# CS-532 Web Science: Assignment #1

Due on Thursday, January 28, 2016

*Dr. Michael L. Nelson*

**Plinio Vargas**

pvargas@cs.odu.edu

# Contents

# List of Figures

# Problem 1

Demonstrate that you know how to use "curl" well enough to correctly POST data to a form. Show that the HTML response that is returned is "correct". That is, the server should take the arguments you POSTed and build a response accordingly. Save the HTML response to a file and then view that file in a browser and take a screen shot.
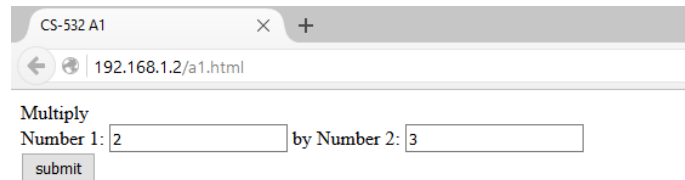
SOLUTION

http://curl.haxx.se[1] was an excellent resource to unveil the power of **curl** command. The are many options available with *curl*, including passing of cookies, which are intensely used by current servers for Cross-Site Request Forgery (CSRF). Then, in order to simplify our demonstration I created a simple html page in my local server that takes two fields: *number1* and *number2*. I named this resource *a1.html*:

a1.html

```
1    <html>
2    <head>
3    <meta http-equiv="Content-Type" content="text/html; charset=utf-8" />
4    <title>CS-532 A1</title>
5    </head>
6
7    <body>
8    <form method="post" action="a1.php">
9    <label>Multiply</label><br />
10    <label>Number 1:</label>
11    <input name="number1" type="text" value="2" />
12    <label> by </label>
13    <label>Number 2:</label>
14    <input name="number2" type="text" value= "3" /> <br />
15    <input name="submit" type="submit" value="submit">
16    </form>
17    </body>
18    </html>
```

**Screen Shot of html page**



Below, is the action script *<a1.php>* run by the server when the form is submitted:

a1.php

```
1    <html>
2    <head>
3    <title>CS-532 A1 Result</title>
4    </head>
5
6    <body>
7    <?php
8        $x = $_POST['number1'];
9        $y = $_POST['number2'];
10
11        echo 'The multiplication of '.$x.' by '.$y.' is '.$x * $y;
12   ?>
13   </body>
14   </html>
```
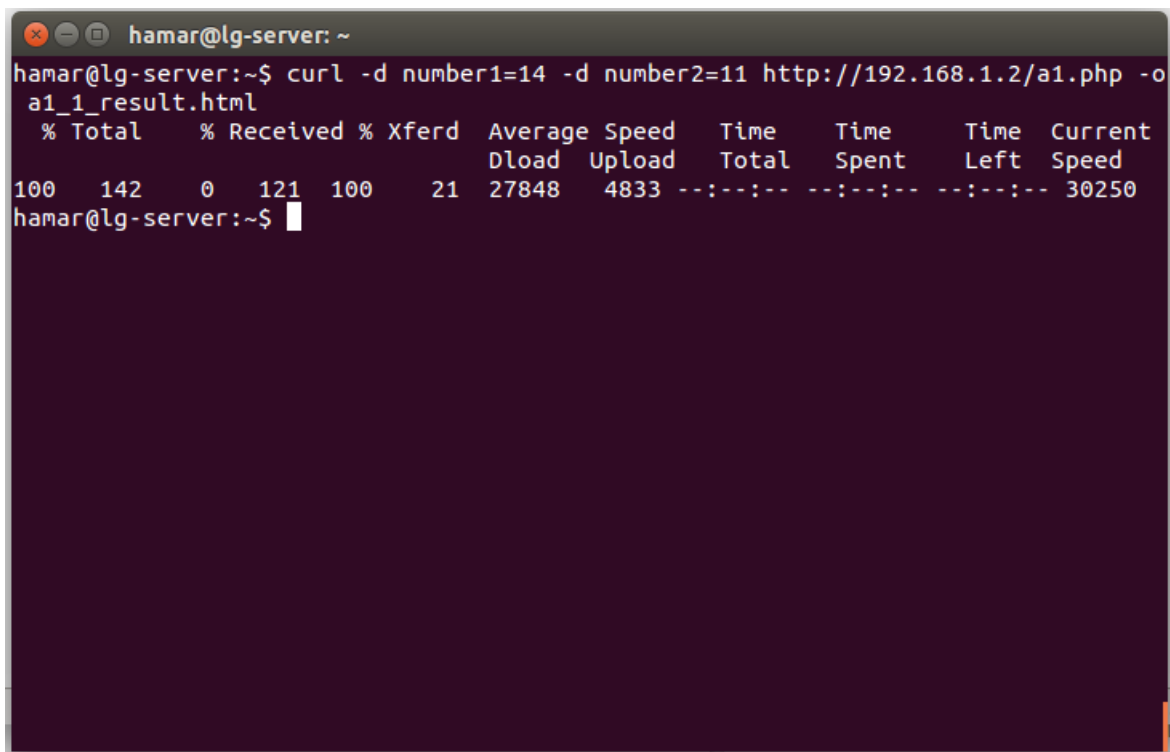
If we type from the command line: **curl -d number1=14 -d number2=11 http://192.168.1.2/a1.php -o a1_1_result.html**, the -d option will combine the field names *number1* and *number2* similar to GET method which is in the format][1]

