

## CS 751: Assignment 1

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

Note that only Twitter API 1.1 is currently available; version 1 code will no longer work.

Reference:

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

Rohit Lambi

Spring 2015

# Contents

1	Q1 . . . . .	2
	1.1 Solution . . . . .	2
	1.2 Number of URIs Encountered: . . . . .	3
	1.3 Code Listing . . . . .	3
2	Q2 . . . . .	9
	2.1 Solution . . . . .	9
	2.2 Code Listing . . . . .	10
	2.3 Summary of Tweet and URI's date delta . . . . .	12

# 1 Q1

1. Save the tweet URIs, and the mapping to the link(s) each tweet contains - Note: tweets can have more than 1 links
2. For each t.co link, record the HTTP headers all the way to a terminal HTTP status (i.e. chase down all the redirects)
3. How many unique final URIs? How many duplicate URIs?
4. Build a histogram of how many redirects (every URI will have at least 1) <http://en.wikipedia.org/wiki/Histogram>
5. Build a histogram of HTTP status codes encountered (you will have at least 20000: 10000 301s, and 10000+ more)

## 1.1 Solution

This assignment has been implemented in Python. Following is the detailed report of how this assignment has been done:

1. To be able to access Twitter's REST API we need the secret keys. I obtained these by registering an app on Twitter's website from my Twitter account.
2. Instead of directly using the Twitter's REST API, I used 'Tweepy' library which is a wrapper to the Twitter's REST API. This makes coding a bit easier.
3. Using the secret keys, get the authentication token from Twitter and use it in the next webservice calls.
4. Twitter has a limitation on the number of times we can access their API within a certain period of time which is 15 minutes. So, to fetch the tweets, I continuously made the webservice calls to fetch the tweets and once the limit got exhausted, I made my program to wait for 15 minutes before making request for the next set of tweets.
5. I saved the Tweet data in a text file in JSON format. The stored information was Tweet id, Tweet text, URLs, Tweet created date
6. To get the final URLs of t.co URLs from Twitter, I used 'Requests' library of Python.
7. For each t.co URL which was stored in file in the previous step, I made a HTTP HEAD request by setting the 'allow redirects' flag to 'True'.
8. From the response, I stored all the HTTP status codes in one text file and in other text file I stored the final URL and number of redirects required to reach the final URI.
9. Moreover, I calculated and stored the count of unique and duplicate URIs in the third text file.

## 1.2 Number of URIs Encountered:

uniqueUrls: 5971

duplicateUrls: 4306

## 1.3 Code Listing

### tweetFetcher.py

```
1  '''
2  Created on Jan 31, 2015
3
4  @author: rlambi
5  '''
6  import tweepy
7  import time
8  import json
9  from datetime import datetime
10
11  CONSUMER_KEY = "iuKUndPfIF5aWnN10Ayq9Ztgt"
12  CONSUMER_SECRET = "QuNsF4gL2LssmbcdtKpyLZGiQctz98T4hXWcAKrBYGh72ZTFC8"
13
14  OAUTH_TOKEN = "549294315-P89swbZzgiP2n9bq6fW2T2jm5etru6Wr6TN08Lg3"
15  OAUTH_TOKEN_SECRET = "NMzDaS5doFtHxXxebE68AunmRHTsFfLxwAkk3LsDN75JH"
16
17  auth = tweepy.OAuthHandler(CONSUMER_KEY, CONSUMER_SECRET)
18  auth.set_access_token(OAUTH_TOKEN, OAUTH_TOKEN_SECRET)
19
20  api = tweepy.API(auth)
21
22  fil = open('tweets.txt', 'w')
23  log = open('log.txt', 'w')
24  data = {}
25  url_list = []
26  urls = []
27
28  seqNum = 0
29  log.write('Start: ' + str(0) + datetime.now().time().__str__() + '\n')
30
31  myCursor = tweepy.Cursor(api.search, q="http:links", rpp=100).items()
32
33  while True:
34      try:
35          tweet = myCursor.next()
36
37          url_list = tweet._json['entities']['urls']
38
39          urls[:] = []
40          if url_list.__len__() > 0:
41              seqNum += 1
42
43              print seqNum, " ", url_list
44
45              for url in url_list:
46                  urls.append(url['url'])
47
48              data['id_str'] = tweet._json['id_str']
49              data['text'] = tweet._json['text']
50              data['created_at'] = tweet._json['created_at']
51              data['urls'] = urls
52              data['seqNum'] = seqNum
53
54              json_tweet = json.dumps(data)
55              fil.write(json_tweet + '\n')
56          except tweepy.TweepError:
57
58              tmp = 'Sleep at ' + str(seqNum) + ' ' + datetime.now().time().__str__() + '\n'
59              print tmp
60              log.write(tmp)
61              time.sleep(900)
62
```

```

63         tmp = 'Wakeup: ' + datetime.now().time().__str__() + '\n'
64         print tmp
65         log.write(tmp)
66         continue
67
68     if seqNum == 10000:
69         break
70
71
72 print "Total tweets: ", seqNum
73 log.write('Stop: ' + datetime.now().time().__str__() + '\n')
74 log.write("Total tweets: " + str(seqNum))
75
76 log.close()
77 fil.close()

