Srividya Majeti

# Assignment 8

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

April 7, 2016

# Contents

# Question 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.**

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

- First I grabbed 100 URIs using the method discussed in class. I stored all the URIs in a set data-structure which has a property of storing only unique values, this is illustrated in the function 'getUrl'. In function 'writeUrlToFile' I have written the 2 URIs mentioned in the question and all the grabbed URIs to a file 'get100URls' . This code is listed in Listing 1.1.

- Later I got the atoms for each blog URI by appending 'feeds/posts/default?max-results=500'. I wrote the atom URIs to a file 'getAtomsFor100Urls'. This function is listed in Listing 1.2.
- I downloaded the 'generateFeedVector.py' from PCI book code and ran it for the atom URIs to get the blog data with word count and words. But I realized that few blogs have more than one page.
- The python code listed in Listing 1.3 results a text file 'noOfPagesForEachBlog.txt' with number of pages and title of the blog.
- To solve the problem of multiple pages I made few modifications to 'generateFeedVector.py'. First I requested each URI and recursively checked if the blog has 'rel="next"' using the library BeautifulSoup and got the links for all the pages for each blog. I stored the resulted links for respective blogs in a list 'blogUrlList'. This is illustrated in function 'getBlogPagesURLs'.
- If a blog has more than 1 page then I ran the 'generateFeedVector' function for the first page and stored the word count in 'wc', then ran the function for the URIs in the 'blogUrlList' and stored the word count in 'nextwc'. Futhermore I consolidated both the dictionaries 'wc' and 'nextwc'. This code is listed in Listing 1.4
- Therefore I generated the output file 'blogdata.txt' which has a blog matrix with blog title as identifier for each blog. This text file is uploaded to github at `https://github.com/majetisiri/cs532-s16/blob/master/a8/q1-blogdata.txt`

**Code Listing**

```
1  def get100BlogUrl():
2      link = "http://www.blogger.com/next-blog?navBar=true&
           blogID=3471633091411211117"
3      set = Set()
4
5      while len(set) < 100:
6          r = requests.get(link, allow_redirects=True)
7          uri= r.url
8
9          if len(uri) > 0:
10             uri = uri.lower()
11
12             parsedUrl = urlparse.urlparse(uri)
13             parsedUrl = parsedUrl.scheme + '://' + parsedUrl.
                   netloc + '/'
14
15             set.add(parsedUrl)
16             print len(set)
17             print parsedUrl
18     return set
19
20 def writeUrlToFile(data):
21     file= open('get100Urls','w')
22     file.write('http://f-measure.blogspot.com/'+'\n')
23     file.write('http://ws-dl.blogspot.com/'+'\n')
24
25     for item in data:
26         file.write(item+'\n')
```

**Listing 1.1.** Function for getting 100 unique blog URIs

**Code Listing**

```
1  def getAtoms():
2      file= open('get100Urls','r')
3      f1= open('getAtomsFor100Urls','w')
4      add= "feeds/posts/default?max-results=500"
5      for uri in file:
6          uri= uri.strip()+add
7          f1.write(uri+"\n")
```

**Listing 1.2.** Function for getting atom URIs

### Code Listing

```
 1  import os
 2  import sys
 3  import urllib
 4  import time
 5  import feedparser
 6
 7  from bs4 import BeautifulSoup
 8
 9  def checkNextPage(url):
10
11      f = urllib.urlopen(url)
12      soup = BeautifulSoup(f.read(), from_encoding=f.info().
            getparam('charset'))
13
14      try:
15          link = soup.find('link', rel='next',href = True)['
                href']
16      except TypeError:
17          link = None
18      return link
19
20
21  def getPages():
22      feedlist    = open('getAtomsFor100Urls').readlines()
23
24      for url in feedlist:
25      d = feedparser.parse(url)
26      title = d['feed']['title']
27      count    = 1
28      nextLink = checkNextPage( url  )
29
30      while nextLink:
31        nextLink = checkNextPage( nextLink )
32        count += 1
33
34      print u'|'.join((str(count),title)).encode('utf-8').
            strip()
35
36  getPages()
```

