

# INTRODUCTION TO WEB SCIENCES: Assignment 4

Babitha Bokka

10 october 2014

# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Question 1</b>                               | <b>2</b>  |
| 1.1      | Approach . . . . .                              | 2         |
| 1.2      | Description of extractHyperLinks.py . . . . .   | 2         |
| 1.3      | Source Code . . . . .                           | 3         |
| 1.3.1    | extractHyperLinks.py . . . . .                  | 3         |
| 1.4      | Input Files . . . . .                           | 5         |
| 1.4.1    | A2FinalOutput.txt . . . . .                     | 5         |
| 1.5      | Output Files . . . . .                          | 6         |
| 1.5.1    | md5Links.txt . . . . .                          | 6         |
| 1.5.2    | a524329028e8dcfc4879b4b453c1040a.link . . . . . | 6         |
| <b>2</b> | <b>Question 2</b>                               | <b>7</b>  |
| 2.1      | Description of createDotFile.py . . . . .       | 7         |
| 2.2      | Source Code . . . . .                           | 8         |
| 2.2.1    | createDotFile.py . . . . .                      | 8         |
| 2.3      | Input Files . . . . .                           | 10        |
| 2.3.1    | md5Links.txt . . . . .                          | 10        |
| 2.4      | Output Files . . . . .                          | 11        |
| 2.4.1    | mapping.txt . . . . .                           | 11        |
| <b>3</b> | <b>Question 3</b>                               | <b>12</b> |
| 3.1      | Approach . . . . .                              | 12        |
| 3.2      | Input Files . . . . .                           | 12        |
| 3.2.1    | mapping.txt . . . . .                           | 12        |
| 3.3      | Visualization with labels . . . . .             | 13        |
| 3.3.1    | Hits . . . . .                                  | 16        |
| 3.3.2    | PageRank . . . . .                              | 17        |
| 3.3.3    | Average . . . . .                               | 18        |
| 3.3.4    | Network Diameter . . . . .                      | 20        |
| 3.3.5    | Connected Components . . . . .                  | 22        |

# 1 Question 1

Choose 100 links from 1000 unique links and extract the outbound links from that page to other URIs.

## 1.1 Approach

To extract the outgoing links from the 100 URIs `extractHyperLinks.py` is loaded with the 1000 unique URIs. The `LIMIT` value in the program selects only 100 URIs, these are URIs which can establish the connection and has no exceptions.

## 1.2 Description of `extractHyperLinks.py`

1. Open the `A2FinalOutput.txt`.
2. Read each URL and extract the outgoing links.
3. Out of 1000 Unique URIs the Program will extract 100 working URIs and get the outgoing links from each URL.
4. Create a md5 file and write all the links to the “.link” file at the same time a log file (`md5Links.txt`) is generated which has the md5 and URL.
5. All the file are saved in `A4/Q1/links`.

## 1.3 Source Code

### 1.3.1 extractHyperLinks.py

```
1 #!/usr/bin/env python
2
3 import os
4 import md5
5 import sys
6 import time
7 import socket
8 import urllib2
9 import unicodedata
10 from bs4 import BeautifulSoup
11
12 LIMIT = 100
13
14 #Main Function
15 def main():
16     # open the 100 links file
17     uniqueLinks = open('A2FinalOutput.txt', 'r')
18     writeFile = open('md5Links.txt', 'w')
19
20     counter = 0
21     for url in uniqueLinks.readlines():
22         hashmd5 = md5.new(url).hexdigest()
23         filename = 'links/' + hashmd5 + '.link'
24         try:
25             request = urllib2.Request(url)
26             response = urllib2.urlopen(url, timeout=30)
27             html_content = response.read()
28
29             # get the html content using beautiful soup
30             soup = BeautifulSoup(html_content)
31             links = soup.find_all('a')
32
33             writeFile.write("{:<20} {}".format(hashmd5, url))
34             counter += 1
35
36             saveFile = open(filename, 'w')
37
38             for tag in links:
39                 try:
40                     link = tag.get('href', None)
41                     if link != None and link.startswith("http"):
42
43                         saveFile.write(link)
44                         saveFile.write('\n')
45
46                 except UnicodeEncodeError:
47                     pass
48             saveFile.close()
49
50         except urllib2.HTTPError:
51             pass
52         except urllib2.URLError:
53             pass
```

```
54         except socket.timeout :
55             pass
56         print filename , counter
57         print "_" * 50
58         # get only the 100 urls data
59         if counter >= LIMIT:
60             break
61
62
63 if __name__ == "__main__":
64     try:
65         main()
66     except KeyboardInterrupt:
67         sys.exit(1)
```

## 1.4 Input Files

### 1.4.1 A2FinalOutput.txt

```
1 Unique URIs
2
3 https://www.facebook.com/MyChickenRun
4 http://www.youtube.com/watch?v=Eubi9YI2dKE
5 http://geladoesntgiveadamn.tumblr.com/
6 http://youtu.be/1BKO2V9EaZ0?a
7 http://www.instagram.com/the_sanging_rebel
8 http://instagram.com/teustimao/
9 http://blogmylunch.com
10 http://Facebook.com/Robbiewhaylez
11 http://ifoodi.blogspot.com/
12 http://www.talkinggoodfood.co.uk
13 http://facebook.com/dimano.sterling
14 http://www.yellowkorner.com
15 http://Emerald.com
16 http://attackontiphan.tumblr.com
17 http://instagram.com/xxmn88
18 http://linggez-network.blogspot.com
19 http://Instagram.com/helloalm_
20 http://www.wagamama.com
21 http://pennyroyaltea.co.vu
22 http://www.oufancyphones.com
23 http://www.facebook.com/rappstartailgate
24 http://thedaintypig.com
25 http://instagram.com/victorparrini
26 http://instagram.com/_lafamiliaonly
27 http://rad-acid.tumblr.com
28 https://www.Youtube.com/iFarLiez
29 http://twitter.com
30 http://melonpatchtv.com
31 http://linkd.in/1qVvTtZ
32 http://www.huffingtonpost.com/alex-palombo/
33 http://youtube.com/user/IvanAlvir
34 http://www.musicsumo.com
35 http://www.KhanaPakana.com
```

## 1.5 Output Files

Aim : To generate the 100 “.link” files which have the outbound links from that page to other URIs. The files are saved in A4/Q1/links.

