries, below.

```
[83]: print(tables[table_index].prettify())
```

```
<table class="wikitable sortable" style="text-align:right">  
<caption>  
  10 most densely populated countries  
<small>  
  (with population above 5 million)  
</small>  
</caption>  
<tbody>  
<tr>  
<th>  
  Rank  
</th>  
<th>  
  Country  
</th>  
<th>  
  Population  
</th>  
<th>  
  Area  
<br/>  
<small>  
  (km  
<sup>  
  2  
</sup>  
)  
</small>  
</th>
```

```

<th>
  Density
  <br/>
  <small>
    (pop/km
    <sup>
      2
    </sup>
    )
  </small>
</th>
</tr>
<tr>
<td>
  1
</td>
<td align="left">
  <span class="flagicon">
    
  </span>
  <a href="/wiki/Singapore" title="Singapore">
    Singapore
  </a>
</td>
<td>
  5,704,000
</td>
<td>
  710
</td>
<td>
  8,033
</td>
</tr>
<tr>
<td>
  2
</td>
<td align="left">
  <span class="flagicon">
    
    </span>
    <a href="/wiki/Bangladesh" title="Bangladesh">
        Bangladesh
    </a>
</td>
<td>
    171,040,000
</td>
<td>
    143,998
</td>
<td>
    1,188
</td>
</tr>
<tr>
<td>
    3
</td>
<td align="left">
    <span class="flagicon">
        
    </span>
    <a href="/wiki/Lebanon" title="Lebanon">
        Lebanon
    </a>
</td>
<td>
    6,856,000
</td>
<td>
    10,452
</td>
<td>
    656
</td>
</tr>

```

```

<tr>
  <td>
    4
  </td>
  <td align="left">
    <span class="flagicon">
      
    </span>
    <a href="/wiki/Taiwan" title="Taiwan">
      Taiwan
    </a>
  </td>
  <td>
    23,604,000
  </td>
  <td>
    36,193
  </td>
  <td>
    652
  </td>
</tr>
<tr>
  <td>
    5
  </td>
  <td align="left">
    <span class="flagicon">
      
    </span>
    <a href="/wiki/South_Korea" title="South Korea">
      South Korea
    </a>
  </td>
  <td>

```



```

51,781,000
</td>
<td>
99,538
</td>
<td>
520
</td>
</tr>
<tr>
<td>
6
</td>
<td align="left">
<span class="flagicon">

</span>
<a href="/wiki/Rwanda" title="Rwanda">
Rwanda
</a>
</td>
<td>
12,374,000
</td>
<td>
26,338
</td>
<td>
470
</td>
</tr>
<tr>
<td>
7
</td>
<td align="left">
<span class="flagicon">

    </span>
    <a href="/wiki/Haiti" title="Haiti">
        Haiti
    </a>
</td>
<td>
    11,578,000
</td>
<td>
    27,065
</td>
<td>
    428
</td>
</tr>
<tr>
<td>
    8
</td>
<td align="left">
    <span class="flagicon">
        
    </span>
    <a href="/wiki/Netherlands" title="Netherlands">
        Netherlands
    </a>
</td>
<td>
    17,620,000
</td>
<td>
    41,526
</td>
<td>
    424
</td>
</tr>
<tr>
<td>
    9

```

```

</td>
<td align="left">
  <span class="flagicon">
    
  </span>
  <a href="/wiki/Israel" title="Israel">
    Israel
  </a>
</td>
<td>
  9,370,000
</td>
<td>
  22,072
</td>
<td>
  425
</td>
</tr>
<tr>
<td>
  10
</td>
<td align="left">
  <span class="flagicon">
    
  </span>
  <a href="/wiki/India" title="India">
    India
  </a>
</td>
<td>
  1,379,660,000
</td>
<td>

```

```

        3,287,240
    </td>
    <td>
        420
    </td>
</tr>
</tbody>
</table>

```

```

[84]: population_data = pd.DataFrame(columns=["Rank", "Country", "Population",
    ↪ "Area", "Density"])

for row in tables[table_index].tbody.find_all("tr"):
    col = row.find_all("td")
    if (col != []):
        rank = col[0].text
        country = col[1].text
        population = col[2].text.strip()
        area = col[3].text.strip()
        density = col[4].text.strip()
        population_data = population_data.append({"Rank":rank, "Country":
    ↪country, "Population":population, "Area":area, "Density":density},
    ↪ignore_index=True)

population_data

