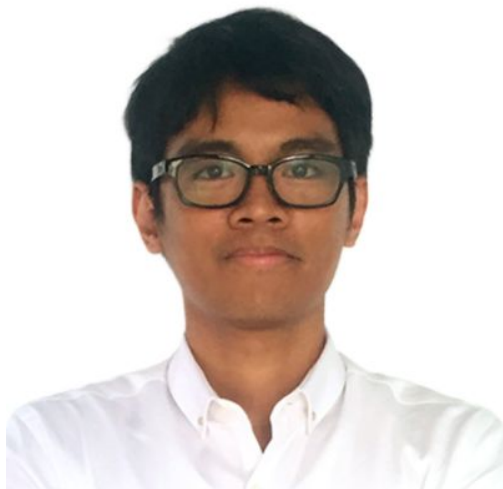


Kaizend: Web Scraping & Automation

April 5, 2020

Who am I?



Joshua “ARVS” Lat

- **Chief Technology Officer** of CBO, Insites, and Jepto
- **AWS Machine Learning Hero**

Goals

- Extract data from webpages
- Learn how to properly use delays to prevent the web scraping jobs from getting blocked or throttled
- Maintain an acceptable level of code quality of the web scraping script(s) as the complexity of the requirement increases

Process

Step # 1: Identify the **target URL**

Step # 2: Download the HTML Page and store the **HTML string** in a variable

Step # 3: Parse the HTML string and **extract the relevant values of the target elements**

DELAY

```
1 from time import sleep
2
3 def delay(seconds):
4     print(f"Sleeping for {seconds} second(s)")
5     sleep(seconds)
6
7
8 def main():
9     print('a')
10    delay(seconds=3)
11    print('b')
12
13
14 if __name__ == "__main__":
15     main()
```

ubuntu:~/environment/labs \$ python3 1.py

```
a
Sleeping for 3 second(s)
b
```

GOAL:

- Loop 5 times
- For each iteration of the loop, sleep for 2 seconds

```
ubuntu:~/environment/labs $ python3 2.py
```

```
0
```

```
Sleeping for 2 second(s)
```

```
1
```

```
Sleeping for 2 second(s)
```

```
2
```

```
Sleeping for 2 second(s)
```

```
3
```

```
Sleeping for 2 second(s)
```

```
4
```

```
Sleeping for 2 second(s)
```

EXTRA CHALLENGE

```
1 from time import sleep
2 from functools import wraps
3
4
5 def delay(seconds):
6     # insert code here
7     def inner_function(function):
8         # insert code here
9         return function
10
11     return inner_function
12
13
14 @delay(seconds=2)
15 def say_something(word):
16     print(word)
17
18
19 def main():
20     say_something("hello")
21
22
23 if __name__ == "__main__":
24     main()
```

```
ubuntu:~/environment/medium $ python3 1.py
```

```
Sleeping for 2 second(s)
```

```
[START]
```

```
hello
```

```
[END]
```


RANDOM NUMBER

```
1 import random
2
3 def get_random_number():
4     return random.randint(1, 5)
5
6 def main():
7     print(get_random_number())
8     print(get_random_number())
9     print(get_random_number())
10    print(get_random_number())
11    print(get_random_number())
12
13
14 if __name__ == "__main__":
15     main()
```

ubuntu:~/environment/labs \$ python3 3.py

```
1
1
4
4
2
```

GOAL:

- Loop 5 times
- For each iteration of the loop, sleep for a random number of seconds

```
ubuntu:~/environment/labs $ python3 4.py
```

```
1
```

```
Sleeping for 2 second(s)
```

```
2
```

```
Sleeping for 2 second(s)
```

```
3
```

```
Sleeping for 3 second(s)
```

```
4
```

```
Sleeping for 1 second(s)
```

