

# **CS-432/532 Introduction to Web Science: Assignment #2**

Due on Thursday, February 11, 2016

*Dr. Michael L. Nelson*

**Plinio Vargas**  
pvargas@cs.odu.edu

# Contents

<b>Preface</b>	<b>1</b>
<b>Problem 1</b>	<b>1</b>
1.1 Approach . . . . .	1
1.2 Solution . . . . .	2
<b>Problem 2</b>	<b>4</b>
2.1 Approach . . . . .	4
2.2 Solution . . . . .	6
<b>Problem 3</b>	<b>7</b>
3.1 Approach . . . . .	7
3.2 Solution Part-1 . . . . .	8
3.3 Solution Part-2 . . . . .	9

## List of Figures

1	Flowchart diagram for <i>TwitterURI.py</i> . . . . .	1
2	Memento Histogram for 1000 URIs Extracted from Twitter . . . . .	6
3	Flowchart diagram for <i>AgeURI.py</i> . . . . .	7
4	Age-Memento Graph . . . . .	8
5	Age-Memento Graph with Various Plot Size . . . . .	9

## List of Tables

1	Histogram Memento Data for 1000 URIs Extracted from Twitter . . . . .	10
2	Age and Memento Data for 1000 URIs Extracted from Twitter . . . . .	11
3	Data Grouped by Age and Memento from Table 2 . . . . .	12

## Preface

Since the course material will become more complex as the semester progresses, and knowing that many solutions to new problems can be built on previous ones, I decided to create a library to access and re-utilize code.

Some of these functions can be really small (4-7 lines), however they solve a specific problem, so it is better to separate them from the rest of the code.

The location, interaction and functionality of this library, and the directory structure for this assignment are presented in more detail in the README file within the same folder hierarchy where this file was found.

## Problem 1

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

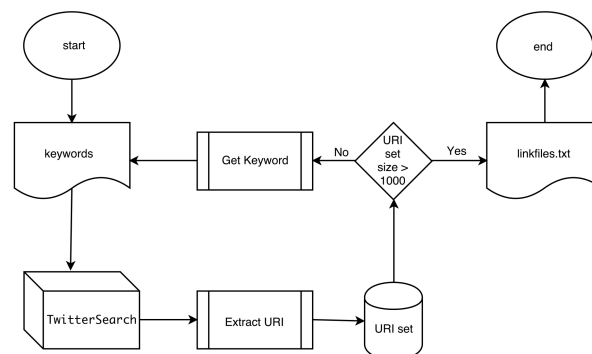
[urlhttp://thomassileo.com/blog/2013/01/25/using-twitter-rest-api-v1-dot-1-with-python/](http://thomassileo.com/blog/2013/01/25/using-twitter-rest-api-v1-dot-1-with-python/)

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.).

### 1.1 Approach

Figure 1: Flowchart diagram for *TwitterURI.py*



*TwitterURI.py* is the Python program solution for problem 1. **keywords** is a text file document which contains subject of interest, such as baseball, “Dominican, Republic”, etc.; **TwitterSearch** is an API that provides a set of tweets from Twitter given a particular keyword; **Extract URI** filters any URI within the tweet; **URI set** is a unique collection of URIs; **Get Keyword** gets the next keyword from the text file; and finally **linkfiles.txt** is a collection of 1000 unique URIs.

Since we do not know the number of tweets returned from a particular keyword search, a file was created

(*keywords.txt*) as input for **TwitterURI.py** (lines 70-75). Then, the API **TwitterSeach** in line (28-34) is used, which returns  $n$  number of tweets. From the tweet object we can extract the URI using regular expression (line 44):

```
url = re.findall(r'(https?://[^\s]+)', tweet['text'])
```

The extracted URI is validated and placed in a set: *all\_links\_set* lines (49-54):

```
# check if url has been already used
if uri not in root_url_array:
    root_url_array.append(uri)

if validators.url(uri):
    AddURI(all_links_set, uri)
```

The last line of the portion code above uses a function *AddURI* imported from the self developed library *WEBlib.py* (a collection of functions that may have use for future applications). The function takes as a parameter a set object and a string (uri). *AddURI* takes care of all details regarding request response, connection timeouts, resource redirects, etc. ONLY if final URI provides a response 200 it gets added to the set.

