Srividya Majeti

# Assignment 3

CS 532: Introduction to Web Science
Dr. Michael Nelson
Spring 2016

February 19, 2016

# Contents

# 1

## Question 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:**

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

- First I got the raw data for all the 1000 unique URIs that I collected in assignment 2, using the following cURL command:
  curl <URI> > < output filename >.
- I stored the raw HTML output generated by the cURL command in separate files for each URI and named the files based on their index. This code is listed in Listing 1.1
- Then I got the processed HTML and stored the output in separate files for each URI using the command:
  lynx -dump -force_html <URI> > < output filename >.
  I named the files with the URI index followed by a hyphen and the word 'processed'. This code is listed in Listing 1.2

### Code Listing

```
1   import commands
2   import os
3
4   count = 0
5   input=open('uri.json','r')
6   t=open('command','w')
7   for line in input:
8           count +=1
9           if count <1001:
10                  filename= str(count)
11                  # print str(line)
12                  command ='curl ' + '"' + str(line).rstrip('\
                        n') + '"'+ '> ./rawData/' + filename
13                  # print command
14                  # t.write(str(command))
15                  output = commands.getoutput(command)
16  input.close()
```

**Listing 1.1.** Python code for getting raw HTML and saving them in files

### Code Listing

```
1   import commands
2   import os
3
4   count = 0
5   input=open('uri.json','r')
6   t=open('command','w')
7   for line in input:
8           count +=1
9           if count <1001:
10                  filename= str(count) + '-processed'
11                  # print str(line)
12                  command ='lynx -dump -force_html ' + '"'+
                        str(line).rstrip('\n') + '"'+'> ./
                        processedData/' + filename
13                  # t.write(str(command))
14                  print command
15                  output = commands.getoutput(command)
16  input.close()
```

**Listing 1.2.** Python code for getting processed HTML without any images or stylesheets and saving them in files

# Question 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.

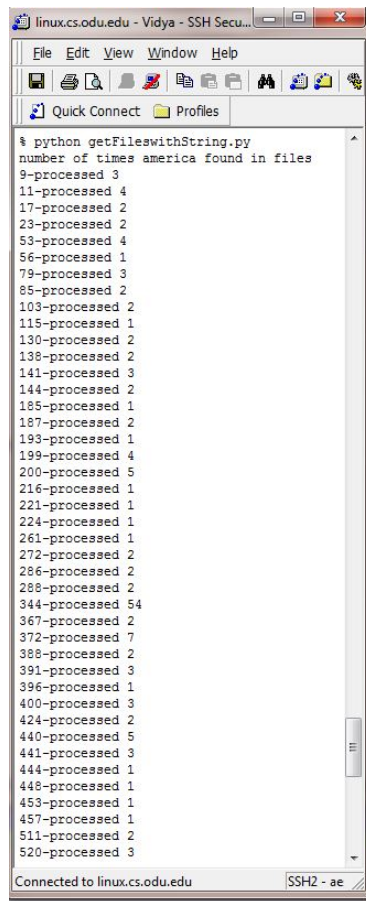**Don't forget the log base 2 for IDF, and mind your significant digits!**

Following are the steps that I have taken to solve this problem:

- I made use of the processed data generated in question 1 and selected a query term.

- Using the following cURL command, I searched for the query term in processed files to find the frequency of the term in each URI.
  cat ./processedData/ <filename> — grep -i -c <queryTerm>
- First I looked for the word 'America' in all the processed files. I got more than 10 URIs, But the number of times the query term appeared in each URI was not satisfactory.The screenshot of the output with filename and number of times the word 'America' appeared in each file is in Figure 2.1.
- Then I changed my query term to 'food' which gave a better frequency. The screenshot for this output is in Figure 2.2
- I also counted the total number of words in each processed file using the command:
  cat ./processedData/ <filename> — wc -w.
- Then I stored the output in a JSON structure with file name, total number of words in the document and frequency of the query term in the document. This code is listed in Listing 2.1 and screenshot of this JSON structure is in Figure 2.3.
- When I searched for the term 'food' I received about 338 URIs but I selected 10 URIs among them based on the versatility of the frequency of the term in each document.
- The selected 10 URIs are as follows:
  ```
  http://www.chilipeppermadness.com/
  http://tcbmag.com/Innovations/January-2016/Why-Chipotle-Is-Closing-
  Its-Stores-Over-The-Lunch?utm_content=27945314&utm_medium=social&
  utm_source=twitter
  http://www.CafeRio.com
  http://www.yelp.com/biz/chipotle-mexican-grill-gainesville-4?
  utm_campaign=CheckIn&utm_medium=twitter&utm_source=ashare
  http://www.nydailynews.com/life-style/eats/chipotle-lettuce-shortage-
  plagues-new-york-article-1.2520586
  http://fieryfork.com/video-chipotle-profit-heats-up-as-it-draws-
  more-diners/
  http://www.eater.com/2016/2/5/10922434/chipotle-e-coli-beef-australia?
  utm_campaign=national&utm_medium=social&utm_source=twitter
  http://www.havingfunsaving.com/2016/01/chipotle-cheese-dip-recipe.
  html
  http://www.ooyuz.com/geturl?aid=10248937
  http://smartmarkradio.com/conspiracy-theory-news/this-chipotle-
  conspiracy-theory-is-the-craziest-thing-weve-read-all-day-grist/
  ```
- To calculate IDF for these URIs, I got the total number of documents in corpus which was 47 billion and documents with the term 'food' which was about 3.52 billion.
- Then I calculated the values of TF, IDF, TFIDF. The code is listed in Listing 2.2
- These values are summarized in the Table 2.1.