**Listing 1.3.** Python code for getting number of pages for each blog

**Code Listing**

```python
def getBlogPagesURLs(url, blogUrlList=[]):
  req = requests.get(url)
  soup = BeautifulSoup(req.text)
  nextLink = soup.find('link', rel='next',href = True)
  if nextLink is not None:
    nextLink = nextLink['href']
    blogUrlList.append(nextLink)
    getBlogPagesURLs(nextLink, blogUrlList)
  return blogUrlList

def getwordcounts(url):
  d=feedparser.parse(url)
  wc={}

  for e in d.entries:
    if 'summary' in e:
      summary=e.summary
    else:
      summary=e.description

    words=getwords(e.title+' '+summary)
    for word in words:
      wc.setdefault(word,0)
      wc[word]+=1
  print d.feed.title
  return d.feed.title,wc

def getwords(html):
  txt=re.compile(r'<[^>]+>').sub('',html)
  words=re.compile(r'[^A-Z^a-z]+').split(txt)
  return [word.lower() for word in words if word!='']

def combineWC(wc, nextwc):
  if len(wc)>0 and len(nextwc)>0:
    for word, wordcount in nextwc.items():
      if word in wc:
        wc[word] = wc[word] + wordcount
      else:
        wc[word] = wordcount
    return wc
  else:
    return {}

def generateFeedVector():
  apcount={}
  wordcounts={}
  feedlist=[line for line in open('getAtomsFor100Urls')]
```

```
48    for feedurl in feedlist:
49      try:
50        blogUrlList = getBlogPagesURLs(feedurl)
51        title ,wc=getwordcounts(feedurl)
52        for url in blogUrlList:
53          title ,nextwc=getwordcounts(feedurl)
54          combineWC(wc,nextwc)
55        wordcounts[title]=wc
56        for word,count in wc.items():
57          apcount.setdefault(word,0)
58          if count >1:
59            apcount[word]+=1
60      except:
61        print 'Failed to parse feed %s' % feedurl
62
63    wordlist=[]
64    countFrequentWords=[]
65    for w,bc in apcount.items():
66      frac=float(bc)/len(feedlist)
67      if frac >0.1 and frac <0.5:
68        countFrequentWords.append((w,bc))
69
70    countFrequentWords=sorted(countFrequentWords,key=lambda x:
          x[1], reverse = True)
71
72    for value in countFrequentWords:
73      value1 = value[0]
74      value2 = value[1]
75      length = len(wordlist)
76      if(length < 500):
77        wordlist.append(value1)
78      else:
79        break
80
81    stop_words_list = [line.rstrip('\r\n') for line in open('
          stopWordList.txt')]
82
83    out=file('blogdata.txt','w')
84    out.write('Blog')
85    for word in wordlist:
86      word1 = word.encode('UTF-8')
87      out.write('\t%s' % word1)
88    out.write('\n')
89    for blog,wc in wordcounts.items():
90      blogName = blog.encode('UTF-8')
91      print blog
92      print blogName
93      out.write(blogName)
94      for word in wordlist:
```

```python
95              if word not in stop_words_list:
96              if word in wc:
97                out.write('\t%d' % wc[word])
98              else:
99                out.write('\t0')
100       out.write('\n')
```

**Listing 1.4.** Function for getting feed vector for each blog

# Question 2

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

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

- I downloaded the python code 'clusters.py' from the 'Programming Collective Intelligence' book by 'Toby Segaran'. I used this for questions 2, 3 and 4.
- I imported the 'clusters.py' and used the code described in 'presentation slide 12' to create an ASCII that clusters the most similar blogs. This code is in Listing 2.1

- The output ascii file is uploaded to github at `https://github.com/majetisiri/cs532-s16/blob/master/a8/q2-AsciiOutput.txt`. The sample output is illustrated in Figure 2.1.



**Fig. 2.1.** Sample ascii output

- Furthermore to get the JPEG dendogram I used 'clusters.py' and the code from 'presentation slide 13'. This code is in Listing 2.2

