CS-432/532 Introduction to Web Science: Assignment #8: Clustering Algorithms

Due on Thursday, April 7, 2016

Dr. Michael L. Nelson

Plinio Vargas pvargas@cs.odu.edu

Contents

Problem 1	1
1.1 Approach	
1.2 Solution	4
Problem 2	4
2.1 Approach	4
2.2 Solution	4
Problem 3	6
3.1 Approach	6
3.2 Solution	6
Problem 4	10
4.1 Approach	10
4.2 Solution	12
Problem 5 - Extra Credit	13
5.1 Approach	13
5.2 Solution	15
List of Figures	
Listings	
1 Creating Blog Matrix: BlogMatrix.py	1
2 Accounting for all Blog Pages: BlogMatrix.py	2
3 Generating ASCII and JPEG Dendrogram: PrintClusters.py	
4 Using K-Means: PrintK-Clusters.py	
5 Creating MDS with: CreateMDS.py	
6 Using TFIDF: CalcTFIDF.py	13

List of Tables

Problem 1

Create a blog-term matrix. Start by grabbing 100 blogs; include:

```
http://f-measure.blogspot.com/
http://ws-dl.blogspot.com/
```

and grab 98 more as per the method shown in class. Note that this method randomly chooses blogs and each student will separately do this process, so it is unlikely that these 98 blogs will be shared among students. In other words, no sharing of blog data. Upload to github your code for grabbing the blogs and provide a list of blog URIs, both in the report and in github..

Use the blog title as the identifier for each blog (and row of the matrix). Use the terms from every item/title (RSS) or entry/title (Atom) for the columns of the matrix. The values are the frequency of occurrence. Essentially you are replicating the format of the "blogdata.txt" file included with the PCI book code. Limit the number of terms to the most "popular" (i.e., frequent) 500 terms, this is *after* the criteria on p. 32 (slide 7) has been satisfied.

1.1 Approach

A great deal of code was obtain from [1] in order to complete this assignment. All the codes in [1] are implemented in Python 2.7, while the implementation in this work is based using Python 3.4. There were minor changes to adapt to and will not be detailed to explain this approach, only major modifications will be pointed out.

BlogMatrix.py is the *Python* code implemented to answer this problem. To form the Matrix, we created a *set* object (*urilist*) containing 98 generated blog titles. Two required blogs were added in lines 42-43.

Listing 1: Creating Blog Matrix: BlogMatrix.py

```
urilist = set()
37
38
      apcount = {}
      wordcounts = {}
      feedlist = []
4(
41
      urilist.add('http://f-measure.blogspot.com/')
      urilist.add('http://ws-dl.blogspot.com/')
42
      for uri in urilist:
48
           title, wc = getallpages(uri)
46
           if title:
47
               wordcounts[title] = wc
48
               print(title, wordcounts[title])
49
               for word, count in wc.items():
                    apcount.setdefault(word, 0)
51
52
                    if count > 1:
                        apcount[word] += 1
               feedlist.append(uri)
```

```
urilist.add(uri)
55
      while len(urilist) < 100:</pre>
57
           lasturi = set()
           AddURI(lasturi, 'http://www.blogger.com/next-blog?navBar=true&blogID
               =3471633091411211117<sup>'</sup>)
           print('Getting word count....')
61
           for uri in lasturi:
               uri = re.match('(^.*\.com\/)(?!expref)', uri).group(0)
63
               if uri not in urilist:
64
                    title, wc = getallpages(uri)
                    if title:
66
                        wordcounts[title] = wc
67
                        print(len(urilist), title, wordcounts[title])
68
                        for word, count in wc.items():
                             apcount.setdefault(word, 0)
                             if count > 1:
71
                                 apcount[word] += 1
                        feedlist.append(uri)
                        urilist.add(uri)
74
75
                    else:
                        print('Discarding URI...')
```

We proceed to find the remaining URL(s) for our matrix. We re-utilized the coded developed on assignment 2 to find what is the ultimate location to the HTTP redirection (line 59).

We removed undesired strings from the returned URL, using regular expression (line 63), discarding any URL without a title (line 66). Finally, we include the URL in a *list* object (line 73), in order to provide all URLs used for this assignment. The remaining of the code was an extraction from [1]

A major difference between this implementation and [1] is that the latest does not account for all the webpages in a blog. Only the first blog page is taken in consideration for word-counts.

