

OLD DOMINION UNIVERSITY

CS 495: INTRODUCTION TO WEB SCIENCE
INSTRUCTOR: MICHAEL L. NELSON, PH.D
FALL 2014 4:20PM - 7:10PM R, ECSB 2120

Assignment # 4

GEORGE C. 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 _____

October 10, 2014

George C. Micros

Written Assignment 4

Fall 2014

CS 495: Introduction to Web Science

Dr. Michael Nelson

October 10, 2014

Contents

1	Written Assignment 4	5
1.1	Question 1	6
1.1.1	The Question	6
1.1.2	The Answer	6
1.2	Question 2	8
1.2.1	The Question	8
1.2.2	The Answer	8
1.3	Question 3	10
1.3.1	The Question	10
1.3.2	The Answer	11
	References	17

Chapter 1

Written Assignment 4

1.1 Question 1

1.1.1 The Question

From your list of 1000 links, choose 100 and extract all of the links from those 100 pages to other pages. We're looking for user navigable links, that is in the form of:

```
<A href="foo">bar</a>
```

We're not looking for embedded images, scripts, `link` elements, etc. You'll probably want to use BeautifulSoup for this.

For each URI, create a text file of all of the outbound links from that page to other URIs (use any syntax that is easy for you). For example:

site:

```
http://www.cs.odu.edu/~mln/
```

links:

```
http://www.cs.odu.edu/
```

```
http://www.odu.edu/
```

```
http://www.cs.odu.edu/~mln/research/
```

```
http://www.cs.odu.edu/~mln/pubs/
```

```
http://ws-dl.blogspot.com/
```

```
http://ws-dl.blogspot.com/2013/09/2013-09-09-ms-thesis-http-mailbox.html
```

etc.

Upload these 100 files to github (they don't have to be in your report).

1.1.2 The Answer

```
1 #!/usr/bin/python
2
3 import urllib2 as ul
4 from bs4 import BeautifulSoup as bs
5 import re
6
7 def getLinks(url):
8     try:
9         req = ul.Request(url);
10        res = ul.urlopen(req);
11        html = res.read();
12        soup = bs(html);
13        links = []
14        for lks in soup.find_all('a'):
15            temp = lks.get('href');
16            #print temp[0:4]
17            if ("http" == (temp[0:4])):
18                g.write(temp+"\n");
19                #h.write(line+" -> "+temp+"\n");
20        except:
21            print "ERROR: "+url+"\n"
22            pass;
23
24 f = open("100","r");
25 #h = open("graph","w");
26
27 #h.write("diagraph test {\n");
28 cnt = 0;
29 for line in f:
30     cnt = cnt +1
31     print cnt;
32
33 g = open("./lks/"+str(cnt).format(str(cnt)),"w+");
34 g.write("site: \n"+str(line)+"links: \n");
35 #print line;
```

```
36| getLinks(line);  
37|  
38| #h.write("}");
```

Listing 1: Bash script that creates a lookup tables of URLs and their hashes

1.2 Question 2

1.2.1 The Question

Using these 100 files, create a single GraphViz “dot” file of the resulting graph. Learn about dot at:

Examples:

<http://www.graphviz.org/content/unix>

<http://www.graphviz.org/Gallery/directed/unix.gv.txt>

Manual:

<http://www.graphviz.org/Documentation/dotguide.pdf>

Reference:

<http://www.graphviz.org/content/dot-language>

<http://www.graphviz.org/Documentation.php>

Note: you’ll have to put explicit labels on the graph, see:

<https://gephi.org/users/supported-graph-formats/graphviz-dot-format/>

(note: actually, I’ll allow any of the formats listed here:

<https://gephi.org/users/supported-graph-formats/>

but “dot” is probably the simplest.)

