# Assignment Nine

Sybil Melton

December 3, 2014

# CONTENTS

# LISTINGS

# 1 BLOG TERM MATRIX

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.

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.

Create a histogram of how many pages each blog has (e.g., 30 blogs with just one page, 27 with two pages, 29 with 3 pages and so on).

## 1.1 SOLUTION

My first step was to grab a list of blogs, using Blogger.com's "next blog" functionality. The blogs were all written to file by net location. After writing the two required blogs, a loop of 200 iterations ran to get additional blogs, using pyCurl to get the net location. [8] Once each blog was retrieved, it was checked to make sure it didn't already exist in the list with the regular expression module. This is the same function used in Assignment One. [2] Once I had the list, I went through it to check for inappropriate content and to remove any picture, video, or foreign language blogs. My final list of blogs was saved as blogs.txt and is included with this report. Listing 1 is the code used to accomplish the first part of this task.

```python
def check(url):
  with open("blogs.txt") as f:
    found = False
    for line in f:  #iterate over the file one line at a time(memory efficient)
      if re.search(url, line):   #if string found is in current line then keep it
        found = True
  return found

blogger = 'https://www.blogger.com/next-blog?navBar=true&blogID=953024975153422094'
bfile = open('blogs.txt', 'w',0)
bfile.write('f-measure.blogspot.com\n')
bfile.write('ws-dl.blogspot.com\n')
bfile.close()
for i in range(1, 200):
  buffer = BytesIO()
  c = pycurl.Curl()
  c.setopt(c.URL, blogger)
  c.setopt(c.WRITEDATA, buffer)
  c.setopt(c.FOLLOWLOCATION, True)
  c.setopt(c.HTTPHEADER, ['Accept-Language: en'])
  bfile = open('blogs.txt', 'a',0)
  try:
    c.perform()
    if c.getinfo(c.RESPONSE_CODE) == 200:
      o = urlparse(c.getinfo(c.EFFECTIVE_URL))
      ch = check(o.netloc)
      if not ch:
        bfile.write(o.netloc + '\n')
        bfile.close()
      else:
        continue
    c.close()
```

```
33   except pycurl.error , error:
34       errno , errstr = error
35       print 'An error occurred: ', errstr
```

Listing 1: Retrieve Blogs

My next step was to retrieve the blogs. Each blog in "blogs.txt" was added the feed path */feeds/posts/default* and encoded with query arguments for Atom, in order to get the raw feed. A feed list was kept, for the stop word hack calculation later on in the program. Each feed was retrieved using the given code to create a getFeed function, using getwordcounts and getwords from the PCI textbook to create the wordcounts dictionary of word counts per blog and apcount dictionary of counts per word. [10] For testing, the Web Science/Digital Libraries blog was downloaded to make sure the feed was retrieved and all BeautifulSoup was used to look for the link to the next page, if available. [3] A counter kept track of the number of pages in the blog, which was used for debugging and stored in a dictionary. The pages per blog and number of blogs was written to blog_pages.txt in order to create the histogram and is included with this report. [7] The apcount dictionary was unsorted, so in order to get the top 500 words, I converted it to a sorted list by word count. A counter was added to the code and stopped the loop if 500 was reached. Listing 2 is the code used to for these tasks.

```
1  def getwordcounts(url):
2      #Returns title and dictionary of word counts for an RSS feed
3      # Parse the feed
4      d = feedparser.parse(url)
5      wc = {}
6      # Loop over all the entries
7      for e in d.entries:
8          if 'summary' in e:
9              summary = e.summary
10             else:
11                 summary = e.description
12             # Extract a list of words
13             words = getwords(e.title + ' ' + summary)
14             for word in words:
15                 wc.setdefault(word, 0)
16                 wc[word] += 1
17      return (d.feed.title , wc)
18
19 def getwords(html):
20     # Remove all the HTML tags
21     txt = re.compile(r'<[^>]+>').sub('', html)
22     # Split words by all non-alpha characters
23     words = re.compile(r'[^A-Z^a-z]+').split(txt)
24     # Convert to lowercase
25     return [word.lower() for word in words if word != '']
26
27 def getFeed(feedurl):
28     global wordcounts
29     global apcount
30     try:
31         (title , wc) = getwordcounts(feedurl)
32         if title in wordcounts:
33             for (w, c) in wc.iteritems():
34                 if w in wordcounts[title].iteritems():
35                     wordcounts[title][w] += c
36                 else:
37                     wordcounts[title][w] = c
38         else:
39             wordcounts[title] = wc
40         for (word, count) in wc.items():
41             apcount.setdefault(word, 0)
42             if count > 1:
43                 apcount[word] += 1
44     except:
45         print 'Failed to parse feed %s' % feedurl
46
47 pages = {}
48 feedlist = []
```

4

```
49  apcount = {}
50  wordcounts = {}
51  wordlist = []
52  scheme = 'http://'
53  path = '/feeds/posts/default'
54  query_arg = {'alt' : 'atom'}
55  udata = urllib.urlencode(query_arg)
56  f = open('blogs.txt', 'r', 0)
57  for line in f:
58    line = line.strip()
59    data = urllib.urlencode(query_arg)
60    full_url = scheme+line+path+"?"+udata
61    feedlist.append(full_url)
62    r = requests.get(full_url)
63    getFeed(full_url)
64    ddata = r.text
65    soup = BeautifulSoup(ddata)
66    next = soup.find('link', rel='next')
67    count = 1
68    while next:
69      n= next.get('href')
70      getFeed(n)
71      r = requests.get(n)
72      data = r.text
73      soup = BeautifulSoup(data)
74      count = count + 1
75      if count % 10 == 0:
76        print ' -- parsing ' + line + ' ' + str(count) + ' pages so far'
77      next = soup.find('link', rel='next')
78    print 'Finished parsing ' + line + ' ' + str(count) + ' pages'
79    if str(count) in pages:
80      pages[str(count)].append(line)
81    else:
82      pages[str(count)] = []
83      pages[str(count)].append(line)
84  f.close()
85
86  p = open('blog_pages.txt', 'w', 0)
87  p.write('Pages\tNumber of Blogs\n')
88  for (pg, num) in pages.iteritems():
89    p.write(pg + '\t' + str(len(num)) + '\n')
90
91  a = sorted(apcount.items(), key=lambda x: x[1], reverse=True)
92  count = 0
93  for (w, bc) in a:
94      frac = float(bc) / len(feedlist)
95      if frac > 0.1 and frac < 0.5:
96          wordlist.append(w)
97          count = count + 1
98          if count >= 500:
99              break
```

