# Assignment Three

Jacob Berlin

CS432 – Spring 2016

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).

```
http://0.tum.news/0UB89
http://1.usa.gov/1ovZzOF
http://1.usa.gov/1SkIALB
http://1000goldschlager.de
http://1000webradios.de/streams/1000goldschlager.m3u
http://1000webradios.de/streams/eurosmoothjazz.m3u
http://1029thehog.com/road-hogs/021216-jim-norton/
http://11noticias.com
http://1ee.me/47v
http://1nayami.xsrv.jp/sp/entry292.html
http://247sports.com/Player/Bradley-Jennings-Jr-87446
http://247wallst.com/investing/2013/03/12/eight-companies-that-will-benefit-from-keystone-pipeline-approval-trp-cnq-cop-de-xom-lyb-pwr-vlo/
http://2jamtt.sakura.ne.jp/lxro/2016/02/13/post-16297/
http://2sk.co/A43cb
http://360musicng.co/soundcloud-responds-rumors-companys-imminent-doom/
http://365diary.net/eDVwMUJ4L2svZG9ucQ--
http://49thshelf.com/Blog/2016/02/11/The-Chat-Trevor-Corkum-Interviews-Damian-Rogers
http://4NN.cx/.98590
http://550909.com/?f5160450
http://550909.com/?f9897728
http://7asnat.com
http://7asnat.com/
http://7czote.com/moc/2016/02/12/post-1934/
http://9gag.com/gag/a57E7MN?ref=tp
http://9gag.com/gag/ajAmbx1?ref=mobile
http://9gag.com/gag/axjM8oY?ref=blackberry
http://9jastreet.com
http://9news.com.au/national/2016/02/12/12/50/stephen-hawking-says-discovery-of-gravitational-waves-provides-new-way-of-looking-at-the-universe
http://a.r10.to/hQvxuO
http://a502.phobos.apple.com/us/r30/Music69/v4/c4/42/54/c44254ae-bd05-e368-6999-414b7f0cc2b8/mzaf_7784649859167295490.plus.aac.p.m4a
http://abc7.com/1187404/
http://abcd2.seesaa.net/
http://abcn.ws/1PForKt
http://abcn.ws/1XoiHdg
http://abizy.com/p/rss.html?user=http://twitrss.me/twitter_user_to_rss/?user=ReallyFreeCams
http://abizy.com/p/view.html?url=http://stackoverflow.com/questions/35373514/how-to-use-pow-function-for-calculating-powers-more-than-232-in-c
http://abr.ai/241C7Jk
http://abuse.sk.211.ca
http://act.credoaction.com/sign/Snyder_Subpoena?sp_ref=173834570.4.154073.o.1.2&referring_akid=.9190838.2IOkH1&source=clickcopy_sp
http://act.democracyforamerica.com/s/254036.45ngAY
http://ad.c-ats.jp/ad/p/r?_site=70&_article=204&_link=320&_image=320
http://adove.top/index.php?no=198103
http://afaee.xyz/qepE2
```

For the first question I took the 1000 different URI's that I had acquired from question one and put them in a new file named 'uriFile.txt' With that information I created a new python program called 'htmlExtractor.py' which covered all of the steps to extract the processed files.

```
#-*- coding: utf-8 -*-

import requests
#import uuid
import os


URIFile = open('uriFile.txt', 'r') #Opens the file w/ all of the uri's to test in read mode
outputFile = open('hashUnedited.txt', 'w')
counter = 0

ScriptPartOne = "echo -n '"
ScriptPartTwo = "' | md5sum >> hashUnedited.txt"
ScriptPartThree = "lynx -source "
ScriptPartFour = " > "
ScriptPartFive = "lynx -dump -force_html "
ScriptPartSix = ".processed"

with open('uriFile.txt') as f:
    lines = f.readlines()
    lines = [x.strip('\n') for x in lines]

for URLS in lines:
    r = os.system(ScriptPartOne + URLS + ScriptPartTwo)
    #outputFile.write(r.text())

HashFile = open('hashUnedited.txt', 'a')
with open('hashUnedited.txt') as z:
    hashName = z.readlines()
    hashName = [x.strip('\n') for x in hashName]
    hashName = [n.strip(' *- ') for n in hashName]

HashFileNames = open('editedFileNames.txt', 'a')
for names in hashName:
    a = os.system(ScriptPartThree + lines[counter] + ScriptPartFour + hashName[counter])
    HashFileNames.write(hashName[counter])
    HashFileNames.write("\n")
    counter += 1
HashFileNames.close()
with open('editedFileNames.txt') as s:
    readFiles = s.readlines()
    readFiles = [i.strip('\n') for i in readFiles]

for loop in readFiles:
    print loop
    q = os.system(ScriptPartFive + loop + ScriptPartFour + loop + ScriptPartSix)
    #outputFile.write(r.text())

URIFile.close()
HashFile.close()
```

