

# Assignment 2

CS532-s16: Web Sciences

Spring 2016

John Berlin

Generated on February 11, 2016



## Question 1

Write a Python program that extracts 1000 unique links from Twitter.  
<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 many different shortened URIs for [www.cnn.com](http://www.cnn.com) (t.co, bit.ly, goo.gl, etc.).

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.

## Answer

Extracting a 1000 unique links from Twitter at first seemed like a daunting task. The first steps I took to solving this problem was to look at the libraries twitter themselves suggested for use [Twitter Libraries](#). I also did search for blog posts on twitter mining and a tutorial that was most helpful.

The [tutorial](#) Mining Twitter Data with Python used the library Tweepy which I was quick to modify for use in this assignment.

To get the uri's use the file `twitter_stream.py`

1. `twitter_stream.py` <-q><the query>
  - Takes a query to look for in the twitter stream
  - Connect to the twitter stream and listen
  - On Tweet that matches query get the json data
  - Extract the uri using the expanded url field
  - Append the uri to the file
  - Code in listing [1](#) starting on page [3](#):

```

1  #!/usr/bin/env python3
2  from tweepy import Stream
3  from tweepy import OAuthHandler
4  from tweepy.streaming import StreamListener
5  import time
6  import argparse
7  import string
8  import config
9  import json
10 import re
11
12 # this code was adapted from the example found at
13 # https://gist.github.com/bonzanini/af0463b927433c73784d
14
15
16 # this regex was used to extract un-expanded urls before thinking
17 # smart was an idea
18 reg = re.compile('http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@.&+]|[*\\(\\)
19 ,]|(?:%[0-9a-fA-F][0-9a-fA-F]))+',re.VERBOSE | re.IGNORECASE)
20
21
22 def get_parser():
23     """Get parser for command line arguments."""
24     parser = argparse.ArgumentParser(description="Twitter
25         Downloader")
26     parser.add_argument("-q",
27         "--query",
28         dest="query",
29         help="Query/Filter",
30         default='')
31
32     return parser
33
34
35 # Subclass StreamListener so that I can handle the data
36 # Normally the StreamListener simply consumes and does nothing
37 # You must subclass for anything to happen
38 class MyListener(StreamListener):
39     """Custom StreamListener for streaming data."""
40
41     def __init__(self, query):
42         query_fname = format_filename(query)
43         self.outfile = "urls%s.dat" % query_fname
44         useragent = 'Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv
45             :44.0) Gecko/20100101 Firefox/44.01 '
46         # self.outfile = "%s/stream-%s.json" % (data_dir,
47             query_fname)
48
49     # when there is a tweet simply append the extracted urls to a
50     # file
51     def on_data(self, data):
52         try:
53             with open(self.outfile, 'a') as f:
54                 map = json.loads(data)
55                 #I found in the json data that there is an
56                 #expanded_url portion
57                 #Which unshortens one level
58                 #Even then some shortened uris get shortened
59                 for it in map['entities']['urls']:

```

```

51         f.write("%s\n"%it['expanded_url'])
52
53         return True
54     except BaseException as e:
55         #if bad data just let me know and move on
56         print("Error on-data: %s" % str(e))
57         time.sleep(5)
58     return True
59
60     def on_error(self, status):
61         print(status)
62         return True
63
64
65     def format_filename(fname):
66         """Convert file name into a safe string.
67         Arguments:
68             fname — the file name to convert
69         Return:
70             String — converted file name
71         """
72         return ''.join(convert_valid(one_char) for one_char in fname)
73
74
75     def convert_valid(one_char):
76         """Convert a character into '_' if invalid.
77         Arguments:
78             one_char — the char to convert
79         Return:
80             Character — converted char
81         """
82         valid_chars = "-_.%s%s" % (string.ascii_letters, string.digits)
83         if one_char in valid_chars:
84             return one_char
85         else:
86             return '_'
87
88
89     if __name__ == '__main__':
90         parser = argparse.ArgumentParser(description="Twitter
91             Downloader")
92         parser.add_argument("-q",
93             "--query",
94             dest="query",
95             help="Query/Filter",
96             default='')
97         args = parser.parse_args()
98         #set up oauth
99         auth = OAuthHandler(config.consumer_key, config.consumer_secret
100             )
101         auth.set_access_token(config.access_token, config.access_secret
102             )
103         # open a stream up and give it my listener
104         twitter_stream = Stream(auth, MyListener(args.query))
105         # filter the open stream for the query
106         # this looks for recent tweets only
107         twitter_stream.filter(track=[args.query])

```

---

Listing 1: Python program to mine urls for a twitter query

## Question 2

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

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

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

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.

