

Assignment 2

Fall 2014

CS595 Web Science

Dr. Michael Nelson

Mathew Chaney

September 26, 2014

Contents

1	Question 1	3
1.1	Question	3
1.2	Resources	3
1.3	Answer	3
2	Question 2	5
2.1	Question	5
2.2	Resources	5
2.3	Answer	5
3	Question 3	7
3.1	Question	7
3.2	Resources	7
3.3	Answer	7

List of Figures

1	Histogram of Site Mementos	6
2	Estimated Age to Memento Count	9

Listings

1	urifinder.py	3
2	mementofinder.py	5
3	build_histogram.r	6
4	carbodate.py	7
5	build_ecd_mementos.py	8
6	ecd_mementos_graph.r	8

1 Question 1

1.1 Question

Write a Python program that extracts 1000 unique links from Twitter. You might want to take a look at:

<http://thomassileo.com/blog/2013/01/25/using-twitter-rest-api-v1-dot-1-with-python/>

But there are many other similar resources available on the web. Note that only Twitter API 1.1 is currently available; version 1 code will no longer work.

Also note that you need to verify that the final target URI (i.e., the one that responds with a 200) is unique. You could have different shortened URIs for www.cnn.com. For example,

<http://cnn.it/1cTNZ3V> <http://t.co/BiYdsGotTd>

Both ultimately redirect to [cnn.com](http://www.cnn.com), so they count as only 1 unique URI. Also note the second URI redirects twice – don't stop at the first redirect.

You might want to use the search feature to find URIs, or you can pull them from the feed of someone famous (e.g., Tim O'Reilly).

Hold on to this collection – we'll use it later throughout the semester.

1.2 Resources

- Getting Started with Twitter API: <http://thomassileo.com/blog/2013/01/25/using-twitter-rest-api-v1-dot-1-with-python/>
- Twitter Search API: <https://dev.twitter.com/rest/public/search>
- Twitter API - Get / Search Tweets: <https://dev.twitter.com/rest/reference/get/search/tweets>

1.3 Answer

Using the python module requests made this task a breeze as well as the initial code provided by Thomas Sileo's blog post.

```
1 # -*- encoding: utf-8 -*-
2 import requests
3 from requests_oauthlib import OAuth1
4 from urllib import quote
5
6 REQUEST_TOKEN_URL = "https://api.twitter.com/oauth/request_token"
7 AUTHORIZE_URL = "https://api.twitter.com/oauth/authorize?oauth_token="
8 ACCESS_TOKEN_URL = "https://api.twitter.com/oauth/access_token"
9
10 CONSUMER_KEY = "PDBekXkvUto4V0XYZrrizcEub"
11 CONSUMER_SECRET = "0E1KNfpNWF8Eh4iwbDFYFKDMBSouni3uRZrpsoGhbJcLZZmnBq"
12
13 OAUTH_TOKEN = "2560074793-g7ES1sQmw13YKAfCJnIBa0lh3wHLjnPqj96XFuV"
14 OAUTH_TOKEN_SECRET = "tGYCQa9LL2i6wmApJzbGzHdVIVA65xiwVffPmbqWJwZPs"
15
16 SEARCH_URI = "https://api.twitter.com/1.1/search/tweets.json?q="
17
18 SEARCH_ITEMS = map(quote, [ 'space x',
19                             'elon musk',
```

```

20         'richard garriott',
21         'starcraft 2',
22         'ebola virus',
23         'world cup',
24         'singularity',
25         'rick and morty',
26         'iphone 6',
27         'android',
28         'robin williams',
29         'tony stewart',
30         'bitcoin',
31         'game of thrones',
32         'facebook',
33         'youtube',
34         'google',
35         'chris roberts',
36         'hyper light drifter',
37         'golang'])
38
39 def get_oauth():
40     return OAuth1(CONSUMER_KEY,
41                   client_secret=CONSUMER_SECRET,
42                   resource_owner_key=OAUTH_TOKEN,
43                   resource_owner_secret=OAUTH_TOKEN_SECRET)
44
45 def find_uris(uris):
46     with open('output', 'a') as outfile:
47         for search_item in SEARCH_ITEMS:
48             result = requests.get(SEARCH_URI + search_item + '&filter%3Alinks&count=1000',
49                                   auth=oauth)
49             for status in result.json()['statuses']:
50                 for url in status['entities']['urls']:
51                     if len(uris) == 1000:
52                         return
53                     if 'expanded_url' in url:
54                         try:
55                             result = requests.get(url['expanded_url'], timeout=4)
56                             # only add expanded uris if they aren't in the list already
57                             if result.status_code == 200 and result.url not in uris:
58                                 add_uri(uris, result.url)
59                                 outfile.write('%s\n' % result.url)
60                         except Exception as e:
61                             print e
62                             continue
63
64 def add_uri(uris, uri):
65     uris.add(uri)
66     print 'added uri #%d: %s' % (len(uris), uri)
67
68 if __name__ == "__main__":
69     oauth = get_oauth()
70     uris = set()
71     # read in previous set of uris
72     try:
73         with open('output', 'r') as infile:
74             for line in infile.readlines():
75                 add_uri(uris, line.strip())
76     except IOError:
77         pass
78     find_uris(uris)