In order to account for all the pages in the blog, we created two functions in **BlogMatrix.py**: getallpages and linkpages.

Listing 2: Accounting for all Blog Pages: BlogMatrix.py

```
def getallpages(uri):
    twc = None
    rss_feed = uri + 'feeds/posts/default/'
    title, wc = getwordcounts(rss_feed)
    if title:
        print('Getting linkpages....', title)
        twc = linkpages(wc, rss_feed)

if not twc:
    return '', ''
```

```
return title, twc
121
  def linkpages(wc, rss_feed):
       twc = {}
122
       for word, count in wc.items():
           twc.setdefault(word, 0)
           twc[word] += count
       flag = True
       while flag:
           page = requests.get(rss_feed).text
129
           soup = BeautifulSoup(page, 'html.parser')
130
           for link in soup.find_all('link'):
131
               t = link.get('rel')
132
                if t and t[0] == 'next':
                    print('Getting %s ...' % link.get('href'))
134
                    title, wc = getwordcounts(link.get('href'))
135
136
                    rss_feed = link.get('href')
137
                    for word, count in wc.items():
138
                        twc.setdefault(word, 0)
139
                         twc[word] += count
141
142
                    print(twc)
                    break
```

For each URL generated a call is made to *getallpages* (see line 46 Listing 1) which has the purpose of passing all the words and title from a particular blog, back to the main function (line 114).

linkpages actually counts all words and their frequencies in a blog. This is accomplished by inspecting the link tag with rel='next'. Using the same scheme as [?], adding the initial values of the first page in the title word count variable *twc* (lines 119-121).

Each page word count will be added to *twc* until there are no more rel='next' links (line 127-142). The iteration will continue until all generated URLs are processed. The remaining of the code is a replica from [?], but is worthy to mention that to find the 500 terms from the blog, we just selected the first 500 column in our

1.2 Solution

- a. File: feedlist.txt in github contains all the URLs feed for the assignment including: http://f-measure.blogspot.com/ http://ws-dl.blogspot.com/
- b. File **blogdata.tex**: contains the matrix with blog title and word counts, similar to [1].

Problem 2

Create an ASCII and JPEG dendrogram that clusters (i.e., HAC)the most similar blogs (see slides 12 & 13). Include the JPEG in your report and upload the ascii file to github (it will be too unwieldy for inclusion in the report).

2.1 Approach

We used **blogdata.tex** file, exact code in [1] to complete this problem. We added a file stream (line 19) in order to print **ascii** file.

Listing 3: Generating ASCII and JPEG Dendrogram: PrintClusters.py

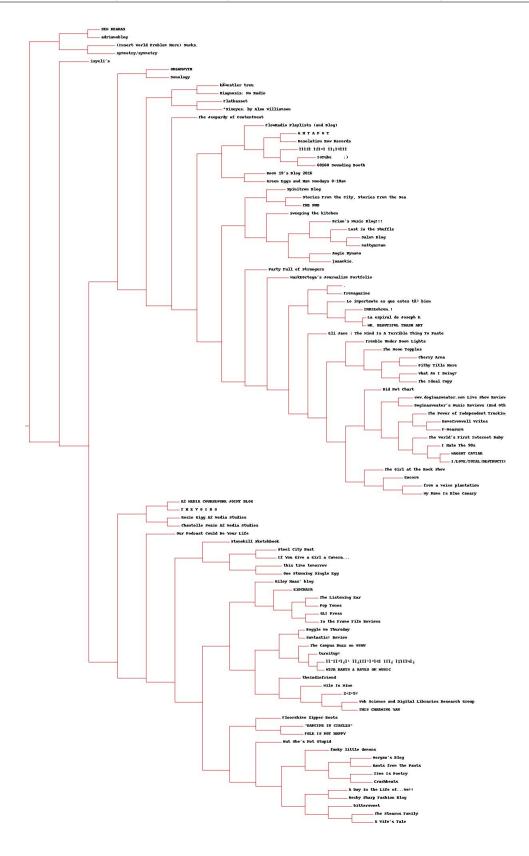
```
import lib.PCI_Code_Folder.chapter3.clusters as clusters

def main():
    blognames,words,data=clusters.readfile('blogdata.txt')
    clust=clusters.hcluster(data)
    outfile = open('ascii-dendogram.txt', 'w')
    clusters.printclust(clust, labels=blognames, file=outfile)
    outfile.close()
    clusters.drawdendrogram(clust,blognames,jpeg='blogclust.jpg')

return
```

2.2 Solution

- a. File: ascii-dendogram.txt in github contains ASCII dendrogram that clusters the most similar blogs.
- b. File **blogclust.jpeg**: contains JPEG dendrogram that clusters the most similar blogs. See example below.