$$< variable1 >=< data1 > \& < variable2 >=< data2 > \&...$$

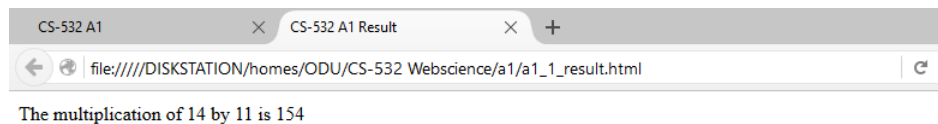The -o option places the response into a given file-name: $< a1\_1\_result.html >$

**Screen Shot of curl POST usage**



**Screen Shot of saved HTML response**

# Problem 2

Write a Python program that:

1. takes as a command line argument a web page

2. extracts all the links from the page

3. lists all the links that result in PDF files, and prints out the bytes for each of the links. (note: be sure to follow all the redirects until the link terminates with a "200 OK".)

4. show that the program works on 3 different URIs, one of which needs to be:
   http://www.cs.odu.edu/ mln/teaching/cs532-s16/test/pdfs.html

   Python program below (***a1.py***) is attached to this file:

a1.py

```
1    import locale
2    import sys
3    import requests
4    import validators
5    from urllib.parse import urlparse
6    from bs4 import BeautifulSoup
7    from time import strftime, localtime, time
8
9    """
10   This Python program
11      1. takes as a command line argument a web page
12      2. extracts all the links from the page
13      3. lists all the links that result in PDF files, and prints out
14        the bytes for each of the links.   (note: be sure to follow
15        all the redirects until the link terminates with a "200 OK".)
16   """
17   __author__ = 'Plinio H. Vargas'
18   __date__ = 'Thu,   Jan 21, 2016 at 22:22:11'
19   __email__ = 'pvargas@cs.odu.edu'
20
21
22   def main(url):
23       locale.setlocale(locale.LC_ALL, 'en_US.utf8')
24       # record running time
25       start = time()
26       print('Starting Time: %s' % strftime("%a,   %b %d, %Y at %H:%M:%S", localtime()))
27       print('Extracting pdf links from: %s\n' % url)
28
29       # get uri status
```
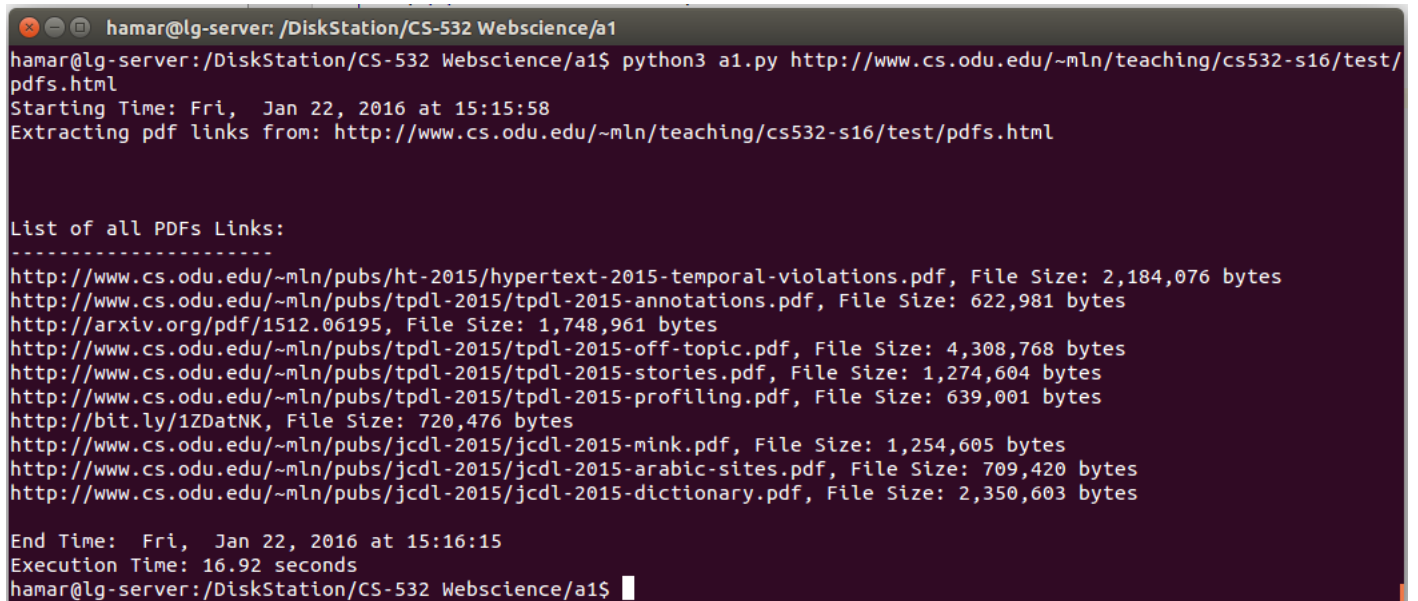
```
30          if requests.get(url).status_code != 200:
31              print('\n\nURI is not available from SERVER. Verify URI.\n')
32              return
33
34          # get source from URI
35          page = requests.get(url).text
36
37          # get parse hostname from URI
38          url = 'http://' + urlparse(url).netloc
39
40          # create BeautifulSoup Object
41          soup = BeautifulSoup(page, 'html.parser')
42
43          # place source link into list
44          all_links = []
45          for link in soup.find_all('a'):