1.2.2 The Answer

```

1 #! /bin/bash
2
3
4
5 echo "digraph G {" > graph.dot
6 #echo "center=true;" >> graph.dot
7 #echo "size=\"6,6\";" >> graph.dot
8 echo "ranksep=8;" >> graph.dot
9 echo "ratio=auto;" >> graph.dot
10 echo "overlap=false;" >> graph.dot
11 for file in lks/*
12 do
13     #cnt=1
14     #file="lks/001"
15     echo $file
16     a=$(head -n 2 $file | tail -n 1)
17     alab=$(echo $a | sed -e 's|^[^/]*//||' -e 's|^www\.||' -e 's|/.*$||')
18     echo "\"$a\" [label=\"$alab\"]" >> graph.dot
19
20     tail -n+4 $file |
21     {
22         while read line
23         do
24             cnt=$((cnt + 1))
25             #echo $cnt
26             line=$line
27             llab=$(echo $line | sed -e 's|^[^/]*//||' -e 's|^www\.||' -e 's|/.*$||')
28             echo "\"$line\" [label=\"$llab\"]" >> graph.dot
29             echo "\"$a\" -> \"$line\";" >> graph.dot
30
31         done
32     }
33 done
34
35 echo "}" >> graph.dot

```

Listing 2: R script to generate a histogram

1.3 Question 3

1.3.1 The Question

Download and install Gephi:

<https://gephi.org/>

Load the dot file created in #2 and use Gephi to:

- visualize the graph (you'll have to turn on labels)
- calculate HITS and PageRank
- avg degree
- network diameter
- connected components

Put the resulting graphs in your report.

You might need to choose the 100 sites with an eye toward creating a graph with at least one component that is nicely connected. You can probably do this by selecting some portion of your links (e.g., 25, 50) from the same site.

