Khaja Masroor Ahmed

# Assignment 1

CS 751: Introduction to Digital Libraries
Dr. Michael Nelson
Spring d2015

February 14, 2015

# Contents

**1**

# Question 1

## 1.1 Question

Write a program that extracts 10000 tweets with links from Twitter.

- Save the tweet URIs, and the mapping to the link(s) each tweet contains.
- For each t.co link, use "curl -I -L" to record the HTTP headers all the way to a terminal HTTP status (i.e. chase down all the redirects)
- How many unique final URI's? How many duplicate URI's?
- Build a histogram of how many redirects (every URI will have at least 1)
- Build a histogram of HTTP status codes encountered (youll have at least 20000: 10000 301s, and 10000+ more)

## 1.2 Solution

For solving the above problem I have used Python as the programming language. Following are the steps I've taken to solve the given problem:

- I first registered a twitter application to generate the consumer key and consumer secret key for using the Twitter API.
- The key's generated in the above step are used for authenticating the application for requesting the tweets.
- I encountered multiple packages while researching for packages for fetching tweet data and decided upon TwitterSearch package.
- I was facing an issue with the TwitterSearch package where I was receiving limited results. In-spite of making modifications to the keyword for search I wasn't able to exceed more than 800 results.
- I eventually moved to tweepy package with which I was receiving approximately 2000 results and it would stop again. Along with the keyword argument I also passed another argument 'since = 2014' to the API search as the tweets fetched would stop at a certain date. This additional argument solved my problem. I was able to retrieve more than 10000 tweets.

- From the tweet data received from the API, I fetched the tweet identifier, created date, user identifier, tweet text and of course the URI.
- The above mentioned data was processed to json format and written to a file 'output.txt'.
- For retrieving the final URI, history of status code's encountered, I found requests package to be the most powerful and appropriate for use in my context.
- I used the requests package in 'requestStatus.py' for fetching the above mentioned information and output this information to 'status.txt' file.
- Now with the raw information available, I processed the data further based on the requirement as specified in the question.
- Based on the tweet data, I found that the number of unique URI's that were encountered were 6431, duplicate links were 3999. The total number of greater than 10000 as some tweets had more than one URI. This information was calculated using script 'getUniqueURI.py' and the result was stored in 'unique.txt'.
- Also, I calculated the number of re-directs for each URI using the script 'getURIRedirect.py' and the result was stored in 'uriRedirect.txt'. This information was processed through R for generating a histogram as shown in the figure below.
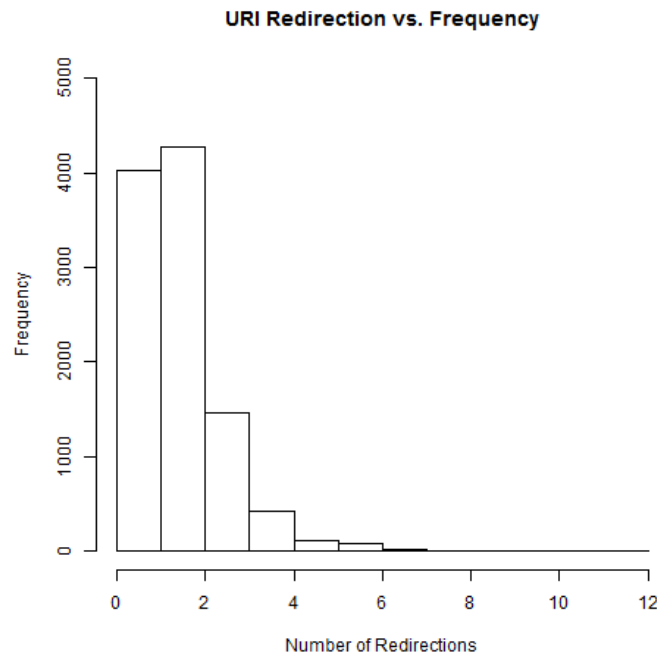