Finally, if the set size is not greater than or not equal to 1000, we continue to get more keywords and repeat the process, otherwise the deliverable (1000 unique URIs) is written in *linkfiles.txt*.

## 1.2 Solution

See attached embedded file: *linkfiles.txt*

Below is the entire code **TwitterURI.py**, which is also attached to this document.

TwitterURI.py

---

```
1 import re
2 import sys
3 import os
4 import requests
5 import validators
6 import urllib3
7
8 PACKAGE_PARENT = '..'
9 SCRIPT_DIR = os.path.dirname(os.path.realpath(os.path.join(os.getcwd(),
10     os.path.expanduser(__file__))))
11 sys.path.append(os.path.normpath(os.path.join(SCRIPT_DIR, PACKAGE_PARENT)))
12
13 from lib.WEBlib import AddURI
14 from TwitterSearch import *
15
16 __author__ = 'Plinio H. Vargas'
17 __date__ = 'Wed, Feb 03, 2016 at 22:27:02'
18 __email__ = 'pvargas@cs.odu.edu'
19
20 """
21 Given a set of key words this program makes use of TwitterSearch API to search for any tweets
```

```
21 related to the key words, and then extract the URI within in the tweet.
22 """
23 root_url_array = []
24 all_links_set = set()
25 final_link = ''
26
27
28 def add_links(ts, key_words):
29     global final_link
30     try:
31         tso = TwitterSearchOrder() # create a TwitterSearchOrder object
32         tso.set_keywords(key_words) # let's define all words we would like to have a look for
33         tso.set_language('en') # we want to see English tweets only
34         tso.set_include_entities(False) # and don't give us all those entity information
35
36         # get all tweet feeds with parameter above
37         for tweet in ts.search_tweets_iterable(tso):
38             if len(all_links_set) >= 1000:
39                 break
40
41             tweet_text = tweet['text']
42
43             # extract url from tweet
44             url = re.findall(r'(https?://[^\s]+)', tweet['text'])
45
46             # analyze url
47             for uri in url:
48
49                 # check if url has been already used
50                 if uri not in root_url_array:
51                     root_url_array.append(uri)
52
53                 if validators.url(uri):
54                     AddURI(all_links_set, uri)
55
56     except TwitterSearchException as e: # take care of all those ugly errors if there are some
57         print(e)
58
59     return
60
61 # it's about time to create a TwitterSearch object with our secret tokens
62 ts = TwitterSearch(
63     consumer_key='mysconsumerkey',
64     consumer_secret='myconsumer_secret',
65     access_token='myaccess_token',
66     access_token_secret='myaccess_token_secret'
67 )
68
69 # find tweets for a particular key word
70 with open('keywords.txt', 'r') as file:
71     for record in file:
72         keys = record.strip().split(',')
73
```

```
74     key_words = []
75     for words in keys:
76         key_words.append(words.strip())
77
78     add_links(ts, key_words)
79     print(len(all_links_set))
80
81 # save unique URIs to a file
82 with open('linkfiles.txt', 'w') as out:
83     for link in all_links_set:
84         out.write(link + '\n')
```

---

## Problem 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>

### 2.1 Approach

A Python code < ***MementosURI.py*** > was developed to generate the histogram data, and the Histogram Graph (Figure 2) was generated using R-code < *histogram.r* >. Both files are attached to this document. Since we will be using the data collected later on in the course, a folder **a2/mementodata** was created to hold all TimeMaps for URIs included in < *linkfiles.txt* > (solution to problem 1).