- The output JPEG that clusters the most similar blogs is illustrated in the Figure 2.2
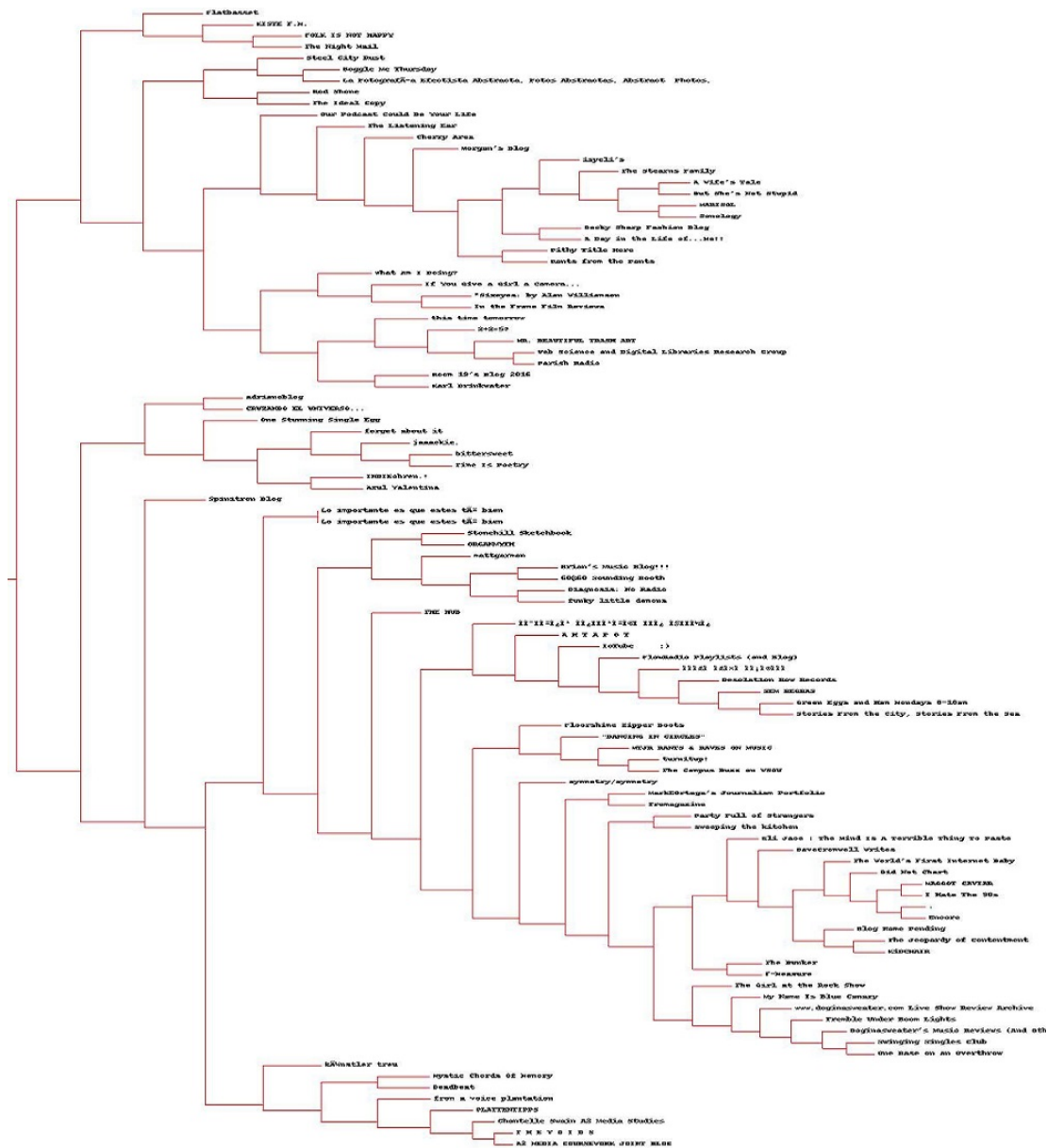


**Fig. 2.2.** JPEG dendogram

### Code Listing

```python
import clusters

def generateAscii():
    blognames,words,data=clusters.readfile('blogdata.txt')
    clust=clusters.hcluster(data)
    clusters.printclust(clust,labels=blognames)

generateAscii()
```

**Listing 2.1.** Python code for generating ASCII

### Code Listing

```python
import clusters

def drawDendogram():
    blognames,words,data=clusters.readfile('blogdata.txt')
    clust=clusters.hcluster(data)
    clusters.drawdendrogram(clust,blognames,jpeg='blogclust.jpg')

drawDendogram()
```

**Listing 2.2.** Python code for generating JPEG dendogram

# 3

## Question 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 interations were required for each value of k?**

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

- I imported the 'clusters.py' mentioned in 'question 2' and used the code described in 'presentation slide 18' to cluster the blogs using K-Means, using k= 5, 10, 20. The output prints the values in each centroid, for each value of k and also the number of iterations required for each value of k. This code is in Listing 3.1

