

CS 532: Assignment 3

Dinesh Kumar Paladhi

Spring 2016

Contents

1	Problem 1	2
	1.1 Solution	2
	1.2 Code Listing	3
	1.3 Results	4
2	Problem 2	7
	2.1 Solution	7
	2.2 Code Listing	9
	2.3 Results	10
3	Problem 3	13
	3.1 Solution	13
	3.2 Results	14
4	Problem 4	15
	4.1 Solution	15
	4.2 Code Listing and Results	16

1 Problem 1

Download the 1000 URIs from assignment #2. ‘curl’, ‘wget’, or ‘lynx’ are all good candidate programs to use. We want just the raw HTML, not the images, stylesheets, etc.

from the command line:

```
% curl http://www.cnn.com/ > www.cnn.com
```

```
% wget -O www.cnn.com http://www.cnn.com/
```

```
% lynx -source http://www.cnn.com/ > www.cnn.com
```

‘www.cnn.com’ is just an example output file name, keep in mind that the shell will not like some of the characters that can occur in URIs (e.g., ‘?’, ‘&’). You might want to hash the URIs, like:

```
% echo -n 'http://www.cs.odu.edu/show_features.shtml?72' | md5
41d5f125d13b4bb554e6e31b6b591eeb
```

(‘md5sum’ on some machines; note the ‘-n’ in echo -- this removes the trailing newline.)

Now use a tool to remove (most) of the HTML markup. ‘lynx’ will do a fair job:

```
% lynx -dump -force_html www.cnn.com > www.cnn.com.processed
```

Use another (better) tool if you know of one. Keep both files for each URI (i.e., raw HTML and processed).

1.1 Solution

1. I started the 1st question by running a sample URI with curl command syntax as given in the question. This created a new file in my directory with html content of the URI that I have mentioned.
2. So I got to know by executing the curl command I will be getting raw data file of the URI.
3. Now the next important point is naming the output files. If we name the output files with URI names, then I will be getting error at one point because it does not accept special characters.
4. This issue can be resolved by creating a hash code which is unique for each URI. I used hashlib library to convert URIs to respective hash codes.
5. Now I have named the file as Raw or processed followed by unique number for each URI followed by the hash code. I add unique number so that it would be useful for me to get URI back from the hash code in future.
6. The same process should be repeated with lynx command (given in question) to get processed file which eliminates all the html content and gives only data.
7. I have written a loop which runs 1000 times and each time it takes a URI and generates a raw file and a processed file. Therefore at the end I will be getting 1000 raw files and 1000 processed files.
8. The code for this program is shown below. Figure 1 and Figure 2 shows a sample Raw file and sample processed file. Figure 3 and Figure 4 shows list of raw and processed files created.

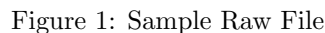
1.2 Code Listing

Here is the Python program for acquiring raw and processed data files for each URI.

```
1 import json
2 import commands
3 import hashlib
4
5 file_1 = open('my_json_data.json','r')
6 count = 0
7 for each_line in file_1.readlines():
8     dummy_line=json.loads(each_line)
9     url = dummy_line['url']
10    hash = hashlib.md5(url.encode())
11    final_hash = hash.hexdigest()
12    count = count +1
13    file_name_1= "Raw"+'-' +str(count) + '-' + final_hash + '.txt'
14    co_1 = 'curl -s -L ' + url + ' > ./Raw-htmldata/' + file_name_1
15    commands.getoutput(co_1)
16    file_name_2= "processed"+'-' +str(count) + '-' + final_hash + '.txt' #Naming a file
17    co_2 = 'lynx -dump -force.html ' + url + ' > ./processed-htmldata/' + file_name_2 #
        writes files into processed-htmldata folder
18    commands.getoutput(co_2)
19    #print url, ' ',count, ' ',file_name_1
```

Listing 1: Python program for acquiring raw and processed data files for each URI

Sample Raw File

[illegible]