URIs in < *linkfiles.txt* > are written sequentially and only separated by a new-line character ('\n'). In order to make a relation TimeMap-URI, the following file name format was established:

*date\_created.relative\_position*

For example if file *20160204.21* is located in folder **a2/mementodata**, then it references the TimeMap for URI located in position 21 out of 1000 in file < *linkfiles.txt* >. If *20160204.1* is not found in the directory, then there is no TimeMap for the first URI in < *linkfiles.txt* >.

Given the schema to solve Problem 2, below is the pseudo-code used for its implementation:

Listing 1: MementosURI.py Pseudo Code

---

```
1  input: (I) file linkfiles.txt
2  output: (O) file histogramdata
3  Let A be an empty array where index  $i$  is number of mementos.
4  Let D be an stack containing all filenames in directory a2/mementodata
5  Let S be an integer containing number of URIs in I
6
7  # Initialize URIs with zero memento to number of URIs in I
8  A[0]  $\leftarrow$  S
9
10 begin
11     while file in D
12         A[0]  $\leftarrow$  A[0] - 1 # decrease number URIs with zero mementos
13          $i \leftarrow$  number of memento in file
14         D.pop()
15         if  $i \in A$  then
16             A[i]  $\leftarrow$  A[i] + 1
17         else
18             A[i]  $\leftarrow$  1
19
20     for  $k := 1 \rightarrow A.size$ 
21         write into O ( $k$ , A[k])
22 end
```

---

The python implementation found in < *MementosURI.py* > and pseudo-code above are very similar, and straightforward if we remove all exception handling required for the implementation. The input file  $I$  is used to get the number of URIs in our sample  $S$  (1000).  $O$  is the output file where histogram data will be saved.  $D$  contains all file names that have TimeMaps. Array  $A$  contains the data that will be written as a tuple in  $O$ .  $A$  starts as an empty array, so ONLY the indexes with values will be saved.

Different from the pseudo-code will be using instead of an array  $A$  a dictionary object [1] in the Python implementation.

Then, the only worthwhile explanation is how to find the value of index  $i$  or the number of mementos in any particular file. This is performed by the function *GetNumberMementos(filename)* located in our developed library *WEBlib.py*:

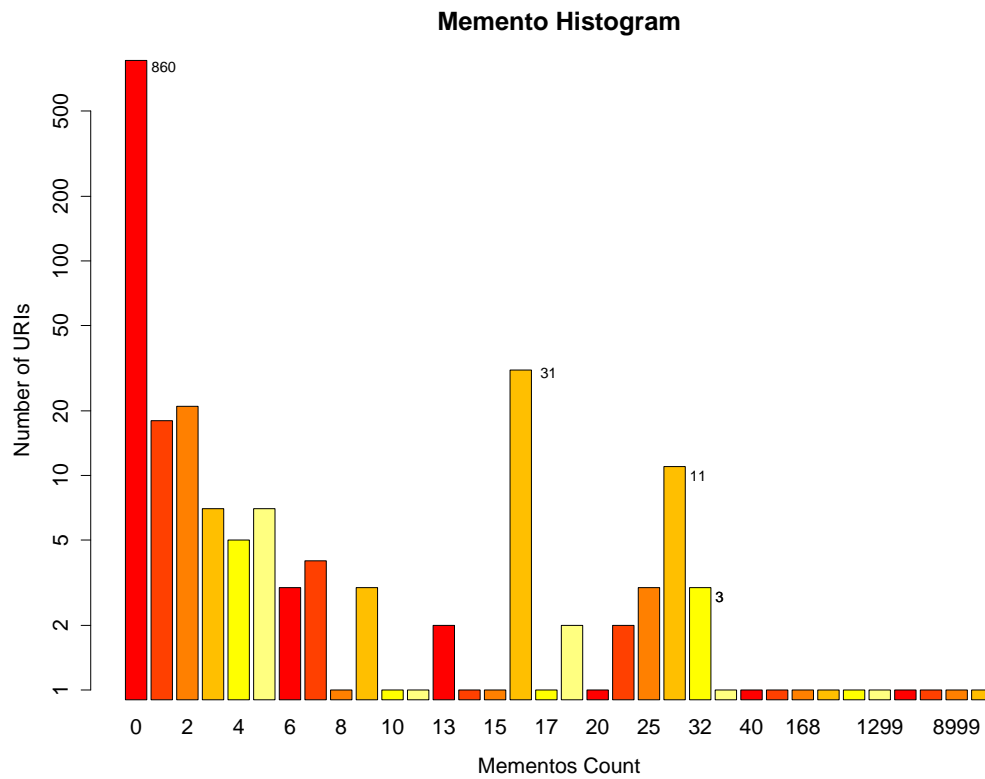
```
def GetNumberMementos(filename):
    """
    :param filename: Name of file to count number of mementos
    :return: number of mementos
    """
    with open(filename, 'r') as file:
        text = file.read()

    return len(re.findall(r'rel="memento"', text))
```