```

Listing 1: urifinder.py

The script was run multiple times to get the desired 1000 unique URIs. It would end prematurely at times, so the data set was initialized with the data of the previous run and then passed on to the `find_uris` function to preserve work performed.

2 Question 2

2.1 Question

Download the TimeMaps for each of the target URIs. We'll use the mementoweb.org Aggregator, so for example:

URI-R = <http://www.cs.odu.edu/>

URI-T = <http://mementoweb.org/timemap/link/http://www.cs.odu.edu/>

You could use the cs.odu.edu aggregator:

URI-T = <http://mementoproxy.cs.odu.edu/aggr/timemap/link/1/http://www.cs.odu.edu/>

But be sure to say which aggregator you use – they are likely to give different answers.

Create a histogram of URIs vs. number of Mementos (as computed from the TimeMaps). For example, 100 URIs with 0 Mementos, 300 URIs with 1 Memento, 400 URIs with 2 Mementos, etc.

See: <http://en.wikipedia.org/wiki/Histogram>

Note that the TimeMaps can span multiple pages. Look for links like:

<<http://mementoweb.org/timemap/link/1000/http://www.cnn.com/>>;rel="timemap"; type="application/link-format"; from = "Sun, 08 Jul 2001 21:30:54 GMT"

This indicates another page of the TimeMap is available. There can be many pages to a TimeMap.

2.2 Resources

- R: <http://www.cs.odu.edu/~sainswor/uploads/Teaching/cs595f13-R.pdf>

2.3 Answer

The python script in Listing 2 was used to retrieve the timemaps and then parse the returned html, traveling down the rabbit hole if the target URI has more than 1000 mementos.

```
1 # -*- encoding: utf-8 -*-
2 #!/usr/bin/python
3
4 import requests
5 import re
6
7 MW_URI = "http://mementoweb.org/timemap/link/"
8
9 if __name__ == '__main__':
10     with open('output', 'r') as f:
11         output = open('results', 'w')
12         mementos = {}
13         for uri in f.read().split('\n'):
14             if uri is '':
15                 continue
16             count = 0
17             target_uri = MW_URI + uri
18             while True:
19                 result = requests.get(target_uri)
20                 if result.ok:
21                     count = count + result.text.count('rel="memento"')
```

```

22         last_line = result.text.split('\n')[-1]
23         if 'rel="timemap"' not in last_line:
24             break
25         sites = re.findall(r'<([^\>]+)>', last_line)
26         target_uri = sites[1]
27         mementos[target_uri] = count
28         print 'found %d mementos for uri: %s' % (count, uri)
29         output.write('%s %d\n' % (uri, count))
30     output.close()

```

Listing 2: mementofinder.py

The dataset created Listing 2. A log scale was used along the y-axis to show more detail among the results. The script in Listing 3 was used to create the histogram in Figure 1, which shows the distribution of mementos per site from the dataset of Question 1.

```

1 #! /usr/bin/Rscript
2
3 data <- read.table("results", header=TRUE, comment.char="")
4 counts <- table(data$Mementos)
5 pdf("hist.pdf")
6 barplot(counts, log="y", ylim=c(.75, nrow(data)), ylab="Sites", xlab="Memento Count", main="
7   Memento Count per Site")
8 dev.off()