**Fig. 1.1.** URI Redirection vs. Frequency

- The raw information from status.txt was also used to calculate the frequency of the status codes encountered for these URI's and the figure below represents the histogram for this data.
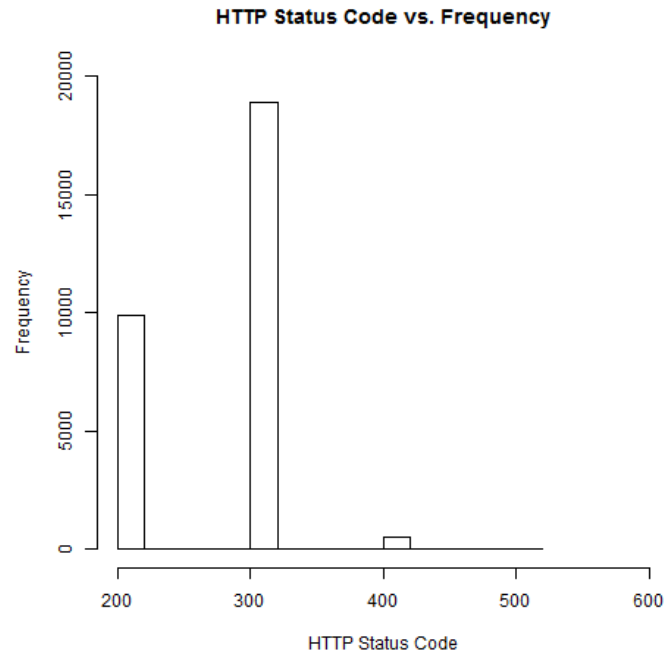


**Fig. 1.2.** HTTP Status Code vs. Frequency

## 1.3 Code Listing

```
1  #CS751 Spring 2015 Offered by Dr. Michael Nelson −
       Assignment 1
2  #Filename: getTweets.py
3  #Author: Khaja Masroor Ahmed
4  #UIN: 00999044
5  #CS Email: kahmed@cs.odu.edu
6
7  import tweepy
8  import json
9  import time
10 import datetime
11
12 consumer_key = "iuKUndPfIF5aWnNl0Ayq9Ztgt"
13 consumer_secret = "
       QuNsF4gL2LssmbcdtKpyLZGiQctz98T4hXWcAKrBYGh72ZTFC8"
14
15 access_token = "549294315−
       P89swbZzgiP2n9bq6fW2T2jm5etru6Wr6TN08Lg3"
16 access_token_secret = "
       NMzDaS5doFtHxXxebE68AunmRHTsFfLxwAkk3LsDN75JH"
17
18 auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
19 auth.set_access_token(access_token, access_token_secret)
20 api = tweepy.API(auth)
21 f = open('output.txt','w')
22
23 tweetCount = 1
24 print('Started execution')
25 cursor = tweepy.Cursor(api.search, q ='http:', since_id =
       2014).items()
26 while True:
27         try:
28                 tweet = cursor.next()
29                 for tweet in cursor:
30                         tweetData = {}
31                         tweetData['tweetId'] = tweet.id
32                         tweetData['createdAt'] = tweet._json
                               ['created_at']
33                         tweetData['userId'] = tweet.user.
                               _json['id']
34                         tweetData['tweetCount'] = tweetCount
35                         tweetData['text'] = tweet.text
36                         urlData = []
37
38                         for url in tweet._json['entities']['
                               urls']:
```

```
39                              urlData.append(url['url'])
40                          tweetData['url'] = urlData
41                          tweetData = json.dumps(tweetData)
42
43                          if tweetCount == 11999:
44                                  print('Complete')
45                                  print(datetime.datetime.now
                                        ())
46                                  f.close()
47                                  break
48                          f.write(tweetData + "\n")
49                          tweetCount = tweetCount + 1
50              except tweepy.TweepError as e:
51                      print(e)
52                      time.sleep(900)
53                      continue
54              except StopIteration:
55                      break
56      print('Program Executed!')
```