### 1.5.1 md5Links.txt

The file acts as a logfile to keep track of URI and corresponding md5 file.

| md5                              | file name | URI   |
|----------------------------------|-----------|---|
| 5414920b1ad98e5a5c2fe91a48c78071 |           | <a href="https://www.facebook.com/MyChickenRun">https://www.facebook.com/MyChickenRun</a>           |
| 8a7aaf3d96587b6edd9381207370a03d |           | <a href="http://www.youtube.com/watch?v=Eubi9YI2dKE">http://www.youtube.com/watch?v=Eubi9YI2dKE</a> |
| a524329028e8dcfc4879b4b453c1040a |           | <a href="http://geladoesntgiveadamn.tumblr.com/">http://geladoesntgiveadamn.tumblr.com/</a>         |
| 96dec514297f45ff7a40f0ec5b527dfc |           | <a href="http://youtu.be/1BKO2V9EaZ0?a">http://youtu.be/1BKO2V9EaZ0?a</a>                           |
| 10d84f5c5226fe47690138352748aa07 |           | <a href="http://www.instagram.com/the_sanging_rebel">http://www.instagram.com/the_sanging_rebel</a> |
| 75c522f8af14203e53d45345fffa79ea |           | <a href="http://instagram.com/teustimao/">http://instagram.com/teustimao/</a>                       |
| c7a1e9146029a56d3fa0acab5d66688b |           | <a href="http://blogmylunch.com">http://blogmylunch.com</a>   |
| 42305950be948963a249dc2e7c82e581 |           | <a href="http://Facebook.com/Robbiewaylez">http://Facebook.com/Robbiewaylez</a>                     |
| aec16fa2999f0e1ed985ba24fc9793d0 |           | <a href="http://ifoodi.blogspot.com/">http://ifoodi.blogspot.com/</a>                               |
| db4488dcef3ed8da9fc7bb6b459c9cde |           | <a href="http://www.talkinggoodfood.co.uk">http://www.talkinggoodfood.co.uk</a>                     |

### 1.5.2 a524329028e8dcfc4879b4b453c1040a.link

A sample “.link” file having the outgoing links.

```
1 URI
2 http://geladoesntgiveadamn.tumblr.com/
3 Links
4 http://www.cherrybam.com/blackandwhite-tumblr-themes.php
5 http://www.cherrybam.com
6 https://www.tumblr.com/reblog/81587066245/kEKZcUHw
7 http://geladoesntgiveadamn.tumblr.com/post/81586972499/toxicgaskarth-all-time-low-so-long-soldier
8 https://www.tumblr.com/reblog/81586972499/CUjrIucX
9 http://geladoesntgiveadamn.tumblr.com/post/81586832033/jackbarakatofficial-all-time-low-inspired-by
10 https://www.tumblr.com/reblog/81586832033/4BCaP7rD
11 http://geladoesntgiveadamn.tumblr.com/post/81586733974/q-whats-the-worst-excuse-you-have-ever-used-for
12 https://www.tumblr.com/reblog/81586733974/KrMsNi8q
13 http://geladoesntgiveadamn.tumblr.com/post/81586630503
14 https://www.tumblr.com/reblog/81586630503/2ZXuDBxA
15 http://www.flickr.com/photos/thenmaysuhsaid/5744956806/
16 http://geladoesntgiveadamn.tumblr.com/post/81585758862/yeahbarakat-the-maine-by-maysa-askar-on-flickr
17 https://www.tumblr.com/reblog/81585758862/jBh5DfiU
18 http://geladoesntgiveadamn.tumblr.com/post/81585751637/john-who-the-fuck-even-sent-you-ocallaghan-x
19 https://www.tumblr.com/reblog/81585751637/zBgMQ9uU
20 http://geladoesntgiveadamn.tumblr.com/post/80128811883
21 https://www.tumblr.com/reblog/80128811883/8GQakRrF
22 http://foreverhalloween.us/
23 http://geladoesntgiveadamn.tumblr.com/post/80128520734/austincxrliles-original-credit
24 https://www.tumblr.com/reblog/80128520734/Y6WamLQe
25 http://geladoesntgiveadamn.tumblr.com/post/77538402738
26 https://www.tumblr.com/reblog/77538402738/2LVznsYP
```

## 2 Question 2

Generate a single DOT file for 100 links.

### 2.1 Description of createDotFile.py

1. Open the md5Links.txt, extract the url and md5.
2. Using glob to find the respective folder.
3. Check whether md5 of the current URL and the file name in the folder are matching if, they match then open the file and append the links inside each file to the URL
4. Write all the URL and respective links to “mapping.dot” file.



## 2.2 Source Code

### 2.2.1 createDotFile.py

```
1 #!/usr/bin/env python
2
3 import os
4 import sys
5 import glob
6
7
8 #Main Function
9 def main():
10
11     path = "/home/bbokka/cs594/A4/Q1/links/*.link"
12     md5_Url = open('md5Links.txt','r')
13     mapping = open('mapping.dot','w')
14     mapping.write('digraph A4_question3 { size = "6,6"; node [color = lightblue2 ,
15 style = filled]; ')
16     for line in md5_Url.readlines():
17         line = line.split(" ")
18         md5 = line[0]
19         url = line[1]
20
21         first_spit = url.split("://")
22         store_first_spit = first_spit[1]
23         second_spit = store_first_spit.split("/")
24         store_second_spit = second_spit[0]
25         url_label_string = store_second_spit
26         #print label_string
27
28     md5_file = glob.glob(path)
29     for each_md5_file in md5_file:
30         #print each_md5_file
31         filename = each_md5_file.split("/links/")
32         md5_file_name = filename[1]
33         md5_code = md5_file_name.split(".link")
34         md5_code_file = md5_code[0]
35         #print md5_code_file
36         if (md5 == md5_code_file):
37             open_md5 = open(each_md5_file, 'r')
38             for link in open_md5.readlines():
39                 first_spit_1 = link.strip().split("://")
40                 store_first_spit_1 = first_spit_1[1]
41                 second_spit_1 = store_first_spit_1.split("/")
42                 store_second_spit_1 = second_spit_1[0]
43                 link_label_string = store_second_spit_1
44                 link = link.strip()
45                 url = url.strip()
46                 #string = ' "' + url.strip() + '"' + '->' + ' "' + links.strip()
47                 string = ' "' + url.strip() + '"' + '->' + ' "' + link.
48 strip() + '"' + ' \n ' + '"' + url.strip() + '"' + ' [label = ' + '"' +
49 url_label_string + '"' + ']; '+' \n ' + '"' + links + '"' + ' [label = ' + '"' +
50 link_label_string + '"' + ']; '+';'
```

```
49     mapping.write ( '}' )
50
51
52
53 if __name__ == "__main__":
54     try:
55         main()
56     except KeyboardInterrupt:
57         sys.exit(1)
```