4

Sample Collection of Raw files

Raw-1-52a180aec13dd535954c7d7e3eae9...	2/16/2016 12:10 AM	Text Document	238 KB
Raw-2-6d6c8abf1e289d6fb9c9a3f69296ed2...	2/16/2016 12:10 AM	Text Document	201 KB
Raw-3-a2685733cd80a265b86886a81a059...	2/16/2016 12:10 AM	Text Document	0 KB
Raw-4-8f6897f748440941226b87bc6aa164...	2/16/2016 12:10 AM	Text Document	248 KB
Raw-5-72b0f1032ff4cf1b0ba422481766eb...	2/16/2016 12:10 AM	Text Document	35 KB
Raw-6-f5e616307918dfcd5d2af430740ca5...	2/16/2016 12:10 AM	Text Document	10 KB
Raw-7-7e6a886b718749606e7f2bc88bd24...	2/16/2016 12:10 AM	Text Document	15 KB
Raw-8-5899ae6829210260535f7a82c3eb0ff4	2/16/2016 12:10 AM	Text Document	95 KB
Raw-9-79398810e32d1f7e706e8316e5a645...	2/16/2016 12:10 AM	Text Document	254 KB
Raw-10-900bd00f2afdad37c599ab139e3a...	2/16/2016 12:10 AM	Text Document	254 KB
Raw-11-3f240d2f9f41b15b126c936b9ec67...	2/16/2016 12:10 AM	Text Document	0 KB
Raw-12-51cfa0314853080f20fd42207a65...	2/16/2016 12:10 AM	Text Document	0 KB
Raw-13-69f93763d8dcd8b17c11d7112fd1...	2/16/2016 12:10 AM	Text Document	0 KB
Raw-14-169545fc867e5e67d58691784d89...	2/16/2016 12:10 AM	Text Document	19 KB
Raw-15-7268a09a67d27bd2c68ac8f1f7fc...	2/16/2016 12:10 AM	Text Document	60 KB
Raw-16-f509b753511eb7b5167066306c95...	2/16/2016 12:10 AM	Text Document	263 KB
Raw-17-34821c1f024e6af5227c922cb8841...	2/16/2016 12:10 AM	Text Document	368 KB
Raw-18-64da98c344a35bd3cd20601e2b85...	2/16/2016 12:10 AM	Text Document	101 KB
Raw-19-97f2743aa4c8c11547914a98f866d...	2/16/2016 12:11 AM	Text Document	111 KB
Raw-20-a6804d0ce21eeac7f334e6753bf29...	2/16/2016 12:11 AM	Text Document	55 KB
Raw-21-ebf2760eb2330c6b6c9603a6d549...	2/16/2016 12:11 AM	Text Document	381 KB
Raw-22-15215b6be13523a786dab4eca90d...	2/16/2016 12:11 AM	Text Document	0 KB
Raw-23-59148e19c1363a38618ede6aab94...	2/16/2016 12:11 AM	Text Document	44 KB
Raw-24-d9cda9477d099afbeea291b6181...	2/16/2016 12:11 AM	Text Document	0 KB
Raw-25-4f219f7dfd101ab6e6e20abd62d...	2/16/2016 12:11 AM	Text Document	93 KB
Raw-26-0c90112179310bf076f950b7eebe3...	2/16/2016 12:11 AM	Text Document	126 KB
Raw-27-8ab5f97b1271fa78a93cfdbbcd6d...	2/16/2016 12:11 AM	Text Document	259 KB
Raw-28-f3180ac97c123cfd9e785b8298620...	2/16/2016 12:11 AM	Text Document	19 KB
Raw-29-6ea33d971b11fbac4e122fab8f0b2...	2/16/2016 12:11 AM	Text Document	0 KB
Raw-30-5dfa1f6d5b287561799172510960...	2/16/2016 12:11 AM	Text Document	111 KB
Raw-31-37ce9c3539dcde981d44094d7d68...	2/16/2016 12:11 AM	Text Document	19 KB
Raw-32-c0a251c1a8ea41fb5fca7923cdb0...	2/16/2016 12:11 AM	Text Document	448 KB
Raw-33-dd261ce7c18915b5a48500d1b4e8...	2/16/2016 12:11 AM	Text Document	236 KB
Raw-34-dcc5f50d74ea411189babe1ad4b8...	2/16/2016 12:11 AM	Text Document	368 KB
Raw-35-3d8fcfb4dec4dd048d0a125669d6...	2/16/2016 12:11 AM	Text Document	0 KB
Raw-36-0e844c6b2bf7f7017d893fe04af4e...	2/16/2016 12:11 AM	Text Document	126 KB
Raw-37-d6600635b246432043c3c626966f...	2/16/2016 12:12 AM	Text Document	203 KB
Raw-38-277b50719331f6efe8d3b11feb2b...	2/16/2016 12:12 AM	Text Document	19 KB
Raw-39-9a97a96dac467e4a6922174bf281...	2/16/2016 12:12 AM	Text Document	70 KB
Raw-40-53324de4c7c4cea7894b0b57b6ad...	2/16/2016 12:12 AM	Text Document	62 KB
Raw-41-b67fb89ab47633cbe26cfc2fd3ef4...	2/16/2016 12:12 AM	Text Document	7 KB
Raw-42-41ae142f1ac0b2636727c4e232a03...	2/16/2016 12:12 AM	Text Document	0 KB
Raw-43-5063c2389b0a5379d83cf0a6e1eb...	2/16/2016 12:12 AM	Text Document	2 KB
Raw-44-ab319d542162787f56cc18ea1d93...	2/16/2016 12:12 AM	Text Document	149 KB
Raw-45-8346bf0dca42f7f1e0a5b3cd8f946...	2/16/2016 12:12 AM	Text Document	10 KB
Raw-46-65fdc07483f97a0cebb1ffb234442...	2/16/2016 12:12 AM	Text Document	40 KB
Raw-47-267926a142b5563ddf6c23aa44d...	2/16/2016 12:12 AM	Text Document	29 KB
Raw-48-a82882855d0f8d307d9b36779236...	2/16/2016 12:12 AM	Text Document	194 KB
Raw-49-cda9c4849cbb8e006e5c250eeb09...	2/16/2016 12:12 AM	Text Document	24 KB