**Listing 1.1.** Python program for retrieving tweets.

```
1   #CS751 Spring 2015 Offered by Dr. Michael Nelson −
        Assignment 1
2   #Filename: requestStatus.py
3   #Author: Khaja Masroor Ahmed
4   #UIN: 00999044
5   #CS Email: kahmed@cs.odu.edu
6
7   import requests
8   import json
9   import string
10
11  f = open('output.txt', 'r+')
12  w = open('status.txt','w')
13  err = open('error.txt','w')
14  count = 1
15  for line in f:
16          data = json.loads(line)
17          i = 0
18          try:
19                  tweetData = {}
20                  tweetData['index'] = count
21                  tweetData['tweetId'] = data['tweetId']
22                  tweetData['createdDate'] = data['createdAt']
23                  tweetURLData = []
24                  for url in data['url']:
25
26                          r = requests.head(data['url'][i],
                                allow_redirects=True, timeout =
                                15)
27                          j = 0
28
29                          URLData = {}
30
31                          URLData['finalURL'] = r.url
32                          statusCode = []
33                          statusCode.append(r.status_code)
34                          for response in r.history:
35                                  statusCode.append(response.
                                        status_code)
36                          URLData['statusCode'] = statusCode
37                          URLData['tweetURL'] = data['url'][i]
38                          tweetURLData.append(URLData)
39
40                          i = i + 1
41                  tweetData['tweetURLData'] = tweetURLData
42
43                  tweetData = json.dumps(tweetData)
44                  print count
```

```
45                          w. write ( tweetData + ”\n” )
46              except :
47                          err . write ( str ( count ) + ’\n’ )
48                          continue
49              count = count + 1
50  f . close ( )
51  w. close ( )
52  err . close ( )
```

**Listing 1.2.** Python program to request the header information for each URI.

```
1   #CS751 Spring 2015 Offered by Dr. Michael Nelson −
        Assignment 1
2   #Filename: getUniqueURI.py
3   #Author: Khaja Masroor Ahmed
4   #UIN: 00999044
5   #CS Email: kahmed@cs.odu.edu
6
7   import json
8   import sets
9
10  f = open('status.txt','r+')
11  w = open('unique.txt','w')
12  uniqueURLs = sets.Set()
13  count = 1
14  for line in f:
15          data = json.loads(line)
16          for urlData in data['tweetURLData']:
17                  uniqueURLs.add(urlData['finalURL'])
18                  count = count + 1
19  w.write('Unique: ' + (str)(len(uniqueURLs)) + '\n')
20  w.write('Duplicates: ' + (str)(count − len(uniqueURLs)) + '\
        n')
21  w.write('Total: ' + (str)(count) + '\n')
22  f.close()
23  w.close()
```

**Listing 1.3.** Python program for retrieving the number of unique, duplicate URI's.

```
1   #CS751 Spring 2015 Offered by Dr. Michael Nelson −
        Assignment 1
2   #Filename: getURIRedirect.py
3   #Author: Khaja Masroor Ahmed
4   #UIN: 00999044
5   #CS Email: kahmed@cs.odu.edu
6
7   import json
8
9   f = open('status.txt','r+')
10  w = open('uriRedirect.txt','w')
11  w2 = open('uriRedirectCount.txt','w')
12  uriRedirect = {}
13  for line in f:
14          data = json.loads(line)
15          count = 0
16          for tweetData in data['tweetURLData']:
17                  for status in tweetData['statusCode']:
18                          if(str(status).startswith('3')):
19                                  count = count + 1
20                  w.write(str(count) + '\t' + tweetData['
                        finalURL'] + '\n')
21                  if count in uriRedirect:
22                          uriRedirect[count] = uriRedirect[
                                count] + 1
23                  else:
24                          uriRedirect[count] = 1
25  w2.write(str(uriRedirect) + '\n')
26  w.close()
27  w2.close()
```