Listing 2: Retrieve Wordcounts

At this point I realized I did not have a the titles or subtitles stored. So I wrote a function, getTitles, in order to accomplish this and were stored as a dictionary. The feed had to be parsed again for both title and subtitle. [5] I found a few still had html markup which caused extra lines and illegal unicode characters which had to be removed. [4], [6] The blog matrix was created and saved as blogdata.txt and is included with this report. The code from the textbook was modified to include the subtitle, or part of the subtitle, as long as the length was less than 100 characters. [11] Listing 3 is the code used to create the blog matrix.

```
1  _illegal_xml_chars_RE = re.compile(u'[\x00-\x08\x0b\xa9\x0c\x0e-\x1F\uD800-\uDFFF\uFFFE\uFFFF]')
2
3  def remove_tags(text):
4      cleanr =re.compile('<.*?>')
5      cleantext = re.sub(cleanr,'', text)
6      cleantext = re.sub("\n", ' ', cleantext)
7      cleantext = cleantext.strip()
```

```
 8      return _illegal_xml_chars_RE.sub(' ', cleantext)
 9
10  def getTitles():
11     titles = {}
12     global scheme, path, udata
13     f = open('blogs.txt', 'r', 0)
14     for line in f:
15        line = line.strip()
16        full_url = scheme+line+path+"?"+udata
17        d = feedparser.parse(full_url)
18        titles[d.feed.title] = remove_tags(d.feed.subtitle)
19     return titles
20
21  titles = {}
22  titles = getTitles()
23  # Blog Term Matrix
24  out = file('blogdata.txt', 'w', 0)
25  out.write('Blog')
26  for word in wordlist:
27     out.write('\t%s' % word)
28  out.write('\n')
29
30  for (blog, wc) in wordcounts.items():
31     if titles[blog] != '':
32        blog = blog + ' - ' + titles[blog]
33     blog = blog.replace(u'\u0144', 'n')
34     blog = _illegal_xml_chars_RE.sub(' ', blog)
35     blog = blog.strip()
36     if len(blog) > 100:
37        blog = blog[:99]
38     print blog
39     out.write(blog)
40     for word in wordlist:
41        if word in wc:
42           out.write('\t%d' % wc[word])
43        else:
44           out.write('\t0')
45     out.write('\n')
46  out.close()
```

Listing 3: Create Blog Matrix

Finally, the histogram was created in R, to show the number of blogs and pages. Figure 1.1 show there were more blogs with fewer pages, as there were only 12 blogs with greater than 18 pages. Listing 4 is the R code used to create it.



Figure 1.1: Pages vs Number of Blogs

```
1  bdata = read.table("C:/Users/sybil/Documents/2014 Fall/Web Science/A9/output/blog_pages.txt", header=T, sep
       ="\t")
2
3  newdata <- bdata[order(pages),]
4  pages <- c(newdata$Pages)
5  blogs <- c(newdata$Number.of.Blogs)
6
7  barplot(blogs, names.arg=pages, xlab="Pages", ylab="Number of Blogs")
```

Listing 4: Pages vs Number of Blogs R

## 2 DENDROGRAMS

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 SOLUTION

First I tried to change the provided function to output to a text file, but it was not successful. So I redirected the standard output temporarily to blog_ascii.txt, which is included with this report. [1] The JPEG dendogram was successful, as shown in Figure 2.1. The python code used is shown in Listing 5.

```
1  old = sys.stdout
2  sys.stdout = open('blog_ascii.txt', 'w', 0)
3  blognames,words,data=clusters.readfile('blogdata.txt')
4  clust = clusters.hcluster(data)
```

```
5  clusters.printclust(clust,labels=blognames)
6  sys.stdout = old
7  # jpeg dendogram
8  clusters.drawdendrogram(clust,blognames,jpeg='blogclust.jpg')
```

Listing 5: Create Dendrograms

Crop 4 kids

Really Simple Sidi (RSS)

User Driven Modelling - This blog is about my PhD research (now finished) at University of the West

ERDAS TITAN - Accessing geospatial data and content from unlimited public and private sources, into

The Ultimate Sales Executive Resource - For B2B sales leaders wanting to measurably improve the pr

All Work No Play - My tech Blog. Symbian Foundation. Mobile phones. Technology. Programming. Trend

Misc - On Software and Science

eclipser-blog - This is a blog of Eclipse enthusiasts from Poznan (Poland) who would like to share

Andrejus Baranovskis Blog - Blog about Oracle technology

MLA Wire

Shige's Research Blog

GNUmed for the masses - This blog deals with the Free and Open So

freegnu - and other stuff too

Luciano Resende - Open source, Apache, SOA, SCA and other good technical tips &amp;

Ruslan's blog

MY NODE - SHARING MY KNOWLEDGE

My life with Android :-)

