# Assignment 1
## CS 532: Introduction to Web Science
Spring 2018
Hrishikesh Gadkari
Finished on January 28, 2018

# 1

## Question

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.

## Answer

To Post data to a form using curl I used the following command:

```
1  curl −i −d "fname=hrishi&lname=gadkari" −X POST http://www.cs.
      odu.edu/~anwala/files/temp/namesEcho.php > output.html
```

Listing 1: Curl command to post data

To approach this problem I referred the `https://gist.github.com/subfuzion/08c5d85437d5d4f00e58` url.

The use of each of the command options mentioned in Listing 1 is as follows:

- -i: To include HTTP headers in the response

- -d: Send specified data in POST request which is followed by the data(parameters) to be sent to the form

- -X: The request method to be used. In this case we have used POST method followed by the php page where data is to be posted

The output in the terminal looks as displayed below:

Figure 1: Response rendered in terminal

The image below displays the view of the output in the browser after saving it in output.html file.



Figure 2: Response rendered in browser

# 2

## Question

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

## Answer

To solve the above problem I wrote the following script using Python3:

```python
1  #!/usr/bin/env python3
2
3  from bs4 import BeautifulSoup
4  import requests
5  import sys
6  import urllib
7  from urllib.parse import urlparse, urljoin
8
9  #Part 1 : Read in command line arguments
10 if len (sys.argv) == 2 :
11         weblink = sys.argv[1]
12
13 else :
14         print ("Usage:python crawl.py url")
15         sys.exit (1)
16
17 # Request data using get
18 try :
19         urllib.request.urlopen(weblink)
20         responsedurl = requests.get(weblink)
21
22 #get http response
23 #if multiple−redirection exists get new URI
24         if responsedurl.history :
25                 print ("Request was redirected to new link")
```

```
26                        print ("New Link: " + responsedurl.url)
27                        response = requests.get(responsedurl.url)
28                        data = response.text
29
30            else :
31                        response = requests.get(weblink)
32                        data = response.text
33
34  except :
35            print ("Enter a valid url")
36            sys.exit (1)
37
38  #Extract all links from the webpage using BeautifulSoup
39  soup = BeautifulSoup(data, 'lxml')
40  print ("Extracting all links from webpage:\n")
41  for link in soup.findAll('a', href=True):
42            webb = link.get('href')
43            print (webb)
44
45  #Extract all pdf links from the webpage if present
46  print ("\nExtracting pdf links from webpage if present:\n")
47  for pdflinks in soup.findAll('a', href=True):
48            web = pdflinks.get('href')
49
50  #check if any relative links present
51            if bool(urlparse(web).netloc) == False :
52                        web = urljoin(weblink, web)
53
54
55            try:
56                        urllib.request.urlopen(web)
57                        req = requests.get(web)
58                        resp = requests.head(web)
59
60  #return pdf links with status_code 200
61                        if req.status_code == 200 and resp.headers.get('
                              content-type') is not None :
62
63                                    if "application/pdf" in resp.headers['
                                        content-type'] :
64                                            print ("First Link: " + web)
65                                            print ("Final Link: " + web)
66                                            print("Bytes: " + resp.headers['
                                                content-length'] + '\n')
67
```

```
68  #return pdf links for the final URI if any redirects
69                  if req.history:
70                          final = req.url
71                          finalresp = requests.head(final)
72                          if finalresp.headers.get('content-type')
                               is not None :
73                                  if "application/pdf" in
                                       finalresp.headers['content-
                                       type'] :
74                                          print ("Request was
                                               redirected")
75                                          print ("First Link: " +
                                               web)
76                                          print ("Final Link: " +
                                               final)
77                                          print("Bytes: " +
                                               finalresp.headers['
                                               content-length'] + '\
                                               n')
78
79          except:
80                  pass
```

Listing 2: Python script that searches for links that end in pdf files

The following libraries were used in my script:

- Beautiful Soup: To parse the HTML and XML documents and extract the data(third party library)

- sys: This module provides access to some variables used or maintained by the interpreter and to functions that interact strongly with the interpreter

- requests: To make HTTP requests

- urllib: This package collects several modules for working with URLs

- urlparse: Focuses on splitting a URL string into its components, or on combining URL components into a URL string

- urljoin: To append a url to the base url

My solution took an iterative approach doing one URI at a time and waiting for each response until moving onto the next URI found.