\* = `https://en.wikipedia.org/wiki/Histogram`

## Answer

To answer this question and the last one I combined the steps to generate the required data into as few files as possible. The following steps give a break down to use the files with detail into what is going on is provided as comments in the code

1. run `local.py` from the carbon date library
  - use method `uniques` and `zeroNoneZero` from `util.py` to ensure the selected uris are unique, also to generate the nonzero uri file
  - copy the file `nonzero.txt` to the directory containing the Carbon Dating library with modified `local.py`
  - The file used to provide the timemap support was originally written by Scott Ainsworth which I modified to be compatible with python3. It is found in listing 3 starting on page 8. The number of mementos is the length of the keys to memento map in TimeMap class contained in this file
  - run and copy dumped json file called `dated.json` back to location where `carbonDating.py` is
  - Code in listing 5 starting on page 12
2. `carbonDating.py`
  - Count the number of Mementos per uri if there is an exception there are no mementos for the uri so associate a count of 0 with that uri
  - Parse the json file `dated.json` with the uri's in contained in the nonzero file produced by running the `zeroNoneZero` from `util.py`
  - the age of the memento is the python Date delta from today to age gotten from running the modified `local.py` in listing 5 starting on page 12 from the CarbonDating library

- Code in listing 6 starting on page 15

The R script to generate the figure 1 shows the histogram of the mementos counts is found is shown below.

```

1 # set the working directory
2 setwd(getwd())
3 # get dataframe site ,count
4 mementoCount <- read.csv('counted.csv')
5 # make histogram
6 h <- hist(mementoCount$count, main="Histogram of Memento Counts",
7           xlab = "Count")
8 # add count labels
9 text(h$mids, h$counts, labels=h$counts, adj=c(0.5, -0.5))

