Assignment 2

Fall 2014 CS595 Web Science Dr. Michael Nelson

Mathew Chaney

September 23, 2014

Contents

1	Que	estion 1
	1.1	Question
	1.2	Resources
	1.3	Answer
2	Que	estion 2
	2.1	Question
	2.2	Resources
	2.3	Answer
L	ist (of Figures
	1	Histogram of Site Mementos
L	istiı	ngs
	1 2	urifinder.py

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, 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.

```
# -*- encoding: utf-8 -*-
import requests
from requests_oauthlib import OAuth1
from urllib import quote

REQUEST_TOKEN_URL = "https://api.twitter.com/oauth/request_token"
AUTHORIZE_URL = "https://api.twitter.com/oauth/authorize?oauth_token="
ACCESS_TOKEN_URL = "https://api.twitter.com/oauth/access_token"

CONSUMER_KEY = "PDBekXkvUto4V0XYZrrizcEub"
CONSUMER_SECRET = "OElKNfpNWF8Eh4iwbDFYFKDMBSouni3uRZrpsoGhbJcLZZmnBq"

OAUTH_TOKEN = "2560074793 - g7ESlsQmwl3YKAfCJnIBaolh3wHLjmPqj96XFuV"
AUTHORIZE_URL = "https://api.twitter.com/oauth/access_token"

SEARCH_URI = "https://api.twitter.com/oauth/access_token"

SEARCH_URI = "https://api.twitter.com/oauth/access_token"

SEARCH_URI = "https://api.twitter.com/oauth/access_token"
```

```
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():
         return OAuth1 (CONSUMER KEY,
40
                        client_secret=CONSUMER_SECRET,
41
                         resource_owner_key=OAUTH_TOKEN,
resource_owner_secret=OAUTH_TOKEN_SECRET)
42
43
44
45
   def find uris(uris):
         with open ('output', 'a') as outfile:
    for search_item in SEARCH_ITEMS:
        result = requests.get (SEARCH_URI + search_item + '&filter%3Alinks&count=1000',
46
47
48
                        auth=oauth)
                   for status in result.json()['statuses']:
for url in status['entities',]['urls']:
49
50
                              if len(uris) = 1000:
51
52
                                   return
                              if 'expanded_url' in url:
53
54
                                   try:
                                         result \ = \ requests \, . \, get \, (\, url \, [\, \mbox{'expanded\_url'}] \, , \ timeout \, {=} 4)
55
                                         # only add expanded uris if they aren't in the list already if result.status_code = 200 and result.url not in uris:
56
57
                                              add_uri(uris, result.url)
outfile.write('%s\n' % result.url)
58
59
                                    except Exception as e:
60
61
                                         print e
62
                                         continue
63
   def add_uri(uris, uri):
64
65
         uris.add(uri)
         print 'added uri #%d: %s' % (len(uris), uri)
66
67
       __name__ == "__main__":
68
        \overline{oauth} = get_oauth()
69
70
         uris = set(\overline{)}
71
         # read in previous set of uris
72
73
              with open ('output', 'r') as infile:
                    for line in infile.readlines()
74
75
                        add_uri(uris, line.strip())
76
         except IOError:
         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

Using a simple python script to get the timemaps and then parse the results, traveling down the rabbit hole if the target uri has more than 1000 mementos, was done with the following code.

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

Listing 2: mementofinder.py

And the results of running the python code in Listing 2. A log scale was used along the y axis to show more detail among the results.

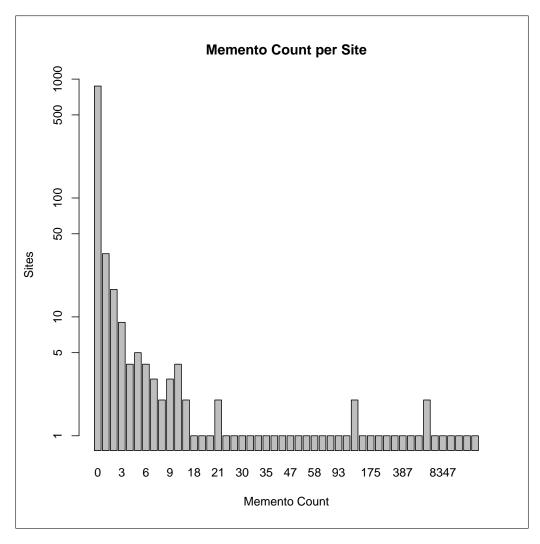


Figure 1: Histogram of Site Mementos