**Listing 1.4.** Python program for calculating the number of redirects for each URI.

```
1   #CS751 Spring 2015 Offered by Dr. Michael Nelson −
         Assignment 1
2   #Filename: getStatusCodeFrequency.py
3   #Author: Khaja Masroor Ahmed
4   #UIN: 00999044
5   #CS Email: kahmed@cs.odu.edu
6
7   import json
8   import string
9
10  f = open('status.txt', 'r+')
11  w = open('statusCode.txt', 'w')
12  statusCode = {}
13  for line in f:
14          data = json.loads(line)
15          if data['index'] == 10000:
16                  print 'Program Executed'
17                  break
18          for tweetData in data['tweetURLData']:
19                  for status in tweetData['statusCode']:
20                          w.write(str(status) + '\n')
21                          if status in statusCode:
22                                  statusCode[status] =
                                          statusCode[status] + 1
23                          else:
24                                  statusCode[status] = 1
```

**Listing 1.5.** Python program for retrieving the different status codes encountered by each URI.

```
1  uriDataset <- read.csv("C:/Users/kahmed/Desktop/uriRedirect.
       txt", stringsAsFactors=F, header=FALSE, sep="\t")
2  data = uriDataset[,1]
3  png("C:/Users/kahmed/Desktop/uriRedirect.png")
4  hist(data, main="URI Redirection vs. Frequency", freq=T, xlab="
       Number of Redirections", ylab="Frequency", xlim=c(0,12),
       ylim=c(0,5000))
```

**Listing 1.6.** R program for generating the histogram for URI Redirection vs. Frequency

```
1  statusCode <- read.csv("C:/Users/kahmed/Desktop/statusCode.
       txt", stringsAsFactors=F, header=FALSE)
2  data = statusCode[,1]
3  png("C:/Users/kahmed/Desktop/statusCode.png")
4  hist(data, main="HTTP Status Code vs. Frequency", freq=T, xlab=
       "HTTP Status Code", ylab="Frequency", xlim=c(200,600), ylim
       =c(0,20000))
```

**Listing 1.7.** R program for generating the histogram for HTTP Status Code vs. Frequency

# Question 2

## 2.1 Question

- Use "Carbon Date" to estimate the age of each link(s) in a tweet.
  - See: http://ws-dl.blogspot.com/2013/04/2013-04-19-carbon-dating-web.html
- For each t.co link, use "curl -I -L" to record the HTTP headers all the way to a terminal HTTP status (i.e. chase down all the redirects)
  - Many (most?) deltas will be 0, but there should be many ¿ 0.
- For these deltas, compute: median, mean, std dev, std err.
- Use wget to download the text for all the links. Hold on to those, well come back to them later. See:
  - http://superuser.com/questions/55040/save-a-single-web-page-with-background-images-with-wget
  - http://stackoverflow.com/questions/6348289/download-a-working-local-copy-of-a-webpage

## 2.2 Solution

- For carbon dating the tweet URI's, I first cloned the github repository.
- As per the instructions of the author of the library, I then created a bitly account and inserted the accesstoken in the configuration file.
- Some of the URI's were taking longer than a minute and some a few hours to retrieve the estimated creation date, I split the input files into 1000 URI's each and created five instances of the local.py file with the necessary updates.
- By splitting these URI's into batches I was able to complete fetching data from each of the URI's in two days.
- Also, some URI's were stuck for over a few hours and got repeatedly stuck at the same location. I deleted these URI's and continued carbon dating the rest in the list.

- In the file 'local.py' I created a JSON structure with the tweet identifier, creation date of tweet and the estimated creation date of the URI.
- I used the 'getTweetAge.py' for fetching the age of each of these URI's. The figure below indicates the histogram of the $Age_{tweet}$ - $Age_{URI}$ vs. Frequency.
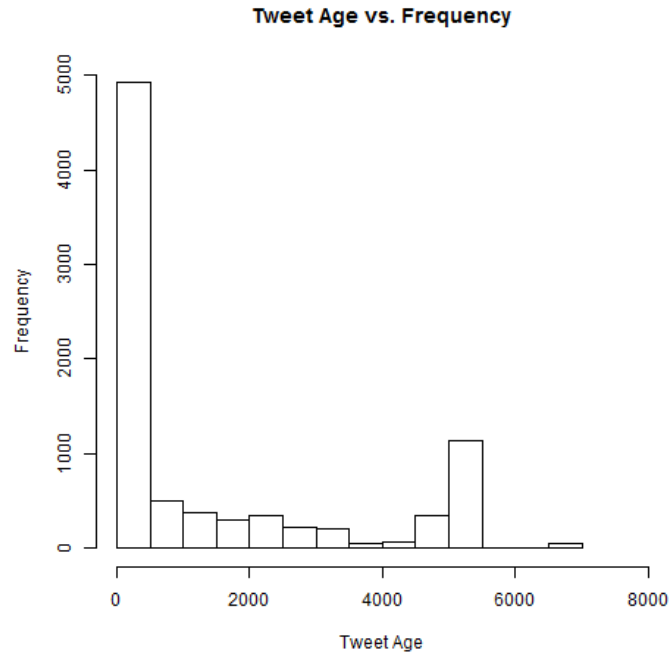