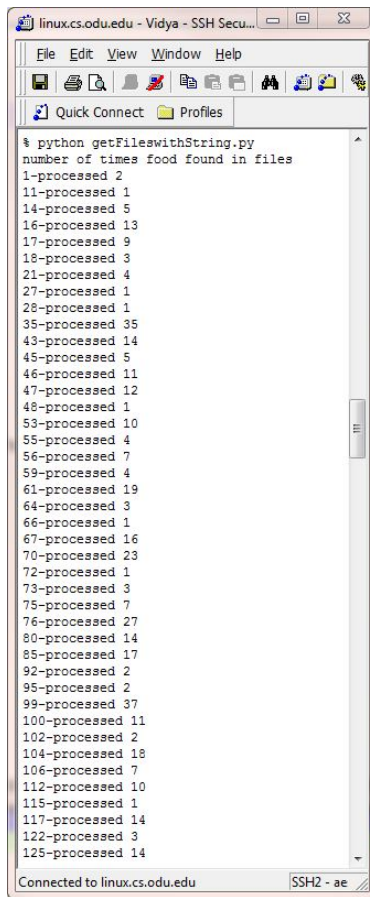**Table 2.1.** Table with calculated TFIDF , TF , IDF and URI

| Rank | TFIDF | TF | IDF | URI |
|---|---|---|---|---|
| 1 | 0.075 | 0.02 | 3.739 | http://fieryfork.com/video-chipotle-profit-heats-up-as-it-draws-more-diners/ |
| 2 | 0.045 | 0.012 | 3.739 | http://www.chilipeppermadness.com/ |
| 3 | 0.026 | 0.007 | 3.739 | http://www.nydailynews.com/life-style/eats/chipotle-lettuce-shortage-plagues-new-york-article-1.2520586 |
| 4 | 0.022 | 0.006 | 3.739 | http://tcbmag.com/Innovations/January-2016/Why-Chipotle-Is-Closing-Its-Stores-Over-The-Lunch?... |
| 5 | 0.019 | 0.005 | 3.739 | http://www.CafeRio.com |
| 6 | 0.019 | 0.005 | 3.739 | http://www.eater.com/2016/2/5/10922434/chipotle-e-coli-beef-australia?utm_... |
| 7 | 0.019 | 0.005 | 3.739 | http://www.ooyuz.com/geturl?aid=10248937 |
| 8 | 0.011 | 0.003 | 3.739 | http://www.yelp.com/biz/chipotle-mexican-grill-gainesville-4?utm_campaign=... |
| 9 | 0.011 | 0.003 | 3.739 | http://smartmarkradio.com/conspiracy-theory-news/this-chipotle-conspiracy-theory-is-the-craziest... |
| 10 | 0.007 | 0.002 | 3.739 | http://www.havingfunsaving.com/2016/01/chipotle-cheese-dip-recipe.html |

**Fig. 2.1.** Number of times "America" appeared in processed files

**Fig. 2.2.** Number of times "food" appeared in processed files

**Fig. 2.3.** JSON structure with file name, total number of words in processed file and number of times the query term occurred in the file