```

Listing 2: Memento Count Histogram R script

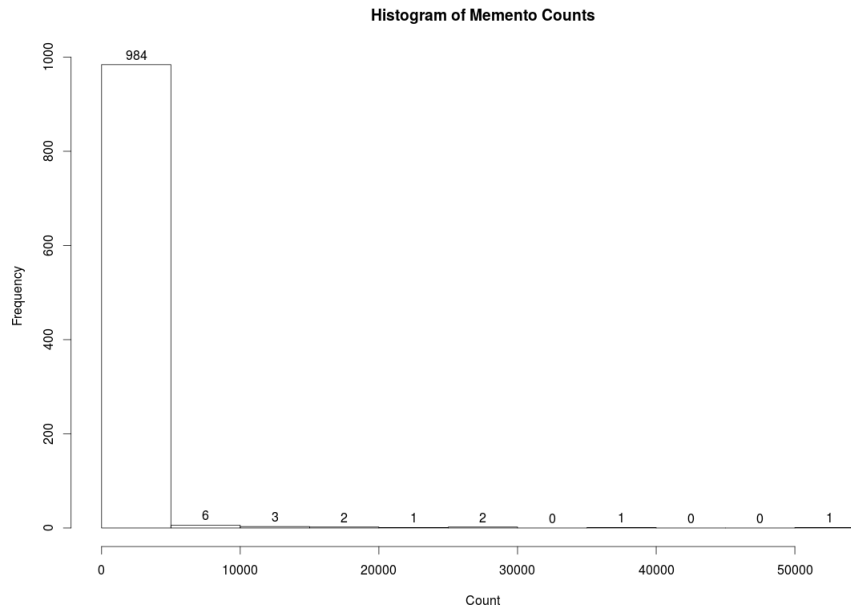


Figure 1: histogramCounts

As seen in the Histogram there were more uris with less than 1000 mementos



```

1  #!/usr/bin/python -B
2
3  import re
4  import urllib.request
5
6  import dateutil.parser
7  import requests
8
9  tokenizer = re.compile('(<[^>]+>|a-zA-Z+= "[^"]*" |[:;])\\s*')
10
11
12  class TimeMap:
13      def __init__(self, timemap_uri):
14          self.original = None
15          self.timebundle = None
16          self.timegate = None
17          self.timemap = None
18          self.first_memento = None
19          self.last_memento = None
20          self.mementos = {}
21          self._tokens = TimeMapTokenizer("http://mementoproxy.cs.
22                                         odu.edu/aggr/timemap/link/1/%s" % timemap_uri)
23          link = self.get_next_link()
24          while link is not None:
25              if link[0] == 'memento':
26                  self.mementos[link[1]] = link[2]
27              elif link[0] == 'original':
28                  self.original = link[2] if link is not None else
29                  None
30              elif link[0] == 'timebundle':
31                  self.timebundle = link[2] if link is not None else
32                  None
33              elif link[0] == 'timegate':
34                  self.timegate = link[2] if link is not None else
35                  None
36              elif link[0] == 'timemap':
37                  self.timemap = link[2] if link is not None else
38                  None
39              elif link[0] == 'first memento':
40                  self.mementos[link[1]] = link[2]
41                  self.first_memento = link[1] if link is not None
42                  else None
43              elif link[0] == 'last memento':
44                  self.mementos[link[1]] = link[2]
45                  self.last_memento = link[1] if link is not None
46                  else None
47              link = self.get_next_link()
48
49      def get_next_link(self):
50          uri = None
51          datetime = None
52          rel = None
53          resource_type = None
54          for token in self._tokens:
55              if token[0] == '<':
56                  uri = token[1:-1]
57              elif token[:9] == 'datetime=':

```

```

51         datetime = token[10:-1]
52     elif token[:4] == 'rel=':
53         rel = token[5:-1]
54     elif token[:5] == 'type=':
55         resource_type = token[6:-1]
56     elif token == ';':
57         None
58     elif token == ',':
59         return (rel, dateutil.parser.parse(datetime)
60                 if datetime is not None else None,
61                 uri, resource_type)
62     else:
63         raise Exception('Unexpected timemap token', token)
64 if uri is None:
65     return None
66 else:
67     return (rel, dateutil.parser.parse(datetime)
68             if datetime is not None else None,
69             uri, resource_type)
70
71 def __getitem__(self, key):
72     return self.mementos[key]
73
74
75 class TimeMapTokenizer:
76     def __init__(self, timemap_uri):
77         self.http = urllib.request.urlopen(timemap_uri)
78         # self._tmfile = requests.get(timemap_uri).iter_lines()
79         self.r = requests.get(timemap_uri)
80         self._tmfile = self.r.iter_lines()
81         # print(self.r.text)
82         self._tokens = []
83         self.lines = []
84         self.size = 0
85         self.cur = 0
86         self.doIt()
87
88     def doIt(self):
89         for line in self._tmfile:
90             self.lines.append(line.decode("utf-8"))
91         self.size = len(self.lines)
92
93     def __next__(self):
94         if len(self._tokens) == 0:
95             if self.cur == self.size:
96                 raise StopIteration
97             line = self.lines[self.cur]
98             self.cur += 1
99             if self.cur == self.size:
100                 raise StopIteration
101             self._tokens = tokenizer.findall(line)
102         return self._tokens.pop(0)
103
104     def __iter__(self):
105         return self