Problem 3

Cluster the blogs using K-Means, using k=5,10,20. (see slide 18). Print the values in each centroid, for each value of k. How many iterations were required for each value of k?

3.1 Approach

For this problem *PrntK-Cluster.py* was implemented, using

Listing 4: Using K-Means: PrintK-Clusters.py

```
import lib.PCI_Code_Folder.chapter3.clusters as clusters
  def main():
      # record running time
      start = time()
17
      print('Starting Time: %s' % strftime("%a, %b %d, %Y at %H:%M:%S", localtime()
          ))
      outfile = open('k-clusters.txt', 'w')
20
      blognames, words, data=clusters.readfile('blogdata.txt')
      k_{means} = [5, 10, 20]
22
      for k in k_means:
          kclust, n_iterations = clusters.kcluster(data, k=k)
25
          counter = 0
26
          for i in range(len(kclust)):
27
              outfile.write('Cluster-%d:\n%s\n' % (i, ('-' * 9)))
              print('Cluster-%d:\n%s' % (i, ('-' * 9)))
29
              for r in kclust[i]:
30
                  counter += 1
                   outfile.write('%s,' % blognames[r])
32
                  print(blognames[r], end=', ')
33
                   if counter % 5 == 0:
34
                       outfile.write('\n')
                       print()
36
               outfile.write('\n')
              print('\n')
39
          outfile.write('Number of iterations for k=\%d is \%d.\n\n' % (k,
40
              n_iterations + 1))
          print('Number of iterations for k=\%d is \%d.\n\n'' \% (k, n_iterations + 1)
      print('\nEnd Time: %s' % strftime("%a, %b %d, %Y at %H:%M:%S", localtime()))
      print('Execution Time: %.2f seconds' % (time()-start))
```