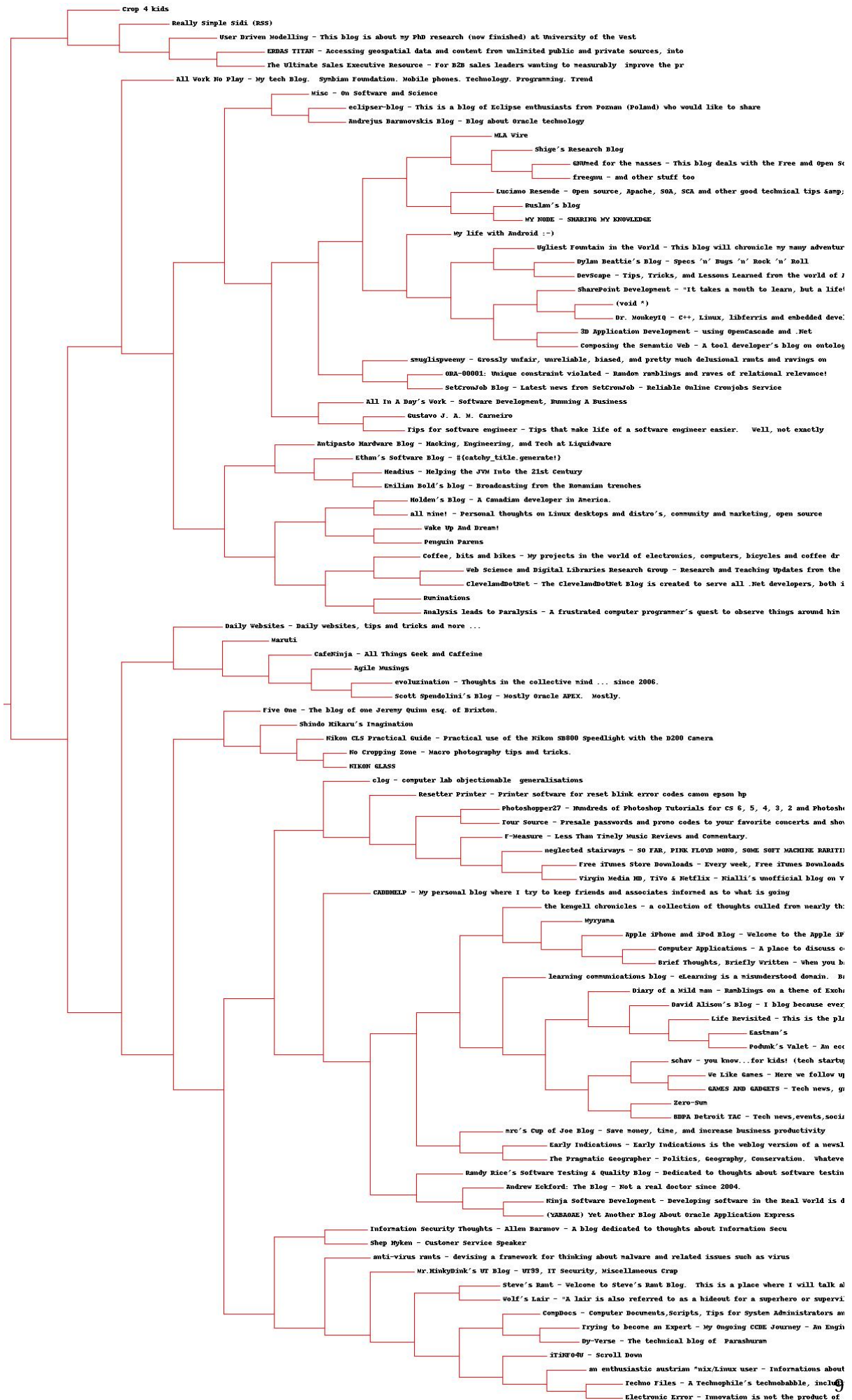Ugliest Fountain in the World - This blog will chronicle my many adventur

Dylan Beattie's Blog - Specs 'n' Bugs 'n' Rock 'n' Roll

DevScape - Tips, Tricks, and Lessons Learned from the world of J

SharePoint Development - "It takes a month to learn, but a lifet

(void *)

Dr. MonkeyIQ - C++, Linux, libferris and embedded devel

3D Application Development - using OpenCascade and .Net

Composing the Semantic Web - A tool developer's blog on ontolog

smuglispweeny - Grossly unfair, unreliable, biased, and pretty much delusional rants and ravings on

ORA-00001: Unique constraint violated - Random ramblings and raves of relational relevance!

SetCronJob Blog - Latest news from SetCronJob - Reliable Online Cronjobs Service

All In A Day's Work - Software Development, Running A Business

Gustavo J. A. M. Carneiro

Tips for software engineer - Tips that make life of a software engineer easier. Well, not exactly

Antipasto Hardware Blog - Hacking, Engineering, and Tech at Liquidware

Ethan's Software Blog - #{catchy_title.generate!}

Headius - Helping the JVM Into the 21st Century

Emilian Bold's blog - Broadcasting from the Romanian trenches

Holden's Blog - A Canadian developer in America.

all mine! - Personal thoughts on Linux desktops and distro's, community and marketing, open source

Wake Up And Dream!

Penguin Parens

Coffee, bits and bikes - My projects in the world of electronics, computers, bicycles and coffee dr

Web Science and Digital Libraries Research Group - Research and Teaching Updates from the

ClevelandDotNet - The ClevelandDotNet Blog is created to serve all .Net developers, both i

Ruminations

Analysis leads to Paralysis - A frustrated computer programmer's quest to observe things around him

Daily Websites - Daily websites, tips and tricks and more ...

Maruti

CafeNinja - All Things Geek and Caffeine

Agile Musings

evoluzination - Thoughts in the collective mind ... since 2006.

Scott Spendolini's Blog - Mostly Oracle APEX. Mostly.

Five One - The blog of one Jeremy Quinn esq. of Brixton.

Shindo Hikaru's Imagination

Nikon CLS Practical Guide - Practical use of the Nikon SB800 Speedlight with the D200 Camera

No Cropping Zone - Macro photography tips and tricks.

