# Web Science: Assignment #1

*Alexander Nwala*

Daniel W. Moore

Sunday, Janury 28, 2018

## Problem 1:

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
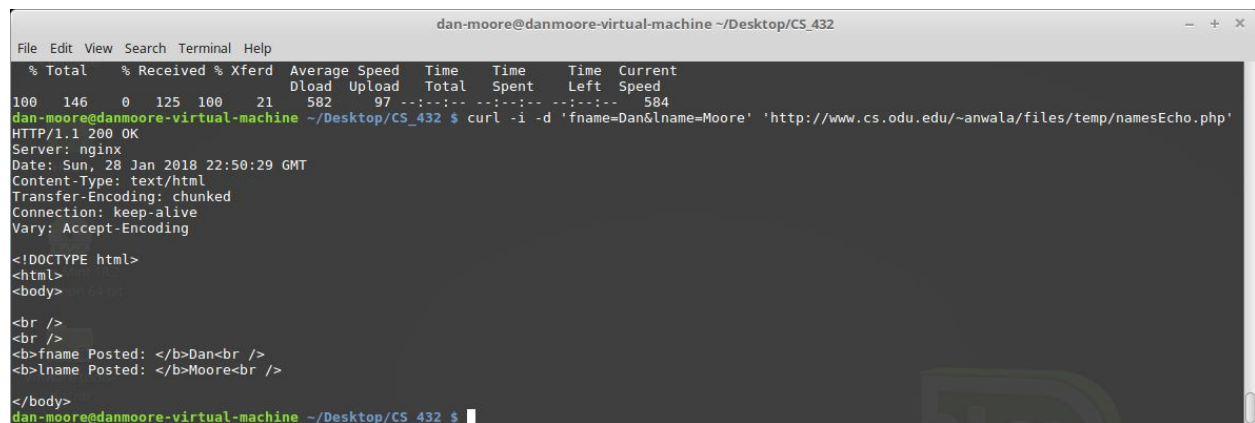
take a screenshot.

## Solution:

The solution for this problem is outlined below:

1. Execute curl command from Linux command line in order to post form data to the
   following url:
   http://www.cs.odu.edu/~anwala/files/temp/namesEcho.php
   See the figure below for command details and output.

**Figure 1**



2.    Redirect output of the above curl command into a html file which may be viewed in a web
browser. The figure below contains a screenshot of the above html-encoded information when
viewed in a web browser.

**Figure 2: Web Browser Representation**



## 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:

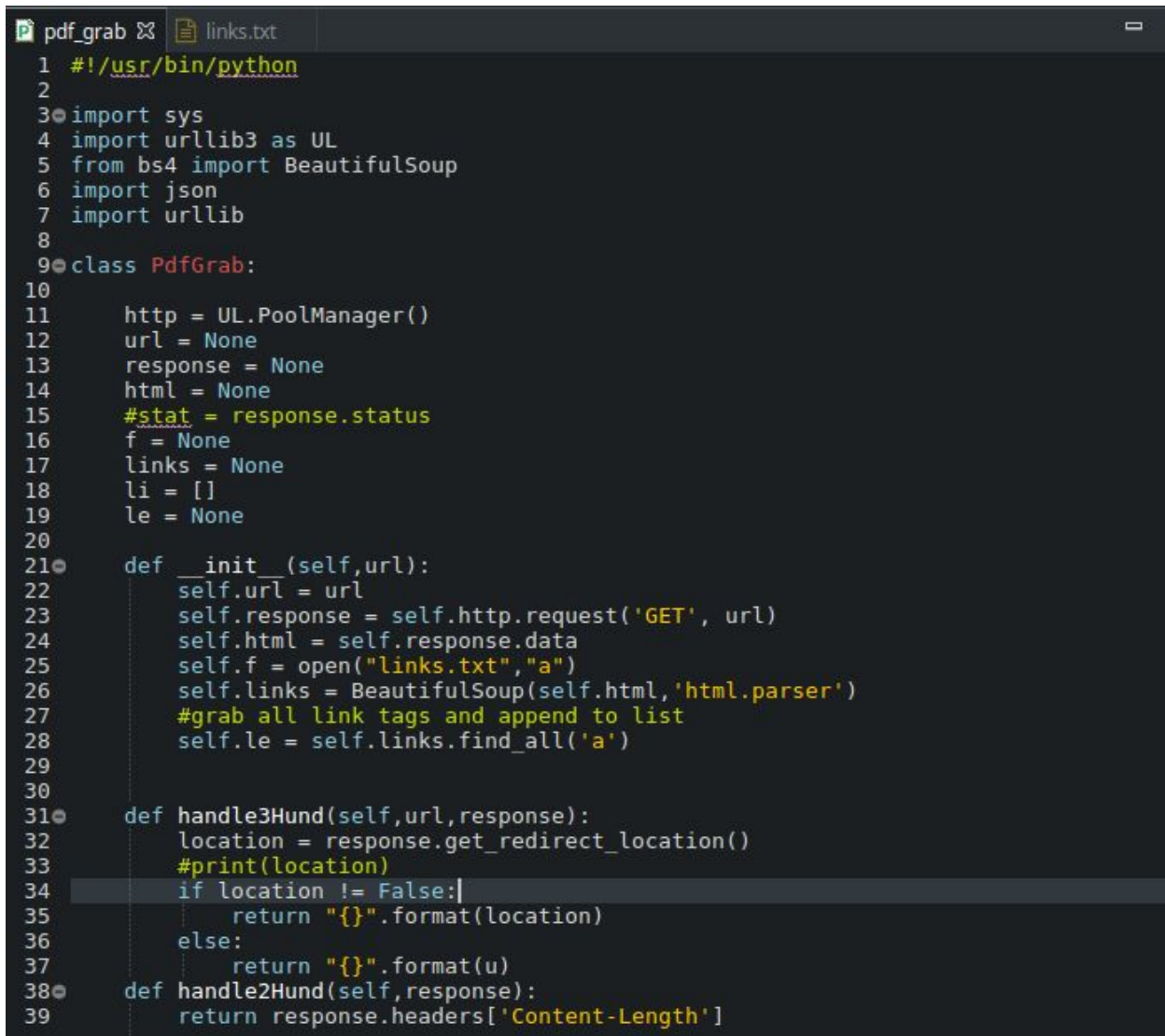        http://www.cs.odu.edu/~mln/teaching/cs532-s17/test/pdfs.html