## 2.3 Input Files

### 2.3.1 md5Links.txt

| md5                              | file name | URI   |
|----------------------------------|-----------|---|
| 5414920b1ad98e5a5c2fe91a48c78071 |           | <a href="https://www.facebook.com/MyChickenRun">https://www.facebook.com/MyChickenRun</a>           |
| 8a7aaf3d96587b6edd9381207370a03d |           | <a href="http://www.youtube.com/watch?v=Eubi9YI2dKE">http://www.youtube.com/watch?v=Eubi9YI2dKE</a> |
| a524329028e8dcfc4879b4b453c1040a |           | <a href="http://geladoesntgiveadamn.tumblr.com/">http://geladoesntgiveadamn.tumblr.com/</a>         |
| 96dec514297f45ff7a40f0ec5b527dfc |           | <a href="http://youtu.be/1BKO2V9EaZ0?a">http://youtu.be/1BKO2V9EaZ0?a</a>                           |
| 10d84f5c5226fe47690138352748aa07 |           | <a href="http://www.instagram.com/the_sanging_rebel">http://www.instagram.com/the_sanging_rebel</a> |
| 75c522f8af14203e53d45345ffa79ea  |           | <a href="http://instagram.com/teustimao/">http://instagram.com/teustimao/</a>                       |
| c7a1e9146029a56d3fa0acab5d66688b |           | <a href="http://blogmylunch.com">http://blogmylunch.com</a>   |
| 42305950be948963a249dc2e7c82e581 |           | <a href="http://Facebook.com/Robbiewhaylez">http://Facebook.com/Robbiewhaylez</a>                   |
| aec16fa2999f0e1ed985ba24fc9793d0 |           | <a href="http://ifoodi.blogspot.com/">http://ifoodi.blogspot.com/</a>                               |
| db4488dcef3ed8da9fc7bb6b459c9cde |           | <a href="http://www.talkinggoodfood.co.uk">http://www.talkinggoodfood.co.uk</a>                     |

## 2.4 Output Files

### 2.4.1 mapping.txt

```
1 digraph A4_question3 {
2   size = "6,6";
3   node [color = lightblue2 ,style = filled];
4   "https://www.facebook.com/MyChickenRun" -> "https://www.facebook.com/recover/
   initiate"
5   "https://www.facebook.com/MyChickenRun" [label = "www.facebook.com"];
6   "https://www.facebook.com/recover/initiate" [label = "www.facebook.com"];
7   "http://www.youtube.com/watch?v=Eubi9YI2dKE" -> "https://plus.google.com/+youtube"
8   "http://www.youtube.com/watch?v=Eubi9YI2dKE" [label = "www.youtube.com"];
9   "https://plus.google.com/+youtube" [label = "plus.google.com"];
10  "http://www.youtube.com/watch?v=Eubi9YI2dKE" -> "https://www.google.com/intl/en/
   policies/privacy/"
11  "http://www.youtube.com/watch?v=Eubi9YI2dKE" [label = "www.youtube.com"];
12  "https://www.google.com/intl/en/policies/privacy/" [label = "www.google.com"];
13  "http://www.youtube.com/watch?v=Eubi9YI2dKE" -> "https://accounts.google.com/
   ServiceLogin?continue=https%3A%2F%2Fwww.youtube.com%2Fsignin%3Fhl%3Den%26app%3
   Ddesktop%26next%3D%252Fwatch%253Fv%253DEubi9YI2dKE%26feature%3Dplaylist%26
   action_handle_signin%3Dtrue&hl=en&service=youtube&passive=true&uilel=3"
14  "http://www.youtube.com/watch?v=Eubi9YI2dKE" [label = "www.youtube.com"]; "https://
   accounts.google.com/ServiceLogin?continue=https%3A%2F%2Fwww.youtube.com%2Fsignin%3
   Fhl%3Den%26app%3Ddesktop%26next%3D%252Fwatch%253Fv%253DEubi9YI2dKE%26feature%3
   Dplaylist%26action_handle_signin%3Dtrue&hl=en&service=youtube&passive=true&uilel
   =3" [label = "accounts.google.com"];
15 }
```

## 3 Question 3

Download and install Gephi then load the dot file to visualize the graph.

### 3.1 Approach

Downloaded gephi and loaded the mapping.dot file into the software. In order to get more insight of the visualization layouts Fruchterman, Reingold and Yifan Hu and Yifan Hu Proportional are used. Run the layouts couple of time to get more insight and export them to a pdf or png format.

### 3.2 Input Files

The input is the output of the Question 2.

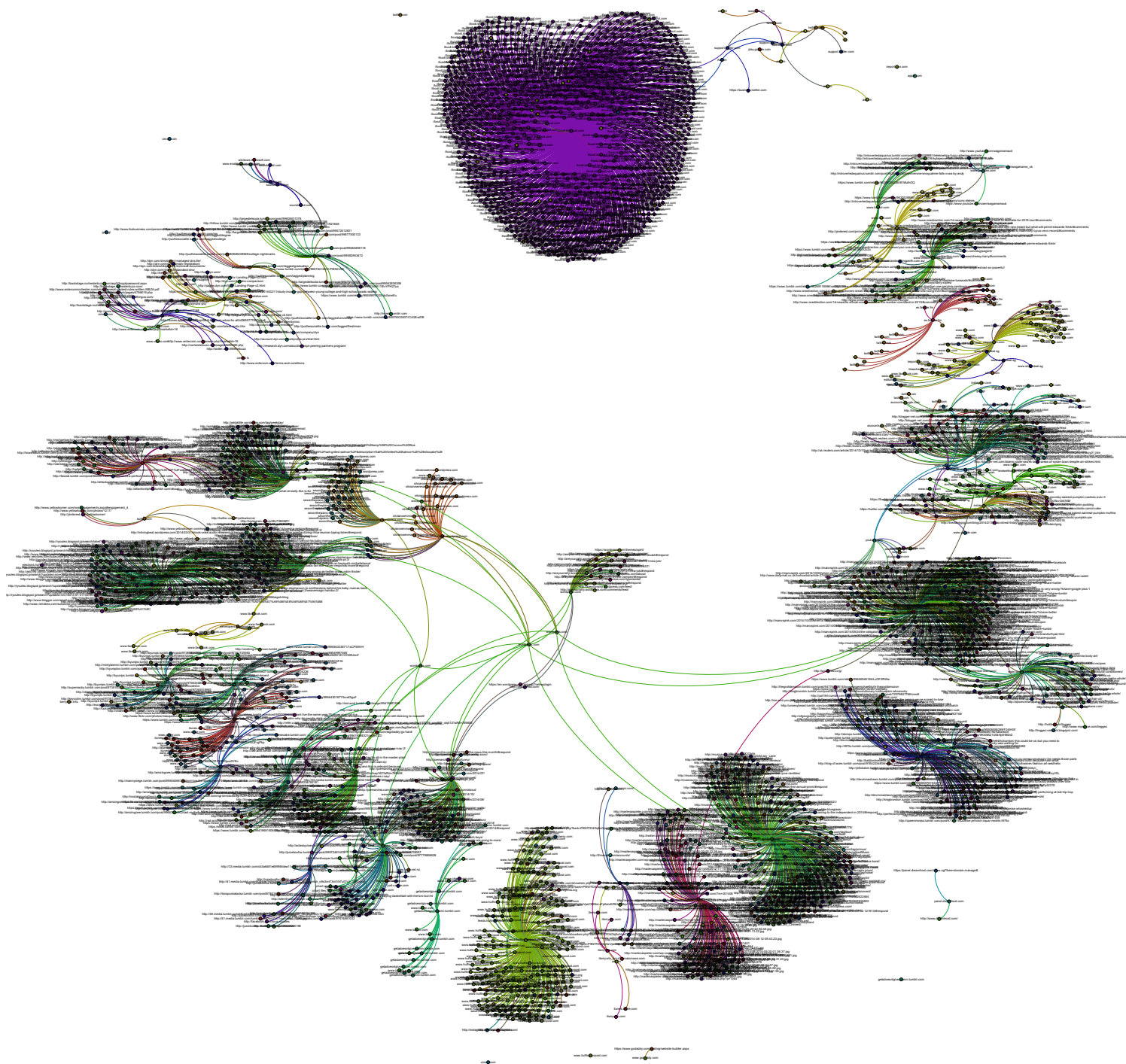
#### 3.2.1 mapping.txt

```
1 digraph A4_question3 {
2     size = "6,6";
3     node [color = lightblue2 ,style = filled];
4     "https://www.facebook.com/MyChickenRun" -> "https://www.facebook.com/recover/
    initiate"
5     "https://www.facebook.com/MyChickenRun" [label = "www.facebook.com"];
6     "https://www.facebook.com/recover/initiate" [label = "www.facebook.com"];
7     "http://www.youtube.com/watch?v=Eubi9YI2dKE" -> "https://plus.google.com/+youtube"
8     "http://www.youtube.com/watch?v=Eubi9YI2dKE" [label = "www.youtube.com"];
9     "https://plus.google.com/+youtube" [label = "plus.google.com"];
10    "http://www.youtube.com/watch?v=Eubi9YI2dKE" -> "https://www.google.com/intl/en/
    policies/privacy/"
11    "http://www.youtube.com/watch?v=Eubi9YI2dKE" [label = "www.youtube.com"];
12    "https://www.google.com/intl/en/policies/privacy/" [label = "www.google.com"];
13    "http://www.youtube.com/watch?v=Eubi9YI2dKE" -> "https://accounts.google.com/
    ServiceLogin?continue=https%3A%2F%2Fwww.youtube.com%2Fsignin%3Fhl%3Den%26app%3
    Ddesktop%26next%3D%252Fwatch%253Fv%253DEubi9YI2dKE%26feature%3Dplaylist%26
    action_handle_signin%3Dtrue&hl=en&service=youtube&passive=true&uilel=3"
14    "http://www.youtube.com/watch?v=Eubi9YI2dKE" [label = "www.youtube.com"]; "https://
    accounts.google.com/ServiceLogin?continue=https%3A%2F%2Fwww.youtube.com%2Fsignin%3
    Fhl%3Den%26app%3Ddesktop%26next%3D%252Fwatch%253Fv%253DEubi9YI2dKE%26feature%3
    Dplaylist%26action_handle_signin%3Dtrue&hl=en&service=youtube&passive=true&uilel
    =3" [label = "accounts.google.com"];
15 }
```