- The output file is uploaded to github at `https://github.com/majetisiri/cs532-s16/blob/master/a8/q3-numberOfIterationsAndValuesInEachCentroid.txt`. The sample output is illustrated in Figure 3.1.

```
K value is 5
    Iteration 0
    Iteration 1
    Iteration 2
    Iteration 3
    Iteration 4
    Iteration 5
    Iteration 6
    Iteration 7
    ['Flatbasset', 'mattgarman', "Brian's Music Blog!!!", '60@60 Sounding Booth', 'The Ideal Copy', 'Rod Shone', 'The Campus Buzz on WSOU',
    'k\xc3\xbcnstler treu', 'Time Is Poetry', '\xce\x94\xce\xaf\xcf\x83\xce\xba\xce\xbf\xce\xb9
    \xce\x9c\xce\xbf\xcf\x85\xcf\x83\xce\xb9\xce\xba\xce\xae\xcf\x82 \xcf\x83\xcf\x84\xce\xbf \xce\xa7\xcf\x81\xcf\x8c\xce\xbd\xce\xbf']
    ['MTJR RANTS & RAVES ON MUSIC', 'turnitup!', 'Floorshime Zipper Boots', 'Did Not Chart', 'DaveCromwell Writes', 'Swinging Singles Club', 'FOLK IS NOT
    HAPPY', 'MAGGOT CAVIAR', 'Spinitron Blog', "The World's First Internet Baby", 'Tremble Under Boom Lights', 'www.doginasweater.com Live Show Review
    Archive', 'One Base on an Overthrow', "Doginasweater's Music Reviews (And Other Horseshit)", 'Eli Jace | The Mind Is A Terrible Thing To Paste', 'The
    Jeopardy of Contentment', 'Boggle Me Thursday', '.', 'KiDCHAIR', '*Sixeyes: by Alan Williamson', 'Blog Name Pending', 'Our Podcast Could Be Your
    Life', 'F-Measure', 'I Hate The 90s']
    ['SEM REGRAS', 'MARISOL', 'THE HUB', 'MR. BEAUTIFUL TRASH ART', 'Green Eggs and Ham Mondays 8-10am', 'Stories From the City, Stories From the Sea', 'A
    H T A P O T', 'adrianoblog', 'CRUZANDO EL UNIVERSO...', 'FlowRadio Playlists (and Blog)', 'INDIEohren.!', 'Desolation Row Records', 'IoTube       :)',
    'Parish Radio', 'If You Give a Girl a Camera...', 'Lo importante es que estes t\xc3\xba bien', 'La Fotograf\xc3\xada Efectista Abstracta. Fotos
    Abstractas. Abstract  Photos.', 'sweeping the kitchen', 'One Stunning Single Egg', 'this time tomorrow', '\xce\x9c\xce\x95\xce\xa3\xce\x91
    \xce\xa3\xce\xa4\xce\x97 \xce\x92\xce\xa1\xce\xa9\xce\x9c\xce\x99\xce\x91', 'KISTE F.M.']
    ['Party Full of Strangers', "MarkEOrtega's Journalism Portfolio", 'jaaackie.', 'The Girl at the Rock Show', 'Mystic Chords Of Memory', 'PLATTENTIPPS',
    'Tremagazine', 'symmetry/symmetry', 'bittersweet', 'Encore']
    ["Riley Haas' blog", 'Pithy Title Here', 'Web Science and Digital Libraries Research Group', 'Steel City Rust', 'ORGANMYTH', 'Diagnosis: No Radio',
    'funky little demons', '2+2=5?', 'Stonehill Sketchbook', 'forget about it', 'T H E V O I D S', 'Chantelle Swain A2 Media Studies', 'A2 MEDIA
    COURSEWORK JOINT BLOG', 'The Listening Ear', "Morgan's Blog", 'My Name Is Blue Canary', 'The Bunker', 'Deadbeat', 'Becky Sharp Fashion Blog', 'from a
    voice plantation', 'The Stearns Family', 'Sonology', 'The Night Mail', 'What Am I Doing?', 'Rants from the Pants', 'A Day in the Life of...Me!!',
    "Room 19's Blog 2016", "A Wife's Tale", 'Karl Drinkwater', 'In the Frame Film Reviews', 'Cherry Area', 'Azul Valentina', "isyeli's", "But She's Not
    Stupid"]
```