IPYTHON EMBED

```
1 from IPython import embed
2
3
4 def main():
5     a = 5
6     b = 3
7
8     embed()
9
10    print(a)
11    print(b)
12
13
14
15 if __name__ == "__main__":
16     main()
```

```
ubuntu:~/environment/labs $ python3 5.py
Python 3.6.9 (default, Nov  7 2019, 10:44:02)
Type 'copyright', 'credits' or 'license' for more information
IPython 7.13.0 -- An enhanced Interactive Python. Type '?' for help.
```

```
In [1]: a
Out[1]: 5
```

```
In [2]: b
Out[2]: 3
```

```
In [3]: █
```

```
1 def debug_input_and_output(func):
2     # Insert code here
3
4     return func
5
6
7 @debug_input_and_output
8 def say_something(word):
9     print(word)
10
11
12 def main():
13     say_something("hello")
14
15
16 if __name__ == "__main__":
17     main()
```

```
ubuntu:~/environment/medium $ python3 2.py
```

```
[INPUT] ARGS: ('hello',)
```

```
[INPUT] KWARGS: {}
```

```
hello
```

```
[OUTPUT]: None
```

BEAUTIFUL SOUP

```
1 from bs4 import BeautifulSoup
2 from IPython import embed
3
4
5 def generate_html():
6     return """
7         <html>
8             <head></head>
9             <body>
10                 <div id="target">Hello World</div>
11             </body>
12         </html>
13     """
14
15
16 def main():
17     html_doc = generate_html()
18     soup = BeautifulSoup(html_doc, 'html.parser')
19     target = soup.find(id="target")
20
21     print(target.get_text())
22
23
24
25 if __name__ == "__main__":
26     main()
```

```
ubuntu:~/environment/labs $ python3 6.py
Hello World
```



```

1 from bs4 import BeautifulSoup
2 from IPython import embed
3
4
5 def generate_html():
6     return """
7         <html>
8             <head></head>
9             <body>
10                 <div>A</div>
11                 <div>B</div>
12                 <div>C</div>
13                 <div>D</div>
14             </body>
15         </html>
16     """
17
18
19 def main():
20     html_doc = generate_html()
21     soup = BeautifulSoup(html_doc, 'html.parser')
22
23     div_elements = soup.find_all('div')
24
25     for div_element in div_elements:
26         print(div_element.get_text())
27
28
29 if __name__ == "__main__":
30     main()

```

ubuntu:~/environment/labs \$ python3 7.py

A
B
C
D

```
1 from bs4 import BeautifulSoup
2 from IPython import embed
3
4
5 def generate_html():
6     return """
7         <html>
8             <head></head>
9             <body>
10                 <a href="/a.html">A</a>
11                 <a href="/b.html">B</a>
12                 <a href="/c.html">C</a>
13                 <a href="/d.html">D</a>
14             </body>
15         </html>
16     """
17
18
19 def main():
20     html_doc = generate_html()
21     embed()
22
23
24 if __name__ == "__main__":
25     main()
```

```
1 from bs4 import BeautifulSoup
2 from IPython import embed
3 import re
4
5
6 def generate_html():
7     return """
8         <html>
9             <head></head>
10            <body>
11                <a href="/a.html">A</a>
12                <a href="/b.html">B</a>
13                <a href="/c.html">C</a>
14                <a href="/d.html">D</a>
15
16                <script>
17                    var hello = "yoh";
18                    alert(hello);
19                </script>
20            </body>
21        </html>
22    """
23
24
25 def main():
26     html_doc = generate_html()
27     soup = BeautifulSoup(html_doc, 'html.parser')
28
29
30 if __name__ == "__main__":
31     main()
```

REQUESTS

```
1 import requests
2
3 from IPython import embed
4
5
6 BASE_URL = "http://sample-target-bucket-with-html-pages.s3-website-ap-southeast-1.amazonaws.com"
7
8
9 def extract_html_content(target_url):
10     response = requests.get(target_url)
11     return response.text
12
13
14 def main():
15     target_page = BASE_URL + "/index.html"
16     print(target_page)
17
18     html_content = extract_html_content(target_page)
19
20     print(html_content)
21
22     embed()
23
24
25 if __name__ == "__main__":
26     main()
```