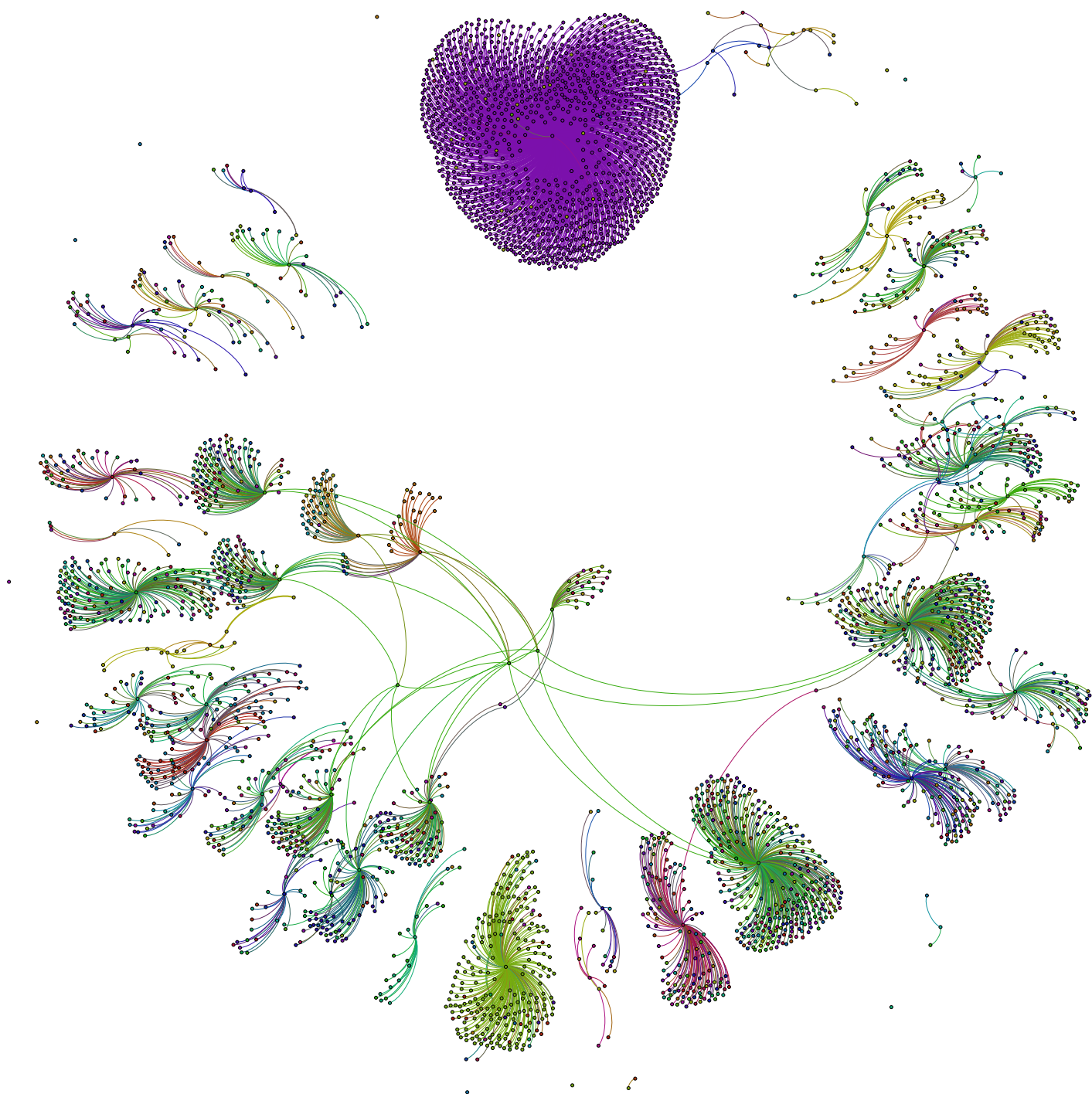
### 3.3 Visualization with labels

The below graph depicts that there are not much connected components very few of the clusters are connected and few of them does not have outgoing links.

The second figure shows you how the nodes and edges are connected and it clearly depicts whether the clusters are connected or not connected when compared to the visualization with labels. A clear insight can be obtained from `connected.pdf` by enlarging.