```

Listing 3: build_histogram.r

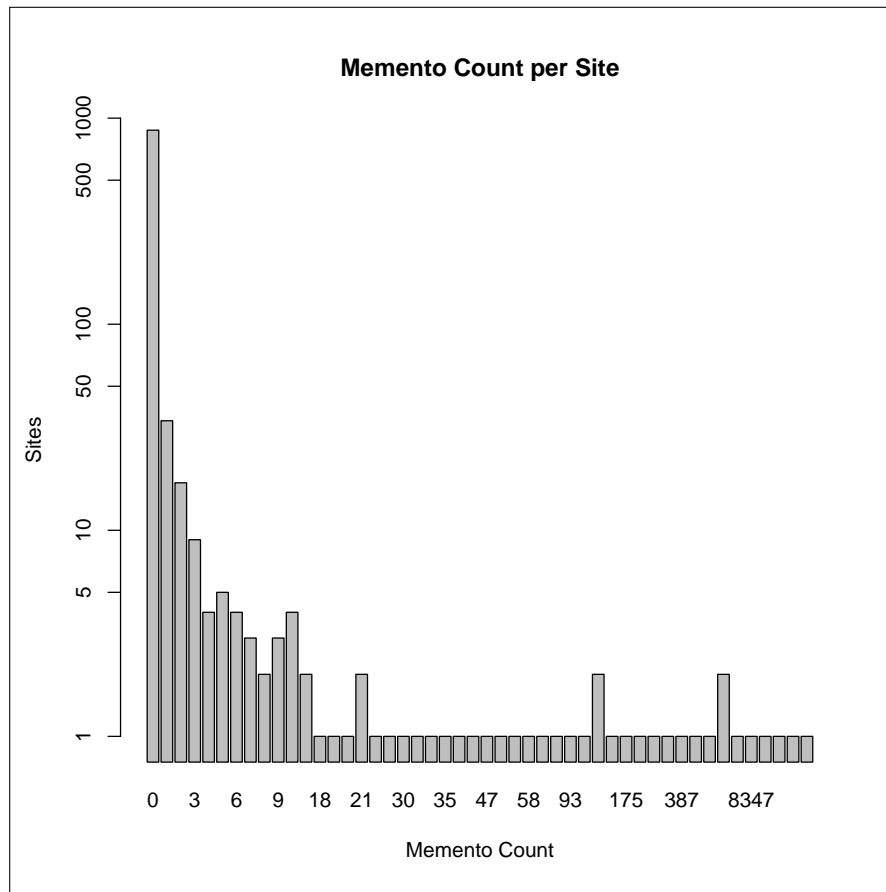


Figure 1: Histogram of Site Mementos

3 Question 3

3.1 Question

Estimate the age of each of the 1000 URIs using the “Carbon Date” tool:

<http://ws-dl.blogspot.com/2013/04/2013-04-19-carbon-dating-web.html>

Note: you’ll have better luck downloading and installing the tool rather than using the web service (which will run slowly and likely be unreliable).

For URIs that have > 0 Mementos and an estimated creation date, create a graph with age (in days) on one axis and number of mementos on the other.

3.2 Resources

- R functions: <http://stackoverflow.com/questions/16860200/row-by-row-operations-and-updates-in-data-table>
- R lapply: <http://stackoverflow.com/questions/20766617/using-apply-functions-in-place-of-for-loop-in-r>
- R strptime: <https://stat.ethz.ch/R-manual/R-devel/library/base/html/strptime.html>

3.3 Answer

Using the python script in Listing 4 to utilize the tools from the “Carbon Date” tool provided by H.S. Eldeen, estimated creation dates were found for 843 of the original URIs. These were stored in the file `site_ecd_all`.

```
1  #!/usr/bin/python
2
3  from local import cd
4  import json
5  import futures
6
7  ECD = 'Estimated Creation Date'
8
9  def getdate(uri):
10     uri = uri.strip()
11     print 'Searching for uri: {}'.format(uri)
12     uri_json = json.loads(cd(uri))
13     if uri_json[ECD]:
14         print 'Found creation date: {}'.format(uri_json[ECD])
15         return (uri, uri_json[ECD])
16     else:
17         print 'Found no ECD'
18         return None
19
20 if __name__ == '__main__':
21     # Remove already completed links
22     with open('site_mementos') as infile:
23         input_uris = [line.split(' ')[0] for line in infile if line.rstrip('\n').split(' ')[1] != '0']
24     with open('site_ecd_all') as prevfile:
25         prev = [line.split(' ')[0] for line in prevfile]
26     uris = [uri for uri in input_uris if uri not in prev]
27     print 'Starting on uri {}'.format(len(uris))
28
29     # Work on the rest
30     with open('site_ecd_all','a') as outfile:
31         with futures.ThreadPoolExecutor(max_workers=8) as executor:
32             urifutures = [executor.submit(getdate, uri) for uri in uris]
33             for future in futures.as_completed(urifutures):
```

```

34         try:
35             data = future.result()
36         except Exception as exc:
37             print '{} generated an exception: {}'.format(uri, exc)
38     if len(data) == 2:
39         print 'Writing data: {}'.format(data)
40         outfile.write('{} {}{}\n'.format(data[0], data[1]))
41     else:
42         print 'Found no data'