```

Listing 3: TimeMaps

```

1 import requests
2 import csv
3
4 # get set of unique urls
5 def unique():
6     s = set()
7     with open("urls.dat", "r") as o:
8         for line in o:
9             url = line.rstrip("\n")
10            s.add(url)
11    with open("urls.dat", "w+") as write:
12        for url in sorted(s):
13            write.write("%s\n"%url)
14
15 # get the real urls for shortened ones
16 def uniqueShortened():
17     useragent = 'Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:44.0)
18                Gecko/20100101 Firefox/44.01'
19     session = requests.Session()
20     session.headers.update({'User-Agent': useragent})
21     s = set()
22     c = 0
23
24     with open("shortenedURLS.txt", "r") as o:
25         for line in o:
26             url = line.rstrip("\n")
27             try:
28                 r = session.get(url)
29                 """ :type : requests.Response """
30                 # if status code is 200 then we have a good link
31                 if r.status_code == 200:
32                     # add it our set of urls
33                     s.add(r.url)
34                     c += 1
35                     print(c)
36             except:
37                 continue
38
39     print(len(s))
40     print(
41         "
42         ++++++
43     ")
44
45     session.close()
46     with open("urls2.dat", "w+") as write:
47         for url in sorted(s):
48             print(url)
49             write.write("%s\n" % url)
50
51 # method to write out the zero and non-zero memento count urls
52 def zeroNoneZero():
53     # count()
54     zero = []
55     nonZero = []
56
57     with open("counted.csv", "r") as csvFile:

```

```

55     reader = csv.DictReader(csvFile)
56     for row in reader:
57         print(row)
58         print(row["site"], row["count"])
59         if int(row["count"]) > 0:
60             nonZero.append(row["site"])
61         else:
62             zero.append(row["site"])
63
64     with open("nonZero.txt", "w+") as onz:
65         for site in nonZero:
66             onz.write("%s\n" % site)
67
68     with open("zero.txt", "w+") as zo:
69         for site in zero:
70             zo.write("%s\n" % site)

```

Listing 4: util

```

1 from checkForModules import checkForModules
2 import json
3 from ordereddict import OrderedDict
4 #import simplejson
5 import urlparse
6 import re
7
8 from getBitly import getBitlyCreationDate
9 from getArchives import getArchivesCreationDate
10 from getGoogle import getGoogleCreationDate
11 from getBacklinks import *
12 from getLowest import getLowest
13
14 from getLastModified import getLastModifiedDate
15 #Topsy service is no longer available
16 #from getTopsyScraper import getTopsyCreationDate
17 from htmlMessages import *
18 from pprint import pprint
19
20 from threading import Thread
21 import Queue
22 import datetime
23
24 import os,sys, traceback
25
26
27
28
29 def cd(url, backlinksFlag = False):
30
31     #print 'Getting Creation dates for: ' + url
32     #scheme missing?
33     parsedUrl = urlparse.urlparse(url)
34     if( len(parsedUrl.scheme)<1 ):
35         url = 'http://' + url
36     threads = []
37     outputArray = ['', '', '', '', '', '']
38     now0 = datetime.datetime.now()
39
40
41     lastmodifiedThread = Thread(target=getLastModifiedDate, args=(
42         url, outputArray, 0))
43     bitlyThread = Thread(target=getBitlyCreationDate, args=(url,
44         outputArray, 1))
45     googleThread = Thread(target=getGoogleCreationDate, args=(url,
46         outputArray, 2))
47     archivesThread = Thread(target=getArchivesCreationDate, args=(
48         url, outputArray, 3))
49
50     if( backlinksFlag ):
51         backlinkThread = Thread(target=
52             getBacklinksFirstAppearanceDates, args=(url,
53                 outputArray, 4))
54
55     #topsyThread = Thread(target=getTopsyCreationDate, args=(url,
56         outputArray, 5))