3.2 Solution

```
Cluster-0:
```

```
A2 MEDIA COURSEWORK JOINT BLOG, A Day in the Life of...Me!!, La espiral de Joseph
    K, theindiefriend, GLI Press,
symmetry/symmetry, MR. BEAUTIFUL TRASH ART, What Am I Doing?, A Wife's Tale, Morgan's Blog,
., The Girl at the Rock Show, KiDCHAIR, 2+2=5?, MarkEOrtega's Journalism Portfolio,
isyeli's, Steel City Rust, Tremagazine, jaaackie., The Ideal Copy,
from a voice plantation, In the Frame Film Reviews, The Jeopardy of Contentment, My Name Is Blue
    Canary, The Listening Ear,
Cherry Area, Web Science and Digital Libraries Research Group, Becky Sharp Fashion Blog, Our Podcast
    Could Be Your Life, If You Give a Girl a Camera...,
THIS CHARMING YAN, Mile In Mine, Pop Tones,
Cluster-1:
Riley Haas' blog,
Cluster-2:
(Insert World Problem Here) Sucks.,
ORGANMYTH, The Stearns Family, One Stunning Single Egg, Time Is Poetry, Sonology,
funky little demons, But She's Not Stupid, Floorshime Zipper Boots, Room 19's Blog 2016, Crashbeats,
The Moon Topples, Pithy Title Here, Rants from the Pants, Tremble Under Boom Lights, bittersweet,
Cluster-3:
-----
Salem Blog, this time tomorrow, , A H T A P O T, IoTube :),
Lo importante es que estes tu bien, The Campus Buzz on WSOU, Angie Dynamo, FlowRadio Playlists (and
    Blog), ,
Chantelle Swain A2 Media Studies, Stories From the City, Stories From the Sea, 60\060 Sounding
    Booth, Party Full of Strangers, THE HUB,
Stonehill Sketchbook, Lost in the Shuffle, Spinitron Blog, SEM REGRAS, INDIE ohren.!,
Desolation Row Records, sweeping the kitchen, mattgarman, Brian's Music Blog!!!, Green Eggs and Ham
    Mondays 8-10am,
adrianoblog,
Cluster-4:
MAGGOT CAVIAR, The World's First Internet Baby, Rosie Gigg A2 Media Studies, "DANCING IN CIRCLES",
T H E V O I D S, kunstler treu, Boggle Me Thursday, DaveCromwell Writes, turnitup!,
Encore, Flatbasset, FOLK IS NOT HAPPY, MTJR RANTS & RAVES ON MUSIC, I Hate The 90s,
Diagnosis: No Radio, Samtastic! Review, *Sixeyes: by Alan Williamson, www.doginasweater.com Live Show
    Review Archive, F-Measure,
I/LOVE/TOTAL/DESTRUCTION, The Power of Independent Trucking, Did Not Chart, Doginasweater's Music
    Reviews (And Other Horseshit), Eli Jace | The Mind Is A Terrible Thing To Paste,
Number of iterations for k=5 is 5.
Cluster-0:
-----
Cluster-1:
MAGGOT CAVIAR, The World's First Internet Baby, symmetry, "DANCING IN CIRCLES", Sonology,
DaveCromwell Writes, Encore, from a voice plantation, My Name Is Blue Canary, I Hate The 90s,
```

```
*Sixeyes: by Alan Williamson, www.doginasweater.com Live Show Review
    Archive, F-Measure, I/LOVE/TOTAL/DESTRUCTION, The Power of Independent Trucking,
Did Not Chart, Tremble Under Boom Lights, Doginasweater's Music Reviews (And Other
    Horseshit), mattgarman, Eli Jace | The Mind Is A Terrible Thing To Paste,
Cluster-2:
Riley Haas' blog,
Cluster-3:
Unicode Character, The Campus Buzz on WSOU, Boggle Me Thursday, turnitup!,
MTJR RANTS & RAVES ON MUSIC, Samtastic! Review,
Cluster-4:
La espiral de Joseph K, MR. BEAUTIFUL TRASH ART, Lo importante es que estes tu bien,
., But She's Not Stupid, MarkEOrtega's Journalism Portfolio, Tremagazine, FOLK IS NOT HAPPY,
SEM REGRAS, INDIEohren.!, adrianoblog,
Cluster-5:
GLI Press, kunstler treu,
KiDCHAIR, Flatbasset, In the Frame Film Reviews, The Jeopardy of Contentment, The Listening Ear,
Diagnosis: No Radio, Our Podcast Could Be Your Life, Pop Tones,
Cluster-6:
_____
A2 MEDIA COURSEWORK JOINT BLOG, ORGANMYTH,
theindiefriend, Rosie Gigg A2 Media Studies, T H E V O I D S, Floorshime Zipper Boots, 2+2=5?,
Chantelle Swain A2 Media Studies, Web Science and Digital Libraries Research Group, THIS CHARMING