For the part 1 of the problem the program first checks whether the number of arguments are correct using the command line arguments and if correct, it will pass the first argument after the script name and considers only a properly formatted URI and performs a HTTP get request using the requests library. To check whether the URI is redirected to a new address, response.history [4] was used, after which get request is made to the new URI.

For the part 2 of the problem, that is to extract all links (absolute as well as relative) from the website, Beautiful Soup [2] as a third party library was used to find all the html a elements that contained href tags using the specified lxml parser.

For the part 3 of the problem, the program iterates through each of the URIs found on the page and request again each of those URIs to determine if the URI would point to pdf file. Also if the program encounters a relative link it appends the base url with the relative url. The flow behind this is that it first checks whether the link is absolute using urlparse(web).netloc [1]. If False it changes the link into absolute link using urljoin() [3]. For this I referred stackoverflow website [5]. After this using head and get requests, the content-type, content-length and status code are obtained for each of the URIs and checks whether it returns status-code [1] as 200 or if the URI has redirects [4] and also if content-type [1] is PDF, thus giving the PDF link for the Final URI.

The script is ran with the command as shown below:

```
python crawl.py URI
```

The URIs I used for this problem were:

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

- `http://www.cs.odu.edu/~mweigle/`

- `http://www.cs.odu.edu/~yaohang/`(the output displays a relative link which is converted to absolute and eventually giving a pdf link)

The output for each of the URIs is displayed below:

```
1  Extracting  all  links  from  webpage :
2
3  http :// twitter .com/ webscidl
```

```
 4  http://www.dlib.org/dlib/november15/vandesompel/11vandesompel.
        html
 5  http://arxiv.org/abs/1508.02315
 6  http://arxiv.org/abs/1508.02315
 7  http://www.cs.odu.edu/~mln/pubs/ht-2015/hypertext-2015-temporal-
        violations.pdf
 8  http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-annotations.
        pdf
 9  http://arxiv.org/pdf/1512.06195
10  http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-off-topic.
        pdf
11  http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-stories.pdf
12  http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-profiling.
        pdf
13  http://dx.doi.org/10.1007/s00799-015-0150-6
14  http://www.cs.odu.edu/~mln/pubs/jcdl-2014/jcdl-2014-brunelle-
        damage.pdf
15  http://arxiv.org/abs/1506.06279
16  http://dx.doi.org/10.1007/s00799-015-0155-1
17  http://bit.ly/1ZDatNK
18  http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-mink.pdf
19  http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-arabic-sites
        .pdf
20  http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-dictionary.
        pdf
21  http://bit.ly/jcdl-pdf
22  http://dx.doi.org/10.1007/s00799-015-0140-8
23
24  Extracting pdf links from webpage if present:
25
26  First Link: http://www.cs.odu.edu/~mln/pubs/ht-2015/hypertext
        -2015-temporal-violations.pdf
27  Final Link: http://www.cs.odu.edu/~mln/pubs/ht-2015/hypertext
        -2015-temporal-violations.pdf
28  Bytes: 2184076
29
30  First Link: http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-
        annotations.pdf
31  Final Link: http://www.cs.odu.edu/~mln/pubs/tpdl-2015/tpdl-2015-
        annotations.pdf
32  Bytes: 622981
33
34  Request was redirected
35  First Link: http://arxiv.org/pdf/1512.06195
36  Final Link: https://arxiv.org/pdf/1512.06195.pdf
```

```
37  Bytes :  1748961
38
39  First  Link :  http://www.cs.odu.edu/~mln/pubs/tpdl−2015/tpdl−2015−
         off−topic.pdf
40  Final  Link :  http://www.cs.odu.edu/~mln/pubs/tpdl−2015/tpdl−2015−
         off−topic.pdf
41  Bytes :  4308768
42
43  First  Link :  http://www.cs.odu.edu/~mln/pubs/tpdl−2015/tpdl−2015−
         stories.pdf
44  Final  Link :  http://www.cs.odu.edu/~mln/pubs/tpdl−2015/tpdl−2015−
         stories.pdf
45  Bytes :  1274604
46
47  First  Link :  http://www.cs.odu.edu/~mln/pubs/tpdl−2015/tpdl−2015−
         profiling.pdf
48  Final  Link :  http://www.cs.odu.edu/~mln/pubs/tpdl−2015/tpdl−2015−
         profiling.pdf
49  Bytes :  639001
50
51  First  Link :  http://www.cs.odu.edu/~mln/pubs/jcdl−2014/jcdl−2014−
         brunelle−damage.pdf
52  Final  Link :  http://www.cs.odu.edu/~mln/pubs/jcdl−2014/jcdl−2014−
         brunelle−damage.pdf
53  Bytes :  2205546
54
55  Request  was  redirected
56  First  Link :  http://bit.ly/1ZDatNK
57  Final  Link :  http://www.cs.odu.edu/~mln/pubs/jcdl−2015/jcdl−2015−
         temporal−intention.pdf
58  Bytes :  720476
59
60  First  Link :  http://www.cs.odu.edu/~mln/pubs/jcdl−2015/jcdl−2015−
         mink.pdf
61  Final  Link :  http://www.cs.odu.edu/~mln/pubs/jcdl−2015/jcdl−2015−
         mink.pdf
62  Bytes :  1254605
63
64  First  Link :  http://www.cs.odu.edu/~mln/pubs/jcdl−2015/jcdl−2015−
         arabic−sites.pdf
65  Final  Link :  http://www.cs.odu.edu/~mln/pubs/jcdl−2015/jcdl−2015−
         arabic−sites.pdf
66  Bytes :  709420
67
```

```
68  First Link: http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-
        dictionary.pdf
69  Final Link: http://www.cs.odu.edu/~mln/pubs/jcdl-2015/jcdl-2015-
        dictionary.pdf
70  Bytes: 2350603
```