NIKON GLASS

clog - computer lab objectionable generalisations

Resetter Printer - Printer software for reset blink error codes canon epson hp

Photoshopper27 - Hundreds of Photoshop Tutorials for CS 6, 5, 4, 3, 2 and Photosh

Tour Source - Presale passwords and promo codes to your favorite concerts and show

F-Measure - Less Than Timely Music Reviews and Commentary.

neglected stairways - SO FAR, PINK FLOYD MONO, SOME SOFT MACHINE RARITIE

Free iTunes Store Downloads - Every week, Free iTunes Downloads

Virgin Media HD, TiVo & Netflix - Nialli's unofficial blog on V

CADDHELP - My personal blog where I try to keep friends and associates informed as to what is going

the kengell chronicles - a collection of thoughts culled from nearly th

Myxyama

Apple iPhone and iPod Blog - Welcome to the Apple iP

Computer Applications - A place to discuss c

Brief Thoughts, Briefly Written - When you b

learning communications blog - eLearning is a misunderstood domain. Bl

Diary of a Mild man - Ramblings on a theme of Exch

David Alison's Blog - I blog because ever

Life Revisited - This is the pla

Eastman's

Podunk's Valet - An eco

schav - you know...for kids! (tech startup

We Like Games - Here we follow up

GAMES AND GADGETS - Tech news, ga

Zero-Sum

BDPA Detroit TAC - Tech news,events,socia

mrc's Cup of Joe Blog - Save money, time, and increase business productivity

Early Indications - Early Indications is the weblog version of a newsl

The Pragmatic Geographer - Politics, Geography, Conservation. Whateve

Randy Rice's Software Testing & Quality Blog - Dedicated to thoughts about software testin

Andrew Eckford: The Blog - Not a real doctor since 2004.

Ninja Software Development - Developing software in the Real World is d

(YABAOAE) Yet Another Blog About Oracle Application Express

Information Security Thoughts - Allen Baranov - A blog dedicated to thoughts about Information Secu

Shep Hyken - Customer Service Speaker

anti-virus rants - devising a framework for thinking about malware and related issues such as virus

Mr.KinkyDink's UT Blog - UT99, IT Security, Miscellaneous Crap

Steve's Rant - Welcome to Steve's Rant Blog. This is a place where I will talk ab

Wolf's Lair - "A lair is also referred to as a hideout for a superhero or supervi

CompDocs - Computer Documents,Scripts, Tips for System Administrators an

Trying to become an Expert - My Ongoing CCDE Journey - An Engin

Dy-Verse - The technical blog of Parashuram

iTiNFO4U - Scroll Down

an enthusiastic austrian *nix/Linux user - Informations about

Techno Files - A Technophile's technobabble, includ

Electronic Error - Innovation is not the product of

Figure 2.1: jPEG Dendogram

# 3 K-CLUSTERING

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

## 3.1 SOLUTION

To get the k-clustering, I used the provided functions in clusters.py from the PCI textbook, however, I also printed the resulting cluster to a file, blog_k.txt. For K = 5, the code took three iterations, but when K = 10 and 20, it took four iterations. With five clusters, there was a group of 'tips and tricks' blogs, three groups of technology blogs, for gadgets and software development, and a group of hobbies - photoshop, photography, concerts, etc. When K was larger, the technology blogs were split up into smaller clusters, while the other clusters seemed to have stuck together. The python code used to accomplish this is in Listing 6 The screen output from the code was saved as k_run.txt

```python
def printClusters(kclust, out):
  for i in range(0,len(kclust)):
    klist = [blognames[r] for r in kclust[i]]
    for x in range(0, len(klist)):
      out.write(klist[x] + '\n')
    out.write('\n')

k = open('blog_k.txt', 'w', 0)
k.write('K = 5\n')
kclust=clusters.kcluster(data,k=5)
printClusters(kclust, k)
k.write('K = 10\n')
kclust=clusters.kcluster(data,k=10)
printClusters(kclust, k)
k.write('K = 20\n')
kclust=clusters.kcluster(data,k=20)
printClusters(kclust, k)
k.close()
```

Listing 6: Create K-clusters

# 4 MDS

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

## 4.1 SOLUTION

In order to create the JPEG, the provided code from the PCI textbook was used. The code used is shown in Listing 7 This section took 265 iterations. Although it wasn't required, the iteration numbers were saved as mds.txt, which also made it easier to count the iterations.

```python
old = sys.stdout
sys.stdout = open('blog_ascii.txt', 'w', 0)
blognames,words,data=clusters.readfile('blogdata.txt')
clust = clusters.hcluster(data)
clusters.printclust(clust,labels=blognames)
sys.stdout = old
# jpeg dendogram
clusters.drawdendrogram(clust,blognames,jpeg='blogclust.jpg')
```

Listing 7: Create MDS

Figure 4.1: MDS

# 5 TFIDF

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 SOLUTION

TF is calculated as the word frequency divided by the total words in the document. The IDF for each term was calculated as:

$$IDF(\text{term}) = log_2(\text{total docs in corpus/docs with term})$$

The total documents in corpus was used as the same as in Assignment Three, which was 42 billion. In order to get the total documents with the term, the python module requests was used to perform the Google search. Beautiful soup found the resultStats, which is where the number of results is stored. My first run was caught by Google's bot detection, so I added a random number sleep element. This took the program a little longer to run, but I was able to get all of the data. The IDF for each term was saved in a dictionary keyed on the term, in order to calculate the TFIDF.

Another blog matrix was created and saved as blog_tfidf.txt. For each word, the TF was calculated using the count from the wordcount dictionary for the term, which was by blog, multiplied by the length of the wordcount entry for the blog, which was the total number of words in the blog. Then TFIDF could be calculated as TF multiplied by IDF, and was rounded to three decimal points. [9]