Figure 3: Sample Collection of Raw files

Sample Collection of Processed files

processed-1-52a180aec13dd535954c7d7e...	2/16/2016 12:10 AM	Text Document	32 KB
processed-2-6d6c8abf1e289d6fbc9a3f692...	2/16/2016 12:10 AM	Text Document	26 KB
processed-3-a2685733cd80a265b86886a8...	2/16/2016 12:10 AM	Text Document	0 KB
processed-4-8f68977f48440941226b87bc6...	2/16/2016 12:10 AM	Text Document	9 KB
processed-5-72b0f1032ff4cf1b0ba422481...	2/16/2016 12:10 AM	Text Document	10 KB
processed-6-f5e616307918dfcd5d2af4307...	2/16/2016 12:10 AM	Text Document	3 KB
processed-7-7e6a886b718749606e7f2bc8...	2/16/2016 12:10 AM	Text Document	10 KB
processed-8-5899ae6829210260535f7a82c...	2/16/2016 12:10 AM	Text Document	13 KB
processed-9-79398810e32d1e7e706e8316...	2/16/2016 12:10 AM	Text Document	31 KB
processed-10-900bd00f2afdad37c599ab1...	2/16/2016 12:10 AM	Text Document	34 KB
processed-11-3f240d2f9f41b15b126c936b...	2/16/2016 12:10 AM	Text Document	0 KB
processed-12-51cfa0314853080f20fd422...	2/16/2016 12:10 AM	Text Document	0 KB
processed-13-69f93763d8cddeb17c11d71...	2/16/2016 12:10 AM	Text Document	0 KB
processed-14-169545fc867e5e67d586917...	2/16/2016 12:10 AM	Text Document	3 KB
processed-15-7268a09a67d27bdd2c68ac8...	2/16/2016 12:10 AM	Text Document	26 KB
processed-16-f509b75311eb7b51670663...	2/16/2016 12:10 AM	Text Document	33 KB
processed-17-34821c1f024e6af5227c922c...	2/16/2016 12:10 AM	Text Document	187 KB
processed-18-64da98c344a35bd3cd20601...	2/16/2016 12:10 AM	Text Document	18 KB
processed-19-97f2743aa4c8c11547914a98...	2/16/2016 12:11 AM	Text Document	19 KB
processed-20-a68040ce21eeacf7334e67...	2/16/2016 12:11 AM	Text Document	10 KB
processed-21-ebf2760eb2330c6b6c9603a...	2/16/2016 12:11 AM	Text Document	187 KB
processed-22-15215b6be13523a786dab4e...	2/16/2016 12:11 AM	Text Document	0 KB
processed-23-59148e19c1363a38618ede6...	2/16/2016 12:11 AM	Text Document	11 KB
processed-24-d9cdca9477d099afbee291...	2/16/2016 12:11 AM	Text Document	0 KB
processed-25-4f219f7df1d101ab6e6e20ab...	2/16/2016 12:11 AM	Text Document	41 KB
processed-26-0c90112179310b076f950b7...	2/16/2016 12:11 AM	Text Document	9 KB
processed-27-8ab5f97b1271fa78a93cfdb...	2/16/2016 12:11 AM	Text Document	33 KB
processed-28-f3180ac97c123cfd9e785b82...	2/16/2016 12:11 AM	Text Document	3 KB
processed-29-6ea33d971b11fbac4e122fa...	2/16/2016 12:11 AM	Text Document	0 KB
processed-30-5dfa1f6d5b2875617991725...	2/16/2016 12:11 AM	Text Document	4 KB
processed-31-37ce9c3539dcde981d44094...	2/16/2016 12:11 AM	Text Document	3 KB
processed-32-c0a251c1a8ea4a1fb5fca792...	2/16/2016 12:11 AM	Text Document	4 KB
