



Task3: Automatic login to a website and data parsing using the “Beautiful Soup” package

By Jawad Khalil,

Why we need to do this task:

Suppose we need to extract the specific data from a specific page which we get after logging in to a specific website.

- ❑ One can use the extracted data for many purposes for example, in **software development, statistics and Insights**, etc.
- ❑ But, here we are going to only **test** the Python code to parse some data.



Before running the code:

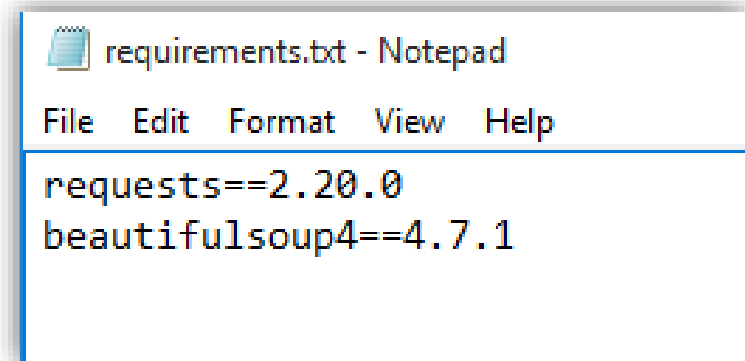
➤ We will need a **website to log** in and parse the data. Here, we will parse the data from <https://the-internet.herokuapp.com> (**Note:** This is a dummy website and it needs a correct username and password to log in.)

➤ Before running the code, we must have installed the following two **packages**

1. **requests** (for pip users, run "**pip install requests**" in Terminal)
2. **beautifulsoup4** (for pip users, run "**pip install beautifulsoup4**" in Terminal.)

Alternatively, we can install the above packages with the "**requirements.txt**" file for this task by running the following command:

"pip install -r requirements.txt"



```
requirements.txt - Notepad
File Edit Format View Help
requests==2.20.0
beautifulsoup4==4.7.1
```

“requests”:

“requests” is an elegant and simple **HTTP library** for Python, built for human beings.

- “requests” allows us to send **HTTP** requests extremely easily
- The HTTP request returns a **Response Object** with all the response data (content, encoding, status, etc).
- It officially supports **Python 3.7+**.
- For the **documentation**, please visit:

<https://requests.readthedocs.io/en/latest/>



<https://requests.readthedocs.io/en/latest/>

“beautifulsoup4”:

BeautifulSoup is a library that makes it easy to **scrape** information from web pages. It creates a parse tree for parsed pages that can be used to extract data from **HTML**, which is useful for **web scraping**.

For the **documentation**, please visit:

<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>



<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>

“BeautifulSoup”:

- “BeautifulSoup” is a class in module “bs4”.
- It is a **data structure** representing a parsed **HTML** or **XML** document.
- Internally, this class defines the **basic interface** called by the tree builders when **converting an HTML/XML** document into a **data structure**.
- Most of the **methods** we call on a BeautifulSoup object are inherited from **PageElement** or **Tag**.

```
BeautifulSoup(markup='', features=None, builder=None, parse_only=None,  
              from_encoding=None, exclude_encodings=None, element_classes=None, **kwargs)
```

Step1: Import the required packages

```
2 import requests
3 from bs4 import BeautifulSoup
4
```

The Difference between bs4 and beautifulsoup4:

"**bs4**" is technically a different package; however, it is a dummy package designed to install the correct package: "**beautifulsoup4**".

We can use either the **short name** or **long name** to install the "beautifulsoup4".

Step2: Login Data

```
36 website_url = 'https://the-internet.herokuapp.com/'
37 login_url = "https://the-internet.herokuapp.com/authenticate"
38
39 login_data = {
40     'username' : 'tomsmith' ,
41     'password' : 'SuperSecretPassword!'
42 }
43
```

1. The “**login_url**” is the page we get after logging in to the required website.

Note: If one visits this URL by just copying it and pasting it in a browser he will get a “Not Found” page.

2. The “**login_data**” contains the required **username** and **password** that we need to pass to the above website to log in.

Step3: Start a session

```
45 with requests.session() as s:  
46
```

By generating a “**session**” object upfront, we can reuse the session; this allows us to store **cookies**, for example, and **re-use** settings that will be utilized for all connections, such as headers and query parameters. Finally, sessions enable **connection pooling**, which allows us to reuse connections to the same host

Step4: Make a request to get a response

```
55 response = s.post(login_url, data=login_data)
56
```

The **post()** is a method that sends a **POST request** to a server.

The **post()** method is used when we need to send some data along with a request to **update** a response.