**Fig. 3.1.** Sample output with number of iterations and values in each centroid for k=5

**Code Listing**

```
1   import clusters
2   def getKmeans():
3     blognames,words,data=clusters.readfile('blogdata.txt')
4     print "K value is 5"
5     kclust=clusters.kcluster(data,k=5)
6     print "\t\t"+str([blognames[r] for r in kclust[0]])
7     print "\t\t"+str([blognames[r] for r in kclust[1]])
8     print "\t\t"+str([blognames[r] for r in kclust[2]])
9     print "\t\t"+str([blognames[r] for r in kclust[3]])
10    print "\t\t"+str([blognames[r] for r in kclust[4]])
11    print "K value is 10"
12    kclust=clusters.kcluster(data,k=10)
13    print "\t\t"+str([blognames[r] for r in kclust[0]])
14    print "\t\t"+str([blognames[r] for r in kclust[1]])
15    print "\t\t"+str([blognames[r] for r in kclust[2]])
16    print "\t\t"+str([blognames[r] for r in kclust[3]])
17    print "\t\t"+str([blognames[r] for r in kclust[4]])
18    print "\t\t"+str([blognames[r] for r in kclust[5]])
19    print "\t\t"+str([blognames[r] for r in kclust[6]])
20    print "\t\t"+str([blognames[r] for r in kclust[7]])
21    print "\t\t"+str([blognames[r] for r in kclust[8]])
22    print "\t\t"+str([blognames[r] for r in kclust[9]])
23    print "K value is 20"
24    kclust=clusters.kcluster(data,k=20)
25    print "\t\t"+str([blognames[r] for r in kclust[0]])
26    print "\t\t"+str([blognames[r] for r in kclust[1]])
27    print "\t\t"+str([blognames[r] for r in kclust[2]])
28    print "\t\t"+str([blognames[r] for r in kclust[3]])
29    print "\t\t"+str([blognames[r] for r in kclust[4]])
30    print "\t\t"+str([blognames[r] for r in kclust[5]])
31    print "\t\t"+str([blognames[r] for r in kclust[6]])
32    print "\t\t"+str([blognames[r] for r in kclust[7]])
33    print "\t\t"+str([blognames[r] for r in kclust[8]])
34    print "\t\t"+str([blognames[r] for r in kclust[9]])
35    print "\t\t"+str([blognames[r] for r in kclust[10]])
36    print "\t\t"+str([blognames[r] for r in kclust[11]])
37    print "\t\t"+str([blognames[r] for r in kclust[12]])
38    print "\t\t"+str([blognames[r] for r in kclust[13]])
39    print "\t\t"+str([blognames[r] for r in kclust[14]])
40    print "\t\t"+str([blognames[r] for r in kclust[15]])
41    print "\t\t"+str([blognames[r] for r in kclust[16]])
42    print "\t\t"+str([blognames[r] for r in kclust[17]])
43    print "\t\t"+str([blognames[r] for r in kclust[18]])
44    print "\t\t"+str([blognames[r] for r in kclust[19]])
45  getKmeans()
```

**Listing 3.1.** Python code for generating ASCII

# 4

# Question 4

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

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

- I imported the 'clusters.py' mentioned in 'question 2' and used the code described in 'presentation slide 29' to create a JPEG of the most similar blogs using MDS. This code is in Listing 4.1

- The output JPEG file is illustrated in 4.1.

**Fig. 4.1.** JPEG of blogs using MDS

- To get the number of iterations I have written a print statement in function 'scaledown(data,distance=pearson,rate=0.01)' of 'clusters.py'. '304' iterations were required for creating the JPEG using MDS. The file 'numberOfIterations.txt' has 'total error' and 'iteration count'.
- The python code 'clusters.py' that I downloaded from the PCI book is illustrated in Listing 4.2

**Code Listing**

```
1  import clusters
2
3  def createMDS():
4
5      blognames,words,data=clusters.readfile('blogdata.txt')
6      coords=clusters.scaledown(data)
7      clusters.draw2d(coords,blognames,jpeg='blogs2d.jpg')
8
9  createMDS()
```