### 3.3.1 Hits

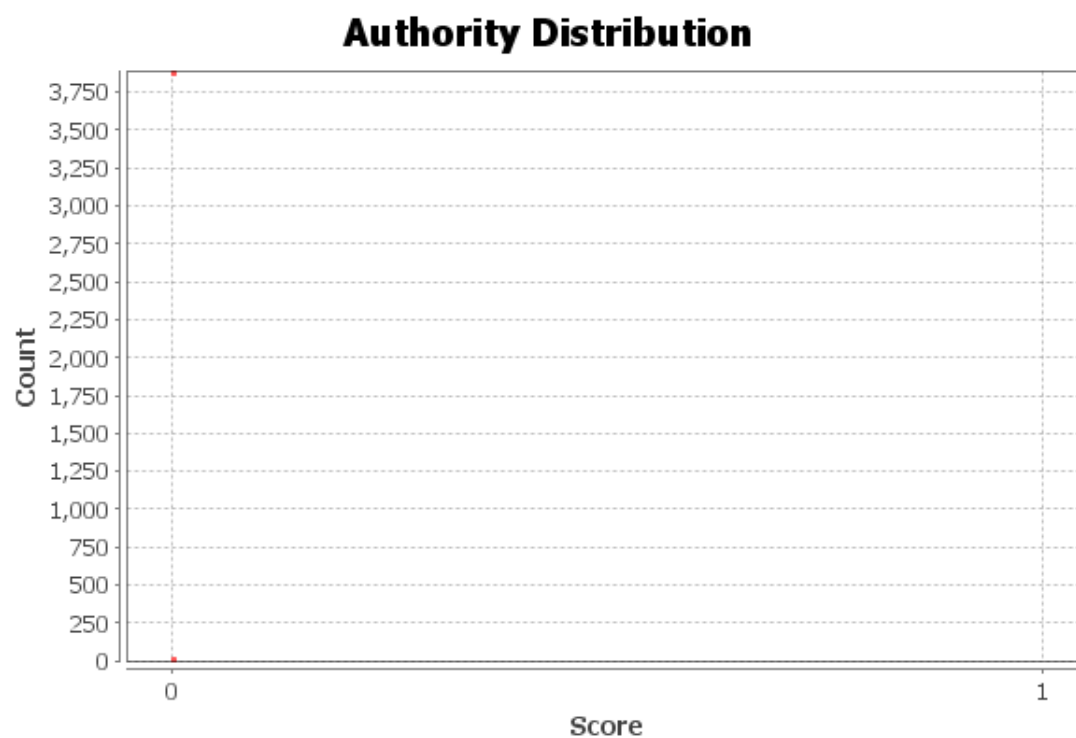


Figure 1: Authorities

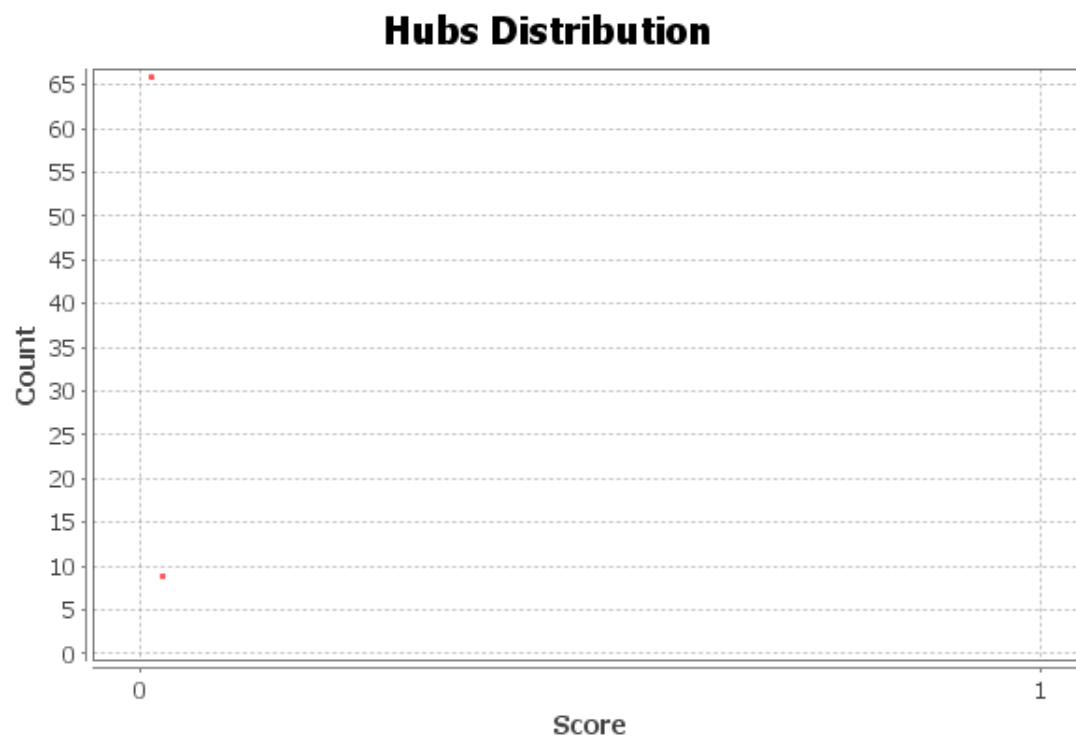


Figure 2: Hubs

### 3.3.2 PageRank

Epsilon = 0.001

Probability = 0.85

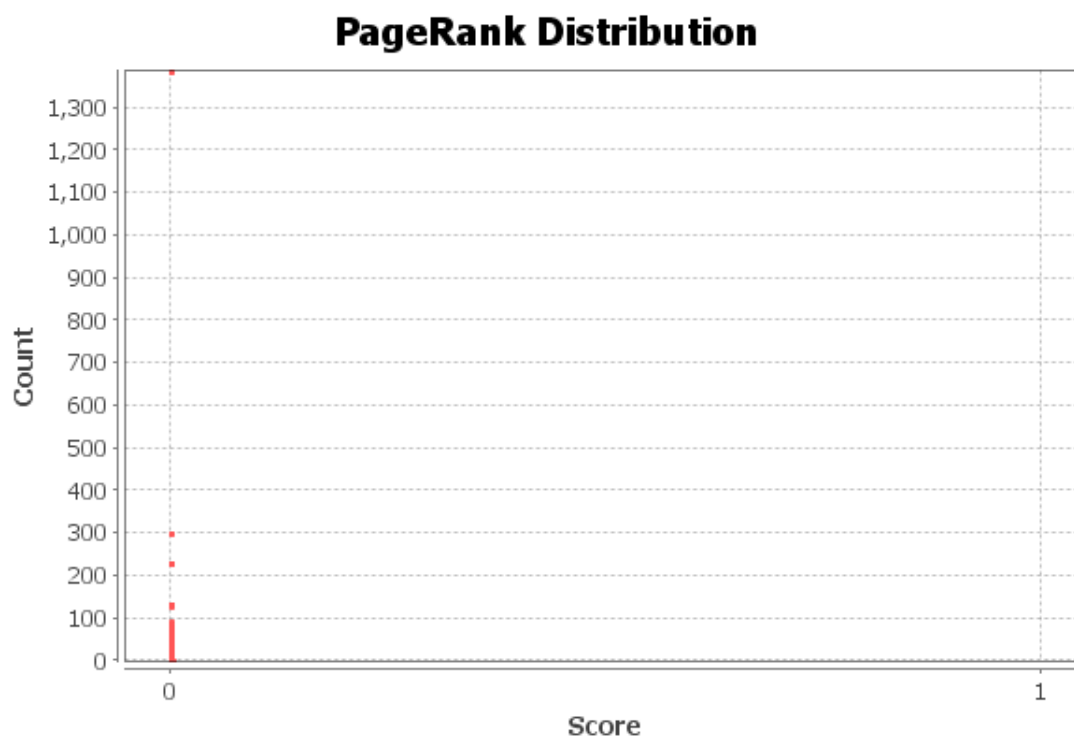


Figure 3: pagerank of links

### 3.3.3 Average

Average Degree: 0.997

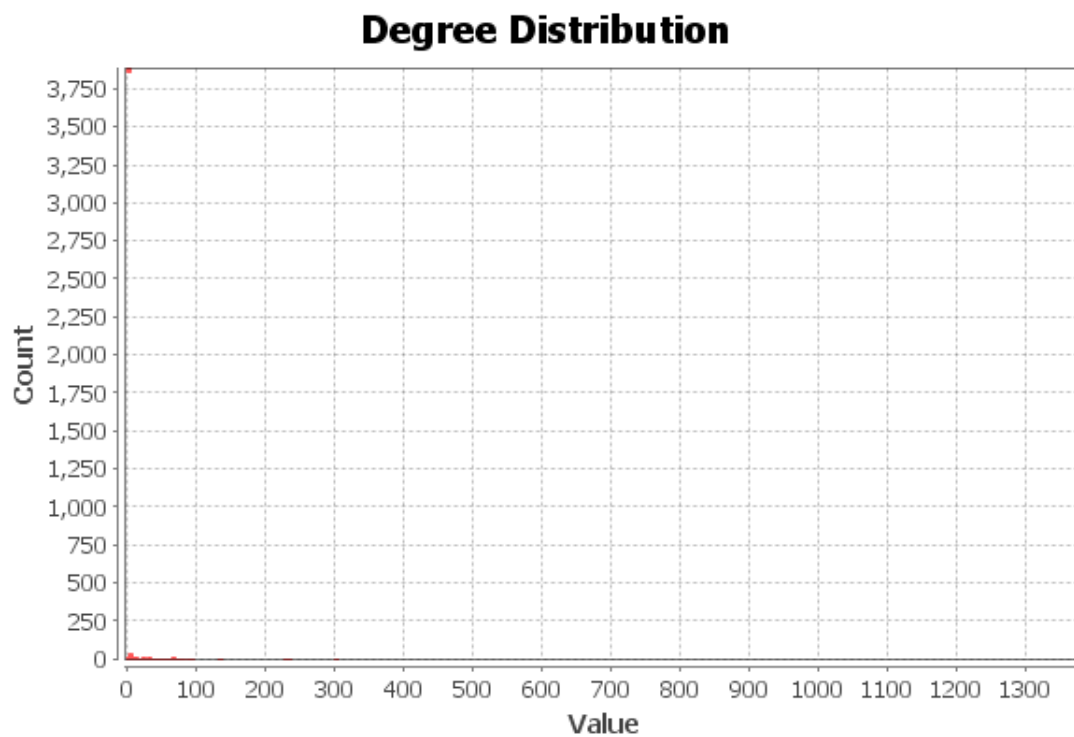


Figure 4: Average degree-distribution