```

```

51
52     # Add threads to thread list
53     threads.append(lastmodifiedThread)
54     threads.append(bitlyThread)
55     threads.append(googleThread)
56     threads.append(archivesThread)
57
58     if( backlinksFlag ):
59         threads.append(backlinkThread)
60
61     #threads.append(topsyThread)
62
63
64     # Start new Threads
65     lastmodifiedThread.start()
66     bitlyThread.start()
67     googleThread.start()
68     archivesThread.start()
69
70     if( backlinksFlag ):
71         backlinkThread.start()
72
73     #topsyThread.start()
74
75
76     # Wait for all threads to complete
77     for t in threads:
78         t.join()
79
80     # For threads
81     lastmodified = outputArray[0]
82     bitly = outputArray[1]
83     google = outputArray[2]
84     archives = outputArray[3]
85
86     if( backlinksFlag ):
87         backlink = outputArray[4]
88     else:
89         backlink = ''
90
91     #topsy = outputArray[5]
92
93     #note that archives["Earliest"] = archives[0][1]
94     try:
95         #lowest = getLowest([lastmodified, bitly, google, archives
96         #                    [0][1], backlink, topsy]) #for thread
97         lowest = getLowest([lastmodified, bitly, google, archives
98         #                    [0][1], backlink]) #for thread
99     except:
100         print sys.exc_type, sys.exc_value, sys.exc_traceback
101
102
103     result = []
104
105     result.append(("URI", url))
106     result.append(("Estimated Creation Date", lowest))

```

```

106     values = OrderedDict(result)
107     r = json.dumps(values, sort_keys=False, indent=2, separators=(
108         ', ', ': '))
109
110     now1 = datetime.datetime.now() - now0
111
112     return r
113
114
115
116 if len(sys.argv) == 1:
117     print "Usage: ", sys.argv[0] + " url backlinksOnOffFlag ( e.g:
118         " + sys.argv[0] + " http://www.cs.odu.edu [--compute-
119             backlinks] )"
120 elif len(sys.argv) == 2:
121     time.strptime("1995-01-01T12:00:00", '%Y-%m-%dT%H:%M:%S')
122     file = open("nonZero.txt", "r")
123     r = []
124     for url in file:
125         print url.rstrip("\n")
126         r.append(cd(url.rstrip("\n")))
127     file.close()
128     f = open("dated.json", "w+");
129     out = ", ".join(r)
130     f.write(out)
131     print out
132     f.close()
133 elif len(sys.argv) == 3:
134     time.strptime("1995-01-01T12:00:00", '%Y-%m-%dT%H:%M:%S')
135     file = open("nonZero2.txt", "r")
136
137     for url in file:
138         print url
139     file.close()
140     if (sys.argv[2] == '--compute-backlinks'):
141         cd(sys.argv[1], True)
142     else:
143         cd(sys.argv[1])

```

Listing 5: Modified local.py

```

1 import csv
2 import json
3 from datetime import date
4
5 from dateutil.parser import parse
6
7 import timemaps
8
9
10 # method to count the number of timemaps for a uri
11 def count():
12     counted = {}
13     # for each url in the 1000 urls
14     with open("urls.dat", "r") as o:
15         for line in o:
16             # remove the newline character
17             url = line.rstrip("\n")
18             try:
19                 # get the time map
20                 tm = timemaps.TimeMap(url)
21                 # the number of mementos is the number of keys
22                 counted[url] = len(tm.mementos.keys())
23                 print("%s, %i\n" % (url, len(tm.mementos.keys())))
24             except:
25                 # if an exception is raised it means no mementos
26                 counted[url] = 0
27     # write the counts to file
28     with open("counted.csv", "w+") as oo:
29         oo.write("site,count\n")
30         for url, count in counted.items():
31             print("%s, %i\n" % (url, count))
32             oo.write("%s, %i\n" % (url, count))
33
34 # method to count and carbon date the uri mementos
35 def countAndDate():
36     count()
37     zero = []
38     nonZero = []
39     all = {}
40     # get todays date
41     now = date.today()
42
43     with open("counted.csv", "r") as csvFile:
44         # read the count file
45         reader = csv.DictReader(csvFile)
46         for row in reader:
47             print(row)
48             print(row["site"], row["count"])
49             # get two collections the nonZero and the zero counts
50             if int(row["count"]) > 0:
51                 nonZero.append(row["site"])
52             else:
53                 zero.append(row["site"])
54             all[row["site"]] = row["count"]
55
56     # open the json data from carbondate
57     jdata = open("dated.json", "r")