The last step was to create the new dendrograms. It was similar to the previous dendrogram, if the blogs were not exactly next to each other after the TFIDF calculation, they were in close proximity. This tells me that the "hack" provided by the book was a pretty good estimation for the TFIDF for each blog. Figure 5.1 is the new JPEG dendrogram. Listing 8 is the python code used to accomplish this task.

```
1  google = 'http://www.google.com/search'
2  corp = 42000000000
3  idf = {}
4  for word in wordlist:
5    print 'Searching for ' + word
6    query_arg = {'q' : word}
7    sdata = urllib.urlencode(query_arg)
8    full_url = google+'?'+sdata
9    r = requests.get(full_url)
10   ddata = r.text
11   soup = BeautifulSoup(ddata)
12   results = soup.find('div', id='resultStats')
13   res = (results.text).rsplit()
14   term = float(res[1].replace(',',''))
15   idf[word] = log((corp/term),2)
16   sl = random.randint(1, 100)
17   time.sleep(sl)
18
19  ti = file('blog_tfidf.txt', 'w', 0)
20  ti.write('Blog')
21  for word in wordlist:
22    ti.write('\t%s' % word)
23
24  ti.write('\n')
25
26  for (blog, wc) in wordcounts.items():
```

```
27    if titles[blog] != '':
28      blog = blog + ' - ' + titles[blog]
29    blog = blog.replace(u'\u0144', 'n')
30    blog = _illegal_xml_chars_RE.sub(' ', blog)
31    blog = blog.strip()
32    if len(blog) > 100:
33      blog = blog[:99]
34    ti.write(blog)
35    if len(blog) > 100:
36      blog = blog[0:99]
37    for word in wordlist:
38      if word in wc:
39        tf = float(wc[word]) / len(wc)
40        tfidf = round(tf * idf[word], 3)
41        ti.write('\t' + str(tfidf))
42      else:
43        ti.write('\t0')
44    ti.write('\n')
45
46  old = sys.stdout
47  sys.stdout = open('blog_tfidf_ascii.txt', 'w', 0)
48  blognames,words,data=clusters.readfile('blog_tfidf.txt')
49  clust = clusters.hcluster(data)
50  clusters.printclust(clust,labels=blognames)
51  sys.stdout = old
52  clusters.drawdendrogram(clust,blognames,jpeg='tfidfclust.jpg')
```
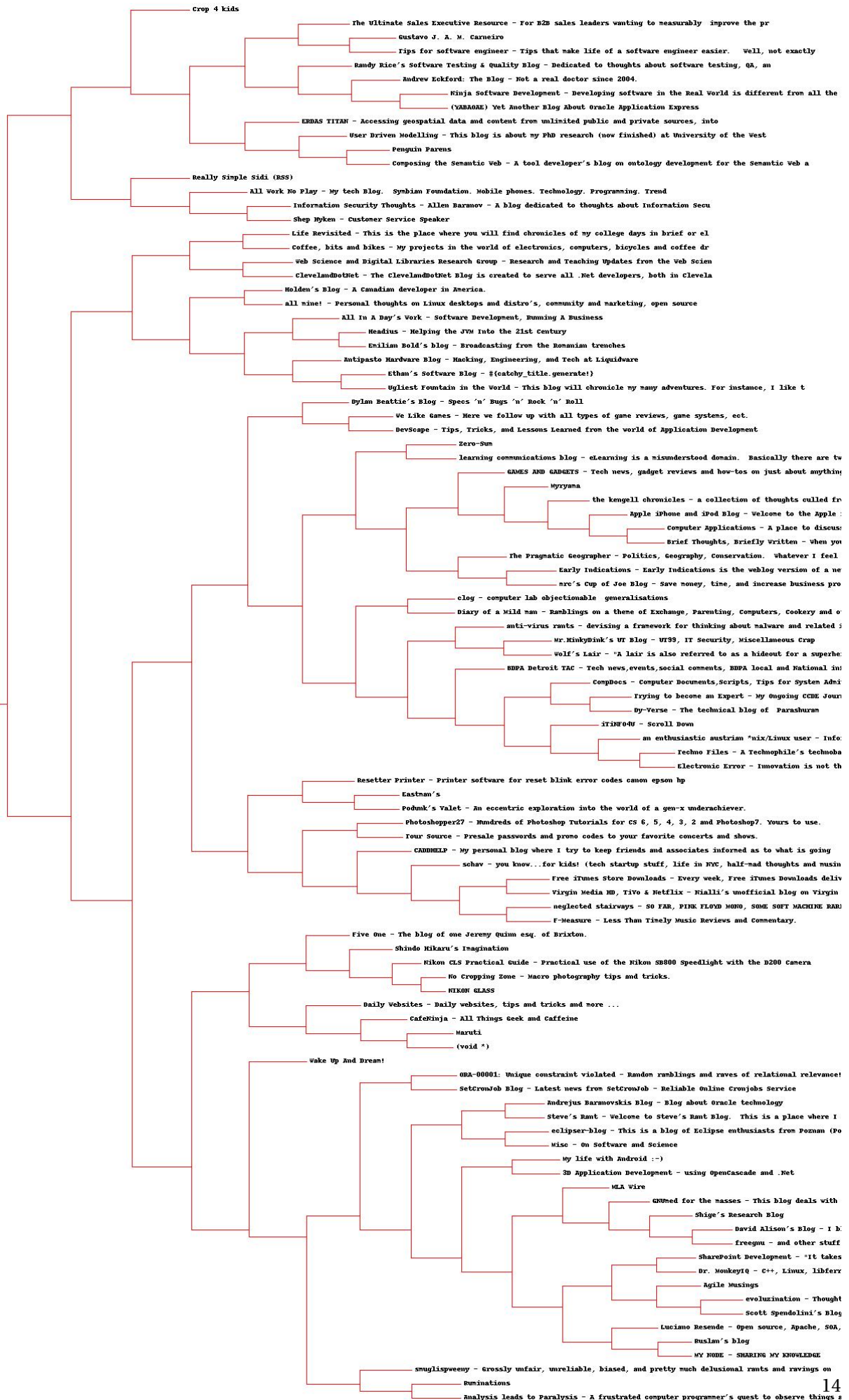
Listing 8: TFIDF Calculation