```

```

[84]:
   Rank  Country      Population      Area  Density
0     1  Singapore    5,704,000      710    8,033
1     2  Bangladesh  171,040,000  143,998    1,188
2     3   Lebanon    6,856,000    10,452     656
3     4   Taiwan   23,604,000    36,193     652
4     5 South Korea  51,781,000    99,538     520
5     6   Rwanda   12,374,000    26,338     470
6     7   Haiti    11,578,000    27,065     428
7     8 Netherlands  17,620,000    41,526     424
8     9   Israel    9,370,000    22,072     425
9    10    India   1,379,660,000  3,287,240     420

```

1.9 Scrape data from HTML tables into a DataFrame using BeautifulSoup and read_html

Using the same `url`, `data`, `soup`, and `tables` object as in the last section we can use the `read_html` function to create a DataFrame.

Remember the table we need is located in `tables[table_index]`

We can now use the `pandas` function `read_html` and give it the string version of the table as well

as the flavor which is the parsing engine bs4.

```
[88]: pd.read_html(str(tables[5]), flavor='bs4')
```

```
[88]: [   Rank   Country  Population  Area(km2)  Density(pop/km2)
      0      1  Singapore    5704000         710           8033
      1      2  Bangladesh   17104000        143998          1188
      2      3   Lebanon     6856000         10452           656
      3      4    Taiwan    23604000         36193           652
      4      5  South Korea   51781000         99538           520
      5      6    Rwanda    12374000         26338           470
      6      7     Haiti    11578000         27065           428
      7      8  Netherlands   17620000         41526           424
      8      9     Israel     9370000         22072           425
      9     10     India   1379660000        3287240          420]
```

The function `read_html` always returns a list of DataFrames so we must pick the one we want out of the list.

```
[89]: population_data_read_html = pd.read_html(str(tables[5]), flavor='bs4')[0]

population_data_read_html
```

```
[89]:   Rank   Country  Population  Area(km2)  Density(pop/km2)
      0      1  Singapore    5704000         710           8033
      1      2  Bangladesh   17104000        143998          1188
      2      3   Lebanon     6856000         10452           656
      3      4    Taiwan    23604000         36193           652
      4      5  South Korea   51781000         99538           520
      5      6    Rwanda    12374000         26338           470
      6      7     Haiti    11578000         27065           428
      7      8  Netherlands   17620000         41526           424
      8      9     Israel     9370000         22072           425
      9     10     India   1379660000        3287240          420
```

1.10 Scrape data from HTML tables into a DataFrame using `read_html`

We can also use the `read_html` function to directly get DataFrames from a url.

```
[90]: dataframe_list = pd.read_html(url, flavor='bs4')
```

We can see there are 25 DataFrames just like when we used `find_all` on the soup object.

```
[91]: len(dataframe_list)
```

```
[91]: 26
```

Finally we can pick the DataFrame we need out of the list.

```
[92]: dataframe_list[5]
```

```
[92]:
```

	Rank	Country	Population	Area(km2)	Density(pop/km2)
0	1	Singapore	5704000	710	8033
1	2	Bangladesh	171040000	143998	1188
2	3	Lebanon	6856000	10452	656
3	4	Taiwan	23604000	36193	652
4	5	South Korea	51781000	99538	520
5	6	Rwanda	12374000	26338	470
6	7	Haiti	11578000	27065	428
7	8	Netherlands	17620000	41526	424
8	9	Israel	9370000	22072	425
9	10	India	1379660000	3287240	420

We can also use the `match` parameter to select the specific table we want. If the table contains a string matching the text it will be read.

```
[93]: pd.read_html(url, match="10 most densely populated countries", flavor='bs4')[0]
```

```
[93]:
```

	Rank	Country	Population	Area(km2)	Density(pop/km2)
0	1	Singapore	5704000	710	8033
1	2	Bangladesh	171040000	143998	1188
2	3	Lebanon	6856000	10452	656
3	4	Taiwan	23604000	36193	652
4	5	South Korea	51781000	99538	520
5	6	Rwanda	12374000	26338	470
6	7	Haiti	11578000	27065	428
7	8	Netherlands	17620000	41526	424
8	9	Israel	9370000	22072	425
9	10	India	1379660000	3287240	420

1.11 Authors

Ramesh Sannareddy

1.11.1 Other Contributors

Rav Ahuja

1.12 Change Log

Date (YYYY-MM-DD)	Version	Changed By	Change Description
2020-10-17	0.1	Joseph Santarcangelo	Created initial version of the lab

Copyright © 2020 IBM Corporation. This notebook and its source code are released under the terms of the [MIT License](#).