46              uri = link.get('href')
47              # include hostname if url is provided by reference
48              if ((len(uri) > 6 and uri[:7].lower() != 'http://') or len(uri) < 7) and
                      uri[:8].lower() != 'https://':
49                  if uri[:2] == '//':      # if url has double backslash then url is not provided by reference
50                      uri = 'http:' + uri
51                  elif uri[0] != '/':      # include backslash if it was not include by reference
52                      uri = url + '/' + uri
53                  else:
54                      uri = url + uri
55
56              # for debugging
57              # print(uri)
58
59              try:
60                  r = requests.get(uri)
61                  if 'content-type' in r.headers and r.headers['content-type'] == 'application/pdf':
62                      if r.status_code == 200:
63                          all_links.append((uri, r.headers['content-length']))
64              except requests.exceptions.SSLError:
65                  print('Couldn\'t open: %s. URL requires authentication.' % uri)
66              except requests.exceptions.ConnectionError:
67                  print('Couldn\'t open: %s. Connection refused.' % uri)
68
69          print('\n\nList of all PDFs Links:')
70          print('-' * len('List of all PDFs Links'))
71
72          pdf_links = set(all_links)
73          all_links = list(pdf_links)
74          if len(all_links) > 0:
75              for i in range(len(pdf_links)):
76                  print('%s, File Size: %s bytes' % (all_links[i][0],
77                                                      locale.format("%d", int(all_links[i][1]), grouping=True)))
```

```
78          else:
79              print('No PDFs links for above URI.')
80
81          print('\nEnd Time:   %s' % strftime("%a,   %b %d, %Y at %H:%M:%S", localtime()))
82          print('Execution Time: %.2f seconds' % (time()-start))
83          return
84
85    if __name__ == '__main__':
86        # checks for argument
87        if len(sys.argv) != 2:
88            print('Please, provide url\nUsage: python3 a1.py [url]')
89            sys.exit(-1)
90        if not validators.url(sys.argv[1]):
91            print('URL is invalid, please correct url and try again')
92            sys.exit(1)
93
94        # call main
95        main(sys.argv[1])
96
97        sys.exit(0)
```