**Code Listing**

```python
1   import os
2   import commands
3   import json
4
5   f = open('filenamesWithCount','w')
6   queryTerm = "food"
7   def getFilesWithString(queryTerm):
8           filenamesCountJson ={}
9           iterateThroughFiles = next(os.walk('./processedData/
                '))[2]
10          print 'number of times %s found in files '%queryTerm
11          for filename in iterateThroughFiles:
12                  command = 'cat ' + './processedData/'+
                        filename + ' | grep -i -c ' + queryTerm
13                  numberOfOccurrencesInDoc = commands.
                        getoutput(command)
14                  if numberOfOccurrencesInDoc > '0':
15                          filenamesCountJson['
                                numberOfOccurrencesOfStringInDoc
                                '] =numberOfOccurrencesInDoc
16                          filenamesCountJson['filename'] =
                                filename
17                          command = 'cat ' + './processedData/
                                '+ filename + ' | wc -w  '
18                          totalWords = commands.getoutput(
                                command)
19                          filenamesCountJson['totalWordsInDoc'
                                ] =totalWords
20                          f.write(json.dumps(
                                filenamesCountJson) + "\n")
21  getFilesWithString(queryTerm)
```

**Listing 2.1.** "Python code for searching a query term and counting the number of times the term appears in each processed file. Also counted the total number of words in each processed file and have written the output in a JSON structure."

### Code Listing

```
1   import json
2   import math
3
4   totalNumberOfDocumentInCorpus  =47000000000
5   documentWithTerm = 3520000000
6
7   def calculate():
8           input  = open('selected10URIsWithSrtingAndWordCount.
                   json','r')
9           for line in input:
10                  data= json.loads(line)
11                  TF = round((float(data['
                       numberOfOccurrencesOfStringInDoc'])/
                       float(data['totalWordsInDoc'])),3)
12                  print 'TF:', TF
13                  IDF = round(logBase2( float(
                       totalNumberOfDocumentInCorpus) / float(
                       documentWithTerm) ),3)
14                  print 'IDF:',IDF
15                  TFIDF = round((TF*IDF),3)
16                  print 'TFIDF:',TFIDF
17
18  def logBase2(number):
19          return math.log(number) / math.log(2)
20
21  calculate()
```

**Listing 2.2.** "Python code for calculating TF, IDF and TFIDF for 10 selected URIs"

# 3

# Question 3

**Now rank the same 10 URIs from question 2, but this time by their PageRank. Use any of the free PR estimaters 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.**

- To find the page rank for all the 10 URIs I used a free PR estimator on the web which is located at:
  `http://www.seocentro.com/tools/search-engines/pagerank.html`.
- The URI and its page rank is summarized in the Table 3.1.
- While finding the page rank for each URI I took a screenshot. The screenshits are in the Figures 3.1, 3.2, 3.3, 3.4, 3.5, 3.6, 3.7, 3.8 3.9, 3.10

Comparing the rankings in question 2 and 3 I observed that their is a lot of difference in rankings. The reason behind this difference is, when we rank the URIs based on TF, IDF and TFIDF, we look for the frequency of the term. But the PR estimator ranks the URIs based on the quality of the document and also the quantity of backlinks. Calculating pagerank based on frequency of the term is not a good idea because even a poor quality document with high query term frequency will be ranked high.

**Table 3.1.** PageRank for 10 selected URIs using `http://www.seocentro.com/tools/search-engines/pagerank.html`

| Rank | PageRank | URI |
|------|----------|-----|
| 3 | 0.7 | http://www.nydailynews.com/ |
| 8 | 0.7 | http://www.yelp.com/ |
| 6 | 0.6 | http://www.eater.com/ |
| 4 | 0.5 | http://tcbmag.com/ |
| 5 | 0.5 | http://www.CafeRio.com |
| 2 | 0.3 | http://www.chilipeppermadness.com/ |
| 10 | 0.3 | http://www.havingfunsaving.com/ |
| 9 | 0.2 | http://smartmarkradio.com/ |
| 1 | 0 | http://fieryfork.com/ |
| 7 | 0 | http://www.ooyuz.com/ |

| Pagerank Results | |
|---|---|
| **Url** | **Pagerank** |
| http://www.chilipeppermadness.com/ | 3 / 10 |

| Pagerank History | |
|---|---|
| **Month** | **Pagerank** |
| 2016-02 | 3 / 10 |

**Fig. 3.1.** page rank for `http://www.chilipeppermadness.com/`

| Pagerank Results | |
|---|---|
| **Url** | **Pagerank** |
| http://tcbmag.com/ | 5 / 10 |

| Pagerank History | |
|---|---|
| **Month** | **Pagerank** |
| 2016-02 | 5 / 10 |

**Fig. 3.2.** page rank for `http://tcbmag.com/Innovations/January-2016/Why-Chipotle-Is-Closing-Its-Stores-Over-The-Lunch?utm_content=27945314&utm_medium=social&utm_source=twitter`

| Pagerank Results | | |
|---|---|---|
| **Url** | | **Pagerank** |
| http://www.CafeRio.com | 5 / 10 | |

| Pagerank History | | |
|---|---|---|
| **Month** | | **Pagerank** |
| 2016-02 | 5 / 10 | |