processed-33-dd261ce7c18915b5a48500d...	2/16/2016 12:11 AM	Text Document	30 KB
processed-34-dcc5f50d74ea411189babel...	2/16/2016 12:11 AM	Text Document	186 KB
processed-35-3d8fcfb4dec4dd048d0a125...	2/16/2016 12:11 AM	Text Document	0 KB
processed-36-0e844c6b2bf7f701d893fe0...	2/16/2016 12:12 AM	Text Document	20 KB
processed-37-d6600635b246432043c3c62...	2/16/2016 12:12 AM	Text Document	27 KB
processed-38-277b50719331f6efe8d3b11f...	2/16/2016 12:12 AM	Text Document	3 KB
processed-39-9a97a96dac467e4a6922174...	2/16/2016 12:12 AM	Text Document	13 KB
processed-40-53324de4c7c4cea7894b0b5...	2/16/2016 12:12 AM	Text Document	43 KB
processed-41-b67fb89ab47633cbe26cfc2f...	2/16/2016 12:12 AM	Text Document	1 KB
processed-42-41ae142f1ac0b2636727c4e...	2/16/2016 12:12 AM	Text Document	0 KB
processed-43-5063c2389b0a5379d83cf0a...	2/16/2016 12:12 AM	Text Document	1 KB
processed-44-ab319d54216278f756cc18e...	2/16/2016 12:12 AM	Text Document	22 KB
processed-45-8346bf0dca42f7f1e0a5b3cd...	2/16/2016 12:12 AM	Text Document	3 KB
processed-46-65fdc07483f97a0cebb1ffb2...	2/16/2016 12:12 AM	Text Document	12 KB
processed-47-267926a142b5563ddfec23...	2/16/2016 12:12 AM	Text Document	3 KB
processed-48-a82882855d0f8d307d9b367...	2/16/2016 12:12 AM	Text Document	27 KB

Figure 4: Sample Collection of Processed files

2 Problem 2

Choose a query term (e.g., ‘shadow’) that is not a stop word (see week 5 slides) and not HTML markup from step 1 (e.g., ‘http’) that matches at least 10 documents (hint: use ‘grep’ on the processed files). If the term is present in more than 10 documents, choose any 10 from your list. (If you do not end up with a list of 10 URIs, you’ve done something wrong).

As per the example in the week 5 slides, compute TFIDF values for the term in each of the 10 documents and create a table with the TF, IDF, and TFIDF values, as well as the corresponding URIs. The URIs will be ranked in decreasing order by TFIDF values. For example:

Table 1. 10 Hits for the term ‘shadow’, ranked by TFIDF.

TFIDF	TF	IDF	URI
0.150	0.014	10.680	http://foo.com/
0.044	0.008	5.510	http://bar.com/

You can use Google or Bing for the DF estimation. To count the number of words in the processed document (i.e., the denominator for TF), you can use ‘wc’:

```
% wc -w www.cnn.com.processed
2370 www.cnn.com.processed
```

It won’t be completely accurate, but it will be probably be consistently inaccurate across all files. You can use more accurate methods if you’d like.

Don’t forget the log base 2 for IDF, and mind your significant digits!