**Fig. 2.1.** Age vs. Frequency

- Using the age of each of these URI's I found the mean as 1409.735 days, median as 235 days, standard deviation as 1948.022 days, standard error as 21.10941 days.

## 2.3 Code Listing

```
1   from checkForModules import checkForModules
2   import json
3   from ordereddict import OrderedDict
4   import urlparse
5   import re
6
7   from getBitly import getBitlyCreationDate
8   from getArchives import getArchivesCreationDate
9   from getGoogle import getGoogleCreationDate
10  from getBacklinks import *
11  from getLowest import getLowest
12
13  from getLastModified import getLastModifiedDate
14  from getTopsyScrapper import getTopsyCreationDate
15  from htmlMessages import *
16  from pprint import pprint
17
18  from threading import Thread
19  import Queue
20  import datetime
21
22  import os, sys, traceback
23
24
25
26
27  def cd(url, backlinksFlag = False):
28
29      #print 'Getting Creation dates for: ' + url
30
31
32      #scheme missing?
33      parsedUrl = urlparse.urlparse(url)
34      if( len(parsedUrl.scheme)<1 ):
35          url = 'http://'+url
36
37
38      threads = []
39      outputArray =['','','','','','']
40      now0 = datetime.datetime.now()
41
42
43      lastmodifiedThread = Thread(target=getLastModifiedDate,
              args=(url, outputArray, 0))
44      bitlyThread = Thread(target=getBitlyCreationDate, args=(
              url, outputArray, 1))
```

```
45        googleThread = Thread(target=getGoogleCreationDate, args
              =(url, outputArray, 2))
46        archivesThread = Thread(target=getArchivesCreationDate,
              args=(url, outputArray, 3))
47
48        if( backlinksFlag ):
49            backlinkThread = Thread(target=
                  getBacklinksFirstAppearanceDates, args=(url,
                  outputArray, 4))
50
51        topsyThread = Thread(target=getTopsyCreationDate, args=(
              url, outputArray, 5))
52
53
54        # Add threads to thread list
55        threads.append(lastmodifiedThread)
56        threads.append(bitlyThread)
57        threads.append(googleThread)
58        threads.append(archivesThread)
59
60        if( backlinksFlag ):
61            threads.append(backlinkThread)
62
63        threads.append(topsyThread)
64
65
66        # Start new Threads
67        lastmodifiedThread.start()
68        bitlyThread.start()
69        googleThread.start()
70        archivesThread.start()
71
72        if( backlinksFlag ):
73            backlinkThread.start()
74
75        topsyThread.start()
76
77
78        # Wait for all threads to complete
79        for t in threads:
80            t.join()
81
82        # For threads
83        lastmodified = outputArray[0]
84        bitly = outputArray[1]
85        google = outputArray[2]
86        archives = outputArray[3]
87
88        if( backlinksFlag ):
```

```
89                backlink = outputArray[4]
90          else:
91                backlink = ''
92
93          topsy = outputArray[5]
94
95          #note that archives["Earliest"] = archives[0][1]
96          try:
97                lowest = getLowest([lastmodified, bitly, google,
                        archives[0][1], backlink, topsy]) #for thread
98          except:
99              print sys.exc_type, sys.exc_value , sys.exc_traceback
100
101
102
103          result = []
104
105          result.append(("URI", url))
106          result.append(("Estimated Creation Date", lowest))
107          result.append(("Last Modified", lastmodified))
108          result.append(("Bitly.com", bitly))
109          result.append(("Topsy.com", topsy))
110          result.append(("Backlinks", backlink))
111          result.append(("Google.com", google))
112          result.append(("Archives", archives))
113          values = OrderedDict(result)
114          r = json.dumps(values, sort_keys=False, indent=2,
                    separators=(',', ': '))
115
116          now1 = datetime.datetime.now() - now0
117
118          print 'runtime in seconds:   ' +  str(now1.seconds) + '\n
                ' + r + '\n'
119
120          return lowest
121
122  f = open('status.txt', 'r+')
123  w = open('carbonDate.txt', 'w')
124
125  for line in f:
126          data = json.loads(line)
127
128          for tweetData in data['tweetURLData']:
129                  print data['index']
130                  print datetime.datetime.now()
131                  carbonDate = {}
132                  carbonDate['oldestDate'] = cd(tweetData['
                        finalURL'])
```

```
133                      carbonDate['createdDate'] = data['
                             createdDate']
134                      carbonDate['tweetId'] = data['tweetId']
135                      carbonDate = json.dumps(carbonDate)
136                      w.write(carbonDate + "\n")
137                      print carbonDate
```

