# Old Dominion Univeristy

CS 495: Introduction to Web Science
Instructor: Michael L. Nelson, Ph.D
Fall 2014 4:20pm - 7:10pm R, ECSB 2120

Assignment # 1

GeorgeC. Micros UIN: 00757376

**Honor Pledge**
I pledge to support the Honor System of Old Dominion University. I will
refrain from any form of academic dishonesty or deception, such as cheating
or plagiarism. I am aware that as a member of the academic community it
is my responsibility to turn in all suspected violations of the Honor Code. I
will report to a hearing if summoned.

Signed _____
September 11, 2014

George C. Micros

# Written Assignment 1

Fall 2014
CS 495: Introduction to Web Science
Dr. Michael Nelson

September 11, 2014

# Contents

# Chapter 1
# Written Assignment 1

## 1.1 Question 1

### 1.1.1 The Question

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.

### 1.1.2 The Answer

```
1  #! /bin/bash
2  rm -f *.htm
3  #curl -v -X POST -d "name=george" -d @gmicrosTest.txt http://www.posttestserver.com/post.php >
        postTest.htm
4  #curl -v -X POST -d "name=george" -d @gmicrosTest.txt http://httpbin.org/post > httpBin.htm
5  curl -v -X POST -d "name=george" -d @gmicrosTest.txt http://greensuisse.zzl.org/product/dump/dump.
        php > green.htm
```

Listing 1: Bash script to call cURL

```
1  This is a test file
2  George C Micros
3  CS495: Intro to Web Science
4  Assignment 1
```

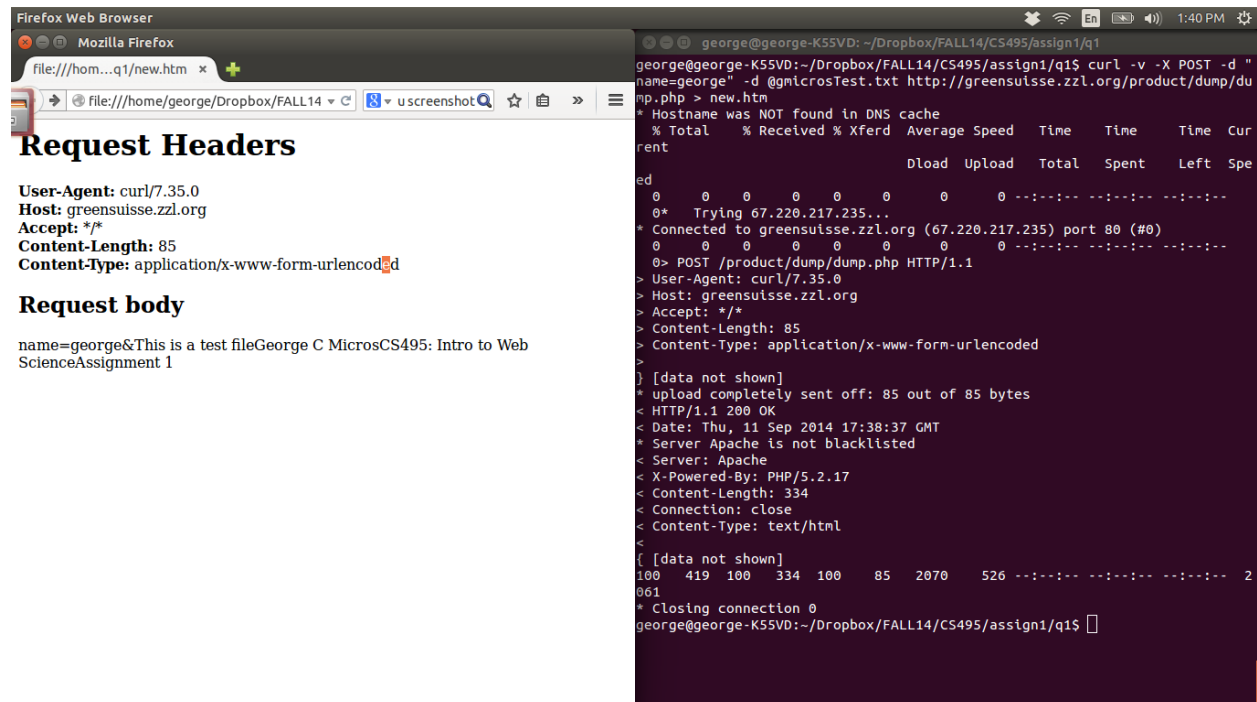Listing 2: Test file to post on server



Fig. 1.1: Screenshot of curl request and return

## 1.2 Question 2

### *1.2.1 The Question*

Write a Python program that:

- takes one argument, like "Old Dominion" or "Virginia Tech"
- takes another argument specified in seconds (e.g., "60" for one minute).
- takes a URI as a third argument:
  `http://sports.yahoo.com/college-football/scoreboard/`
  or
  `http://sports.yahoo.com/college-football/scoreboard/?week=2&conf=all`
  or
  `http://sports.yahoo.com/college-football/scoreboard/?week=1&conf=72`
  etc.
- dereferences the URI, finds the game corresponding to the team argument, prints out the current score (e.g., "Old Dominion 27, East Carolina 17), sleeps for the specified seconds, and then repeats (until control-C is hit).

### *1.2.2 The Answer*

The beautiful soup library converts the complex HTML file tree to a python object tree that contains methods that permit the easy traversal and search of the tree. Initially the structure site of interest was understood by inspecting the page in a web client, i.e. Chrome. This way the specific tags of interest can be easily known without exploring the entire tree in python. Then by targeting these tags we can find all of their instances and find the one that matched our search. Therefore all the away team tags are searched for the team of interest, then all the home teams. Once the tag of interest is found we can use it as a frame of reference to find the other items we need.

Of course there is a condition for the possibility that the desired team was not available. Assuming it was located its sibling and parents can be used to easy find the location of the opponent team and score without searching the entire tree again. The time library allows for a sleep to be inserted at any point in the code.

```bash
#! /bin/bash

./urlCode.py "Old Dominion" 5 "http://sports.yahoo.com/college-football/scoreboard/?week=2&conf=all"
```

Listing 3: Command line input

```python
#! /usr/bin/python
from bs4 import BeautifulSoup
import bs4
import urllib
import urllib2
import sys
import re
import time
from datetime import datetime

def oops():
  print "Didnt find it";
  exit()

print "URI program"
team = sys.argv[1];
sec = float(sys.argv[2]);
site = sys.argv[3];

while(1):
```

```python
21  # fetching the page from the internet
22  request = urllib2.Request(site);
23  file = urllib2.urlopen(request);
24  # getting the soup from the html
25  html = file.read();
26  soup = BeautifulSoup(html);
27
28  team1 = [];
29  reverse = 0;
30  # check the away tsd first
31  tds = soup.find_all("td", "away");
32  for td in tds:
33   for kids in td.children:
34    if isinstance(kids, bs4.element.Tag):
35     for kid in kids.children:
36      if isinstance(kid, bs4.element.Tag) and kid.string == team:
37       there = kid;
38       team1 = there.string;
39  # the team was found in aways
40  if team1 != []:
41   # get parent tag
42   away = there.parent.parent
43   # hop over to score
44   score =  away.next_sibling.next_sibling
45   # weird cntr to store values in list
46   cnt = 0
47   for kids in score.h4.a.children:
48    if isinstance(kids, bs4.element.Tag):
49     score[cnt] = kids.string
50     cnt = cnt +1;
51   # hop over to home
52   home = score.next_sibling.next_sibling;
53   for kids in home.children:
54    if isinstance(kids, bs4.element.Tag):
55     for kid in kids.children:
56      if isinstance(kid, bs4.element.Tag) and kid.string != None:
57       team2 = kid.string;
58
59  # check for team in home
60  else:
61   reverse = 1
62   tds = soup.find_all("td", "home");
63   for td in tds:
64    for kids in td.children:
65     if isinstance(kids, bs4.element.Tag):
66      for kid in kids.children:
67       if isinstance(kid, bs4.element.Tag) and kid.string == team:
68        there = kid;
69        team1 = there.string;
70   if team1 == []:
71    oops();
72   else:
73    home = there.parent.parent
74    score = home.previous_sibling.previous_sibling;
75    cnt = 1;
76    for kids in score.h4.a.children:
77     if isinstance(kids, bs4.element.Tag):
78      score[cnt] = kids.string;
79      cnt = cnt - 1;
80    away = score.previous_sibling.previous_sibling;
81    for kids in away.children:
82     if isinstance(kids, bs4.element.Tag):
83      for kid in kids.children:
84       if isinstance(kid, bs4.element.Tag) and kid.string != None:
85        team2 = kid.string;
86  if team1 != [] and team !=[]:
87   if reverse == 0:
88    print team1, score[0],"-", score[1], team2,"\t", datetime.now().strftime('%H:%M:%S')
89   elif reverse == 1:
90    print team2, score[1],"-", score[0], team1, "\t", datetime.now().strftime('%H:%M:%S')
91   else:
92    print "Something went wrong"
93  else:
94                                    print "Sorry cant find that team. check your spelling";i
95
96  time.sleep(sec);
```