2.1 Solution

1. Aim of this question is to first get 10 URIs which contains a specific key word in it.
2. The keyword should not be a stop word like a,the,and etc and not html mark-up .I have chosen “burger” keyword.
3. So my first task was to write a program which takes 1000 URIs as input and gives all the list of URIs which have “burger” in it.
4. This can be done by using “grep -lr” command followed by “burger”. This command is executed directly in putty in the folder which has all the 1000 processed files.
5. This gave me more than 10 files and that result is showed in Figure 6. I have randomly taken 10 URIs from this list and this list is shown in Figure 7.
6. Now in order to calculate TFIDF I need to calculate TF and IDF separately and then multiply both of them.
7. IDF is common for all the 10 URIs because it is based on the corpus. I have taken “google” as my search engine and I have found the total number of document in the corpus from reference 6. rel=mementos,rel=first memento and rel=memento last from the time maps..

8. Now I queried for “burger” in google which gave 263 million results which is the number of documents with burger in it. This is shown in Figure 5.
9. So my IDF will be total documents in corpus/263 million. This is same for all the 10 URIs.
10. IDF will be growing too large as corpus grows so in order to limit it I used log base 2 function of the IDF.
11. Now TF can be calculated if we get the total number of words in each document and the number of occurrences of the word “burger” in each document.
12. I got the total number of words in each document using “wc -w” command followed by the file name of which I want the word count.
13. I got the number of times the word occurred in each document by using “grep -lr” command followed by my keyword and followed by the file name in which the keyword should be searched.
14. By dividing occurrences/word count,I got TF.
15. Now this TF is multiplied with IDF which gave me TFIDF value.
16. These TFIDF values for each URI are arranged in ascending order and these are shown in Figure 8.
17. The URI having high TFIDF value is first shown to the user when searched in a search engine.

2.2 Code Listing

Python Code for getting URIs with “burger” and TFIDF values for each URI

```
1 import commands
2 import re
3 import math
4 import json
5 # def find_urls (URI_num):
6     # URI_number1= URI_num
7     # count =0
8     # for each_line in file_2.readlines():
9         # count=count+1
10        # print count
11        # if count == URI_number1:
12            # print each_line
13            # count =0
14            # break
15    # return each_line
16 file_1=open ('10_urls.txt','r')
17 file_2 = open('only_urls.txt','r+') .readlines()
18 Tot_doc_Google= 4700000000.00
19 burger_in_google= 263000000.00
20 temp_IDF= Tot_doc_Google/burger_in_google
21 t_IDF= math.log(temp_IDF)/math.log(2)
22 IDF= round(t_IDF,4)
23 #print IDF
24 print "TFIDF", ' ', "Rounded_TF", ' ', "IDF", ' ', "URI"
25 for f_name in file_1.readlines(): #To calculate TF and TFIDF
26     #print f_name
27     #print count
28     count=0
29     w_doc1= 'wc -w ' + f_name.strip()
30     oc_doc1= 'grep -c ' + "burger" + f_name.strip()
31     #print oc_doc1
32     w_doc=commands.getoutput(w_doc1)
33     want_word_count=commands.getoutput(oc_doc1)
34     tot_word=w_doc.split(' ')[0]
35     #print tot_word, ' ', want_word_count
36     TF= float(want_word_count)/float(tot_word)
37     Rounded_TF= round(TF,4)
38     TFIDF= round(Rounded_TF*IDF,4)
39     URI_find=f_name.split('-')
40     URI_num=int (URI_find[1])
41     i=0
42     for each in file_2:
43         count =count +1
44         if URI_num==count:
45             URI=each
46             print TFIDF, ' ', Rounded_TF, ' ', IDF, ' ', URI
47             i=i+1
48             break
49
50     # for i in a:
51         # print i
52         # print count
53         # if count == i:
54             # print each
55             # break
56
57     #print URI_num
58     #URI_1=find_urls (URI_num)
59     #print URI_1
60 # print TF_array
61 # print TFIDF_array
62 # print URI_array
```

Listing 2: Python Code for getting URIs with “burger” and TFIDF values for each URI

