

CS 432

## Assignment 3

Daniel Moore

## Part 1:

Part one of this assignment involved capturing the content of our 1000 unique URLs from assignment 2 and place this content in files. Once this was completed I was to remove the HTML markup from each file and place the processed content in separate files. I accomplished these tasks with the following script:

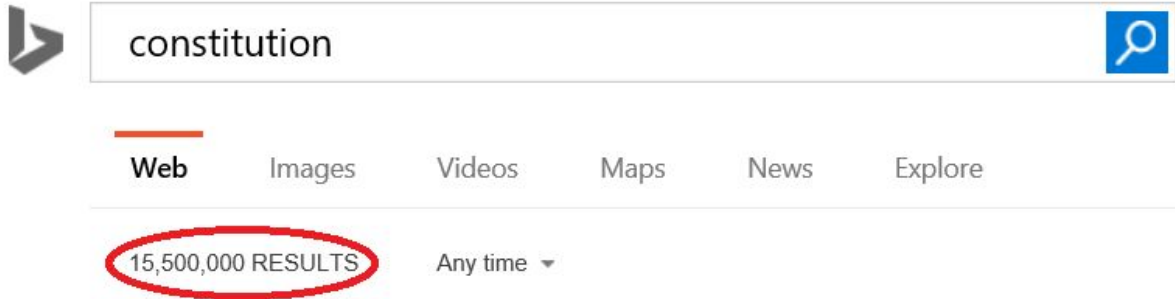
```
1  #!/bin/bash
2
3  # test-echo-urls.sh
4
5  while read -r full_url
6  do
7      echo $full_url
8      wget -O temp.txt "${full_url}"
9      temp=`echo -n $full_url | md5sum | awk '{ print $1 }'`
10     cp temp.txt ./outFiles/"${temp}.raw"
11     echo "${full_url}" >> ./outFiles/"${temp}.raw"
12     echo "${full_url}" > ./outFiles/"${temp}.processed"
13     lynx -dump -force_html temp.txt >> ./outFiles/"${temp}.processed"
14 done < links.txt
```

The first portion of phase two involved finding a word which was contained within at least ten of the processed files(I chose the word "constitution")and selecting ten of those files for later processing. I did this with the following command script :

```
1  #!/bin/bash
2
3  for file in $(grep -lir "constitution" ./outFiles/*.processed | tail -10)
4  do
5      cp "${file}" ./selection
6  done
```

I then used the following script to compute the TF values and print them to a file along with their URI:

```
1  #!/bin/bash
2
3  for file in ./selection/*.processed
4  do
5      oc=$(grep -o 'constitution' "${file}" | wc -l)
6      count=$(wc -w "${file}" | awk '{print $1}')
7      tf=$(echo "scale=4; $oc/$count" | bc)
8      echo "${tf}" >> tf.txt
9      head -1 "${file}" >> tf.txt
10 done
```



Using Bing, a search of the word “constitution” returned 15,500,000 results. Using an estimated index size for Bing of 20,000,000,000 this gives an IDF value of 10.3. For each entry URI I calculated the TFIDF value by multiplying the TF and IDF values. I did this with the following script:

```
1  from math import *
2
3  f = open("tf.txt", "r")
4  out = open("out.csv", "a")
5  i = 0
6  tf = []
7  URI = []
8  for i in range(1, 21):
9      if i % 2 == 0:
10         URI.append(f.readline())
11     else:
12         tf.append(f.readline())
13
14     out.write("TFIDF,TF,IDF,URI\n")
15
16 for x in range(0, 10):
17
18     TFIDF = float(tf[x]) * 10.3
19     out.write(str(TFIDF) + "," + str(tf[x]) + "," + str("10.3") + "," + str(URI[x]) + "\n")
20
```

I then placed these values in the following table:

TFIDF	TF	IDF	URI
0.03811	0.0037	10.3	<a href="https://www.nraila.org/issues/terrorist-watchlistno-f...">https://www.nraila.org/issues/terrorist-watchlistno-f...</a>
0.02884	0.0028	10.3	<a href="http://www.breitbart.com/big-journalism/2015/12/0...">http://www.breitbart.com/big-journalism/2015/12/0...</a>
0.02266	0.0022	10.3	<a href="http://www.tucsonsentinel.com/opinion/report/12...">http://www.tucsonsentinel.com/opinion/report/12...</a>
0.01751	0.0017	10.3	<a href="http://thefederalist.com/2015/11/11/chicagos-not-...">http://thefederalist.com/2015/11/11/chicagos-not-...</a>
0.00824	0.0008	10.3	<a href="http://opinion.injo.com/2016/01/251742-dont-gun-i...">http://opinion.injo.com/2016/01/251742-dont-gun-i...</a>
0.00721	0.0007	10.3	<a href="https://www.nraila.org/articles/20151222/nra-cond...">https://www.nraila.org/articles/20151222/nra-cond...</a>
0.00721	0.0007	10.3	<a href="http://www.news-journalonline.com/article/201511...">http://www.news-journalonline.com/article/201511...</a>
0.00721	0.0007	10.3	<a href="http://www.washingtontimes.com/news/2016/jan/...">http://www.washingtontimes.com/news/2016/jan/...</a>
0.00618	0.0006	10.3	<a href="https://www.nraila.org/articles/20151027/wayne-la...">https://www.nraila.org/articles/20151027/wayne-la...</a>
0.00309	0.0003	10.3	<a href="http://www.knoxnews.com/news/local/right-name-...">http://www.knoxnews.com/news/local/right-name-...</a>

I utilized the following source to calculate page rank:

<http://www.seocentro.com/tools/search-engines/pagerank.html>

All of the results for full URIs were “Undefined” I had to shorten the URIs in order to achieve a result. The URI portions for which I was able to obtain a score are in blue text.

URI	Page Rank
<a href="https://www.nraila.org/articles/20151222/nra-condemns-virginia-attorney-general-s-decision-to-play-politics-with-self-defense-rights">https://www.nraila.org/articles/20151222/nra-condemns-virginia-attorney-general-s-decision-to-play-politics-with-self-defense-rights</a>	0.5
<a href="http://thefederalist.com/2015/11/11/chicagos-not-so-brilliant-bullet-tax/">http://thefederalist.com/2015/11/11/chicagos-not-so-brilliant-bullet-tax/</a>	0.0
<a href="https://www.nraila.org/issues/terrorist-watchlistno-fly-list/everytowns-claims-and-the-facts/">https://www.nraila.org/issues/terrorist-watchlistno-fly-list/everytowns-claims-and-the-facts/</a>	0.5
<a href="http://www.tucsonsentinel.com/opinion/report/121215_watch_list_op/watch-lists-should-never-used-deny-constitutional-rights/">http://www.tucsonsentinel.com/opinion/report/121215_watch_list_op/watch-lists-should-never-used-deny-constitutional-rights/</a>	0.4
<a href="http://opinion.injo.com/2016/01/251742-dont-gun-ive-still-got-3-questions-obama-gun-control-orders/">http://opinion.injo.com/2016/01/251742-dont-gun-ive-still-got-3-questions-obama-gun-control-orders/</a>	Undefined

<a href="http://www.knoxnews.com/news/local/right-name-wrong-man-knoxville-veterinarian-cant-get-off-no-fly-list-26f4d248-2e52-111a-e053-0100007-362692771.html">http://www.knoxnews.com/news/local/right-name-wrong-man-knoxville-veterinarian-cant-get-off-no-fly-list-26f4d248-2e52-111a-e053-0100007-362692771.html</a>	0.6
<a href="https://www.nraila.org/articles/20151027/wayne-lapierre-to-president-obama-this-is-how-you-stop-violent-crime">https://www.nraila.org/articles/20151027/wayne-lapierre-to-president-obama-this-is-how-you-stop-violent-crime</a>	0.5
<a href="http://www.news-journalonline.com/article/20151130/opinion/151139997?tc=ar">http://www.news-journalonline.com/article/20151130/opinion/151139997?tc=ar</a>	0.6
<a href="http://www.washingtontimes.com/news/2016/jan/5/gun-rules-wont-make-easier-charges-ex-atf-agents/">http://www.washingtontimes.com/news/2016/jan/5/gun-rules-wont-make-easier-charges-ex-atf-agents/</a>	0.7
<a href="http://www.breitbart.com/big-journalism/2015/12/07/la-times-counters-obama-background-checks-not-include-no-fly-list/">http://www.breitbart.com/big-journalism/2015/12/07/la-times-counters-obama-background-checks-not-include-no-fly-list/</a>	0.6