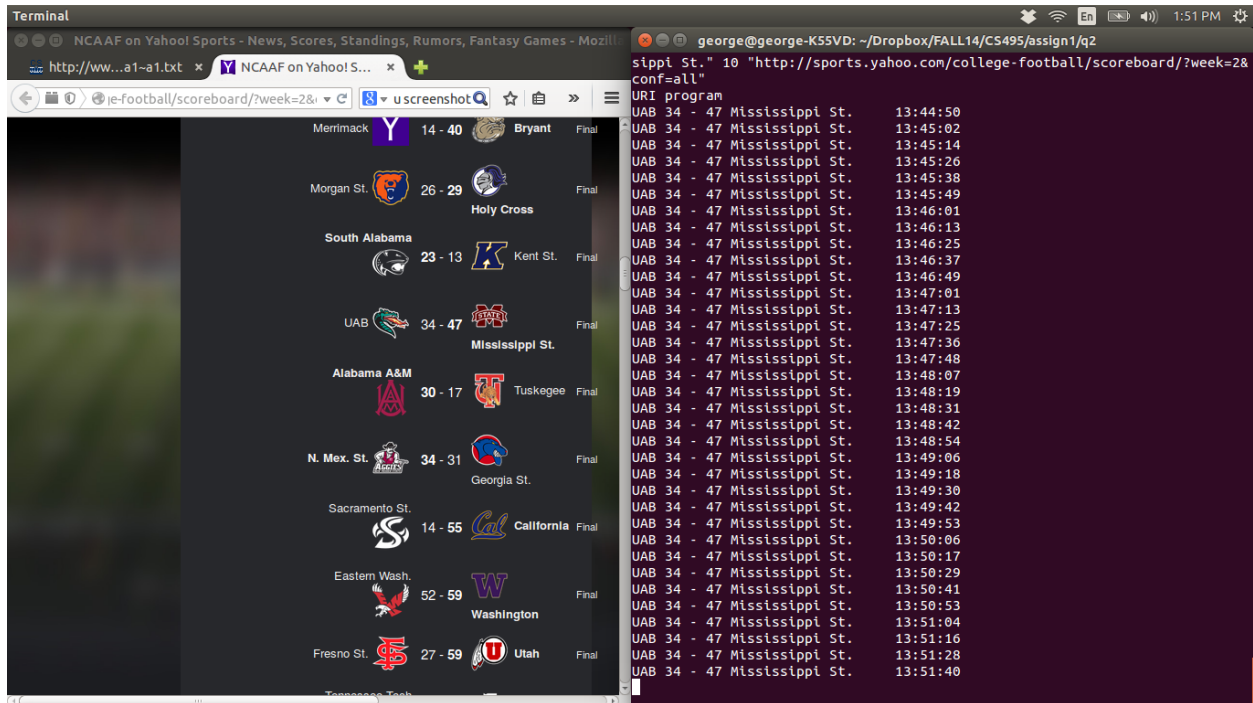Listing 4: Python code implementing the score search