1.3.2 *The Answer*

1.3.2.1 Network Graph

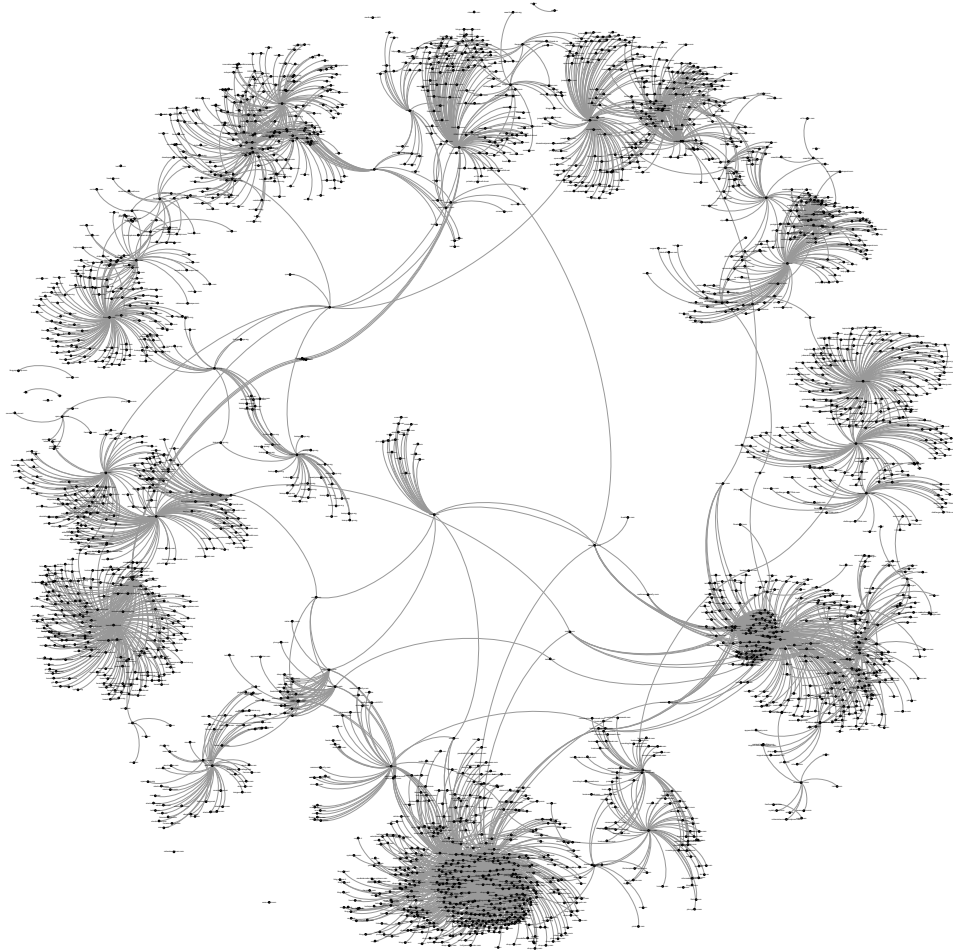


Fig. 1.1: Network Graph

1.3.2.2 HITS

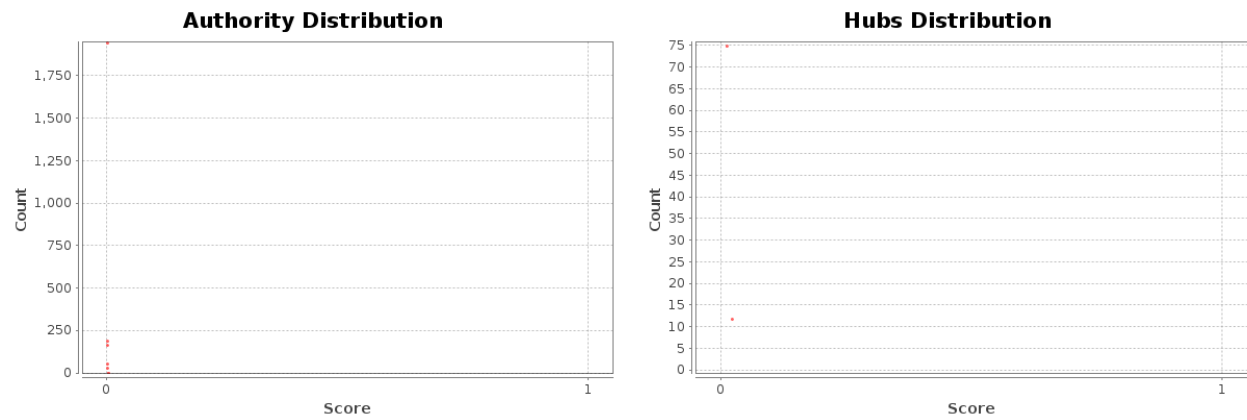


Fig. 1.2: HITS results