    YAN.
Cluster-7:
Cluster-8:
this time tomorrow, A H T A P O T,
        :),FlowRadio Playlists (and Blog),UnicodeCharacters,Stories From the City, Stories From
IoTube
    the Sea,60060 Sounding Booth,
Party Full of Strangers, THE HUB, Lost in the Shuffle, Spinitron Blog, Desolation Row Records,
sweeping the kitchen, Brian's Music Blog!!!, Green Eggs and Ham Mondays 8-10am,
Cluster-9:
Salem Blog, (Insert World Problem Here) Sucks.,
A Day in the Life of...Me!!, The Stearns Family, One Stunning Single Egg, What Am I Doing?, A Wife's
Time Is Poetry, Morgan's Blog, Angie Dynamo, funky little demons, The Girl at the Rock Show,
isyeli's, Steel City Rust, jaaackie., The Ideal Copy, Room 19's Blog 2016,
Crashbeats, The Moon Topples, Cherry Area, Becky Sharp Fashion Blog, Stonehill Sketchbook,
If You Give a Girl a Camera..., Pithy Title Here, Rants from the Pants, Mile In Mine, bittersweet,
Number of iterations for k=10 is 11.
Cluster-0:
-----
```

```
Salem Blog, IoTube :), Stonehill Sketchbook, Lost in the Shuffle, sweeping the kitchen,
mattgarman, Brian's Music Blog!!!,
Cluster-1:
_____
MAGGOT CAVIAR, The World's First Internet Baby, "DANCING IN CIRCLES",
.,Floorshime Zipper Boots,DaveCromwell Writes,FOLK IS NOT HAPPY,I Hate The 90s,
Web Science and Digital Libraries Research Group, F-Measure, I/LOVE/TOTAL/DESTRUCTION, THIS CHARMING
    YAN, The Power of Independent Trucking,
Did Not Chart, Doginasweater's Music Reviews (And Other Horseshit), Eli Jace | The Mind Is A
    Terrible Thing To Paste,
Cluster-2:
A Day in the Life of ... Me!!, ORGANMYTH,
The Stearns Family, A Wife's Tale, Sonology, Angie Dynamo, jaaackie.,
Room 19's Blog 2016, Crashbeats, Diagnosis: No Radio, Cherry Area, Pithy Title Here,
Mile In Mine, Green Eggs and Ham Mondays 8-10am, bittersweet,
Cluster-3:
-----
Riley Haas' blog,
Cluster-4:
Cluster-5:
-----
Cluster-6:
_____
kunstler treu,
Flatbasset, *Sixeyes: by Alan Williamson,
Cluster-7:
(Insert World Problem Here) Sucks.,
Cluster-8:
A2 MEDIA COURSEWORK JOINT BLOG, Rosie Gigg A2 Media Studies,
T H E V O I D S, Chantelle Swain A2 Media Studies,
Cluster-9:
Unicode-Characters, The Campus Buzz on WSOU, turnitup!,
MTJR RANTS & RAVES ON MUSIC,
Cluster-10:
this time tomorrow, the indiefriend, One Stunning Single Egg, What Am I Doing?,
funky little demons,2+2=5?, The Ideal Copy, The Moon Topples, If You Give a Girl a Camera...,
Spinitron Blog, Rants from the Pants,
Cluster-11:
-----
Cluster-12:
La espiral de Joseph K, MR. BEAUTIFUL TRASH ART, Lo importante es que estes tu bien,
MarkEOrtega's Journalism Portfolio, The Jeopardy of Contentment, SEM REGRAS, INDIE ohren.!, adrianoblog,
```

```
Cluster-13:
Cluster-14:
A H T A P O T, FlowRadio Playlists (and Blog), Unicode-Characters, Stories From the City, Stories
    From the Sea,60060 Sounding Booth,
Party Full of Strangers, THE HUB, Desolation Row Records,
Cluster-15:
-----
Boggle Me Thursday, Samtastic! Review,
Cluster-16:
symmetry/symmetry, The Girl at the Rock Show, Encore, from a voice plantation, My Name Is Blue Canary,
www.doginasweater.com Live Show Review Archive, Becky Sharp Fashion Blog, Tremble Under Boom Lights,
Cluster-17:
-----
Time Is Poetry, But She's Not Stupid,
isyeli's, Tremagazine,
Cluster-18:
_____
Cluster-19:
_____
GLI Press, Morgan's Blog, KiDCHAIR,
Steel City Rust, In the Frame Film Reviews, The Listening Ear, Our Podcast Could Be Your Life, Pop
    Tones,
Number of iterations for k=20 is 6.
```

```
Number of iterations for k=5 is 5.

Number of iterations for k=10 is 11.

Number of iterations for k=20 is 6.
```