Number of searches for “burger” in google

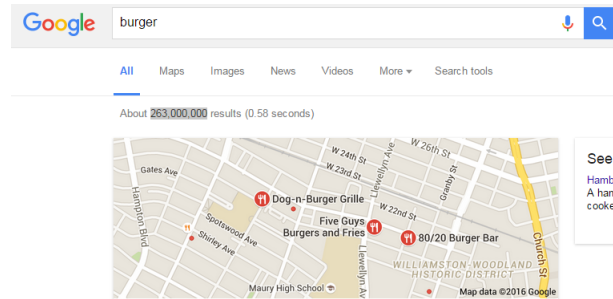


Figure 5: Number of searches for “burger” in google

Processed URI files- having “burger” keyword in it

```
processed-102-d3403063f8ebba384ca21f8614d0dbd80.txt
processed-156-219f0bd3f3c58c58ce84d21cf7ae07827.txt
processed-183-e3ba61a9e4c84074604480a354083001.txt
processed-209-012c0bd346f601eba30854677b1d1d96.txt
processed-284-a9d69418362909039fa5828bc69ee913.txt
processed-340-eac2bbab72a0b86a213cfa2d668723d3.txt
processed-342-87f3c109e3deb8e83a25334d756bd732.txt
processed-393-9fc057fca4226b8b203e5155168b9686.txt
processed-479-8cd74ea8fb1c84fae98d3693f4cd4ace.txt
processed-609-e27c5af8226e59287f861b190b28ae72.txt
processed-759-e41c8c33636eb40a1fc56246cd3312d8.txt
processed-845-5224c720d259569b949f89f9bfbdb5a45.txt
processed-846-8644cb3e6d47c5d068fce2deccc62e2fa.txt
processed-849-7e29c780d4394ca02c7842d5d7035d3.txt
processed-850-a3b35a81d7f9bdb9b62ea8189f429d18.txt
processed-851-b4104274a1174ac62b4d7f8f5d4645494.txt
processed-864-eae02d599cc73271ce8d779efa91d6ab.txt
processed-866-a36fb014fe877089b7292d043e1fd0a5.txt
processed-872-e12c98abc5643fb15caab8b814328b3db.txt
processed-886-c97c0b9c014677c50d8afbf22c0e4e27b.txt
processed-897-7cb4fb9237319e0aa0430d09062cd9da.txt
processed-899-1857a08321cd0eb8b2f83663cbd0ee51.txt
processed-901-af7cdd6003e50b720abec11e7dad24e7.txt
processed-906-3ac0e5a1248a543bfc9d91d1261e38aa.txt
processed-907-9d53e0d3174ce73c5cb8ef49cdcf363f.txt
processed-908-aff77ca85f2175022dae23a793b35f0d4.txt
processed-920-ce5ef25ce196643c430577cfbde84d35.txt
processed-942-0ffe7f8aa13096be354cdc54f1c45f9b.txt
processed-951-2ba63a1375102c33a93b3a2447729728.txt
processed-954-2e7a7645dc2925ef1268c75e6c4101e12.txt
processed-962-f88162fbc92acf4463106b0cfcf082d6.txt
processed-964-296109e7fac3e1cbf6e19e632c2e443b.txt
processed-968-cc5e1afede0efc44bc400958fc91e090.txt
processed-969-9776da6853b90cd39cd5d5a326f32972.txt
processed-986-2800e2119587df434a92eaab941e63e9.txt
```

Figure 6: Processed URI files- having “burger” keyword in it

10 random URIs from Processed URI files- having “burger” keyword in it

```

processed-897-7cb4fb9237319e0aa4030d09062cd9da.txt
processed-899-1857a08321cd0eb8b2f83663cbd0ee51.txt
processed-901-af7cdd6003e50b720abec11e7dad24e7.txt
processed-906-3ac0e5a1248a543bfcd9d1d1261e38aa.txt
processed-907-9d53e0d3174ce73c5cb8ef49cdcf363f.txt
processed-908-aff7ca85f2175022dae23a793b35f0d4.txt
processed-920-ce5ef25ce196643c430577cfbde84d35.txt
processed-942-0ffe7f8aa13096be354cdc54f1c45f9b.txt
processed-951-2ba63a1375102c33a93b3a2447729728.txt
processed-102-d3403063f8ebba384ca21f8614d0bd80.txt