1.3.2.3 Page Rank

Epsilon = 0.001 Probability = 0.85

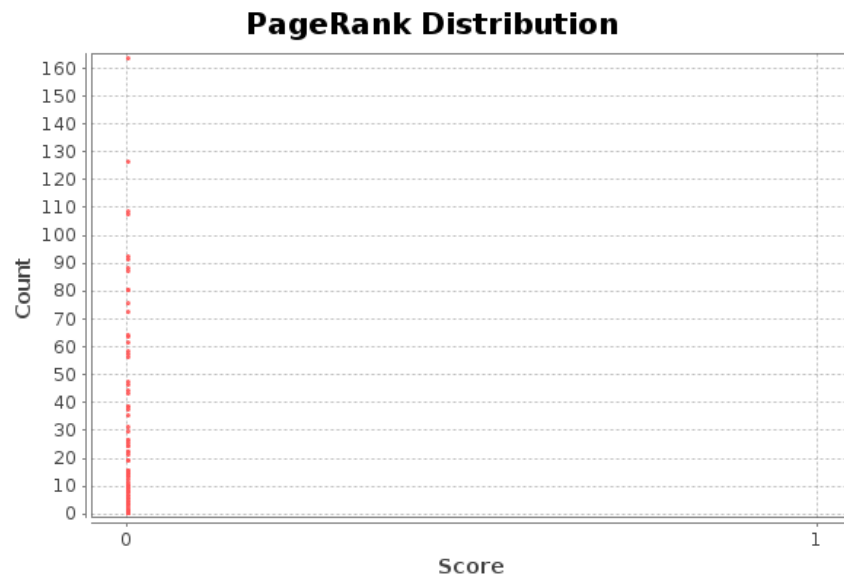
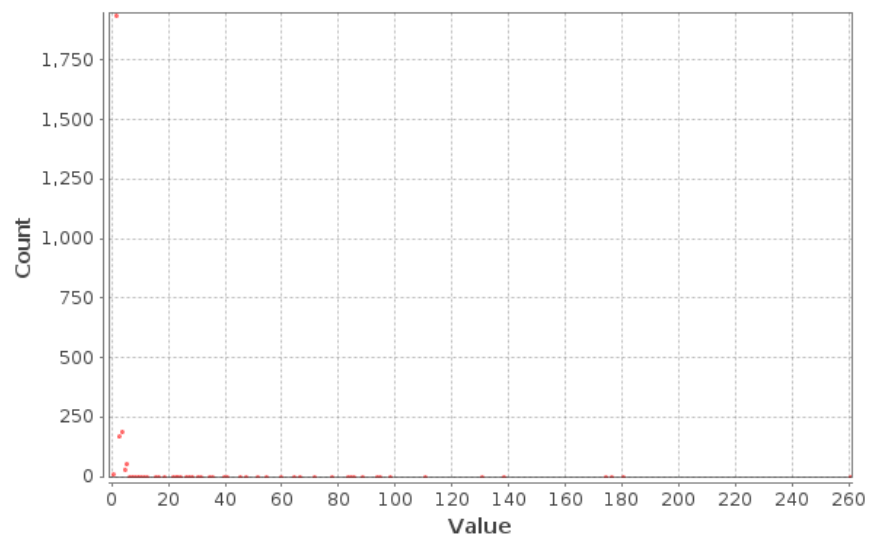
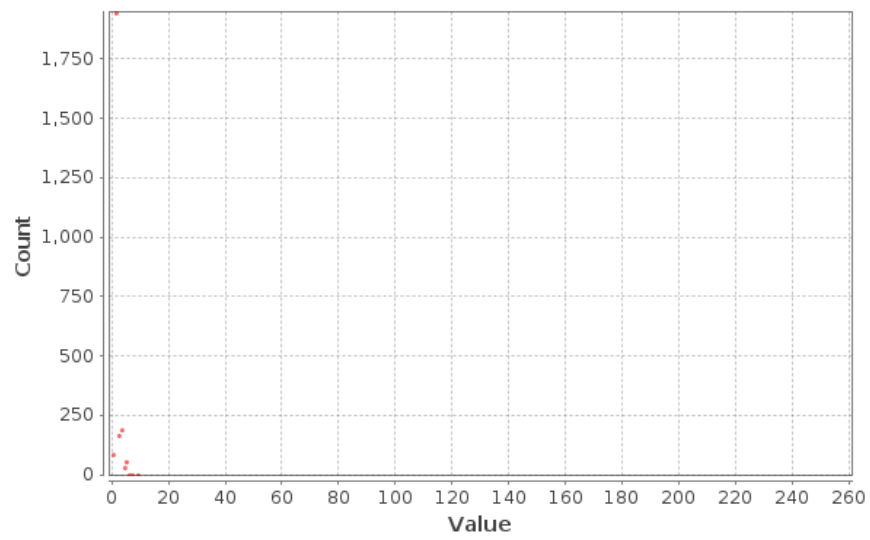