```

Listing 1: Python program to fetch 10,000 Tweets with links

## fetchFinalURI.py

```
1  '''
2  Created on Feb 5, 2015
3
4  @author: rlambi
5  '''
6  import requests
7  import requests.exceptions
8  import json
9  import csv
10 import os
11
12 filePath = '/home/rlambi/rohit/Courses/Digital-Libraries/Assignments/'
13 outputDir = filePath + 'output_2/'
14
15 if not os.path.exists(outputDir):
16     os.makedirs(outputDir)
17
18 finalURIs = []
19 duplicateUrlsCount = 0
20 uniqueUrlsCount = 0
21 skippedUrlsCount = 0
22 serialNo = 1
23 fileTweets = open(filePath + 'tweets.txt', 'r')
24 fileHttpStatuses = open(outputDir + 'http-statuses.txt', 'w')
25 fileUrls = open(outputDir + 'tweets-processed-1.txt', 'w')
26 skippedUrls = open(outputDir + 'skippedUrls.txt', 'w')
27
28 xx = 1
29 for line in fileTweets:
30
31     tweet = json.loads(line)
32
33     tweetedOn = tweet['created_at']
34     urlList = tweet['urls']
35
36     for url in urlList:
37
38         resp = requests.Response()
39         try:
40             resp = requests.head(url, allow_redirects=True, timeout=10)
41
42             redirectCount = len(resp.history)
43             terminalURI = resp.url
44             print serialNo, '\t', tweet['seqNum'], '\t' + url + '\t' + terminalURI + '\t' +
45                 tweetedOn + '\t' + str(redirectCount)
46             fileUrls.write(str(serialNo) + '\t' + str(tweet['seqNum']) + '\t' + url + '\t' +
47                 terminalURI + '\t' + tweetedOn + '\t' + str(redirectCount) + '\n')
48             resp.history.append(resp)
49
50             # Check in HttpStatusMap if this statusCode is already present and add/update
51             # its count accordingly
52             for respHist in resp.history:
53                 fileHttpStatuses.write(str(respHist.status_code) + '\n')
54
55             if finalURIs.count(terminalURI) == 0:
56                 finalURIs.append(terminalURI)
57             else:
58                 duplicateUrlsCount += 1
59
60             serialNo += 1
61
62         except requests.exceptions.Timeout:
63             print 'Timeout ' + str(tweet['seqNum']) + '\t' + url
64             skippedUrls.write(str(tweet['seqNum']) + '\t' + url)
65             skippedUrlsCount += 1
66
67         except:
68             print 'Skipped ' + str(tweet['seqNum']) + '\t' + url
69             skippedUrls.write('Skipped ' + str(tweet['seqNum']) + '\t' + url)
70             skippedUrlsCount += 1
```

```

68 #         if xx == 50:
69 #             break
70 #         xx += 1
71
72
73
74 uniqueUrlsCount = len(finalURIs)
75 print 'uniqueUrls: ', uniqueUrlsCount
76 print 'duplicateUrls: ', duplicateUrlsCount
77 print 'skippedUrls: ', skippedUrlsCount
78
79
80 fileUrlCount = open(filePath + 'url-count.txt', 'w')
81 fileUrlCount.write('uniqueUrls: ', '\t', str(uniqueUrlsCount))
82 fileUrlCount.write('duplicateUrls: ', '\t', str(duplicateUrlsCount))
83 fileUrlCount.close()
84
85 skippedUrls.close()
86 fileUrls.close()
87 fileTweets.close()
88 fileHttpStatuses.close()