The highest page rank was obtained by [www.washingtontimes.com](http://www.washingtontimes.com) while this URI had a TFIDF value of only .007 the highest TFIDF value was associated with: <https://www.nraila.org/issues/terrorist-watchlistno-fly-list/everytowns-claims-and-the-facts/> whereas [www.nra.org](http://www.nra.org) possessed the median page rank of 5/10.

URI	Page Rank	
<a href="https://www.nraila.org/articles/20151222/nra-condemns-virginia-attorney-general-s-decision-to-play-politics-with-self-defense-rights">https://www.nraila.org/articles/20151222/nra-condemns-virginia-attorney-general-s-decision-to-play-politics-with-self-defense-rights</a>	0.5	0.0072
<a href="http://thefederalist.com/2015/11/11/chicagos-not-so-brilliant-bullet-tax/">http://thefederalist.com/2015/11/11/chicagos-not-so-brilliant-bullet-tax/</a>	0.0	0.018
<a href="https://www.nraila.org/issues/terrorist-watchlistno-fly-list/everytowns-claims-and-the-facts/">https://www.nraila.org/issues/terrorist-watchlistno-fly-list/everytowns-claims-and-the-facts/</a>	0.5	0.038

<a href="http://www.tucsonsentinel.com/opinion/report/121215_watch_list_op/watch-lists-should-never-used-deny-constitutional-rights/">http://www.tucsonsentinel.com/opinion/report/121215_watch_list_op/watch-lists-should-never-used-deny-constitutional-rights/</a>	0.4	0.023
<a href="http://opinion.injo.com/2016/01/251742-dont-gun-ive-still-got-3-questions-obama-gun-control-orders/">http://opinion.injo.com/2016/01/251742-dont-gun-ive-still-got-3-questions-obama-gun-control-orders/</a>	Undefined	0.0082
<a href="http://www.knoxnews.com/news/local/right-name-wrong-man-knoxville-veterinarian-cant-get-off-no-fly-list-26f4d248-2e52-111a-e053-0100007-362692771.html">http://www.knoxnews.com/news/local/right-name-wrong-man-knoxville-veterinarian-cant-get-off-no-fly-list-26f4d248-2e52-111a-e053-0100007-362692771.html</a>	0.6	.0031
<a href="https://www.nraila.org/articles/20151027/wayne-lapierre-to-president-obama-this-is-how-you-stop-violent-crime">https://www.nraila.org/articles/20151027/wayne-lapierre-to-president-obama-this-is-how-you-stop-violent-crime</a>	0.5	0.0062
<a href="http://www.news-journalonline.com/article/20151130/opinion/151139997?tc=ar">http://www.news-journalonline.com/article/20151130/opinion/151139997?tc=ar</a>	0.6	0.0072
<a href="http://www.washingtontimes.com/news/2016/jan/5/gun-rules-wont-make-easier-charges-ex-atf-agents/">http://www.washingtontimes.com/news/2016/jan/5/gun-rules-wont-make-easier-charges-ex-atf-agents/</a>	0.7	0.0072
<a href="http://www.breitbart.com/big-journalism/2015/12/07/lap-times-counters-obama-background-checks-not-include-no-fly-list/">http://www.breitbart.com/big-journalism/2015/12/07/lap-times-counters-obama-background-checks-not-include-no-fly-list/</a>	0.6	0.029