```

Listing 4: carbondate.py

To prepare the dataset for graphing, the Python script in Listing 5 was used to capture the desired subset of data from the `site_ecd_all` file; i.e, those site-memento pairs where the memento count was greater than zero. The results were stored in the file `ecd_mementos`.

```

1  #!/usr/bin/python
2
3  if __name__ == '__main__':
4      with open('site_ecd_all') as ecdfile:
5          ecdmap = dict([line.rstrip('\n').split() for line in ecdfile])
6      with open('site_mementos') as memfile:
7          mementos = dict([line.rstrip('\n').split() for line in memfile if line.rstrip('\n').
8                          split()[1] != '0'])
9      with open('ecd_mementos', 'w') as outfile:
10         outfile.write('ECD Mementos\n')
11         for uri, mem in mementos.iteritems():
12             try:
13                 outfile.write('{} {}{}\n'.format(ecdmap[uri], mem))
14             except KeyError as e:
15                 print '{}\n'.format(e)

```

Listing 5: build_ecd_mementos.py

Using the R script in Listing 6, with the dataset obtained from Listing 5 (the `ecd_mementos` file), the graph in Figure 2 was created. This graph shows that the older a site is the higher that site's memento count tends to be. It also shows that there are a far greater number of new sites than old ones.

```

1  #!/usr/bin/Rscript
2
3  data <- read.table("ecd_mementos", header=TRUE)
4  strtodate = function(x) {
5      result = Sys.time() - strptime(x, "%Y-%m-%dT%H:%M:%S")
6      return(result)
7  }
8  data$ECD <- lapply(data$ECD, strtodate)
9  pdf("ecd_mementos.pdf")
10 plot(data, log="y", xlab="Estimated Site Age in days", ylab="Number of Mementos", main="
    Estimated Site Age to Memento Count")
11 dev.off()

```

Listing 6: ecd_mementos_graph.r

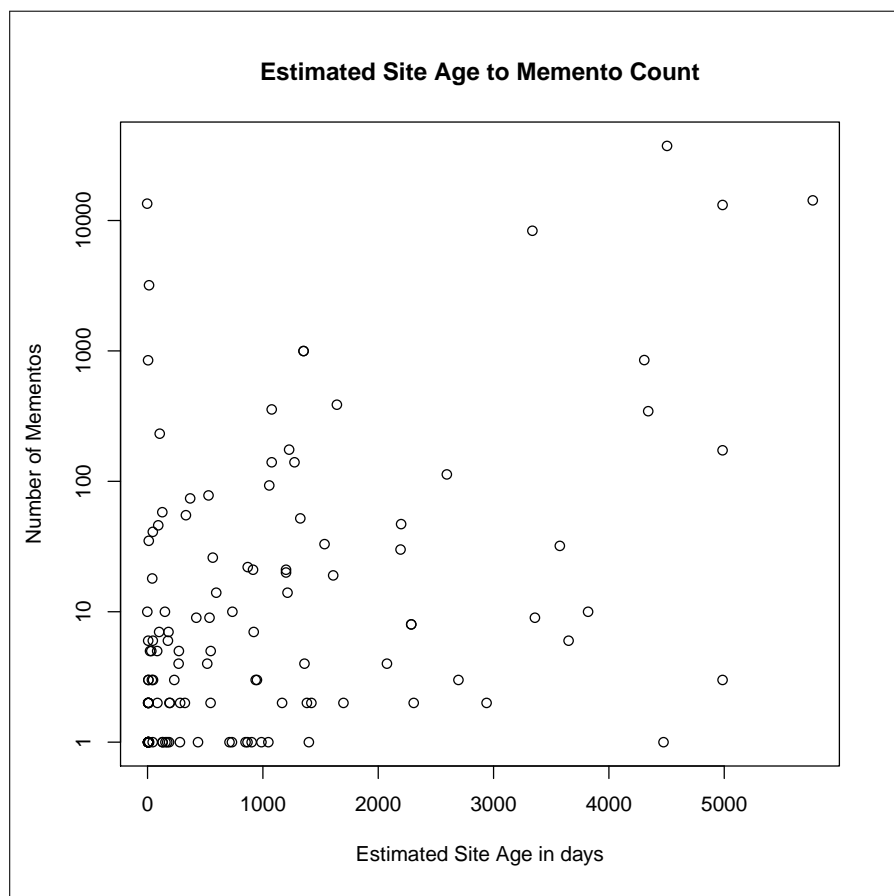


Figure 2: Estimated Age to Memento Count