```

Listing 2: Python program to fetch the final URI for each of the t.co URIs from the Listing 1

## Histograms

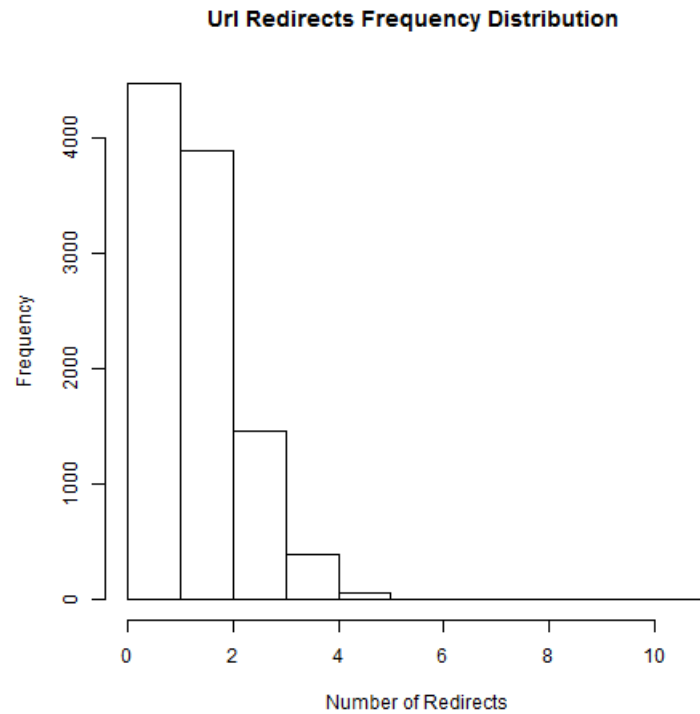


Figure 1: Number of Redirects

### url-redirect-histogram.R

```
1 redirects=read.csv("tweets-processed-1.txt",stringsAsFactors=F,header=FALSE,sep="\t")
2 redirects
3 data=redirects[,1]
4 png("url-redirect-histogram.png")
5 hist(data,main="Url Redirects Frequency Distribution",freq=T,xlab="Number of Redirects",
6     ylab="Frequency")
7 dev.off()
```

Listing 3: R program to plot histogram for the number of redirects encountered



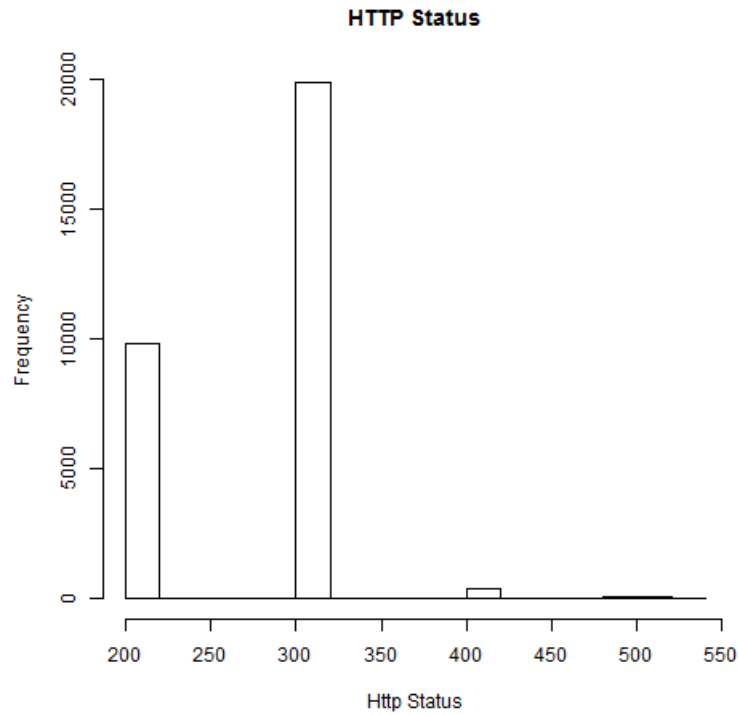


Figure 2: HTTP statuses encountered

#### http-statuses.R

```

1 inputFile=read.csv("http-statuses.txt",stringsAsFactors=F,header=FALSE,sep="\t")
2 inputFile
3 data=inputFile[,1]
4 png("http-statuses.png")
5 hist(data,main="HTTP Status",freq=T,xlab="Http Status",ylab="Frequency")
6 #barData <- table(data)
7 #png("http-statuses-bar.png")
8 #barplot(barData, main="HTTP Status", xlab="Http Status",ylab="Frequency")
9 dev.off()

```

Listing 4: R program to plot histogram for the HTTP statuses encountered

## 2 Q2

1. 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>
2. Create a histogram of (Agetweet - Agelink) - Many (most?) deltas will be 0, but there should be many more than 0
3. For these deltas, compute: median, mean, std dev, std err
4. Use wget to download the text for all the links. Hold on to those, we will come back to them later.