Using the md5sum commands I was able to traverse through my URI file and create a hash value for each to be used as a naming convention. This was all stored initially inside of a file named 'hashUnedited' due to md5sum creating a trailing '*-' after each has value.

```
b1667f94d89d2eff4af4a1b1cbdb3e47 *-       b03710141f1f1bec1fa7533c0d865f56
da6b60b4b4978f3b5b978f1754f3c457 *-       9a95b8419da3cb23dc515070ab10a514
8867e76abff84f75ddc053902662a60b *-       e0c44cf2483467a95da398bb56b31080
ae17e1966d43b1d062be59fd01f0dc03 *-       28e738d8cf355578e73dda7b7ac882e8
ed78788da45ba6db1e0bd0ddad5a73f9 *-       6735d49db53561c5500aca73d01d353d
6533c991a236bce3c4df6215577a5921 *-       3d32c1fde4f142817476aafde9bd7145
52026916d5b2c2b679c3a38fcfbda2e7 *-       75384c9180dc62c30ec35e0ad263ce82
23d911ba989e5f57ae6aa6088fb0be13 *-       8f707f857351e9f39bee81278e223339
fcb458cd93d86564f0e2f1fa997b1292 *-       d9d6d392381ce2d2cc76feea6258d9c1
7a6b7b5ca1c09b92c24f39d37be1e4b8 *-       ba42b2621b210d1fe843281be2c96ecf
c8cd2423db653a2a77820cba738a9c48 *-       09a69936db9872bb22fa0038ce16c9f7
0957f31e48544a30f36b7e9ce04cf438 *-       d8a1ab30cc1b5a8141830199c65da664
582cf4e4cd48badcc4753cbfcca179ae *-       bf538dfd9a3cb9c7f032507d6337deae
b81617650d7cc6c4ec5d721e30a8d789 *-       74b1aa2dbd110c32bd4ab3fc01c58966
01d27434e0d841039e2323f980302c87 *-       c0922e9d8ac29440ee64fc5edfebf2f5
82ff7e215a38ba29fe02431237cee2a4 *-       c70f02a1f214d789f04bc867108ef873
fe800b15bb93a444d4a0fc105f81f437 *-       1814922a0433ee3005bb83cb95268842
fa36cdbbaa1195feb620a44474fbab1c *-       035fe57f901f2a3142698450a313e5f2
98f1f47bfefb03f1383619ea173477fe *-       dbd36b13233f374a4065df14770d9b0d
4faa22f822bc80b0492352eb28fb52c0 *-       cf7023f9da5563aef5dbbfe47aca2dbd
636e7eff34ea92e3d3c890151619c39f *-       2b6d9be961e7eb30c2bc2f029220bb1b
d226bdcfa2f60ff5ab90bad47da1bfad *-       a49ece1824d9bd31c11140d87e30d6bd
dc05aff872329c28730798bff26084df *-       f249832f5b1074487eadd3705b017269
e03db9af018355d695fdb603d3041d99 *-       70ec34751961c47dbb54666c208d66a0
c1b749306badc0e3f750cec21e482d02 *-       5aa838c2cdb61cbf5edd14047acbc340
9888c7f9740003c38069951128a14366 *-       e8b7770d51ca3b4a21d7cca0fba40787
47a4b852fe5690f06d70cba90c91142d *-       ab8c33550b7ced323ba9a509c77f742a
0fdabb027e6e4543d8a525abe3a378dd *-       2c24e05b81f99145605dcd5d3ce002a8
d34d3042c5a4171e19a9a0bdac544452 *-       abe490a1901491e3f724052586a80afc
86f13c295cdc558b111d7c3c590cc0e9 *-       54d9ad5df2b2fba00217412cc7eefbe1
45c374d37ce33c2b0b5318f280dcb788 *-       f50ea2ab5587b1400ac23909f0278c13
f0c2945d964ed6eeda1ca043c95d56bc *-       b429026d1edefe2954efa345830b1f0a
176e3c6b9b8078ea45f86451606cab51 *-       32f5857a73942cec182309ddbcf7a109
0b0e30c8dcc345f617fb10472a30f074 *-       1b94649a322002d1edcff00886fc9958
193ee5097813e35b16d3d06893b3cf1c *-       ac72c29c0d371b93011ad80d408dc37a
659cf7d31c9c9023389afe28a2887710 *-       4480d14ca804575295f0998096f899d0
b8ba666b724df9c0232d00cfb98a2f4f *-       d1c787bfdcc780b6408f9d0379747b32
f404cbf3aed73a6c341d59f1cc69971e *-       25e68a231817e72119e18d81d48d57ab
033e8b9934af50a191897d6c66d650d9 *-       6ee2df41c8aa1ae32690fa6f2b3e249a
a25d76bb91630613d052a6fb29e72706 *-       a94cb74f75165750cd246ff33bc04d14
2361024e354eb0b6116e5c296e0fdec6 *-       d08fc04230348f56147bc8c04034040f
9ac54005fd43faee97052d875f036424 *-       ed00c76032cfa33f3db85a1f0f140a8d
2c2c76b8eb444326d6d0a22411476517 *-       e033b19ac936c42bf26fdc3a5aef7fd2
aa3ef4a9e6a819d143be9a342b9d487c *-       f2a84a78b2ba142ec1e7b4a22b58a5a5
84252e412cdc059c4237916176ef171b *-       34f63696ea3a8acb944f98bbebf5c888
ec755aa88e22873d58265c831e43c419 *-       63d3ff8b1193b161697f7b0772f58b0b
4c3b270e02fea91f2c626c4e790922a4 *-       ed6082d2dea413702de51928473fb35b
cfbebe4995a68366cc1712f6213f0785 *-       ba62dafc0cb6d79543046453233a2328
c85d3d6b62a0264092c7a2e6522b9a42 *-       8e686c201cbe93d22dab974c675cfdd3
c2ff0d27eace68341b904b03b44de176 *-       b0651082f7decb1b9471082e3ad11340
ae7c73443ad5c6bdf6348d968b46161b *-       4ba0f7d55055745c60055d4a9b45ba15
                                          2e502ad228ab7ae2138517e63df771d9
```