Crop 4 kids

The Ultimate Sales Executive Resource - For B2B sales leaders wanting to measurably  improve the pr
Gustavo J. A. M. Carneiro
Tips for software engineer - Tips that make life of a software engineer easier.   Well, not exactly
Randy Rice's Software Testing & Quality Blog - Dedicated to thoughts about software testing, QA, an
Andrew Eckford: The Blog - Not a real doctor since 2004.
Ninja Software Development - Developing software in the Real World is different from all the
(YABAOAE) Yet Another Blog About Oracle Application Express
ERDAS TITAN - Accessing geospatial data and content from unlimited public and private sources, into
User Driven Modelling - This blog is about my PhD research (now finished) at University of the West
Penguin Parens
Composing the Semantic Web - A tool developer's blog on ontology development for the Semantic Web a

Really Simple Sidi (RSS)
All Work No Play - My tech Blog.  Symbian Foundation. Mobile phones. Technology. Programming. Trend
Information Security Thoughts - Allen Baranov - A blog dedicated to thoughts about Information Secu
Shep Hyken - Customer Service Speaker
Life Revisited - This is the place where you will find chronicles of my college days in brief or el
Coffee, bits and bikes - My projects in the world of electronics, computers, bicycles and coffee dr
Web Science and Digital Libraries Research Group - Research and Teaching Updates from the Web Scien
ClevelandDotNet - The ClevelandDotNet Blog is created to serve all .Net developers, both in Clevela
Holden's Blog - A Canadian developer in America.
all mine! - Personal thoughts on Linux desktops and distro's, community and marketing, open source
All In A Day's Work - Software Development, Running A Business
Headius - Helping the JVM Into the 21st Century
Emilian Bold's blog - Broadcasting from the Romanian trenches
Antipasto Hardware Blog - Hacking, Engineering, and Tech at Liquidware
Ethan's Software Blog - #{catchy_title.generate!}
Ugliest Fountain in the World - This blog will chronicle my many adventures. For instance, I like t
Dylan Beattie's Blog - Specs 'n' Bugs 'n' Rock 'n' Roll
We Like Games - Here we follow up with all types of game reviews, game systems, ect.
DevScape - Tips, Tricks, and Lessons Learned from the world of Application Development
Zero-Sum
learning communications blog - eLearning is a misunderstood domain.  Basically there are tw
GAMES AND GADGETS - Tech news, gadget reviews and how-tos on just about anything
Myxyama
the kengell chronicles - a collection of thoughts culled fr
Apple iPhone and iPod Blog - Welcome to the Apple
Computer Applications - A place to discuss
Brief Thoughts, Briefly Written - When you
The Pragmatic Geographer - Politics, Geography, Conservation.  Whatever I feel
Early Indications - Early Indications is the weblog version of a ne
mrc's Cup of Joe Blog - Save money, time, and increase business pro
clog - computer lab objectionable  generalisations
Diary of a Wild man - Ramblings on a theme of Exchange, Parenting, Computers, Cookery and o
anti-virus rants - devising a framework for thinking about malware and related i
Mr.HinkyDink's VT Blog - VT99, IT Security, Miscellaneous Crap
Wolf's Lair - "A lair is also referred to as a hideout for a superhe
BDPA Detroit TAC - Tech news,events,social comments, BDPA local and National ini
CompDocs - Computer Documents,Scripts, Tips for System Admi
Trying to become an Expert - My Ongoing CCDE Jour
Dy-Verse - The technical blog of  Parashuram
iTiNFO4U - Scroll Down
an enthusiastic austrian *nix/Linux user - Info
Techno Files - A Technophile's technoba
Electronic Error - Innovation is not th
Resetter Printer - Printer software for reset blink error codes canon epson hp
Eastman's
Podunk's Valet - An eccentric exploration into the world of a gen-x underachiever.
Photoshopper27 - Hundreds of Photoshop Tutorials for CS 6, 5, 4, 3, 2 and Photoshop7. Yours to use.
Tour Source - Presale passwords and promo codes to your favorite concerts and shows.
CADDHELP - My personal blog where I try to keep friends and associates informed as to what is going
schav - you know...for kids! (tech startup stuff, life in NYC, half-mad thoughts and musin
Free iTunes Store Downloads - Every week, Free iTunes Downloads deliv
Virgin Media HD, TiVo & Netflix - Nialli's unofficial blog on Virgin
neglected stairways - SO FAR, PINK FLOYD MONO, SOME SOFT MACHINE RARI
F-Measure - Less Than Timely Music Reviews and Commentary.
Five One - The blog of one Jeremy Quinn esq. of Brixton.
Shindo Hikaru's Imagination
Nikon CLS Practical Guide - Practical use of the Nikon SB800 Speedlight with the D200 Camera
No Cropping Zone - Macro photography tips and tricks.
NIKON GLASS
Daily Websites - Daily websites, tips and tricks and more ...
CafeNinja - All Things Geek and Caffeine
Maruti
(void *)

Wake Up And Dream!
ORA-00001: Unique constraint violated - Random ramblings and raves of relational relevance!
SetCronJob Blog - Latest news from SetCronJob - Reliable Online Cronjobs Service
Andrejus Baranovskis Blog - Blog about Oracle technology
Steve's Rant - Welcome to Steve's Rant Blog.  This is a place where I
eclipser-blog - This is a blog of Eclipse enthusiasts from Poznan (Po
Misc - On Software and Science
My life with Android :-)
3D Application Development - using OpenCascade and .Net
MLA Wire
GNUmed for the masses - This blog deals with
Shige's Research Blog
David Alison's Blog - I b
freegnu - and other stuff
SharePoint Development - "It takes
Dr. MonkeyIQ - C++, Linux, libferr
Agile Musings
evoluzination - Thought
Scott Spendolini's Blog
Luciano Resende - Open source, Apache, SOA,
Ruslan's blog
MY NODE - SHARING MY KNOWLEDGE
smuglispweeny - Grossly unfair, unreliable, biased, and pretty much delusional rants and ravings on
Ruminations
Analysis leads to Paralysis - A frustrated computer programmer's quest to observe things a