### 2.1 Solution

1. For Carbon-Dating, I downloaded its source code from github, installed the pre-requisites.
2. I registered an app on Bitly.com and updated the access token in Carbon-Date's 'config' file and made the other necessary configuration changes.
3. I started Carbon-Date server locally and wrote a Python script that fetches the Estimated URI created date for each URI stored in the previous step. At the same time, I calculated the difference between this date and the date on which this URI was Tweeted.
4. Output of this script is a text file which contains URI, Tweet created date(this information is carry forwarded from the previous step), Estimated URI created date and delta between the two dates.
5. Lastly, I created a python program to fetch the text for all the links.

## 2.2 Code Listing

### URIAGE.py

```
1  '''
2  Created on Feb 7, 2015
3
4  @author: rlambi
5  '''
6  import requests
7  import json
8  import csv
9  import os
10 import datetime
11
12
13 def skip(row):
14     print 'Skipped ' + '\t' + row['serialNo'] + '\t' + row['tweetSeqNo'] + '\t' + row['tcoURL'] + '\t' + row['finalURL'] + '\t' + row['tweetedOn'] + '\n'
15     errorLogFile.write('Skipped ' + '\t' + row['serialNo'] + '\t' + row['tweetSeqNo'] + '\t' + row['tcoURL'] + '\t' + row['finalURL'] + '\t' + row['tweetedOn'] + '\n')
16
17 serverURL = 'http://localhost:8080/cd?url='
18
19 tweetAgeMap = {}
20 filePath = '/home/rlambi/rohit/Courses/Digital-Libraries/Assignments/'
21 inputFilePath = filePath + 'output_2/tweets-processed-1.txt'
22 outputDir = filePath + 'output_3/'
23 outputFilePath = outputDir + 'url-dates.csv'
24 xx = 1
25 errorLogFile = open(outputDir + 'errorLog.txt', 'w')
26
27 if not os.path.exists(outputDir):
28     os.makedirs(outputDir)
29
30 with open(inputFilePath) as inputFile:
31     reader = csv.DictReader(inputFile, delimiter='\t', fieldnames=['serialNo', 'tweetSeqNo', 'tcoURL', 'finalURL', 'tweetedOn', 'redirectCount'])
32
33     with open(outputFilePath, 'wb') as outputFile:
34         writer = csv.writer(outputFile, delimiter='\t')
35
36         try:
37             for row in reader:
38                 print 'Processing ', row['serialNo'], row['finalURL']
39                 resp = requests.get(serverURL + row['finalURL'])
40                 print resp.status_code
41
42                 if (resp.status_code == 200):
43                     jsonResponse = json.loads(resp.text)
44
45                     creationDate = datetime.datetime.strptime(jsonResponse['Estimated Creation Date'], '%Y-%m-%dT%H:%M:%S') # 2006-02-22T00:00:00 "%Y-%m-%d"
46                     tweetedDate = datetime.datetime.strptime(row['tweetedOn'], "%a %b %d %H:%M:%S +0000 %Y") # Tue Feb 03 08:16:18 +0000 2015
47
48                     age = abs((creationDate.date() - tweetedDate.date()).days)
49
50                     writer.writerow([row['serialNo'], row['tweetSeqNo'], row['tcoURL'], row['finalURL'], row['tweetedOn'], jsonResponse['Estimated Creation Date'], age])
51
52                     if xx == 5:
53                         break
54                     xx += 1
55                 else:
56                     skip(row)
57
58
59 except requests.exceptions.ConnectionError:
60     print 'Cannot connect to server. Make sure server is running'
```

```

61         except:
62             print 'Exception'
63             skip(row)
64
65
66 errorLogFile.close()