```

Figure 7: 10 random URIs from Processed URI files- having “burger” keyword in it

TFIDF Table

TFIDF	Rounded_TF	IDF	URI
0.0045	0.0006	7.4815	http://www.theguardian.com/us-news/2016/feb/08/mcdonalds-chef-neal-fraser-food-influencers
0.0067	0.0009	7.4815	http://www.eater.com/2016/1/28/10861108/mcdonalds-all-day-breakfast-mcgriddle-testing
0.0075	0.001	7.4815	http://www.rachaelrayshow.com/celebs/22190_rach_surprises_a_new_york_city_subway_hero/
0.0082	0.0011	7.4815	http://www.handelsblatt.com/unternehmen/handel-konsumgueter/schluss-mit-bio-mcdonalds-schmeisst-mcb-aus-dem-sortiment/12936600.html
0.0082	0.0011	7.4815	http://time.com/4206461/mcdonalds-happy-meals-books/?xid=tcoshare
0.015	0.002	7.4815	http://www.welt.de/wirtschaft/article151987009/McDonalds-scheitert-mit-seinem-Bio-Burger.html
0.0352	0.0047	7.4815	http://www.sueddeutsche.de/wirtschaft/fast-food-industrie-mcdonalds-verabschiedet-sich-schon-wieder-vom-bio-burger-1.2855377
0.0449	0.0006	7.4815	http://www.absatzwirtschaft.de/nach-durchwachsenem-markttest-in-deutschland-mcdonalds-beerdigt-den-bio-burger-74903/
0.1167	0.0156	7.4815	http://www.seriousseats.com/recipes/2015/09/the-mcwhopper-burger-king-mcdonalds.html
0.3606	0.0482	7.4815	https://twitter.com/burgerking/status/695760664028119041

Figure 8: TFIDF Table

3 Problem 3

Now rank the same 10 URIs from question #2, but this time by their PageRank. Use any of the free PR estimators on the web, such as:

http://www.prchecker.info/check_page_rank.php
<http://www.seocentro.com/tools/search-engines/pagerank.html>
<http://www.checkpagerank.net/>

If you use these tools, you'll have to do so by hand (they have anti-bot captchas), but there is only 10. Normalize the values they give you to be from 0 to 1.0. Use the same tool on all 10 (again, consistency is more important than accuracy).

Create a table similar to Table 1:

Table 2. 10 hits for the term "shadow", ranked by PageRank.

PageRank	URI
0.9	http://bar.com/
0.5	http://foo.com/

Briefly compare and contrast the rankings produced in questions 2 and 3.

3.1 Solution

1. For this question I need to get the Page rank for each of the 10 URIs.
2. This can be found out by using one of the 3 links provided in the question.
3. I used 2nd link as it did not ask captcha to enter each time I check for a rank. This just reduced my work a little bit.
4. I manually copied each URI and pasted it in the PR estimator website and it gave me the rank.
5. First when I tried to get the rank from my URIs all my 10 URIs gave n/a output.
6. Later I read the comment in the google groups and found that it can use top-level URIs rather than my entire URI.
7. This is because deep links into a site does not have PR.
8. So I took top level of each URI and got the Page rank for each one of them. This can be seen in Figure 9.
9. Now each Page rank is normalized as asked in the question.
10. The list of Page ranks arranged in ascending order and the URIs are shown in Figure 10.
11. In the question I was asked to compare Page rank and TFIDF, so both of them are shown in a single table which can be seen in Figure 11.
12. When we compare the rankings based on TFIDF and page rank there is not relation between them. Page rank considers back link structures, so it gives the rank of the document with respect to other documents in the web, but TFIDF gives rank to the documents based on the level of similarity to the term. So If a document has high TFIDF value then it ranks it high even though if it is unpopular.

3.2 Results

Sample Rank for one URI

Pagerank Results		
Uri	Pagerank	
http://www.rachaelrayshow.com	6 / 10	<div><div></div></div>

Figure 9: Sample Rank for one URI

Page rank along with 10 URIs

Page Rank	URI
0	http://www.seriousseats.com/recipes/2015/09
0	http://time.com
0	http://www.welt.de
0.6	http://www.sueddeutsche.de/wirtschaft
0.6	http://www.handelsblatt.com/unternehmen
0.6	http://twitter.com/burgerking
0.6	http://www.absatzwirtschaft.de
0.6	http://www.eater.com
0.6	http://www.rachaelrayshow.com
0.7	http://www.theguardian.com

Figure 10: Page rank along with 10 URIs

Comparison between Pagerank and TFIDF

TFIDF	Page Rank	URI
0.1167	0	http://www.seriousseats.com/recipes/2015/09
0.0082	0	http://time.com
0.015	0	http://www.welt.de
0.0352	0.6	http://www.sueddeutsche.de/wirtschaft
0.0082	0.6	http://www.handelsblatt.com/unternehmen
0.3606	0.6	http://twitter.com/burgerking
0.0449	0.6	http://www.absatzwirtschaft.de
0.0067	0.6	http://www.eater.com
0.0075	0.6	http://www.rachaelrayshow.com
0.0045	0.7	http://www.theguardian.com

Figure 11: Comparison between Pagerank and TFIDF

4 Problem 4

4. Compute the Kendall Tau_b score for both lists (use ‘b’ because there will likely be tie values in the rankings). Report both the Tau value and the ‘p’ value.