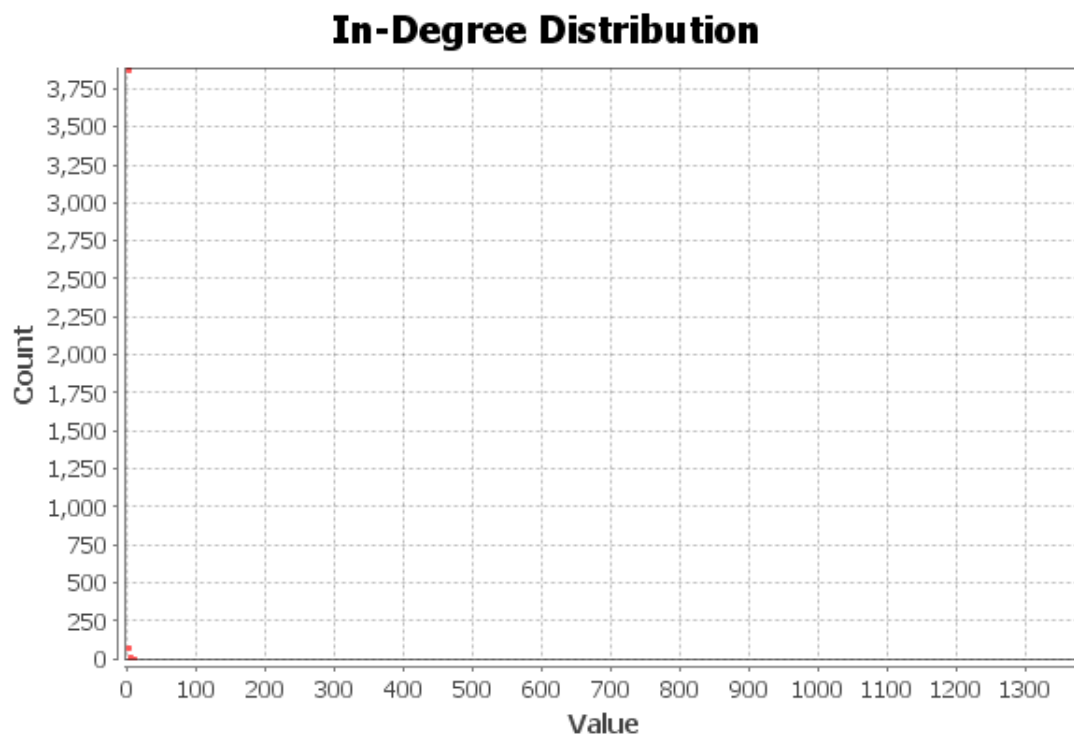


Figure 5: Average indegree-distribution

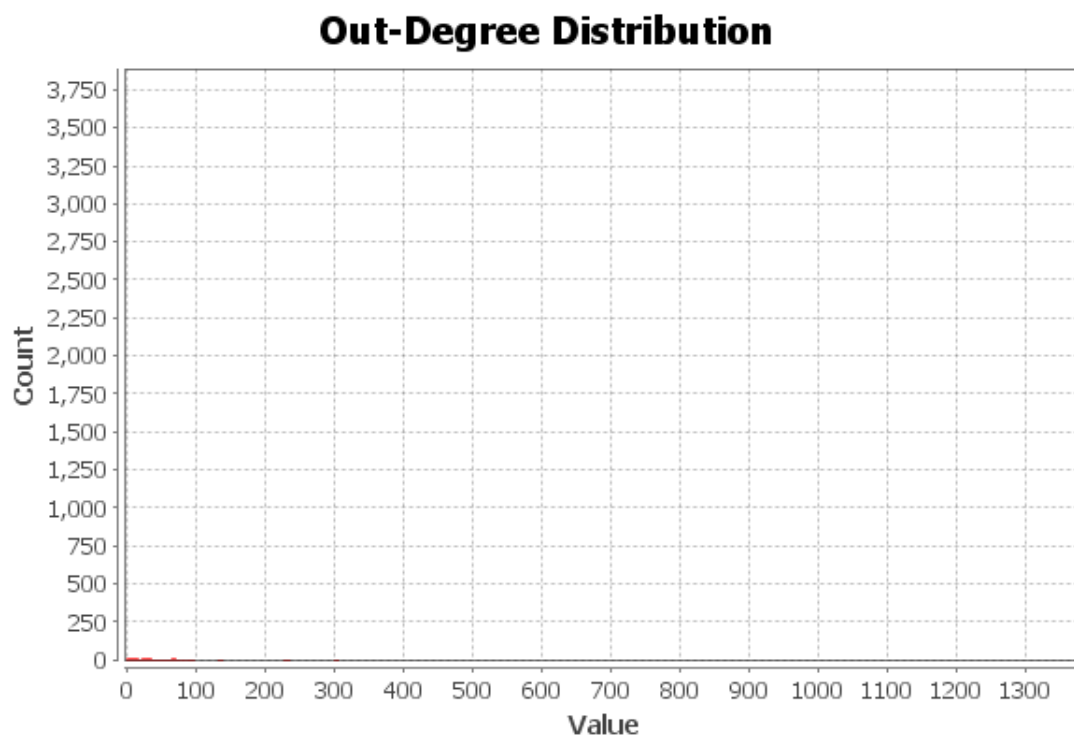


Figure 6: average outdegree-distribution

### 3.3.4 Network Diameter

Diameter: 1

Radius: 0

Average Path length: 1.0

Number of shortest paths: 3969

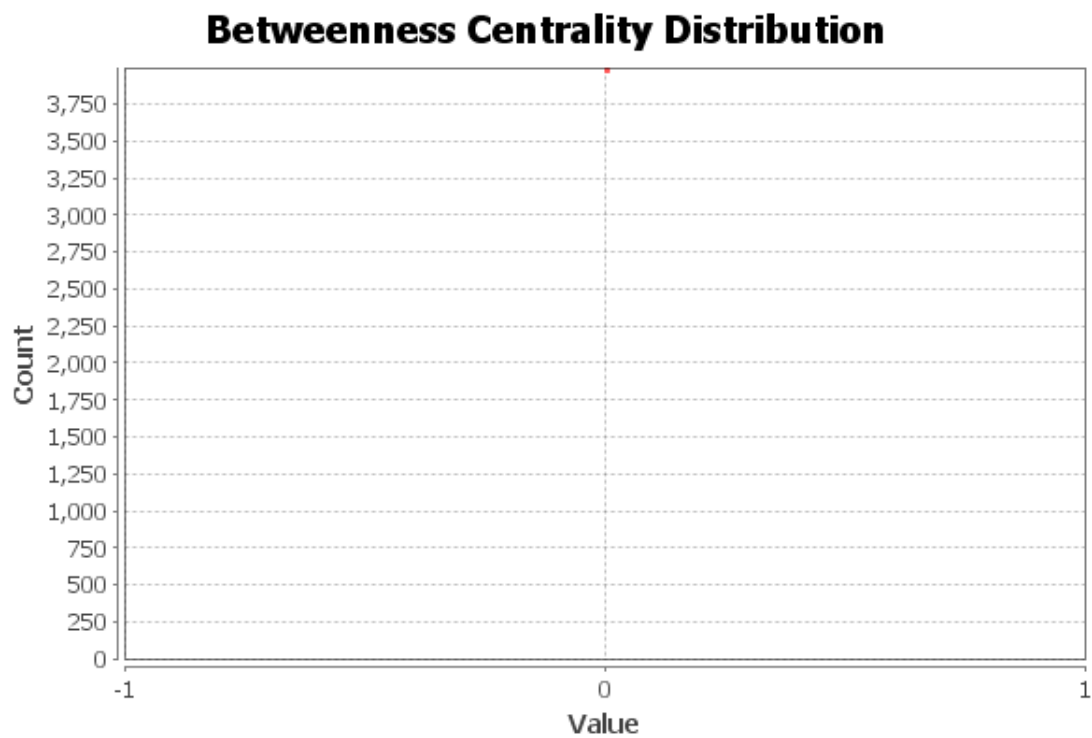


Figure 7: Betweenness

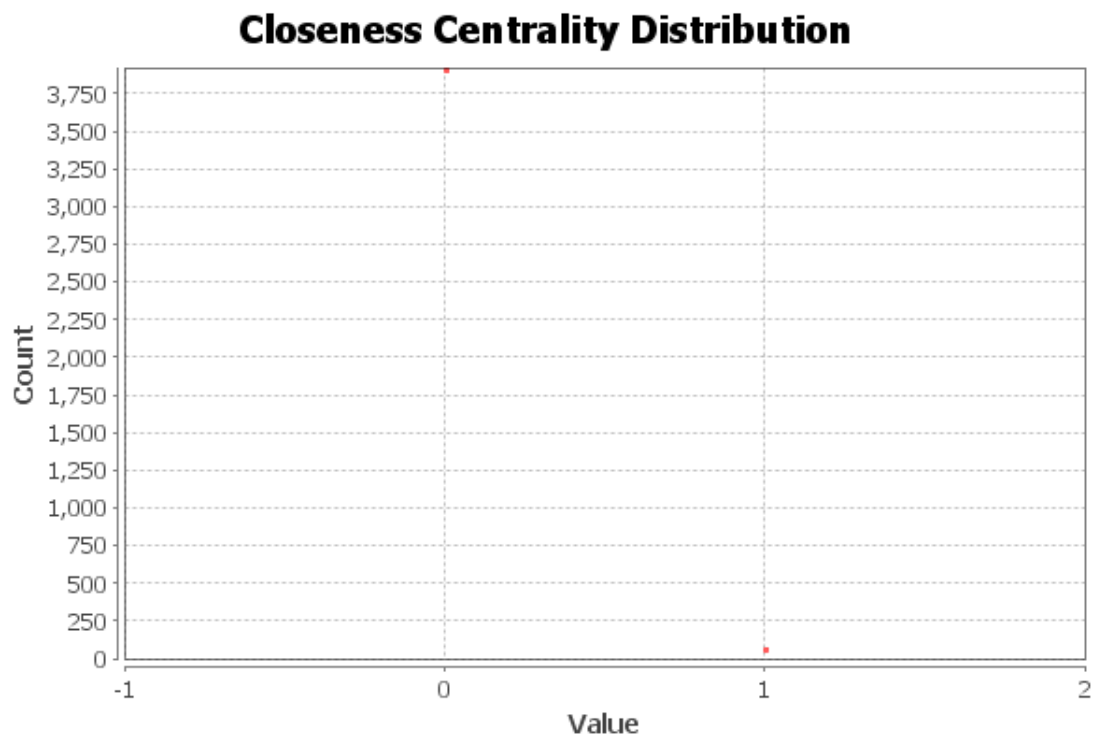