Listing 3: Output from `http://www.cs.odu.edu/~mln/teaching/cs532-s17/test/pdfs.html`

```
1   Extracting all links from webpage:
2
3   http://www.cs.odu.edu
4   http://www.odu.edu
5   http://www.cs.odu.edu/~mweigle/Main/Home?action=login
6   http://www.cs.odu.edu/~mweigle/Main/Home
7   http://www.cs.odu.edu/~mweigle/Main/Research
8   http://www.cs.odu.edu/~mweigle/Main/PubsByYear
9   http://www.cs.odu.edu/~mweigle/Main/Students
10  http://www.cs.odu.edu/~mweigle/files/CV.pdf
11  http://www.cs.odu.edu/~mweigle/Resources/WorkingWithMe
12  http://www.cs.odu.edu/~mweigle/Resources/ResearchMethods
13  http://www.cs.odu.edu/~mweigle/Resources/InfoVis
14  http://www.cs.odu.edu/~mweigle/Main/Teaching
15  http://www.cs.odu.edu/~mweigle/Main/Sched
16  http://www.cs.odu.edu/~mweigle/Main/Personal
17  http://www.cs.odu.edu/~mweigle/CS725-S18/Home
18  https://graduate.cs.odu.edu/
19  http://www.cs.odu.edu/~yaohang
20  https://securegrants.neh.gov/publicquery/main.aspx?f=1&gn=HAA
        -256368-17
21  https://www.neh.gov/divisions/odh/grant-news/announcing-new
        -2017-odh-grant-awards
22  http://ws-dl.cs.odu.edu
23  http://ws-dl.blogspot.com
24  http://ws-dl.blogspot.com/2018/01/2018-01-08-introducing-
        reconstructive.html
25  http://ws-dl.blogspot.com/2018/01/2018-01-07-review-of-ws-dls
        -2017.html
26  http://ws-dl.blogspot.com/2018/01/2018-01-06-two-wsdl-classes-
        offered-for.html
27  http://ws-dl.blogspot.com/2018/01/2018-01-02-link-to-web-
        archives-not.html
28  http://ws-dl.blogspot.com/2017/12/2017-12-31-digital-blackness-
        in-archive.html
29  http://www.cs.odu.edu/~mweigle/Research/InfoVis-Gallery
```

```
30  https://arxiv.org/abs/1712.03140
31  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=aturban−arxiv17
32  https://arxiv.org/abs/1708.05790
33  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=mccoy−arxiv17
34  http://dx.doi.org/10.1145/3041656
35  http://www.cs.odu.edu/~mweigle/papers/alkwai−tois17−preprint.pdf
36  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=alkwai−tois17
37  http://www.cs.odu.edu/~mweigle/papers/alam−jcdl17.pdf
38  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=alam−jcdl17
39  http://www.cs.odu.edu/~anwala/files/publications/NwalaJCDL_LMP.
       pdf
40  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=nwala−jcdl17
41  http://www.cs.odu.edu/~mweigle/papers/alnoamany−websci17.pdf
42  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=alnoamany−websci17
43  http://www.cs.odu.edu/~mweigle/papers/brunelle−jcdl17.pdf
44  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=brunelle−jcdl17
45  http://dx.doi.org/10.1109/JCDL.2017.7991619
46  http://www.cs.odu.edu/~mkelly/papers/2017_jcdl_wail.pdf
47  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=berlin−jcdl17
48  http://dx.doi.org/10.1109/JCDL.2017.7991601
49  http://www.cs.odu.edu/~mkelly/papers/2017_jcdl_countingMementos.
       pdf
50  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=kelly−jcdl17
51  http://arxiv.org/abs/1705.06218
52  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=alnoamany−arxiv17
53  http://dx.doi.org/10.1109/JCDL.2017.7991601
54  http://www.cs.odu.edu/~mkelly/papers/2017_jcdl_countingMementos.
       pdf
55  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=kelly−jcdl17
56  http://www.cs.odu.edu/~mln/pubs/jcdl−2015/jcdl−2015−mink.pdf
57  http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile
       =Main.bibtex&bibref=jordan−jcdl15
58  http://www.cs.odu.edu/~mln/pubs/jcdl−2015/jcdl−2015−arabic−sites
       .pdf
```