See:

<http://stackoverflow.com/questions/2557863/measures-of-association-in-r-kendalls-tau-b-and-tau-c>

http://en.wikipedia.org/wiki/Kendall_tau_rank_correlation_coefficient#Tau-b

http://en.wikipedia.org/wiki/Correlation_and_dependence

4.1 Solution

1. Here I need to calculate the kendall Tau b value and p-value.
2. I used Tau b because there are tie values in the rankings.
3. Tau b score is used to measure the association between the ranks produced by TFIDF and Page rank.
4. I have both the TFIDF and page rank values into two vectors and “cor.test” command is used in “R” giving method as “Kendall” which returns Tau b value and p-value.
5. Output :- tau = 0.7833, z= 2.8088, p-value = 0.0024
6. Here 3 URIs have 0 rank and 6 of them have 0.6 rank. So, there are a lot of ties and Kendall’s tau coefficient shows a high correlation.

4.2 Code Listing and Results

R code and results for Kendal's Tau b value and p value

```
1
2 R version 3.1.2 (2014-10-31) — "Pumpkin Helmet"
3 Copyright (C) 2014 The R Foundation for Statistical Computing
4 Platform: i386-w64-mingw32/i386 (32-bit)
5
6 R is free software and comes with ABSOLUTELY NO WARRANTY.
7 You are welcome to redistribute it under certain conditions.
8 Type 'license()' or 'licence()' for distribution details.
9
10 Natural language support but running in an English locale
11
12 R is a collaborative project with many contributors.
13 Type 'contributors()' for more information and
14 'citation()' on how to cite R or R packages in publications.
15
16 Type 'demo()' for some demos, 'help()' for on-line help, or
17 'help.start()' for an HTML browser interface to help.
18 Type 'q()' to quit R.
19
20 > Pagerank <- c(0,0,0,0.6,0.6,0.6,0.6,0.6,0.6,0.7)
21 >
22 > TFIDF <- c(0.0045,0.0067,0.0075,0.0082,0.0082,0.015,0.0352,0.0449,0.1167,0.3606)
23 >
24 > cor.test(TFIDF,Pagerank, method = "kendall", alternative = "greater")
25
26 Kendall's rank correlation tau
27
28 data: TFIDF and Pagerank
29 z = 2.8088, p-value = 0.002486
30 alternative hypothesis: true tau is greater than 0
31 sample estimates:
32 tau
33 0.7833495
34
35 Warning message:
36 In cor.test.default(TFIDF, Pagerank, method = "kendall", alternative = "greater") :
37 Cannot compute exact p-value with ties
38 >
```

Listing 3: R code for Kendal's Tau b value

Bibliography

- [1] Hashlib Secure Hashes and message digests <https://docs.python.org/2/library/hashlib.html> , Python Software Foundation, 2016
- [2] Measures of association in R Kendall's tau-b and tau-c <http://stackoverflow.com/questions/2557863/measures-of-association-in-r-kendalls-tau-b-and-tau-c>, Developers from Stack Overflow, 2008
- [3] Searching For A String In Multiple Files <http://blamcast.net/articles/recursive-grep-find-strings-files>, John, 2008
- [4] getting syntax error near unexpected token ';' in python <http://stackoverflow.com/questions/10900617/getting-syntax-error-near-unexpected-token-in-python>, Developers from Stack Overflow, 2008
- [5] math Mathematical functions <https://docs.python.org/2/library/math.html>, Python Software Foundation, 2016
- [6] The size of the World Wide Web (The Internet) <http://www.worldwidewebsize.com/>, Maurice de Kunder, 2006
- [7] Google PageRank Checker <http://www.seocentro.com/tools/search-engines/pagerank.html>, Developers at SEOCentro, 2016