Here in this code, we are sending the “**login_data**” to a server to get some data back.

Step5: Call the BeautifulSoup() for the response and then inspect the HTML and text

```
58     soup = BeautifulSoup(response.content, "html.parser")
59     print("type(soup)=", type(soup))
60
```

Output: type(soup)= <class 'bs4.BeautifulSoup'>

Step5: Call the BeautifulSoup() for the response and then inspect the HTML and text

```

66
67     print("##### soup.pretty() started #####")
68     print(soup.pretty())
69     print("##### soup.pretty() ended #####")
70

```

Output:

```

##### soup.pretty() started #####
<!DOCTYPE html>
<!--[if IE 8]>         <html class="no-js lt-ie9" lang="en" > <![endif]-->
<!--[if gt IE 8]><!-->
<html class="no-js" lang="en">
<!--<![endif]-->
<html>
  <head>
    <script src="/js/vendor/298279967.js">
    </script>
    <meta charset="utf-8"/>
    <meta content="width=device-width" name="viewport"/>
    <title>
      The Internet
    </title>
    <link href="/css/app.css" rel="stylesheet"/>
    <link href="/css/font-awesome.css" rel="stylesheet"/>
    <script src="/js/vendor/jquery-1.11.3.min.js">
    </script>
    <script src="/js/vendor/jquery-ui-1.11.4/jquery-ui.js">
    </script>
    <script src="/js/foundation/foundation.js">
    </script>
    <script src="/js/foundation/foundation.alerts.js">
    </script>
    <script>
      $(document).foundation();
    </script>
  </head>
  <body>
    <div class="row">
      <div class="large-12 columns" id="flash-messages">
        <div class="flash success" data-alert="" id="flash">
          You logged into a secure area!
          <a class="close" href="#">
            x
          </a>
        </div>
      </div>
    </div>
  </div>
</div>

```

```

<div class="row">
  <a href="https://github.com/touredave/the-internet">
    
  </a>
  <div class="large-12 columns" id="content">
    <div class="example">
      <h2>
        <i class="icon-lock">
        </i>
        Secure Area
      </h2>
      <h4 class="subheader">
        Welcome to the Secure Area. When you are done click logout below.
      </h4>
      <a class="button secondary radius" href="/logout">
        <i class="icon-2x icon-signout">
        Logout
        </i>
      </a>
    </div>
  </div>
</div>
<div class="row" id="page-footer">
  <div class="large-4 large-centered columns">
    <hr/>
    <div style="text-align: center;">
      Powered by
      <a href="http://elementalselenium.com/" target="_blank">
        Elemental Selenium
      </a>
    </div>
  </div>
</div>
</body>
</html>
##### soup.pretty() ended #####

```

Step5: Call the BeautifulSoup() for the response and then inspect the HTML and text

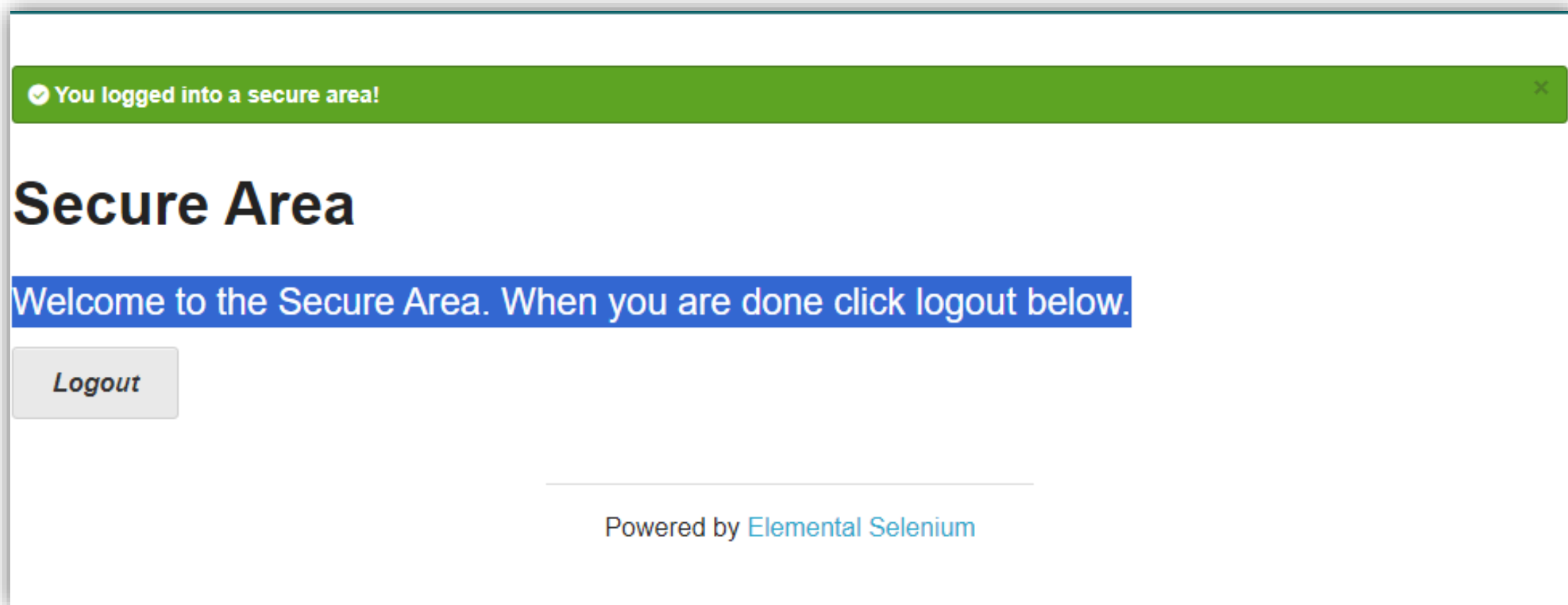
```
71  
72 print("##### Soup.text started #####")  
73 print(soup.text)  
74 print("##### Soup.text ended #####")  
75
```