## Solution:

The solution for this problem is outlined below:

1. Wrote and executed the following Python program which accepts a url as a command line parameter.

**Figure 3:pdf_grab.py**

```python
#!/usr/bin/python

import sys
import urllib3 as UL
from bs4 import BeautifulSoup
import json
import urllib

class PdfGrab:

    http = UL.PoolManager()
    url = None
    response = None
    html = None
    #stat = response.status
    f = None
    links = None
    li = []
    le = None

    def __init__(self,url):
        self.url = url
        self.response = self.http.request('GET', url)
        self.html = self.response.data
        self.f = open("links.txt","a")
        self.links = BeautifulSoup(self.html,'html.parser')
        #grab all link tags and append to list
        self.le = self.links.find_all('a')


    def handle3Hund(self,url,response):
        location = response.get_redirect_location()
        #print(location)
        if location != False:
            return "{}".format(location)
        else:
            return "{}".format(u)
    def handle2Hund(self,response):
        return response.headers['Content-Length']
```

(This space left intentionally blank)

**Figure 3 Continued**

```
40
41  if __name__ == "__main__":
42      grab = PdfGrab(sys.argv[1])
43      for link in grab.le:
44          #url of extrcted link
45          u = "{}".format(link.get('href'))
46          r = grab.http.request('GET', u,redirect=False)
47          #print(r.get_redirect_location())
48          while r.status in range (300,399):
49              u = "{}".format(grab.handle3Hund(u,r))
50              #print (u)
51              r = grab.http.request('GET', u,redirect=False)
52○         #if r.status == 200:
53              #print(r.headers)
54          if r.headers['Content-Type'] == 'application/pdf':
55              size = grab.handle2Hund(r)
56              grab.f.write("{}    {}\n".format(u,size))
57          else:
58              pass
```

2.  The above program extracts all links from the url supplied at command line and prints the url and file size of any links which result in a pdf file to a file.  The program will output the redirect location of any links which respond with status code in the 3xx range. The output of this program using the following url as a command line parameter is shown in the figure below:

    http://www.cs.odu.edu/~mln/teaching/cs532-s17/test/pdfs.html

**Figure 4 links.txt**

```
 1 http://www.cs.odu.edu/~mln/pubs/ht-2015/hypertext-2015-temporal-violations.pdf    218
 2 http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-annotations.pdf    622981
 3 http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-off-topic.pdf    4308768
 4 http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-stories.pdf    1274604
 5 http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-profiling.pdf    639001
 6 http://www.cs.odu.edu/~mln/pubs/jcdl-2014/jcdl-2014-brunelle-damage.pdf    2205546
 7 http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-temporal-intention.pdf    720476
 8 http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-mink.pdf    1254605
 9 http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-arabic-sites.pdf    709420
10 http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-dictionary.pdf    2350603
11
```

## 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 --> K

L --> D

M --> A

M --> N

N --> D

O --> A

P --> G

For the above graph, give the values for:

IN:

SCC:

OUT:

Tendrils:

Tubes:

Disconnected:

## Solution:

1. Placed provided graph data into a text file to be read by Python program.
2. Wrote and executed the following program which utilizes the networkx and matplotlib libraries to produce a visual representation of the provided graph data.
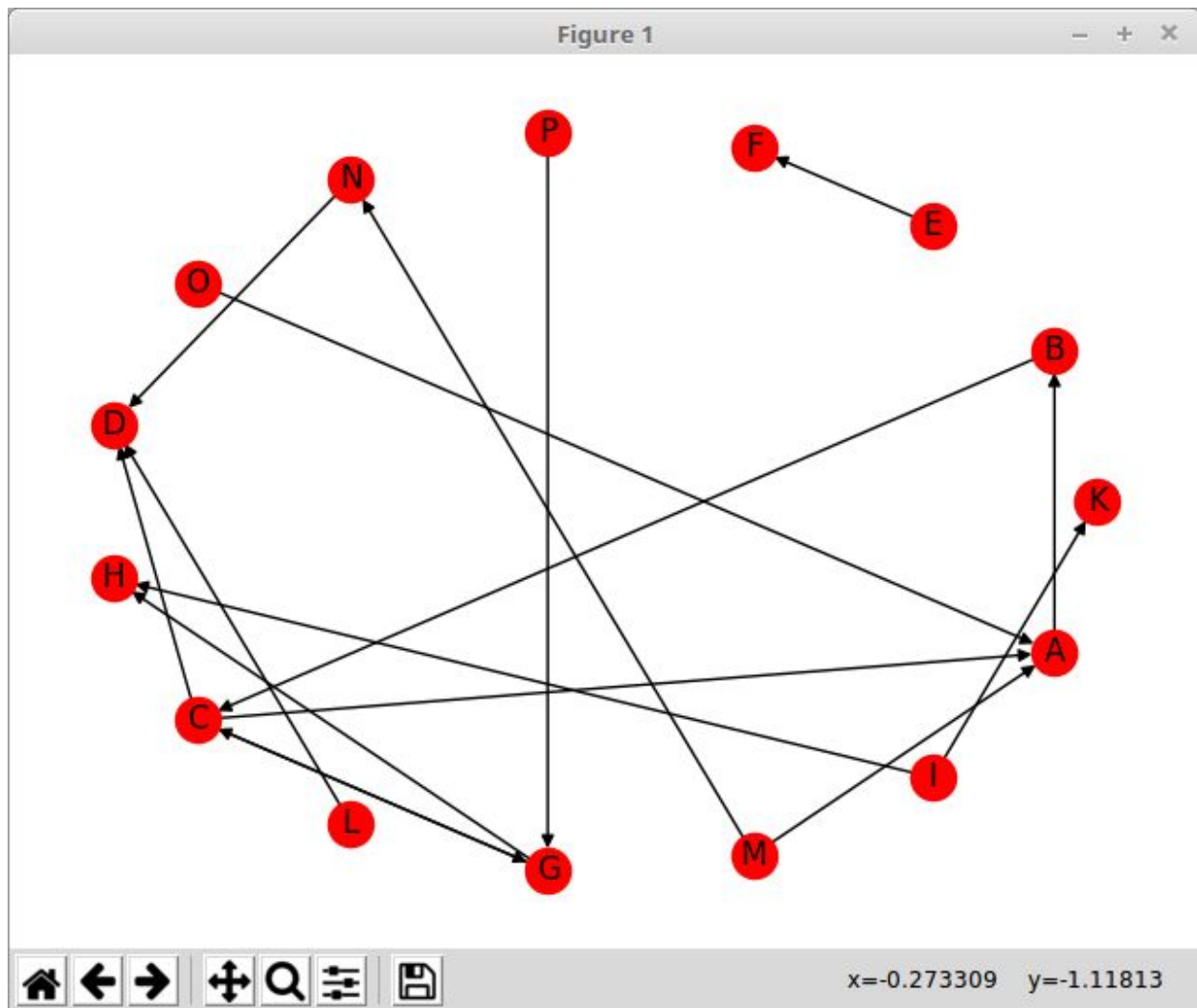
**Figure 5 web_graph.py**

```python
1  #!/usr/bin/python
2
3  import networkx as nx
4  import matplotlib.pyplot as plt
5
6  graph = nx.DiGraph()
7  f = open("nodes.txt","r")
8
9  for line in f:
10     n1 = line[0]
11     n2 = line[6]
12     graph.add_node(n1)
13     graph.add_node(n2)
14     graph.add_edge(n1,n2)
15
16  #print(graph.nodes)
17
18  nx.draw_circular(graph)
19  nx.draw_networkx_labels(graph, pos=nx.circular_layout(graph))
20  plt.show()
```

3. Visually examined the graph below to determine to which group each node belongs

(This space left intentionally blank)

**Figure 6 Directional Node Graph**



SCC:  A, B, C, G

IN: M, O, P

OUT: D, H, K

TENDRILS: L, I

TUBES: N

DISCONNECTED: E, F