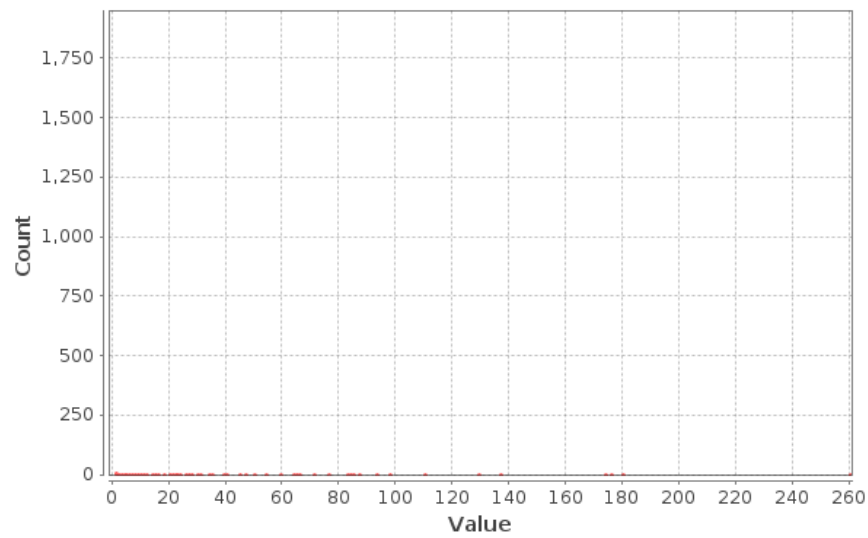
Fig. 1.3: Page Rank Results

1.3.2.4 Average Degrees

Average Degree: 1.347

Degree Distribution**In-Degree Distribution**

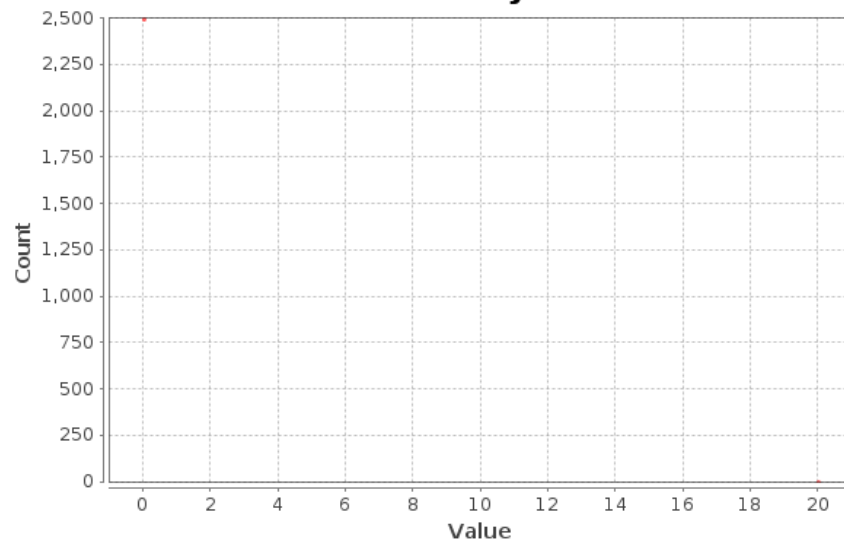
Out-Degree Distribution

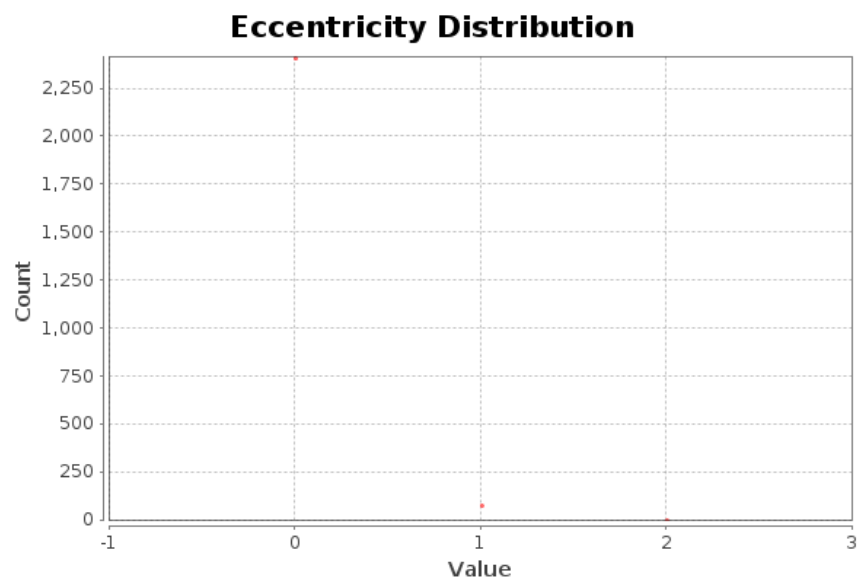
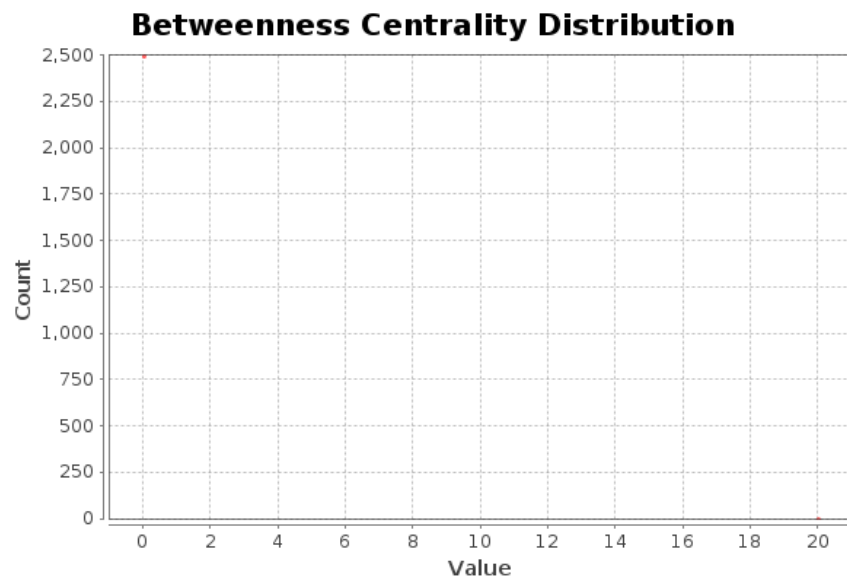