After obtaining the initial unedited hash values, I put all of the original URIs through the 'lynx –source foo.bar > hashValue' meanwhile trimming off the ends of the hash values. This gave me both the list of edited hash values and the 1000 files of pure html named after the hash codes under the filename 'editedFileNames.txt'.

After that, I ran each generated file through the 'lynx –dump' command to get the 1000 processed files without the HTML structure thus completing question 1.

```
Social Media Marketing for Business


     * MavSocial Secure Login
     * _____
     * _____
     * [_] Remember Me [1]Forgot Your Password?
     *
     * Don't have an account?
     * [2]Sign up here. (Submit) Login


  [mavsocial_login_image.png]


  Software for Social Media Marketing Powered by MavSocial


          (c) Copyright 2013-2016 Maventus Group Incorporated


  (BUTTON) ×

Getting Started with MavSocial


  1
  [3]Social Networks
  Add any of your Social Networks and start using MavSocial.
  2
  [4]Digital Library
  Manage your digital content effectively with MavSocial Digital Library
  .
  3
  [5]Post Manager
  Share your content with the powerful MavSocial Post Manager.

References

  Visible links
  1. file://localhost/cygdrive/c/Users/jacob_000/Desktop/CS432/assignmentThree/HTMLWork/0a128737e76ec230f4cecb5b01fa8d1f#forgot_password
  2. http://www.mavsocial.com/sign-up/
  3. https://app.maventus.com/account/manage-social-network
  4. https://app.maventus.com/media/image-library
  5. https://app.maventus.com/post-manager/index

  Hidden links:
  6. https://app.maventus.com/home/home
```

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 10URIs, 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.

All of my work can be shown in the 'calculations.xlsx' file. My keyword was 'smartphone'.