As we may notice, function *GetNumberMementos()* simple takes as an argument a file name. It opens the file for reading and it uses regular expression to count and return the number of rel="memento" occurrences in the file.

## 2.2 Solution

Figure 2: Memento Histogram for 1000 URIs Extracted from Twitter



The y-axis or Number of URIs is in a log scale. The data used to plot this histogram is contained in Table 1. The histogram resembles characteristics of a log curve (80/20) values. In our case 860 URIs out of a 1,000 had 0 mementos. At the tail end of the histogram we can notice that the memento counts gets higher as the number of URIs get smaller.



## Problem 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.

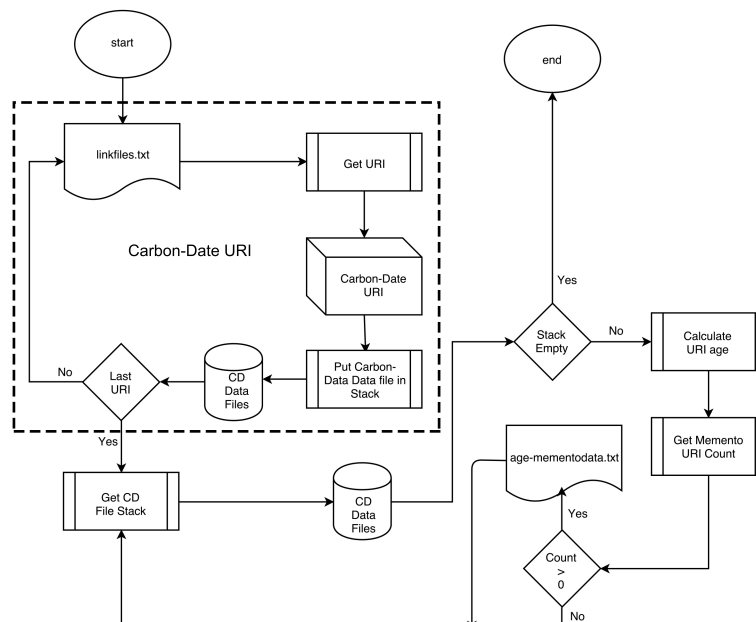
Not all URIs will have Mementos, and not all URIs will have an estimated creation date. State how many fall into either categories.

### 3.1 Approach

We will be using the same nomenclature file schema to associate the Carbon Date data with an URI's age. However, the Carbon Date files will be stored in a different directory: **a2/uriagedata**. Similar to problem 2 approach, python was used to generate the data, and R was utilized to create the graphs.