Figure 8: Closeness

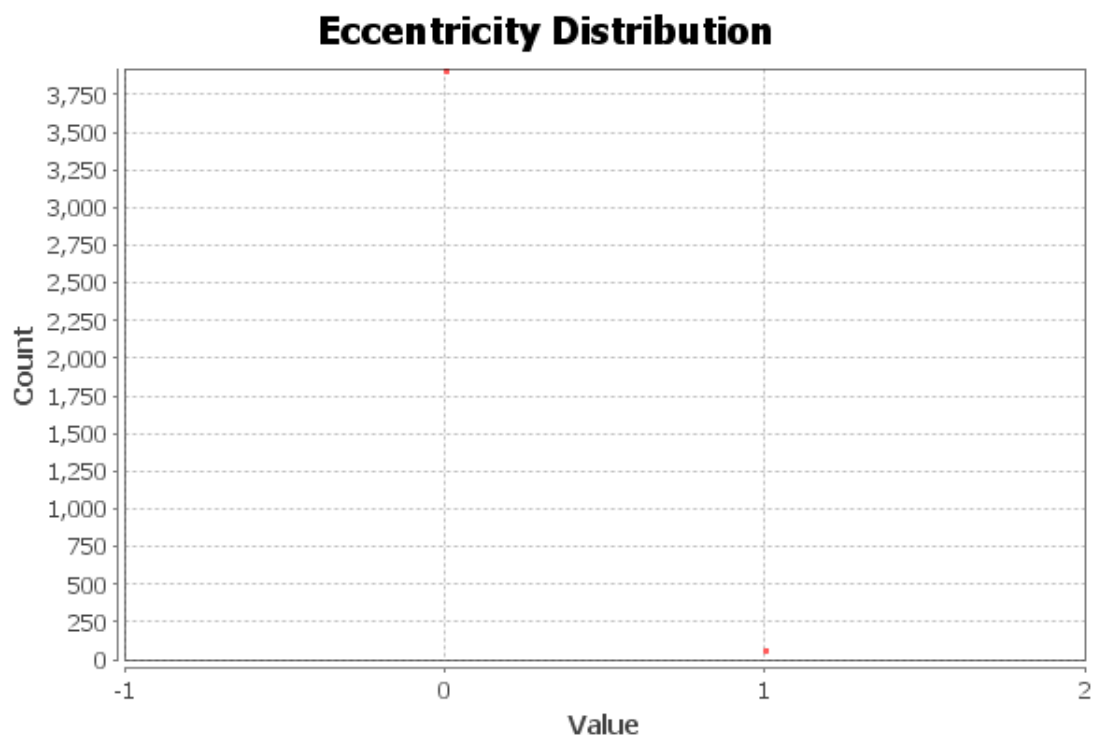


Figure 9: EccentricityDistribution

### 3.3.5 Connected Components

Number of Weakly Connected Components: 58

Number of Strongly Connected Components: 3990

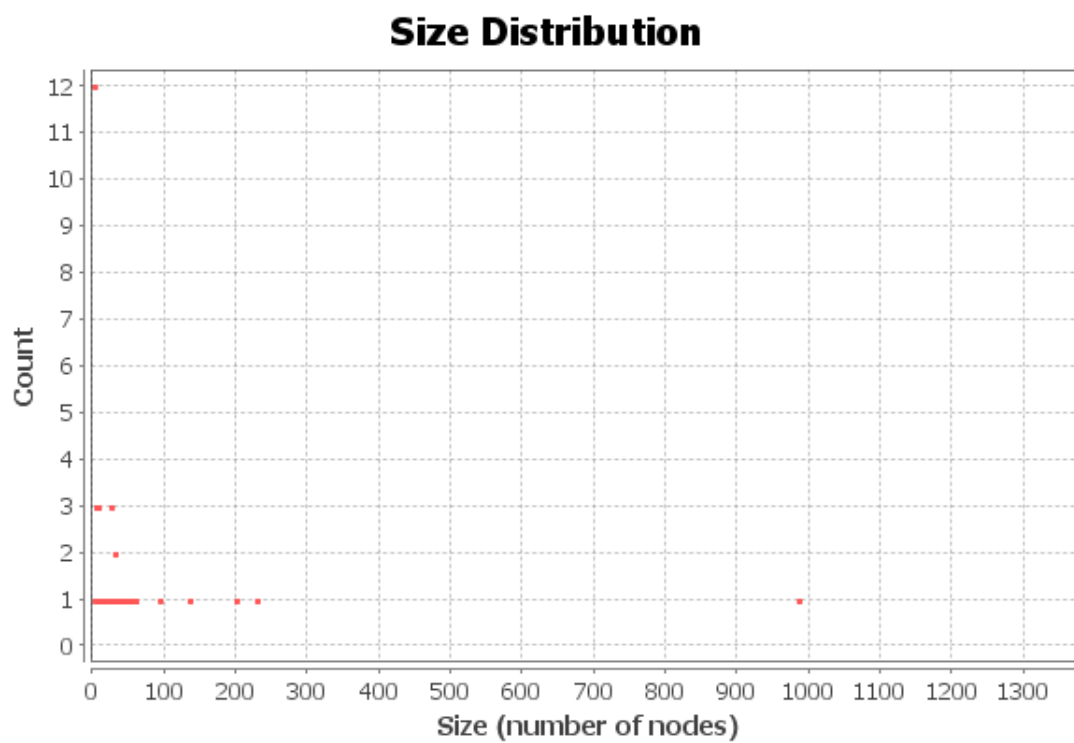


Figure 10: Graph Connected Components

## References

- [1] python. <http://docs.python-requests.org/en/latest/user/quickstart/>.
- [2] python. <http://askubuntu.com/questions/352198/reading-all-files-from-a-directory>.
- [3] python. <http://stackoverflow.com/questions/5815747/beautifulsoup-getting-href>.
- [4] Urlunicode. <https://docs.python.org/2/howto/unicode.html>.
- [5] Urlunicode. <http://stackoverflow.com/questions/19212306/difference-between-ascii-and-unicode>.

□