1.3.2.5 Network Diameter

Diameter: 2 Radius: 0 Average Path length: 1.0059294396679515 Number of shortest paths: 3373

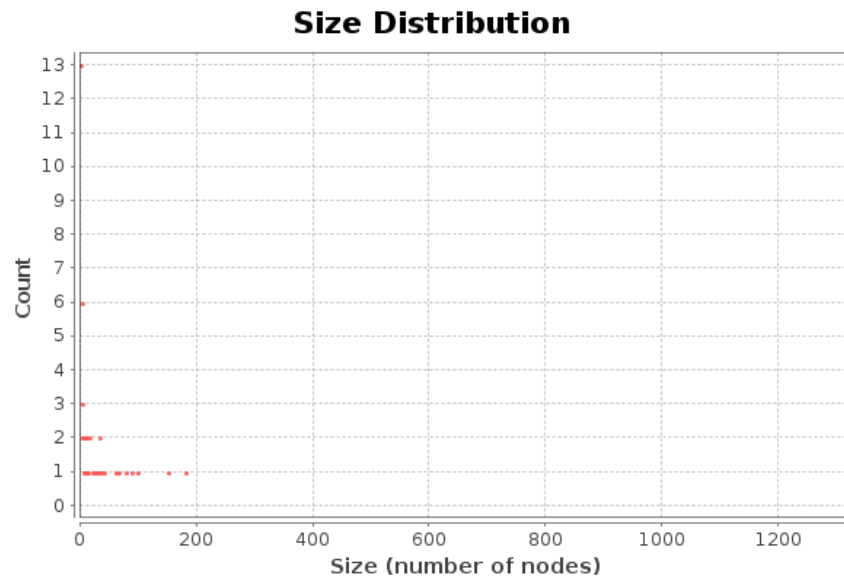
Betweenness Centrality Distribution





1.3.2.6 Connected Components

Number of Weakly Connected Components: 54 Number of Strongly Connected Components: 2498



References

1. <https://stat.ethz.ch/pipermail/r-help//2012-August/333656.html>
2. <http://digitalpbk.com/perl/perl-script-check-google-pagerank>
3. <http://www.google.com>
4. <http://jakevdp.github.io/blog/2012/10/14/scipy-sparse-graph-module-word-ladders/>
5. <http://curl.haxx.se/docs/httpscripting.html>
6. <http://www.crummy.com/software/BeautifulSoup/bs4/doc/>
7. <http://www.rmi.net/~lutz/>
8. <http://www.cs.cornell.edu/home/kleinber/networks-book/>
9. <http://thomassileo.com/blog/2013/01/25/using-twitter-rest-api-v1-dot-1-with-python/>
10. <http://www.cs.odu.edu/~mklein/cs796/lecture/>