Fig. 1.2: Screenshot of the script execution

## 1.3 Question 3

### 1.3.1 The Question

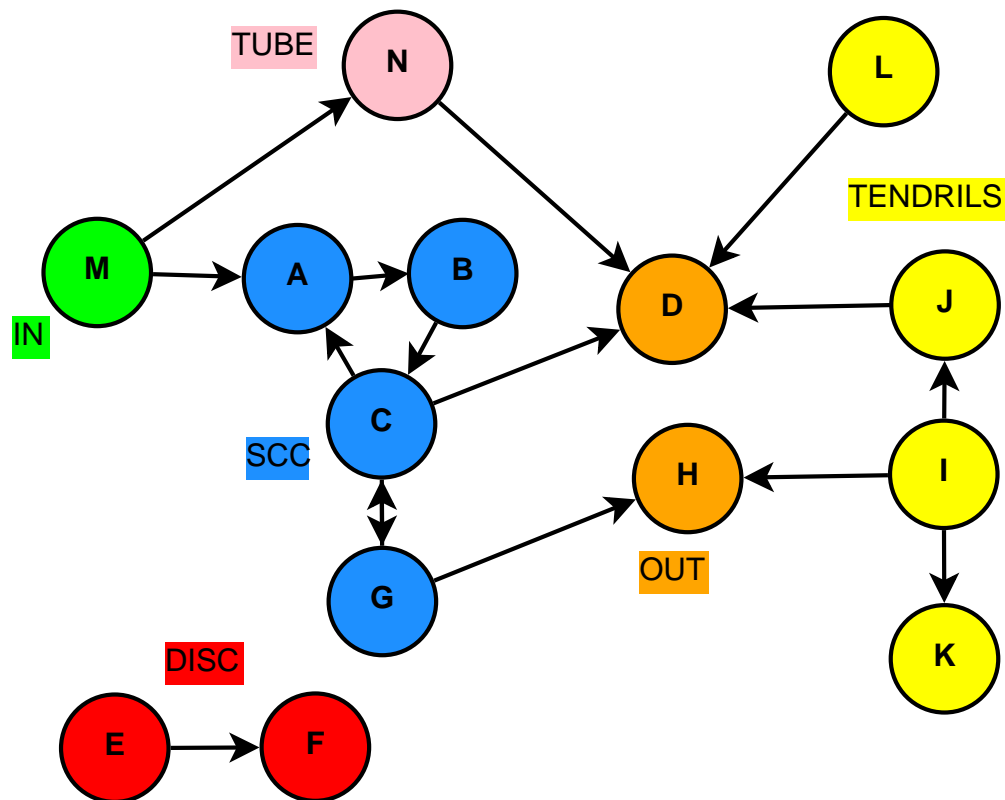Consider the "bow-tie" graph in the Broder et al. paper (fig 9): `http://www9.org/w9cdrom/160/160.html`
Now consider the following graph:

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

For the above graph, give the values for: IN, SCC, OUT, Tendrils, Tubes, Disconnected

### 1.3.2 The Answer

The most critical part of analyzing a graph and identifying the constituent components is to find the strongly connected component. This can be done visually or heuristically in this case. However, in the case of complex and large networks this is not an efficient method. Thankfully algorithms have been developed to solve this problem and they are implemented in python libraries. This way we can verify the result we obtain manually.

```python
1  #! /usr/bin/python
2
3  # script that imports graph data and finds the strongly connected groups
4
5  from collections import defaultdict
6  import numpy as np
7  from scipy.sparse import csgraph
8
9
10 print "Graph Script\n"
11 f = open("graph.txt", 'r');
12 graph = defaultdict(list);
13 array = np.zeros((14,14))
14
15 for line in f:
16   line = line.replace(" ", "");
17    print line
18   line = line.replace("-->","");
19   key = line[0];
20   value = line[1];
21   i = ord(key) - ord('A');
22   j = ord(value) - ord('A')
23   array[i][j] = 1;
24   graph[key].append(value);
25 f.close();
26
27 graph = sorted(graph.items())
28 comp, comp_list =  csgraph.connected_components(array, True, "strong")
29 a =  [list((np.where(comp_list == i))) for i in range(comp)]
30
31 final = [];
32 for i in range(len(a)):
33    b =  (a[i][0].tolist());
34    out = []
35    for j in range(len(b)):
36      out.append( chr(b[j] + ord('A')))
37    final.append(out)
38
39 print "The strongly connected groups are:"
40 print final
```

Listing 5: Code Listing the relevance file corrector

# References

1. http://www.google.com
2. http://jakevdp.github.io/blog/2012/10/14/scipy-sparse-graph-module-word-ladders/
3. http://curl.haxx.se/docs/httpscripting.html
4. http://www.crummy.com/software/BeautifulSoup/bs4/doc/
5. http://www.rmi.net/~lutz/
6. http://www.cs.cornell.edu/home/kleinber/networks-book/