```



```

58 data = json.load(jdata)
59 jdata.close()
60 cds = {}
61
62 for jd in data:
63     # create mapping for the uri to date
64     cds[jd['URI']] = parse(jd['Estimated Creation Date'])
65     print(jd, type(parse(jd['Estimated Creation Date'])))
66
67 with open("dated.csv", "w+") as out:
68     out.write("age,mementos\n")
69     # for all nonZero memento write the estimated age in days
70     for it in nonZero:
71         # python dates have subtraction for dates
72         # to mean the days between two daes
73         age = now - date(cds[it].year, cds[it].month, cds[it].
74                         day)
75         """ :type datetime.timedelta """
76         out.write("%s,%s\n" % (age.days, all[it]))
77         print(it, age.days, all[it])
78
79
80
81
82 if __name__ == "__main__":
83     countAndDate()

```

Listing 6: Python program to mine urls for a twitter query

### Question 3

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

<http://ws-dl.blogspot.com/2014/11/2014-11-14-carbon-dating-web-version-20.html>

Note: you'll should download the library and run it locally; don't try to use the web service.

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.

Not all URIs will have Mementos, and not all URIs will have an estimated creation date. State how many fall into either categories.

### Answer

Out of the 1000 uri mementos 28 of them were found to have an age of 0 but none of the uri were found to have a memento count of 0.

Using the generated files from the python files I used R code found in listing 7 starting on page 17 to generate figure 2.

In the r script I found the median memento count and partitioned the data frame into two halves those that fell below the median and those the fall about it. This was to have two graphs that clearly show the entire graph.

The number of mementos that were above the median are fare more dispersed than those below it. As seen in the all age count graph it is difficult to see the distribution due to outliers.