59 | http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile=Main.bibtex&bibref=alkwai−jcdl15
60 | http://dx.doi.org/10.1109/MASS.2014.91
61 | http://www.cs.odu.edu/~mweigle/papers/mohrehkesh−misenet14.pdf
62 | http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile=Main.bibtex&bibref=mohrehkesh−misenet14
63 | http://dx.doi.org/10.1109/JCDL.2014.6970187
64 | http://www.cs.odu.edu/~mln/pubs/jcdl−2014/jcdl−2014−brunelle−damage.pdf
65 | http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile=Main.bibtex&bibref=brunelle−jcdl14
66 | http://dx.doi.org/10.1145/2509338.2509340
67 | http://www.cs.odu.edu/~mweigle/papers/olariu−misenet13.pdf
68 | http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile=Main.bibtex&bibref=olariu−misenet13
69 | http://dx.doi.org/10.1007/978−3−642−40501−3_35
70 | http://www.cs.odu.edu/~mln/pubs/tpdl−2013/paper_149.pdf
71 | http://arxiv.org/abs/1309.4016
72 | http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile=Main.bibtex&bibref=alnoamany−tpdl13
73 | http://dx.doi.org/10.1109/SOLI.2009.5203967
74 | http://www.cs.odu.edu/~mweigle/papers/yan−soli09.pdf
75 | http://www.cs.odu.edu/~mweigle/Main/Home?action=bibentry&bibfile=Main.bibtex&bibref=yan−soli09
76 | https://securegrants.neh.gov/publicquery/main.aspx?f=1&gn=HAA−256368−17
77 | https://mellon.org/grants/grants−database/grants/old−dominion−university/11600663/
78 | http://www.nsf.gov/awardsearch/showAward?AWD_ID=1526700
79 | http://www.imls.gov/grants/awarded/lg−71−15−0077−15
80 | http://www.cs.odu.edu/~mweigle/files/CV.pdf
81 | http://www.cs.odu.edu/~acmw
82 | http://www.ncwit.org/alliances/aa
83 | http://www.cs.odu.edu/
84 | http://www.odu.edu/
85 | http://www.clemson.edu/ces/departments/computing/
86 | http://www.clemson.edu
87 | http://www.cs.unc.edu
88 | http://www.unc.edu
89 | http://www.ulm.edu/cba/computerscience/index.html
90 | http://www.ulm.edu/honors
91 | http://www.ulm.edu
92 | http://www.cs.odu.edu/~mweigle/Main/Home?action=print
93 | http://www.cs.odu.edu/~mweigle/Site/Search
94 | http://www.cs.odu.edu/~mweigle/Main/Home?action=login

```
 95
 96  Extracting  pdf  links  from  webpage  if  present:
 97
 98  First  Link:  http://www.cs.odu.edu/~mweigle/files/CV.pdf
 99  Final  Link:  http://www.cs.odu.edu/~mweigle/files/CV.pdf
100  Bytes:  101583
101
102  First  Link:  http://www.cs.odu.edu/~mweigle/papers/alkwai−tois17−
         preprint.pdf
103  Final  Link:  http://www.cs.odu.edu/~mweigle/papers/alkwai−tois17−
         preprint.pdf
104  Bytes:  1430568
105
106  First  Link:  http://www.cs.odu.edu/~mweigle/papers/alam−jcdl17.
         pdf
107  Final  Link:  http://www.cs.odu.edu/~mweigle/papers/alam−jcdl17.
         pdf
108  Bytes:  1600140
109
110  First  Link:  http://www.cs.odu.edu/~anwala/files/publications/