MANAGING COMPLEXITY

Best Practices

- Making debugging easier using verbose logging and usage of tools such as IPython.embed
- Set a limit to the number of lines for each function (e.g., 15 lines)
- Avoid the usage of try-catch statements
- Add a delay when downloading / extracting pages and files from online sources
- Make function and variable names clear and consistent

```
1 import random
2
3 from IPython import embed
4 from time import sleep
5
6
7 BASE_URL = "http://sample-target-bucket-with-html-pages.s3-website-ap-southeast-1.amazonaws.com"
8
9 PAGES = ["1.html",
10         "2.html",
11         "3.html",
12         "4.html",
13         "5.html"]
14
15
16 def delay(seconds):
17     print(f"Sleeping for {seconds} second(s)")
18     sleep(seconds)
19
20
21 def get_random_number():
22     return random.randint(1, 5)
23
24
25 def main():
26     for page in PAGES:
27         delay(get_random_number())
28         print(BASE_URL + "/" + page)
29
30
31 if __name__ == "__main__":
32     main()
```


CHALLENGE # 4

```
1 import random
2 import requests
3
4 from IPython import embed
5 from time import sleep
6 from bs4 import BeautifulSoup
7
8
9 BASE_URL = "http://sample-target-bucket-with-html-pages.s3-website-ap-southeast-1.amazonaws.com"
10
11 PAGES = ["group1/1.html", "group1/2.html", "group1/3.html", "group1/4.html", "group1/5.html"]
12
13
14 def delay(seconds):
15     print(f"Sleeping for {seconds} second(s)")
16     sleep(seconds)
17
18
19 def get_random_number():
20     return random.randint(1, 3)
21
22
23 def extract_html_content(target_url):
24     return ""
25
26
27 def extract_target_value_from_page(html_string):
28     return ""
29
30
31 def main():
32     for page in PAGES:
33         target_page = BASE_URL + "/" + page
34         print(target_page)
35
36
37 if __name__ == "__main__":
38     main()
```

```
ubuntu:~/environment/medium $ python3 4.py
[START: extract_html_content]
Downloading HTML content of http://sample-target-bucket-with-html-pages.s3-website-ap-southeast-1.amazonaws.com/group1/1.html
[END: extract_html_content]]
<html>
  <body>
    <div id="target">Apple</div>
  </body>
</html>
[START: extract_target_value_from_page]
[END: extract_target_value_from_page]]
Apple
[START: delay]
Sleeping for 3 second(s)
[END: delay]]
[START: extract_html_content]
Downloading HTML content of http://sample-target-bucket-with-html-pages.s3-website-ap-southeast-1.amazonaws.com/group1/2.html
[END: extract_html_content]]
<html>
  <body>
    <div id="target">Banana</div>
  </body>
</html>
[START: extract_target_value_from_page]
[END: extract_target_value_from_page]]
Banana
[START: delay]
Sleeping for 3 second(s)
[END: delay]]
[START: extract_html_content]
```

LINK: <https://sample-target-bucket-with-html-pages.s3-ap-southeast-1.amazonaws.com/group2/index.html>

```
<html>
  <body>
    <a href="/group2/oipuyouiyouo.html">1</a>
    <a href="/group2/qwerqasdfadsfa.html">2</a>
    <a href="/group2/zxvzcxvzxasdqasf.html">3</a>
  </body>
</html>
```

LINK: <https://sample-target-bucket-with-html-pages.s3-ap-southeast-1.amazonaws.com/group2/oipuyouiyouo.html>

```
<html>
  <body>
    <div id="target">42</div>
  </body>
</html>
```

It works!

This is the default welcome page used to test the correct operation of the Apache2 server after installation on Ubuntu systems. It is based on the equivalent page on Debian, from which the Ubuntu Apache packaging is derived. If you can read this page, it means that the Apache HTTP server installed at this site is working properly. You should **replace this file** (located at `/var/www/html/index.html`) before continuing to operate your HTTP server.