Output:

```
Soup.text started #####  
  
The Internet  
  
$(document).foundation();  
  
You logged into a secure area!  
×  
  
Secure Area  
Welcome to the Secure Area. When you are done click logout below.  
Logout  
  
Powered by Elemental Selenium  
  
Soup.text ended #####
```

Step6: Get the required data by inspecting the HTML tree

The screenshot of the logged in page:



Task: We need to print the highlighted text

Step6: Get the required data by inspecting the HTML tree

To get the **position(element)** of the required data we need to inspect the login page.

One can inspect the HTML tree in the terminal or output **but...**

```
</head>
<body>
  <div class="row">
    <div class="large-12 columns" id="flash-messages">
      <div class="flash success" data-alert="" id="flash">
        You logged into a secure area!
        <a class="close" href="#">
          ×
        </a>
      </div>
    </div>
  </div>
  <div class="row">
    <a href="https://github.com/tourdedave/the-internet">
      
    <div class="large-12 columns" id="content">
      <div class="example">
        <h2>
          <i class="icon-lock">
          </i>
          Secure Area
        </h2>
        <h4 class="subheader">
          Welcome to the Secure Area. When you are done click logout below.
        </h4>
        <a class="button secondary radius" href="/logout">
          <i class="icon-2x icon-signout">
            Logout
          </i>
        </a>
      </div>
    </div>
  </div>
  <div class="row" id="page-footer">
    <div class="large-4 large-centered columns">
      <hr/>
      <div style="text-align: center;">
        Powered by
        <a href="http://elementalselenium.com/" target="_blank">
          Elemental Selenium
        </a>
      </div>
    </div>
  </div>
</body>
</html>
</html>
##### soup.prettify() ended #####
```


Step6: Get the required data by inspecting the HTML tree

the **easiest way** is to inspect the login page in a browser using the **developer's tools (inspect)**

✓ You logged into a secure area!

Secure Area

h4.subheader 970 × 32.19

Welcome to the Secure Area. When you are done click logout below.

Logout

Powered by [Elemental Selenium](#)

Fork me on GitHub

The screenshot shows the browser's developer tools with the 'Elements' panel open. The HTML tree on the right shows the structure of the page, with the 'h4.subheader' element selected. The 'Styles' panel below shows the default styles for the 'h4.subheader' element, including 'line-height: 1.4' and 'color: #6f6f6f'. The 'Console' panel at the bottom shows a message about the Chrome 121 update.

```
lang= en > <![endif]-->
<!--[if gt IE 8]><!-->
<html class="no-js fbcfaltn idc0_350" lang="en">
<!--<![endif]-->
<head> </head>
<body data-new-gr-c-s-check-loaded="14.1152.0" data-gr-ext-installed>
  <div class="row"> </div>
  <div class="row">
    ::before
    <a href="https://github.com/tourdedave/the-inter-net"> </a>
    <div id="content" class="large-12 columns">
      <div class="example">
        <h2> </h2>
        ...
        <h4 class="subheader">Welcome to the Secure Area. When you are done click logout below.
        </h4> == $0
  </div>
</body>
</html>
```

div#content.large-12.columns div.example h4.subheader

Styles Computed Layout Event Listeners >>

Filter :hov .cls + - []

```
element.style {
}

.subheader {
  line-height: 1.4;
  color: #6f6f6f;
} app.css:1995
```

Console What's new X

Highlights from the Chrome 121 update

@font-palette-values in Elements

The Elements panel now supports and shows a new section in Styles for the @font-palette-values at-rule.