         NwalaJCDL_LMP.pdf
111  Final  Link:  http://www.cs.odu.edu/~anwala/files/publications/
         NwalaJCDL_LMP.pdf
112  Bytes:  17623699
113
114  First  Link:  http://www.cs.odu.edu/~mweigle/papers/alnoamany−
         websci17.pdf
115  Final  Link:  http://www.cs.odu.edu/~mweigle/papers/alnoamany−
         websci17.pdf
116  Bytes:  6962016
117
118  First  Link:  http://www.cs.odu.edu/~mweigle/papers/brunelle−
         jcdl17.pdf
119  Final  Link:  http://www.cs.odu.edu/~mweigle/papers/brunelle−
         jcdl17.pdf
120  Bytes:  1276346
121
122  First  Link:  http://www.cs.odu.edu/~mkelly/papers/2017_jcdl_wail.
         pdf
123  Final  Link:  http://www.cs.odu.edu/~mkelly/papers/2017_jcdl_wail.
         pdf
124  Bytes:  412476
125
126  First  Link:  http://www.cs.odu.edu/~mkelly/papers/2017
         _jcdl_countingMementos.pdf
```

```
127   Final Link: http://www.cs.odu.edu/~mkelly/papers/2017
          _jcdl_countingMementos.pdf
128   Bytes: 274265
129
130   First Link: http://www.cs.odu.edu/~mkelly/papers/2017
          _jcdl_countingMementos.pdf
131   Final Link: http://www.cs.odu.edu/~mkelly/papers/2017
          _jcdl_countingMementos.pdf
132   Bytes: 274265
133
134   First Link: http://www.cs.odu.edu/~mln/pubs/jcdl−2015/jcdl−2015−
          mink.pdf
135   Final Link: http://www.cs.odu.edu/~mln/pubs/jcdl−2015/jcdl−2015−
          mink.pdf
136   Bytes: 1254605
137
138   First Link: http://www.cs.odu.edu/~mln/pubs/jcdl−2015/jcdl−2015−
          arabic−sites.pdf
139   Final Link: http://www.cs.odu.edu/~mln/pubs/jcdl−2015/jcdl−2015−
          arabic−sites.pdf
140   Bytes: 709420
141
142   First Link: http://www.cs.odu.edu/~mweigle/papers/mohrehkesh−
          misenet14.pdf
143   Final Link: http://www.cs.odu.edu/~mweigle/papers/mohrehkesh−
          misenet14.pdf
144   Bytes: 1231147
145
146   First Link: http://www.cs.odu.edu/~mln/pubs/jcdl−2014/jcdl−2014−
          brunelle−damage.pdf
147   Final Link: http://www.cs.odu.edu/~mln/pubs/jcdl−2014/jcdl−2014−
          brunelle−damage.pdf
148   Bytes: 2205546
149
150   First Link: http://www.cs.odu.edu/~mweigle/papers/olariu−
          misenet13.pdf
151   Final Link: http://www.cs.odu.edu/~mweigle/papers/olariu−
          misenet13.pdf
152   Bytes: 405542
153
154   First Link: http://www.cs.odu.edu/~mln/pubs/tpdl−2013/paper_149.
          pdf
155   Final Link: http://www.cs.odu.edu/~mln/pubs/tpdl−2013/paper_149.
          pdf
156   Bytes: 813692
```

```
157
158  First Link: http://www.cs.odu.edu/~mweigle/papers/yan−soli09.pdf
159  Final Link: http://www.cs.odu.edu/~mweigle/papers/yan−soli09.pdf
160  Bytes: 278184
161
162  First Link: http://www.cs.odu.edu/~mweigle/files/CV.pdf
163  Final Link: http://www.cs.odu.edu/~mweigle/files/CV.pdf
164  Bytes: 101583
```