**Listing 2.1.** Python program for carbon dating the URI's.

```
 1  import json
 2  import datetime
 3
 4  f = open('carbonData.txt','r+')
 5  w = open('age.txt','w')
 6  count = 0
 7  for line in f:
 8          data = json.loads(line)
 9          if data['oldestDate'] != '':
10                  count = count + 1
11                  oldestDate = datetime.datetime.strptime(data
                        ['oldestDate'], '%Y-%m-%dT%H:%M:%S')
12                  createdDate = datetime.datetime.strptime(
                        data['createdDate'], "%a %b %d %H:%M:%S
                        +0000 %Y")
13                  age = abs((oldestDate.date() - createdDate.
                        date()).days)
14                  w.write(str(age) +'\n')
15  print('Execution Complete')
```

**Listing 2.2.** Python program for calculating the Age$_{\text{tweet}}$ - Age$_{\text{URI}}$.

```
1  age <- read.csv("C:/Users/kahmed/Desktop/age.txt",
        stringsAsFactors=F, header=FALSE)
2  data = age[,1]
3
4  png("C:/Users/kahmed/Desktop/tweetAge.png")
5  hist(data,main="Tweet Age vs. Frequency",freq=T,xlab="Tweet
        Age",ylab="Frequency",xlim=c(1,8000),ylim=c(0,5000))
6
7  se <- function(x) sqrt(var(x)/length(x))
8  write("Mean", file = "C:/Users/kahmed/Desktop/meanMedMode.
        txt")
9  write(mean(data), file = "C:/Users/kahmed/Desktop/
        meanMedMode.txt",
10         append = TRUE, sep = "\t")
11  write("Median", file = "C:/Users/kahmed/Desktop/meanMedMode.
        txt",
12         append = TRUE, sep = "\t")
13  write(median(data, na.rm = FALSE), file = "C:/Users/kahmed/
        Desktop/meanMedMode.txt",
14         append = TRUE, sep = "\t")
15  write("Standard Deviation", file = "C:/Users/kahmed/Desktop/
        meanMedMode.txt",
16         append = TRUE, sep = "\t")
17  write(sd(data, na.rm = FALSE), file = "C:/Users/kahmed/
        Desktop/meanMedMode.txt",
18         append = TRUE, sep = "\t")
19  se <- function(data) sd(data)/sqrt(length(data))
20  write("Standard Error", file = "C:/Users/kahmed/Desktop/
        meanMedMode.txt",
21         append = TRUE, sep = "\t")
22  write(se(data), file = "C:/Users/kahmed/Desktop/meanMedMode.
        txt",
23         append = TRUE, sep = "\t")
```

**Listing 2.3.** R program for generating the histogram for Age vs. Frequency and calculating mean, median, mode and standard error