| TFIDF | TF | IDF | URI |
|---|---|---|---|
| 0.0095 | 0.00153 | 6.2143 | http://wrld.bg/Y8Rg2 |
| 0.0243 | 0.00392 | 6.2143 | http://www.elmejormovil.net/#Ranking_mejores_moviles |
| 0.0059 | 0.00096 | 6.2143 | http://www.news.com.au/technology/online/social/latest-financial-report-for-streaming-service-soundcloud-is-far-worse-than-anybody-imagined/news-story/771fdf1f88079550f14c0ba07d09bbb9 |
| 0.0028 | 0.00045 | 6.2143 | http://host.madison.com/business/article_0dd4070e-d278-5293-a1ec-82e59e54351a.html |
| 0.0369 | 0.00594 | 6.2143 | http://www.gizmag.com/laser-weld-neurons/41807/ |
| 0.0049 | 0.0008 | 6.2143 | http://www.jagran.com/spiritual/religion-when-why-and-how-the-worship-of-goddess-saraswati-vasant-panchami-13571414.html |
| 0.0059 | 0.00096 | 6.2143 | http://www.news.com.au/technology/online/social/latest-financial-report-for-streaming-service-soundcloud-is-far-worse-than-anybody-imagined/news-story/771fdf1f88079550f14c0ba07d09bbb9 |
| 0.0037 | 0.0006 | 6.2143 | http://www.examiner.com/review/wild-game-cookbook-fried-moose-ribs-with-poached-pears-and-candied-chestnuts |
| 0.0097 | 0.00156 | 6.2143 | http://www.macrumors.com/2016/02/12/eddy-cue-craig-federighi-bloated-software/ |
| 0.0053 | 0.00085 | 6.2143 | http://lt.cl/Yh311 |

**IDF(t) = log_2(Total number of documents / Number of documents with term t in it)**

Total Docs with term in it: 633 Total Docs in Corpus = 47Billion m

Total Documents

IDF(t) = log_2(47,000,000,000 / 633,000,000)

IDF(t)=log_2(74.2496050553)

IDF(t)= 6.214311446912

IDF(t) = 6.2143

TF(t) = (Number of times term t appears in a document) / (Total number of terms in the document)

| Word Count | Total words | TF | URI |
|---|---|---|---|
| 2 | 1305 | 0.00153 | http://wrld.bg/Y8Rg2 |
| 21 | 5363 | 0.00392 | http://www.elmejormovil.net/#Ranking_mejores_moviles |
| 1 | 1038 | 0.00096 | http://www.news.com.au/technology/online/social/latest-financial-report-for-streaming-service-soundcloud-is-far-worse-than-anybody-imagined/news-story/771fdf1f88079550f14c0ba07d09bbb9 |
| 1 | 2233 | 0.00045 | http://host.madison.com/business/article_0dd4070e-d278-5293-a1ec-82e59e54351a.html |
| 8 | 1346 | 0.00594 | http://www.gizmag.com/laser-weld-neurons/41807/ |
| 4 | 4993 | 0.0008 | http://www.jagran.com/spiritual/religion-when-why-and-how-the-worship-of-goddess-saraswati-vasant-panchami-13571414.html |
| 5 | 5203 | 0.00096 | http://www.news.com.au/technology/online/social/latest-financial-report-for-streaming-service-soundcloud-is-far-worse-than-anybody-imagined/news-story/771fdf1f88079550f14c0ba07d09bbb9 |
| 1 | 1672 | 0.0006 | http://www.examiner.com/review/wild-game-cookbook-fried-moose-ribs-with-poached-pears-and-candied-chestnuts |
| 2 | 1281 | 0.00156 | http://www.macrumors.com/2016/02/12/eddy-cue-craig-federighi-bloated-software/ |
| 4 | 4697 | 0.00085 | http://lt.cl/Yh311 |

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.

| PageRank Table | USING  http://www.seocentro.com/tools/search-engines/pagerank.html | |
|---|---|---|
| TFIDF Rank | PageRank | URI |
| 1 | 0.4 | http://www.movilzona.es/2016/02/12/mejora-el-uso-de-iphone-con-esta-curiosa-organizacion-de-las-aplicaciones/ |
| 2 | 0.6 | http://www.gizmag.com/laser-weld-neurons/41807/ |
| 3 | 0.1 | http://www.elmejormovil.net/#Ranking_mejores_moviles |
| 4 | 0.7 | http://www.macrumors.com/2016/02/12/eddy-cue-craig-federighi-bloated-software/ |
| 5 | 0.2 | http://wrld.bg/Y8Rg2 |
| 6 | 0.8 | http://www.news.com.au/technology/online/social/latest-financial-report-for-streaming-service-soundcloud-is-far-worse-than-anybody-imagined/news-story/771fdf1f88079550f14c0ba07d09bbb9 |
| 7 | 0.1 | http://lt.cl/Yh311 |
| 8 | 0.6 | http://www.jagran.com/spiritual/religion-when-why-and-how-the-worship-of-goddess-saraswati-vasant-panchami-13571414.html |
| 9 | 0.8 | http://www.examiner.com/review/wild-game-cookbook-fried-moose-ribs-with-poached-pears-and-candied-chestnuts |
| 10 | 0.6 | http://host.madison.com/business/article_0dd4070e-d278-5293-a1ec-82e59e54351a.html |