Listing 4: Output from `http://www.cs.odu.edu/~mweigle/`

```
 1  Extracting all links from webpage:
 2
 3  http://www.cs.odu.edu/~yaohang
 4  teaching.html
 5  projects.html
 6  Tools.html
 7  education.html
 8  pubs.html
 9  people.html
10  contact.html
11  YaohangLiCV2010.pdf
12  Services.html
13  software.html
14  links.html
15  http://www.cs.odu.edu/
16  http://www.odu.edu/
17  mailto:yaohang@cs.odu.edu
18  http://www.cs.odu.edu/
19  http://www.odu.edu/
20  http://www.shodor.org
21  http://www.cs.appstate.edu/nc−hpc/
22  http://www.sura.org
23  postdoctoral_researcher_position.htm
24  phdstudent_positions.htm
25  http://www.orau.org
26  http://www.nsf.gov/awardsearch/showAward.do?AwardNumber=0845702
27  http://www.fsu.edu
28  http://www.scut.edu.cn
29  http://www.ornl.gov
30  http://fellowships.ncsa.uiuc.edu/summer/past.html
31  http://www.ncsa.uiuc.edu/
32  http://www.ornl.gov
33  http://www.ncat.edu/
34
35  Extracting pdf links from webpage if present:
```

```
36
37   First  Link:  http://www.cs.odu.edu/~yaohang/YaohangLiCV2010.pdf
38   Final  Link:  http://www.cs.odu.edu/~yaohang/YaohangLiCV2010.pdf
39   Bytes:  39166
```