If you are a normal user of this web site and don't know what this page is about, this probably means that the site is currently unavailable due to maintenance. If the problem persists, please contact the site's administrator.

Configuration Overview

Ubuntu's Apache2 default configuration is different from the upstream default configuration, and split into several files optimized for interaction with Ubuntu tools. The configuration system is **fully documented in `/usr/share/doc/apache2/README.Debian.gz`**. Refer to this for the full documentation. Documentation for the web server itself can be found by accessing the **manual** if the `apache2-doc` package was installed on this server.

The configuration layout for an Apache2 web server installation on Ubuntu systems is as follows:

```
/etc/apache2/
|-- apache2.conf
/   |-- ports.conf
|-- mods-enabled
/   |-- *.load
/   |-- *.conf
|-- conf-enabled
/   |-- *.conf
|-- sites-enabled
/   |-- *.conf
```

- `apache2.conf` is the main configuration file. It puts the pieces together by including all remaining configuration files when starting up the web server.
- `ports.conf` is always included from the main configuration file. It is used to determine the listening ports for incoming connections, and this file can be customized anytime.
- Configuration files in the `mods-enabled/`, `conf-enabled/` and `sites-enabled/` directories contain particular configuration snippets which manage modules, global configuration fragments, or virtual host configurations, respectively.
- They are activated by symlinking available configuration files from their respective `*-available/` counterparts. These should be managed by using our helpers `a2enmod`, `a2dismod`, `a2ensite`, `a2dissite`, and `a2enconf`, `a2disconf`. See their respective man pages for detailed information.
- The binary is called `apache2`. Due to the use of environment variables, in the default configuration, `apache2` needs to be started/stopped with `/etc/init.d/apache2` or `apache2ctl`.
Calling `/usr/bin/apache2` directly will not work with the default configuration.

Document Roots

By default, Ubuntu does not allow access through the web browser to *any* file apart of those located in `/var/www`, **public_html** directories (when enabled) and `/usr/share` (for web applications). If your site is using a web document root located elsewhere (such as in `/srv`) you may need to whitelist your document root directory in `/etc/apache2/apache2.conf`.

The default Ubuntu document root is `/var/www/html`. You can make your own virtual hosts under `/var/www`. This is different to previous releases which provides better security out of the box.

Reporting Problems

Please use the `ubuntu-bug` tool to report bugs in the Apache2 package with Ubuntu. However, check **`href="https://bugs.launchpad.net/ubuntu/+source/apache2" rel="nofollow">existing bug reports`** before reporting a new bug.

Please report bugs specific to modules (such as PHP and others) to respective packages, not to the web

LINK: <https://sample-target-bucket-with-html-pages.s3-ap-southeast-1.amazonaws.com/group3/target.html>

```
ubuntu:~/environment/labs $ python3 final-challenge.py
```

```
Downloading HTML content of http://sample-target-bucket-with-html-pages.s3-website-ap-southeast-1.amazonaws.com/group3/target.html
```

```
apache2.conf is the main configuration file. It puts the pieces together by including all remaining configuration files when starting up the web server.
```

ports.conf is always included from the main configuration file. It is used to determine the listening ports for incoming connections, and this file can be customized anytime.

Configuration files in the mods-enabled/, conf-enabled/ and sites-enabled/ directories contain particular configuration snippets which manage modules, global configuration fragments, or virtual host configurations, respectively.

They are activated by symlinking available configuration files from their respective *-available/ counterparts. These should be managed by using our helpers a2enmod, a2dismod, a2ensite, a2dissite, and a2enconf, a2disconf. See their respective man pages for detailed information.

The binary is called apache2. Due to the use of environment variables, in the default configuration, apache2 needs to be started/stopped with /etc/init.d/apache2 or apache2ctl. Calling /usr/bin/apache2 directly will not work with the default configuration.

Kaizend: Web Scraping & Automation

THE END