**Fig. 3.3.** page rank for `http://www.CafeRio.com`

| Pagerank Results | | |
|---|---|---|
| **Url** | | **Pagerank** |
| http://www.yelp.com/ | 7 / 10 | |

| Pagerank History | | |
|---|---|---|
| **Month** | | **Pagerank** |
| 2016-02 | 7 / 10 | |
| 2015-11 | 7 / 10 | |
| 2015-10 | 7 / 10 | |
| 2015-05 | 7 / 10 | |
| 2015-04 | 7 / 10 | |
| 2014-08 | 7 / 10 | |
| 2014-06 | 7 / 10 | |
| 2014-05 | 7 / 10 | |
| 2014-04 | 7 / 10 | |

**Fig. 3.4.** page rank for `http://www.yelp.com/biz/chipotle-mexican-grill-gainesville-4?utm_campaign=CheckIn&utm_medium=twitter&utm_source=ashare`

| Pagerank Results | | |
|---|---|---|
| **Url** | | **Pagerank** |
| http://www.nydailynews.com/ | 7 / 10 | |

| Pagerank History | | |
|---|---|---|
| **Month** | | **Pagerank** |
| 2016-02 | 7 / 10 | |
| 2015-09 | 7 / 10 | |
| 2015-06 | 7 / 10 | |

**Fig. 3.5.** page rank for `http://www.nydailynews.com/life-style/eats/chipotle-lettuce-shortage-plagues-new-york-article-1.2520586`

| Pagerank Results | | |
|---|---|---|
| **Url** | | **Pagerank** |
| http://fieryfork.com/ | 0 / 10 | |

| Pagerank History | | |
|---|---|---|
| **Month** | | **Pagerank** |
| 2016-02 | 0 / 10 | |

**Fig. 3.6.** page rank for `http://fieryfork.com/video-chipotle-profit-heats-up-as-it-draws-more-diners/`

**Pagerank Results**

| Url | Pagerank | |
|-----|----------|--|
| http://www.eater.com | 6 / 10 | ▬ |

**Pagerank History**

| Month | Pagerank | |
|-------|----------|--|
| 2016-02 | 6 / 10 | ▬ |

**Fig. 3.7.** page rank for `http://www.eater.com/2016/2/5/10922434/chipotle-e-coli-beef-australia?utm_campaign=national&utm_medium=social&utm_source=twitter`

**Pagerank Results**

| Url | Pagerank | |
|-----|----------|--|
| http://www.havingfunsaving.com/ | 3 / 10 | ▬ |

**Pagerank History**

| Month | Pagerank | |
|-------|----------|--|
| 2016-02 | 3 / 10 | ▬ |

**Fig. 3.8.** page rank for `http://www.havingfunsaving.com/2016/01/chipotle-cheese-dip-recipe.html`

**Pagerank Results**

| Url | Pagerank | |
|-----|----------|--|
| http://www.ooyuz.com/ | 0 / 10 | ▬ |

**Pagerank History**

| Month | Pagerank | |
|-------|----------|--|
| 2016-02 | 0 / 10 | ▬ |

**Fig. 3.9.** page rank for `http://www.ooyuz.com/geturl?aid=10248937`

**Pagerank Results**

| Url | Pagerank | |
|-----|----------|--|
| http://smartmarkradio.com/ | 2 / 10 | ▬ |

**Pagerank History**

| Month | Pagerank | |
|-------|----------|--|
| 2016-02 | 2 / 10 | ▬ |

**Fig. 3.10.** page rank for `http://smartmarkradio.com/conspiracy-theory-news/this-chipotle-conspiracy-theory-is-the-craziest-thing-weve-read-all-day-grist/`

# 4

## Extra-Credit Question

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

- The Kendall Tau value is '0.926' and p value is '0.000217' for both the lists.
- Below is the R code for calculating the Tau and p values.



**Fig. 4.1.** R code for calculating Tau and p values

# References

1. On what basis PR estimator ranks a URI: http://www.seomastering.com/pagerank-prediction.php, Professional SEO authors, Igor Ivanov
2. How to find author of website: http://www.wikihow.com/Find-the-Author-Of-a-Website, Jack Herrick, 2005
3. How to calculate Tau values: http://stackoverflow.com/questions/2557863/measures-of-association-in-r-kendalls-tau-b-and-tau-c, Joel Spolsky, 2008