Shown below is a flowchart diagram for a python program used to generate the data:

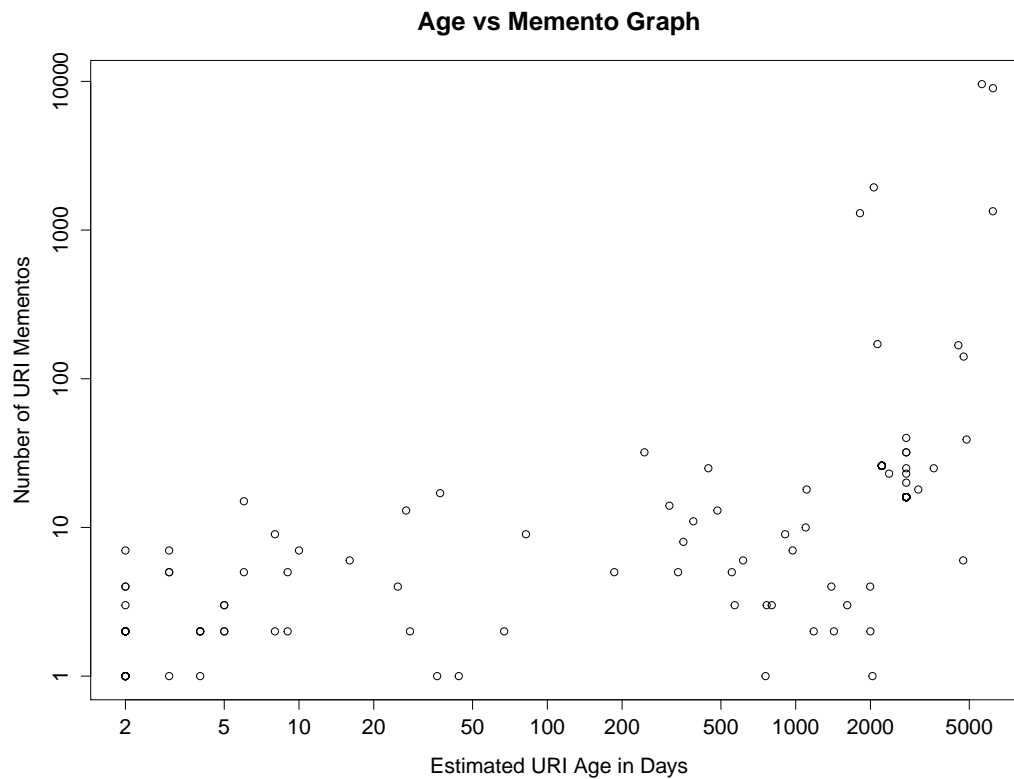
Figure 3: Flowchart diagram for *AgeURI.py*



*AgeURI.py* is the Python implementation for problem 3. The section in the diagram labeled **Carbon-Date URI** takes all URIs generated to solve problem 1 and uses the **Carbon Date Tool** obtained from <http://ws-dl.blogspot.com/2014/11/2014-11-14-carbon-dating-web-version-20.html>. The HTTP request output is saved into **uriagedata** directory. This directory is labeled in the diagram **CD Data Files**. The remaining part of the flowchart explains that **CD Data Files** result gets merged with the result from data file solution from problem 2 (**Get Memento URI count**) to produce our deliverable: **age-memento-data.txt** file.