```

Listing 5: Python program to find the estimated creation date of URI using Carbon Dating and calculating the delta between this date and URI tweeted date

### fetchWebpages.py

```

1  '''
2  Created on Feb 8, 2015
3
4  @author: rlambi
5  '''
6  import subprocess
7  import os
8  import thread
9  import threading
10 import csv
11 import datetime
12
13 def fetchWebPage(url):
14     print 'Fetching ', url
15     subprocess.Popen(["wget", "-E", "-H", "-k", "-K", "-p", url])
16
17 sites = 'sites'
18 if not os.path.exists(sites):
19     os.makedirs(sites)
20
21 os.chdir(sites)
22
23 fieldNames = ['sno', 'seqNum', 'tcoUrl', 'url']
24
25 print datetime.datetime.now()
26 with open('../sample-urls.txt') as csvfile:
27     reader = csv.DictReader(csvfile, fieldnames=fieldNames, delimiter='\t')
28     for row in reader:
29         fetchWebPage(row['url'])
30         thread.start_new_thread(fetchWebPage, (row['url'], ))
31
32 print datetime.datetime.now()
33
34 print 'Waiting for all threads to complete'
35 while threading.activeCount() > 1:
36     print str(threading.activeCount())
37     pass
38
39 print 'Completed fetching all webpages'
40 print datetime.datetime.now()

```

Listing 6: Python program to fetch the text for all the links

## Histogram

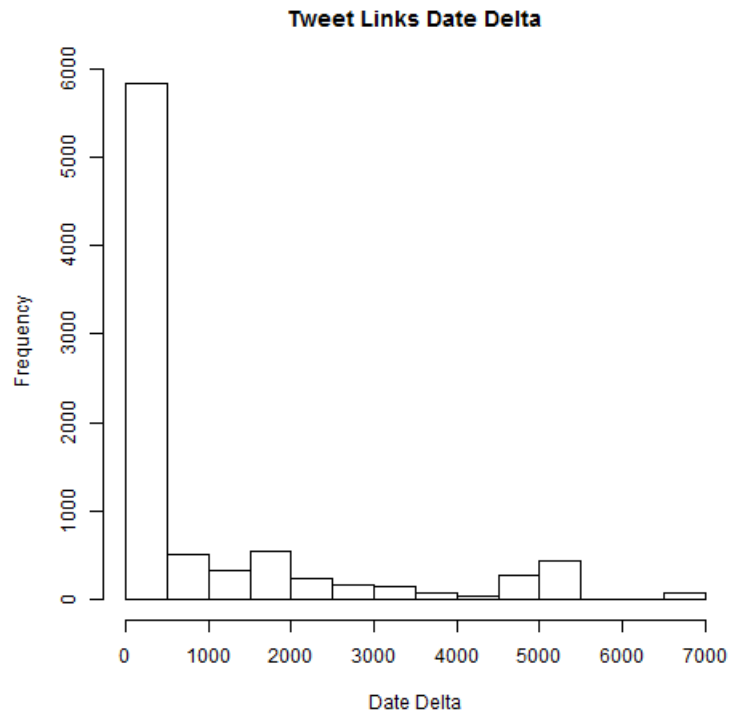


Figure 3: Delta between the Tweet and URI dates and their number of occurrences

### tweet-link-age-histogram.R

```
1 inputFile=read.csv("url-dates.csv",stringsAsFactors=F,header=FALSE,sep="\t")
2 inputFile
3 data=inputFile[,1]
4 png("tweet-link-age-histogram.png")
5 hist(data,main="Tweet Links Date Delta",freq=T,xlab="Date Delta",ylab="Frequency")
6 dev.off()
```

Listing 7: R program to plot the histogram for Tweet and URI dates

## 2.3 Summary of Tweet and URI's date delta

Median: 93

Mean: 952.355

Standard Deviation: 1587.182

Standard Error: 17.03

### 3 References

- [1] Tweepy library. <https://github.com/tweepy/tweepy>.
- [2] Tweepy Documentation. <http://docs.tweepy.org/en/v3.2.0/>.
- [3] Using twitter api keys. "http://thomassileo.com/blog/2013/01/25/using-twitter-rest-api-v1-dot-1-with-python/."
- [4] Github for carbondate. <https://github.com/HanySalahEldeen/CarbonDate>.
- [5] Producing simple graphs in r. <http://www.harding.edu/fmccown/r/>.