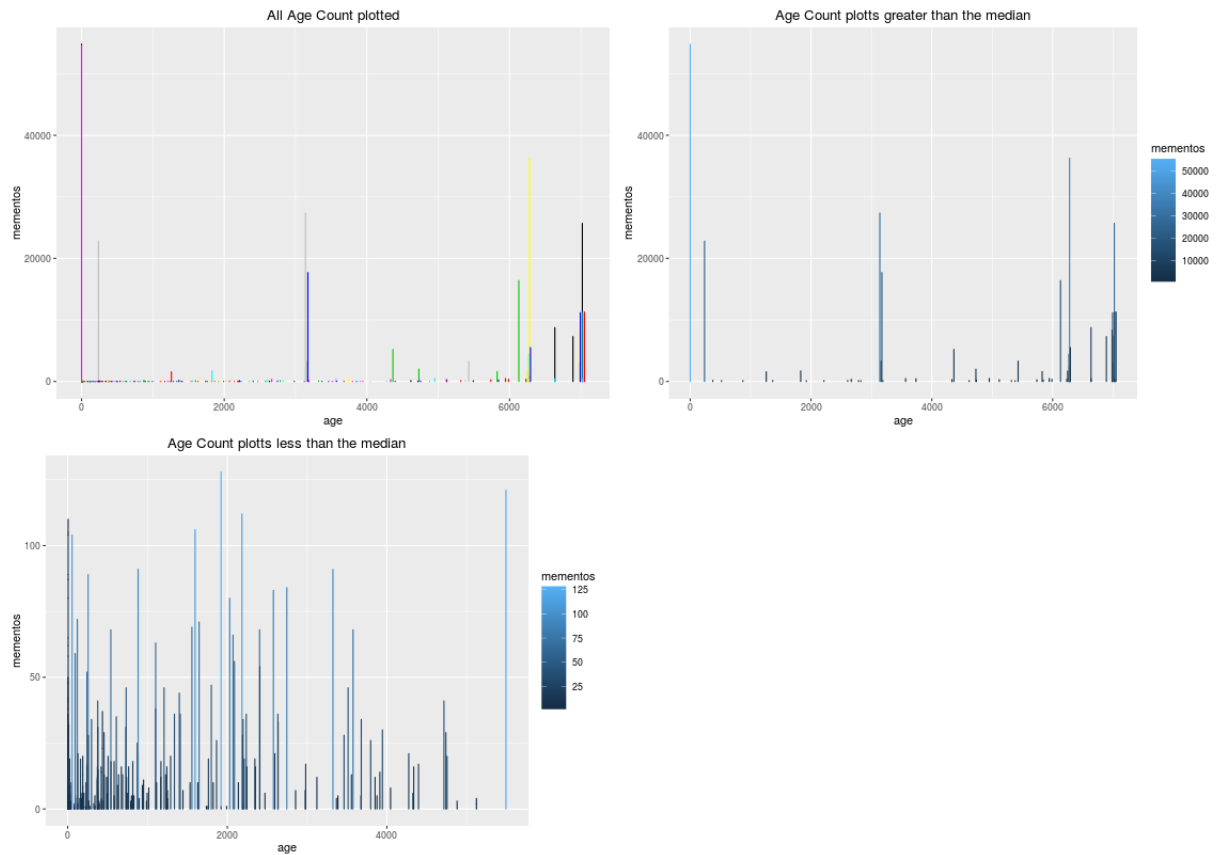


Figure 2: histogramCounts

```
1 library(ggplot2)
2 setwd(getwd())
3
4 #this function wonderfully borrowed from
5 #http://www.cookbook-r.com/Graphs/Multiple_graphs_on_one_page-%28
6 #ggplot2%29/
7 multiplot <- function(..., plotlist=NULL, file, cols=1, layout=NULL) {
```

```

7   library(grid)
8
9   # Make a list from the ... arguments and plotlist
10  plots <- c(list(...), plotlist)
11  numPlots = length(plots)
12
13  # If layout is NULL, then use 'cols' to determine layout
14  if (is.null(layout)) {
15    # Make the panel
16    # ncol: Number of columns of plots
17    # nrow: Number of rows needed, calculated from # of cols
18    layout <- matrix(seq(1, cols * ceiling(numPlots/cols)),
19                      ncol = cols, nrow = ceiling(numPlots/cols))
20  }
21
22  if (numPlots==1) {
23    print(plots[[1]])
24  } else {
25    # Set up the page
26    grid.newpage()
27    pushViewport(viewport(layout = grid.layout(nrow(layout), ncol(
28      layout))))
29
30    # Make each plot, in the correct location
31    for (i in 1:numPlots) {
32      # Get the i,j matrix positions of the regions that contain
33      # this subplot
34      matchidx <- as.data.frame(which(layout == i, arr.ind = TRUE))
35
36      print(plots[[i]], vp = viewport(layout.pos.row = matchidx$row
37      ,
38      layout.pos.col = matchidx$col
39      ))
40    }
41  }
42
43  # get the age,mementos dataframe
44  dateNum <- read.csv('dated.csv')
45
46  # find the median value to split the plots
47  mcount <- sort(unique(dateNum$mementos))
48  midval <- median(mcount)
49
50  # find all values less than the midval
51  lt <- subset(dateNum, dateNum$mementos <= midval)
52  # find all values greater than the midval
53  gt <- subset(dateNum, dateNum$mementos >= midval)
54
55  # plot the entire dataframe
56  # use stat_identity to plot the actual data not a count
57  a<-ggplot(dateNum, aes(y=mementos,x=age,color=mementos))+
58  geom_bar(stat="identity")+scale_color_identity()+
59  ggtitle("All Age Count plotted")
60
61  # plot the lower values, use color values based on memento count

```

```

60 b<-ggplot(lt, aes(y=mementos,x=age,color=mementos))+
61   geom_bar(stat="identity",position = "stack",aes(color=mementos))+
62   ggtitle(as.character("Age Count plotts less than the median"))
63
64 # plot the upper values, use color values based on memento count
65 c<-ggplot(gt, aes(y=mementos,x=age))+
66   geom_bar(stat="identity",position = "stack",aes(color=mementos))+
67   scale_fill_identity(breaks = lt$mementos,guide = "legend")+
68   ggtitle(as.character("Age Count plotts greater than the median"))
69
70 #combine the three plots
71 multiplot(a,b,c,cols=2)

```

Listing 7: Memento Count Age R script