### 3.2 Solution Part-1

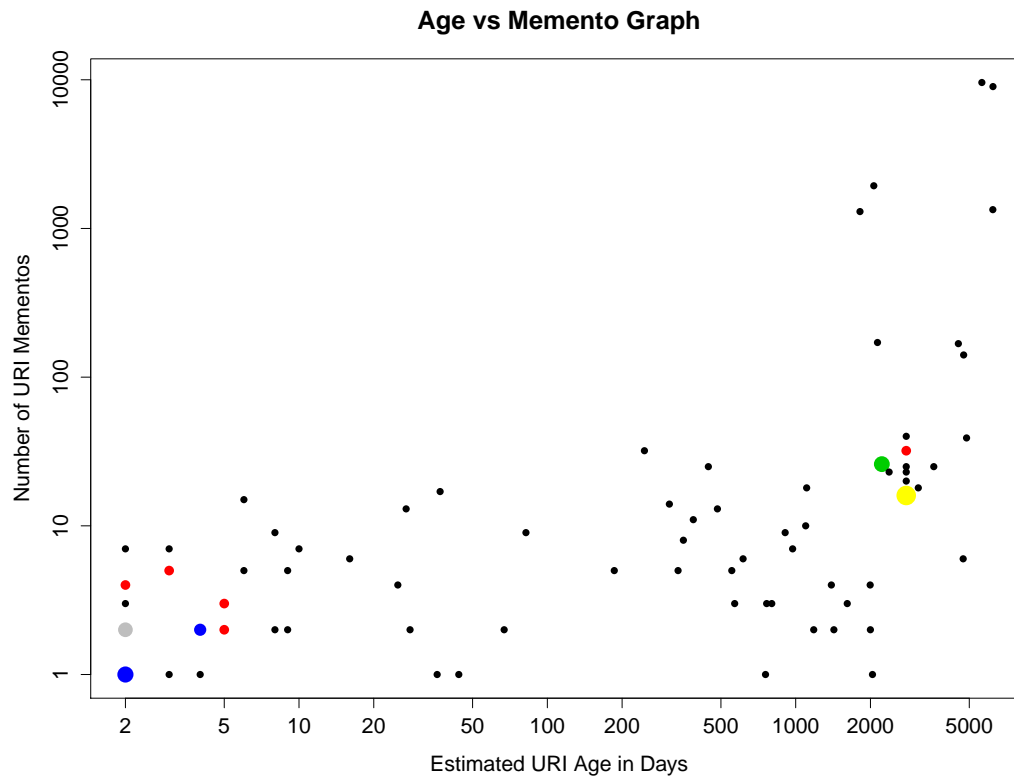
Figure 4: Age-Memento Graph



The x-axis and y-axis are in a log scale. The data used to plot this graph is contained in Table 2. Number of plotting points do not correspond to number of data set in Table 2 due to tuple value repetition with different URIs. It is difficult to appreciate repeating plots with darker plotting points on this graph. Then, it is a need to capture plot point repetition: Figure 5.

From the graph above we can notice a relationship between an URI's age and its Mementos count. In general, the younger the URI the less count of mementos it has, and vice versa. Figure 5 captures much better this inference, since the plot density is more representative.

Figure 5: Age-Memento Graph with Various Plot Size



The x-axis and y-axis are in a log scale. The data used to plot this graph is contained in Table 3. The plot point size corresponds to column labeled *Freq* from the same table utilizing the formula:  $\log Freq/2 + 1$ . The larger the plot point size, the greater the number of occurrences for any particular data set.

The analysis of Figure 5 is similar to Figure 4. In general, the younger the URI the less count of mementos it has, and vice versa. However, the plotting is more representative of its data. For example, we can appreciate that a great deal of pages 2 days old have 1 memento. This information is not capture on Figure 4.

### 3.3 Solution Part-2

Not all URIs will have Mementos, and not all URIs will have an estimated creation date. *AgeURI.py* calculated this information:

Out of 1000 URIs 382 or 38.20% have no estimated creation date.

Out of 1000 URIs 860 or 86.00% have no mementos.

Additionally, the number of URIs with an estimated creation date and a memento is 140 or 14.00%.

Table 1: Histogram Memento Data for 1000 URIs Extracted from Twitter

No-URI	Mementos
860	0
18	1
21	2
7	3
5	4
7	5
3	6
4	7
1	8
3	9
1	10
1	11
2	13
1	14
1	15
31	16
1	17
2	18
1	20
2	23
3	25
11	26
3	32
1	39
1	40
1	141
1	168
1	171
1	240
1	1299
1	1338
1	1937
1	8999
1	9586

Histogram data was generated on 02/04/2016 utilizing timeMaps from <http://mementoproxy.cs.odu.edu/aggr/timemap/link/1/> concatenated to the URIs contained in the attached file *linkfiles.txt*. Any particular tuple (No-URI, Memento) represents number of URIs that have a specific number of mementos. For example the first tuple translate that from our 1,000 URIs sample 850 of them had 0 mementos.