Problem 4

Use MDS to create a JPEG of the blogs similar to slide 29. How many iterations were required?

4.1 Approach

For this problem *CreateMDS.py* was implemented, using

Listing 5: Creating MDS with: CreateMDS.py

```
blognames, words, data=clusters.readfile('blogdata.txt')
coords, n_iterations = clusters.scaledown(data)
clusters.draw2d(coords, blognames, jpeg='blogs2d.jpg')

print("Number of iterations is %d" % n_iterations)

print('\nEnd Time: %s' % strftime("%a, %b %d, %Y at %H:%M:%S", localtime()))
print('Execution Time: %.2f seconds' % (time()-start))
return
```

Above code was extracted from [1]

4.2 Solution

File name: blogclust.jpg



Number of iterations required was 424

Problem 5 - Extra Credit

Re-run question 2, but this time with proper TFIDF calculations instead of the hack discussed on slide 7 (p. 32). Use the same 500 words, but this time replace their frequency count with TFIDF scores as computed in assignment #3. Document the code, techniques, methods, etc. used to generate these TFIDF values. Upload the new data file to github.

Compare and contrast the resulting dendrogram with the dendrogram from question #2.

Note: ideally you would not reuse the same 500 terms and instead come up with TFIDF scores for all the terms and then choose the top 500 from that list, but I'm trying to limit the amount of work necessary.

5.1 Approach

The approach is the same as in problem 2, but we substituted the matrix data calculating with TFIDF values. Our **Total Document Corpus** is accumulated as we read the matrix in line 38. **Document Term** values are accumulated in line 37.

Listing 6: Using TFIDF: CalcTFIDF.py

```
def main():
      n = 500
19
20
      m = 100
      matrix = []
      blognames = []
22
      matrix_total = [0] * n
      matrix_TDIF = [[0] * n for i in range(m)]
24
25
      total_doc_corpus = 0
      counter = 0
26
      with open('blogdata.txt', 'r') as infile:
27
           for line in infile:
               terms = line.strip().split('\t')
29
30
               if counter > 0:
                   row = []
31
                   row.append(terms[0])
32
                   blognames.append(terms[0])
33
34
                   i = 0
                   for size in terms[1:]:
                        row.append(size)
36
                        matrix_total[i] += int(size)
                        total_doc_corpus += int(size)
38
                        i += 1
39
                   matrix.append(row)
40
41
               counter += 1
      total_doc_corpus
                                               # total docs in corpus
      for i in range(m):
44
           for k in range(n):
               matrix_TDIF[i][k] = math.log(total_doc_corpus / matrix_total[k], 2) *
                   int(matrix[i][k + 1])
```

```
clust=clusters.hcluster(matrix_TDIF)

outfile = open('ascii-dendogramP5.txt', 'w')

clusters.printclust(clust, labels=blognames, file=outfile)

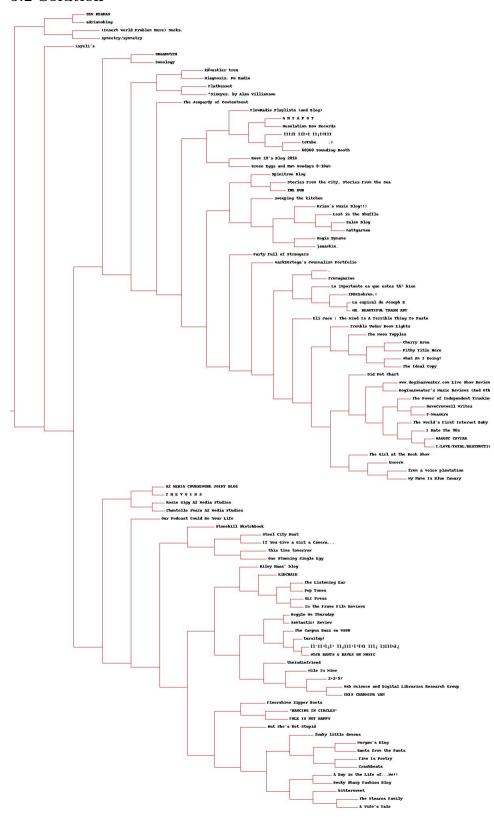
outfile.close()

clusters.drawdendrogram(clust,blognames,jpeg='blogclustP2.jpg')

return
```

Lines 44 through 47 iterates through the matrix to calculate the TFIDF for each blog term. The remaining of the code is a copy of problem 2.

5.2 Solution



The result is very interesting. A great deal of the blogs are clustering with similarity, but they seem to be following on a different categories. For example, blogs F-Measure, DaveCrownell Writes, I hate the 90's,

The girl at the Rock Show in problem 2 are clustering together in a flatter way while using TFIDF makes the relationship much hierarchical.

References

[1] Segarn, Toby. Programming Collective Intelligence. Building Smart Web 2.0 Application. (pp 29-53). Sebastopol, CA: O'Reilly Media.