Listing 5: Output from `http://www.cs.odu.edu/~yaohang/`

# 3

## Question

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:

## Answer

A graph was generated for the above values using webgraphviz [6]. The following snippets show how it got generated:

```
digraph g{
  rankdir=LR;
  "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"

}
```
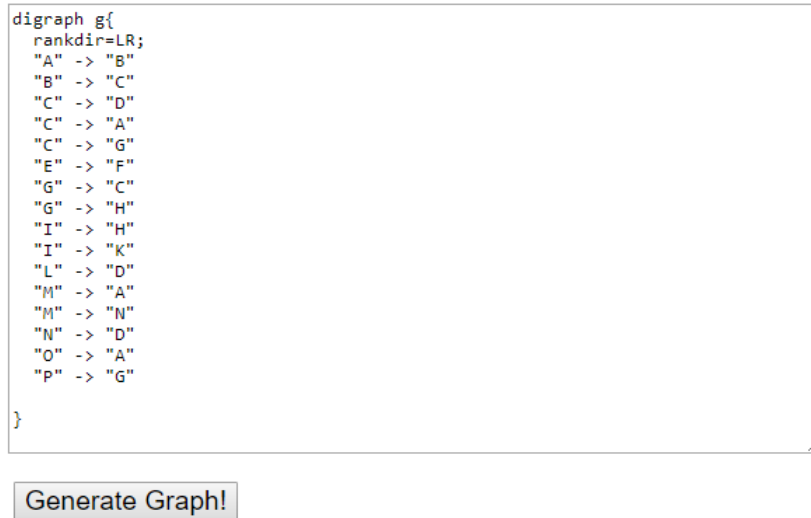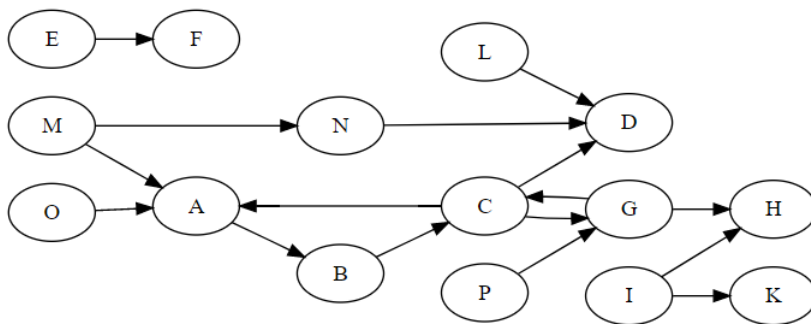
Generate Graph!

Figure 3: Graph generation with WebGraphviz



Figure 4: Graph generated with WebGraphviz

As per the figure 9 in Broder et al. paper and also from the reference to https://www.harding.edu/fmccown/classes/comp475-s13/web-structure-homework.pdf, it made my approach clear to this problem that I should start from SCC. From my observation the values could be as follows:

17

**SCC:** A, B, C, G

Since SCC is at the center of the graph, it will contain values which have directed links amongst each other and also have in-links or out-links to other nodes outside of SCC.

**IN:** M, O, P

Values which have no in-links but can reach the SCC values are considered as IN values.

**OUT:** D, H

Values which have in-links from SCC values but have no out-links are considered as OUT values.

**Tendrils:** I, K, L

Values which have out-links to OUT values and have in-links from IN values but do not refer to SCC values at any point are considered as Tendrils.

**Tubes:** N

Value which is in the middle of IN and OUT value and acts as a direct path from IN to OUT but do not refer to SCC at any point is considered as a Tube.

**Disconnected:** E, F

These values do not connect to anything in the graph, but only amongst them.

# References

[1] "Urllib.requests Documentation." Developer Interface Requests 2.18.4 documentation, Web. 28 Jan, 2018. `https://docs.python.org/3.0/library/urllib.request.html`.

[2] Richardson, Leonard. "Beautiful Soup Documentation." Beautiful Soup Documentation - Beautiful Soup 4.4.0 Documentation. N.p., n.d. Web. 28 Jan. 2018. `https://www.crummy.com/software/BeautifulSoup/bs4/doc/`.

[3] "Urllib.parse Documentation." 21.8. urllib.parse Parse URLs into components Python 3.6.4 documentation, Web. 28 Jan. 2018. `https://docs.python.org/3/library/urllib.parse.html`.

[4] Martijn Pieters. "How to Check Redirects using Python?" Stack Overflow. N.p., n.d. Web. 28 Jan. 2018. `https://stackoverflow.com/questions/20475552/python-requests-library-redirect-new-url`.

[5] Lalinsk, Luk. "How Can I Check If a URL Is Absolute Using Python?" Stack Overflow. N.p., n.d. Web. 28 Jan. 2018. `https://stackoverflow.com/questions/8357098/how-can-i-check-if-a-url-is-absolute-using-python`.

[6] "Webgraphviz." Webgraphviz. N.p., n.d. Web. 28 Jan. 2018. `http://www.webgraphviz.com/`.