Table 2: Age and Memento Data for 1000 URIs Extracted from Twitter

Age (days)	Memento Qty	Age (days)	Memento Qty	Age (days)	Memento Qty
2	2	28	2	2787	16
2	4	36	1	2787	16
2	1	37	17	2787	16
2	2	44	1	2787	25
2	2	67	2	2787	16
2	2	82	9	2787	16
2	3	186	5	2787	32
2	2	246	32	2787	20
2	1	310	14	2787	16
2	1	336	5	2787	16
2	1	353	8	2787	16
2	2	387	11	2787	16
2	2	445	25	2787	16
2	1	484	13	2787	16
2	1	553	5	2787	16
2	1	568	3	2787	32
2	1	614	6	2787	16
2	1	756	1	2787	16
2	1	764	3	2787	16
2	2	801	3	2787	16
2	1	907	9	2787	16
2	4	972	7	2787	16
2	7	1097	10	2787	16
2	1	1107	18	2787	16
3	5	1182	2	2787	16
3	1	1392	4	2787	40
3	5	1425	2	2787	16
3	7	1613	3	2787	23
4	1	1815	1299	2787	16
4	2	1995	4	2787	16
4	2	1999	2	2787	16
4	2	2038	1	2787	16
4	2	2063	1937	2787	16
5	2	2135	171	2787	16
5	3	2223	26	2787	16
5	2	2223	26	2787	16
5	3	2223	26	3117	18
6	15	2223	26	3598	25
6	5	2223	26	4519	168
8	9	2223	26	4725	6
8	2	2223	26	4746	141
9	5	2223	26	4876	39
9	2	2223	26	5619	9586
10	7	2223	26	6221	1338
16	6	2223	26	6229	8999
25	4	2378	23		
27	13	2787	16		

The Age in days for all URIs were calculated on 02/05/2016 utilizing the Carbon-Dating tool found at <http://ws-dl.blogspot.com/2014/11/2014-11-14-carbon-dating-web-version-20.html>. The same methodology utilized on Table 2 to find mementos count was utilize here to find Mementos Qty.

Table 3: Data Grouped by Age and Memento from Table 2

Age (days)	Memento Qty	Freq	Age (days)	Memento Qty	Freq
2	1	12	568	3	1
2	2	8	614	6	1
2	3	1	756	1	1
2	4	2	764	3	1
2	7	1	801	3	1
3	1	1	907	9	1
3	5	2	972	7	1
3	7	1	1097	10	1
4	1	1	1107	18	1
4	2	4	1182	2	1
5	2	2	1392	4	1
5	3	2	1425	2	1
6	5	1	1613	3	1
6	15	1	1815	1299	1
8	2	1	1995	4	1
8	9	1	1999	2	1
9	2	1	2038	1	1
9	5	1	2063	1937	1
10	7	1	2135	171	1
16	6	1	2223	26	11
25	4	1	2378	23	1
27	13	1	2787	16	31
28	2	1	2787	20	1
36	1	1	2787	23	1
37	17	1	2787	25	1
44	1	1	2787	32	2
67	2	1	2787	40	1
82	9	1	3117	18	1
186	5	1	3598	25	1
246	32	1	4519	168	1
310	14	1	4725	6	1
336	5	1	4746	141	1
353	8	1	4876	39	1
387	11	1	5619	9586	1
445	25	1	6221	1338	1
484	13	1	6229	8999	1
553	5	1			

Same data set that resulted in generating the data for Table 2 was utilized to generate values for this table. The grouping of (Age, Memento Qty) from Table 2 added a new column: Freq.

## References

- [1] Lutz, Mark (2013). List and Dictionaries. *Learning Python* (5th ed.). (pp. 262-263). Sebastopol, CA: O'Reilly Media.