Below is the result of testing *a1.py* on 3 different URIs:



hamar@lg-server: /DiskStation/CS-532 Webscience/a1

```
hamar@lg-server:/DiskStation/CS-532 Webscience/a1$ python3 a1.py http://www.cs.odu.edu/~mln/teaching/cs532-s16/test/
pdfs.html
Starting Time: Fri,  Jan 22, 2016 at 15:15:58
Extracting pdf links from: http://www.cs.odu.edu/~mln/teaching/cs532-s16/test/pdfs.html


List of all PDFs Links:
--------------------
http://www.cs.odu.edu/~mln/pubs/ht-2015/hypertext-2015-temporal-violations.pdf, File Size: 2,184,076 bytes
http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-annotations.pdf, File Size: 622,981 bytes
http://arxiv.org/pdf/1512.06195, File Size: 1,748,961 bytes
http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-off-topic.pdf, File Size: 4,308,768 bytes
http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-stories.pdf, File Size: 1,274,604 bytes
http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-profiling.pdf, File Size: 639,001 bytes
http://bit.ly/1ZDatNK, File Size: 720,476 bytes
http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-mink.pdf, File Size: 1,254,605 bytes
http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-arabic-sites.pdf, File Size: 709,420 bytes
http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-dictionary.pdf, File Size: 2,350,603 bytes

End Time:  Fri,  Jan 22, 2016 at 15:16:15
Execution Time: 16.92 seconds
hamar@lg-server:/DiskStation/CS-532 Webscience/a1$
```

Figure 1: Test of *a1.py* at http://www.cs.odu.edu/~mln/teaching/cs532-s16/test/pdfs.html

Figure 2: Test of *a1.py* at `http://www.cs.odu.edu/`



Figure 3: Test of *a1.py* at `http://www.vbschools.com/curriculum/gifted/`

# Problem 3

Consider the "bow-tie" graph in the Broder et al. paper (fig 9):
http://www9.org/w9cdrom/160/160.html

Now consider the following graph:

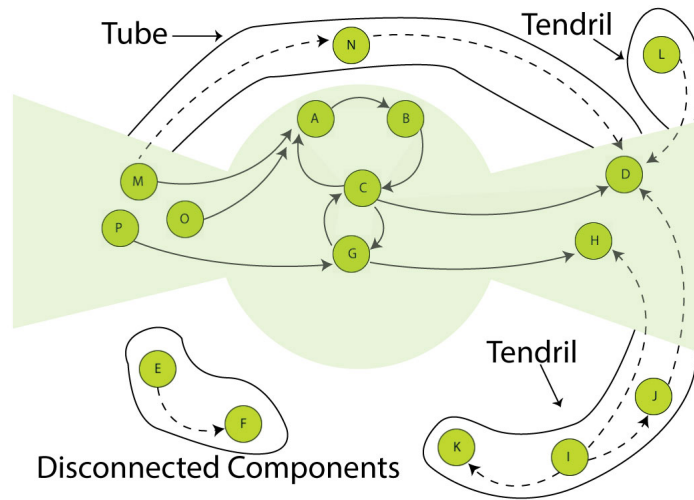| | | |
|---|---|---|
| A | → | B |
| B | → | C |
| C | → | D |
| C | → | A |
| C | → | G |
| E | → | F |
| G | → | C |
| G | → | H |
| I | → | H |
| I | → | J |
| I | → | K |
| J | → | D |
| L | → | D |
| M | → | A |
| M | → | N |
| N | → | D |
| O | → | A |
| P | → | G |



Figure 4: Bow-Tie Graph Representation

For the above graph, give the values for:

**IN**: {M,O,P}        "pages that can reach the SCC, but cannot be reached from it"[2]

**SCC**: {A,B,C,G}        "central core, all of whose pages can reach one another along directed links – this "giant strongly connected component" (SCC) is at the heart of the web". [2]

**OUT**: {D,H}        "pages that are accessible from the SCC, but do not link back to it"[2]

**Tendrils**: {I,J,K,L}        " pages that cannot reach the SCC, and cannot be reached from the SCC"[2]

**Tubes**: {N}        "passage from a portion of IN to a portion of OUT without touching SCC"[2]

**Disconnected**: {E,F}    Everything NOT fitting all criteria above.

# References

[1] Graph structure in the web. (n.d.) Retrieved January 23, 2016, from `http://curl.haxx.se/docs/manual.html`

[2] Graph structure in the web. (n.d.) Retrieved January 23, 2016, from `http://http://www9.org/w9cdrom/160/160.html`