14

Figure 5.1: TFIDF JPEG Dendogram

# 6 Python Code

```python
1   import requests
2   import pycurl
3   import re
4   from urlparse import urlparse
5   import urllib
6   import feedparser
7   from bs4 import BeautifulSoup
8   import clusters
9   import time
10  from math import log
11  import HTMLParser
12  try:
13      from io import BytesIO
14  except ImportError:
15      from StringIO import StringIO as BytesIO
16
17  def check(url):
18      with open("blogs.txt") as f:
19        found = False
20        for line in f:  #iterate over the file one line at a time(memory efficient)
21          if re.search(url, line):    #if string found is in current line then keep it
22            found = True
23      return found
24
25  blogger = 'https://www.blogger.com/next-blog?navBar=true&blogID=953024975153422094'
26  bfile = open('blogs.txt', 'w',0)
27  bfile.write('f-measure.blogspot.com\n')
28  bfile.write('ws-dl.blogspot.com\n')
29  bfile.close()
30  for i in range(1, 200):
31      buffer = BytesIO()
32      c = pycurl.Curl()
33      c.setopt(c.URL, blogger)
34      c.setopt(c.WRITEDATA, buffer)
35      c.setopt(c.FOLLOWLOCATION, True)
36      c.setopt(c.HTTPHEADER, ['Accept-Language: en'])
37      bfile = open('blogs.txt', 'a',0)
38      try:
39        c.perform()
40        if c.getinfo(c.RESPONSE_CODE) == 200:
41          o = urlparse(c.getinfo(c.EFFECTIVE_URL))
42          ch = check(o.netloc)
43          if not ch:
44            bfile.write(o.netloc + '\n')
45            bfile.close()
46          else:
47            continue
48        c.close()
49      except pycurl.error, error:
50        errno, errstr = error
51        print 'An error occurred: ', errstr
52
53
54
55  def getwordcounts(url):
56      #Returns title and dictionary of word counts for an RSS feed
57      # Parse the feed
58      d = feedparser.parse(url)
59      wc = {}
60      # Loop over all the entries
61      for e in d.entries:
62          if 'summary' in e:
63              summary = e.summary
64          else:
65              summary = e.description
66          # Extract a list of words
67          words = getwords(e.title + ' ' + summary)
68          for word in words:
69              wc.setdefault(word, 0)
70              wc[word] += 1
```

```python
71          return (d.feed.title, wc)
72
73  def getwords(html):
74          # Remove all the HTML tags
75          txt = re.compile(r'<[^>]+>').sub('', html)
76          # Split words by all non-alpha characters
77          words = re.compile(r'[^A-Z^a-z]+').split(txt)
78          # Convert to lowercase
79          return [word.lower() for word in words if word != '']
80
81  def getFeed(feedurl):
82          global wordcounts
83          global apcount
84          try:
85                  (title, wc) = getwordcounts(feedurl)
86                  if title in wordcounts:
87                          for (w, c) in wc.iteritems():
88                                  if w in wordcounts[title].iteritems():
89                                          wordcounts[title][w] += c
90                                  else:
91                                          wordcounts[title][w] = c
92                  else:
93                          wordcounts[title] = wc
94                  for (word, count) in wc.items():
95                          apcount.setdefault(word, 0)
96                          if count > 1:
97                                  apcount[word] += 1
98          except:
99                  print 'Failed to parse feed %s' % feedurl
100
101 pages = {}
102 feedlist = []
103 apcount = {}
104 wordcounts = {}
105 wordlist = []
106 scheme = 'http://'
107 path = '/feeds/posts/default'
108 query_arg = {'alt' : 'atom'}
109 udata = urllib.urlencode(query_arg)
110 f = open('blogs.txt', 'r', 0)
111 for line in f:
112   line = line.strip()
113   data = urllib.urlencode(query_arg)
114   full_url = scheme+line+path+"?"+udata
115   feedlist.append(full_url)
116   r = requests.get(full_url)
117   getFeed(full_url)
118   ddata = r.text
119   soup = BeautifulSoup(ddata)
120   next = soup.find('link', rel='next')
121   count = 1
122   while next:
123     n= next.get('href')
124     getFeed(n)
125     r = requests.get(n)
126     data = r.text
127     soup = BeautifulSoup(data)
128     count = count + 1
129     if count % 10 == 0:
130       print ' -- parsing ' + line + ' ' + str(count) + ' pages so far'
131     next = soup.find('link', rel='next')
132   print 'Finished parsing ' + line + ' ' + str(count) + ' pages'
133   if str(count) in pages:
134     pages[str(count)].append(line)
135   else:
136     pages[str(count)] = []
137     pages[str(count)].append(line)
138 f.close()
139
140 p = open('blog_pages.txt', 'w', 0)
141 p.write('Pages\tNumber of Blogs\n')
142 for (pg, num) in pages.iteritems():
143   p.write(pg + '\t' + str(len(num)) + '\n')
```

```python
144
145 a = sorted(apcount.items(), key=lambda x: x[1], reverse=True)
146 count = 0
147 for (w, bc) in a:
148     frac = float(bc) / len(feedlist)
149     if frac > 0.1 and frac < 0.5:
150         wordlist.append(w)
151         count = count + 1
152         if count >= 500:
153             break
154
155 _illegal_xml_chars_RE = re.compile(u'[\x00-\x08\x0b\xa9\x0c\x0e-\x1F\uD800-\uDFFF\uFFFE\uFFFF]')
156
157 def remove_tags(text):
158     cleanr =re.compile('<.*?>')
159     cleantext = re.sub(cleanr,'', text)
160     cleantext = re.sub("\n", ' ', cleantext)
161     cleantext = cleantext.strip()
162     return _illegal_xml_chars_RE.sub(' ', cleantext)
163
164 def getTitles():
165     titles = {}
166     global scheme, path, udata
167     f = open('blogs.txt', 'r', 0)
168     for line in f:
169         line = line.strip()
170         full_url = scheme+line+path+"?"+udata
171         d = feedparser.parse(full_url)
172         titles[d.feed.title] = remove_tags(d.feed.subtitle)
173     return titles
174
175 titles = {}
176 titles = getTitles()
177 # Blog Term Matrix
178 out = file('blogdata.txt', 'w', 0)
179 out.write('Blog')
180 for word in wordlist:
181     out.write('\t%s' % word)
182 out.write('\n')
183
184 for (blog, wc) in wordcounts.items():
185     if titles[blog] != '':
186         blog = blog + ' - ' + titles[blog]
187     blog = blog.replace(u'\u0144', 'n')
188     blog = _illegal_xml_chars_RE.sub(' ', blog)
189     blog = blog.strip()
190     if len(blog) > 100:
191         blog = blog[:99]
192     print blog
193     out.write(blog)
194     for word in wordlist:
195         if word in wc:
196             out.write('\t%d' % wc[word])
197         else:
198             out.write('\t0')
199     out.write('\n')
200 out.close()
201 # Ascii
202 old = sys.stdout
203 sys.stdout = open('blog_ascii.txt', 'w', 0)
204 blognames,words,data=clusters.readfile('blogdata.txt')
205 clust = clusters.hcluster(data)
206 clusters.printclust(clust,labels=blognames)
207 sys.stdout = old
208 # jpeg dendogram
209 clusters.drawdendrogram(clust,blognames,jpeg='blogclust.jpg')
210 # K-clustering
211 def printClusters(kclust, out):
212     for i in range(0,len(kclust)):
213         klist = [blognames[r] for r in kclust[i]]
214         for x in range(0, len(klist)):
215             out.write(klist[x] + '\n')
216         out.write('\n')
```

```
217
218 k = open('blog_k.txt', 'w', 0)
219 k.write('K = 5\n')
220 kclust=clusters.kcluster(data,k=5)
221 printClusters(kclust, k)
222 k.write('K = 10\n')
223 kclust=clusters.kcluster(data,k=10)
224 printClusters(kclust, k)
225 k.write('K = 20\n')
226 kclust=clusters.kcluster(data,k=20)
227 printClusters(kclust, k)
228 k.close()
229 # MDS
230 coords = clusters.scaledown(data)
231 clusters.draw2d(coords,blognames,jpeg='blogs2d.jpg')
232 #
233 #Extra credit - TFIDF
234 #
235 google = 'http://www.google.com/search'
236 corp = 42000000000
237 idf = {}
238 for word in wordlist:
239     print 'Searching for ' + word
240     query_arg = {'q' : word}
241     sdata = urllib.urlencode(query_arg)
242     full_url = google+'?'+sdata
243     r = requests.get(full_url)
244     ddata = r.text
245     soup = BeautifulSoup(ddata)
246     results = soup.find('div', id='resultStats')
247     res = (results.text).rsplit()
248     term = float(res[1].replace(',',''))
249     idf[word] = log((corp/term),2)
250     sl = random.randint(1, 100)
251     time.sleep(sl)
252
253 ti = file('blog_tfidf.txt', 'w', 0)
254 ti.write('Blog')
255 for word in wordlist:
256     ti.write('\t%s' % word)
257
258 ti.write('\n')
259
260 for (blog, wc) in wordcounts.items():
261     if titles[blog] != '':
262         blog = blog + ' - ' + titles[blog]
263     blog = blog.replace(u'\u0144', 'n')
264     blog = _illegal_xml_chars_RE.sub(' ', blog)
265     blog = blog.strip()
266     if len(blog) > 100:
267         blog = blog[:99]
268     ti.write(blog)
269     if len(blog) > 100:
270         blog = blog[0:99]
271     for word in wordlist:
272         if word in wc:
273             tf = float(wc[word]) / len(wc)
274             tfidf = round(tf * idf[word], 3)
275             ti.write('\t' + str(tfidf))
276         else:
277             ti.write('\t0')
278     ti.write('\n')
279
280 old = sys.stdout
281 sys.stdout = open('blog_tfidf_ascii.txt', 'w', 0)
282 blognames,words,data=clusters.readfile('blog_tfidf.txt')
283 clust = clusters.hcluster(data)
284 clusters.printclust(clust,labels=blognames)
285 sys.stdout = old
286 clusters.drawdendrogram(clust,blognames,jpeg='tfidfclust.jpg')
```