Step6: Get the required data by inspecting the HTML tree

```
▼ <div class="example">  
  ▶ <h2> ... </h2>  
***  <h4 class="subheader">Welcome to the Secure Area. When you are done click logout below.  
      </h4> == $0  
  ▶ <a class="button secondary radius" href="/logout"> ... </a>
```

One can clearly see that the required data is in the tag “h4” with **class=“subheader”**. Now having this tag, we can write the **code** that will extract the required data.

Step6: Get the required data by inspecting the HTML tree

We can use `find()` or `find_all()` methods to find the required(unique) data from the html tree.

find(): It takes some arguments and returns the data according to the given criteria

```
find(name, attrs, recursive, string, **kwargs)
```

find_all(): It also takes some arguments and returns the data (**bs4.element.ResultSet** or a **list of results**)

The only **difference** between `find()` and `find_all()` is that `find_all()` returns a **list** containing a single result while `find()` just returns the **result**.

Step6: Get the required data by inspecting the HTML tree

Extracting the Required Data with find() :

```
74 required_data1 = soup.find("h4", class_="subheader").text
75 print('Required Data=', required_data1)
76
```

Output: Required Data= Welcome to the Secure Area. When you are done click logout below.

Here we used the find() function to look for the tag “h4” where **class=“subheader”**.
“**.text**” returns only text from the tag “h4”.

```
▼ <div class="example">
  ▶ <h2> ... </h2>
...   <h4 class="subheader">Welcome to the Secure Area. When you are done click logout below.
      </h4> == $0
  ▶ <a class="button secondary radius" href="/logout"> ... </a>
```

Step6: Get the required data by inspecting the HTML tree

Extracting the Required Data with `find_all()` :

Since there was only one “h4” tag with class=“subheader”, so we got the data by just using the **`find()`**.

Now let's assume there are more than one “h4” tag with class=“subheader”. If we use the `find()` function, it will return only the first “h4” tag with class=“subheader”. So, we can't get the data from other same tags and classes with `find()` function.

In order to extract the data from the other same “h4” tag with same class we'll use the **`find_all()`** function.

Let's see an example in our case...

Step6: Get the required data by inspecting the HTML tree

Extracting the Required Data with find_all() :

```
78 required_data2 = soup.find_all("div", class_="row")
79 required_data2 = required_data2[1].find('h4', class_="subheader").text
80 print("Required Data=", required_data2)
81
```

Output:

Required Data= Welcome to the Secure Area. When you are done click logout below.

```
<!DOCTYPE html>
<!--[if IE 8]>      <html class="no-js lt-ie9" lang="en" > <![endif]-->
<!--[if gt IE 8]><!-->
<html class="no-js rzvtxtnv idc0_350" lang="en">
  <!--<![endif]-->
  <head> ... </head>
  <body data-new-gr-c-s-check-loaded="14.1152.0" data-gr-ext-installed>
    <div class="row"> ... </div>
    <div class="row">
      ::before
      <a href="https://github.com/tourdedave/the-internet"> ... </a>
      <div id="content" class="large-12 columns">
        <div class="example">
          <h2> ... </h2>
          ...
          <h4 class="subheader">Welcome to the Secure Area. When you are done click logout below.</h4> == $0
          <a class="button secondary radius" href="/logout"> ... </a>
        </div>
      </div>
      ::after
    </div>
    <div id="page-footer" class="row"> ... </div>
  </body>
  <grammarly-desktop-integration data-grammarly-shadow-root="true"> ... </grammarly-desktop-integration>
</html>
```

Step6: Get the required data by inspecting the HTML tree

Extracting the Required Data with “attrs”(attribute) search:

```
83 required_data3 = soup.find(attrs={"class": "subheader"}).text
84
```

Extracting the Required Data with “tag” and “attrs”(attribute) search:

```
86 required_data4 = soup.find("h4", attrs={"class": "subheader"}).text
```

```
74 required_data1 = soup.find("h4", class_="subheader").text
```

```
78 required_data2 = soup.find_all("div", class_="row")
79 required_data2 = required_data2[1].find('h4', class_="subheader").text
```

```
88
89 print(required_data1==required_data2==required_data3==required_data4)
90
```

Output: True

Thank you.