Listing 9: Complete Feed Parser

## REFERENCES

[1] Alan Gauld. how to redirect sys.stdout back to screen in IDLE. `https://mail.python.org/pipermail/tutor/2002-November/018792.html`. Accessed: 2014-11-20.

[2] Ashwini Chaudhary. Python: Checking if string is in a text file. `http://stackoverflow.com/questions/16432203/python-checking-if-string-is-in-a-text-file`. Accessed: 2014-9-20.

[3] Beautiful soup documentation. `http://www.crummy.com/software/BeautifulSoup/bs4/doc/`. Accessed: 2014-10-03.

[4] c24b. Python code to remove HTML tags from a string. `http://stackoverflow.com/questions/9662346/python-code-to-remove-html-tags-from-a-string`. Accessed: 2014-11-20.

[5] Kurt McKee. Common Atom Elements. `https://pythonhosted.org/feedparser/common-atom-elements.html`. Accessed: 2014-11-20.

[6] Leo Simons. stripping illegal characters out of xml in python. `http://lsimons.wordpress.com/2011/03/17/stripping-illegal-characters-out-of-xml-in-python/`. Accessed: 2014-11-20.

[7] Philip Southam. 'too many values to unpack', iterating over a dict. key=>string, value=>list. `http://stackoverflow.com/questions/5466618/too-many-values-to-unpack-iterating-over-a-dict-key-string-value-list`. Accessed: 2014-11-20.

[8] HTTP with Python âĂŞ PycURL by Example. `http://www.angryobjects.com/2011/10/15/http-with-python-pycurl-by-example/`. Accessed: 2014-11-20.

[9] Python Software Foundation. Built-in Functions. `https://docs.python.org/2/library/functions.html`. Accessed: 2014-11-20.

[10] Toby Segaran. *Programming Collective Intelligence*. O'Reilly Media, 2007.

[11] tutorialspoint. Python Strings. `http://www.tutorialspoint.com/python/python_strings.htm`. Accessed: 2014-11-20.