**Listing 4.1.** Python code for creating MDS

**Code Listing**

```
1  from PIL import Image,ImageDraw
2
3  def readfile(filename):
4    lines=[line for line in file(filename)]
5
6    colnames=lines[0].strip().split('\t')[1:]
7    rownames=[]
8    data=[]
9    for line in lines[1:]:
10     p=line.strip().split('\t')
11     rownames.append(p[0])
12     data.append([float(x) for x in p[1:]])
13   return rownames,colnames,data
14
15
16  from math import sqrt
17
18  def pearson(v1,v2):
19    sum1=sum(v1)
20    sum2=sum(v2)
21    sum1Sq=sum([pow(v,2) for v in v1])
22    sum2Sq=sum([pow(v,2) for v in v2])
23    pSum=sum([v1[i]*v2[i] for i in range(len(v1))])
24    num=pSum-(sum1*sum2/len(v1))
25    den=sqrt((sum1Sq-pow(sum1,2)/len(v1))*(sum2Sq-pow(sum2,2)/
            len(v1)))
26    if den==0: return 0
```

```
27
28    return 1.0−num/den
29
30 class bicluster:
31    def __init__(self,vec,left=None,right=None,distance=0.0,id
          =None):
32      self.left=left
33      self.right=right
34      self.vec=vec
35      self.id=id
36      self.distance=distance
37
38 def hcluster(rows,distance=pearson):
39    distances={}
40    currentclustid=−1
41    clust=[bicluster(rows[i],id=i) for i in range(len(rows))]
42
43    while len(clust)>1:
44      lowestpair=(0,1)
45      closest=distance(clust[0].vec,clust[1].vec)
46
47      for i in range(len(clust)):
48        for j in range(i+1,len(clust)):
49          if (clust[i].id,clust[j].id) not in distances:
50            distances[(clust[i].id,clust[j].id)]=distance(
                  clust[i].vec,clust[j].vec)
51
52          d=distances[(clust[i].id,clust[j].id)]
53
54          if d<closest:
55            closest=d
56            lowestpair=(i,j)
57
58      mergevec=[
59      (clust[lowestpair[0]].vec[i]+clust[lowestpair[1]].vec[i
          ])/2.0
60      for i in range(len(clust[0].vec))]
61
62      newcluster=bicluster(mergevec,left=clust[lowestpair[0]],
63                           right=clust[lowestpair[1]],
64                           distance=closest,id=currentclustid)
65
66      currentclustid−=1
67      del clust[lowestpair[1]]
68      del clust[lowestpair[0]]
69      clust.append(newcluster)
70
71    return clust[0]
72
```

```python
73  def printclust(clust,labels=None,n=0):
74    for i in range(n): print ' ',
75    if clust.id<0:
76      print '-'
77    else:
78      if labels==None: print clust.id
79      else: print labels[clust.id]
80
81    if clust.left!=None: printclust(clust.left,labels=labels,n
          =n+1)
82    if clust.right!=None: printclust(clust.right,labels=labels
          ,n=n+1)
83
84  def getheight(clust):
85    if clust.left==None and clust.right==None: return 1
86    return getheight(clust.left)+getheight(clust.right)
87
88  def getdepth(clust):
89    if clust.left==None and clust.right==None: return 0
90    return max(getdepth(clust.left),getdepth(clust.right))+
          clust.distance
91
92
93  def drawdendrogram(clust,labels,jpeg='clusters.jpg'):
94    h=getheight(clust)*20
95    w=1200
96    depth=getdepth(clust)
97
98    scaling=float(w-150)/depth
99
100   img=Image.new('RGB',(w,h),(255,255,255))
101   draw=ImageDraw.Draw(img)
102
103   draw.line((0,h/2,10,h/2),fill=(255,0,0))
104
105   drawnode(draw,clust,10,(h/2),scaling,labels)
106   img.save(jpeg,'JPEG')
107
108  def drawnode(draw,clust,x,y,scaling,labels):
109    if clust.id<0:
110      h1=getheight(clust.left)*20
111      h2=getheight(clust.right)*20
112      top=y-(h1+h2)/2
113      bottom=y+(h1+h2)/2
114      ll=clust.distance*scaling
115      draw.line((x,top+h1/2,x,bottom-h2/2),fill=(255,0,0))
116      draw.line((x,top+h1/2,x+ll,top+h1/2),fill=(255,0,0))
117      draw.line((x,bottom-h2/2,x+ll,bottom-h2/2),fill
            =(255,0,0))
```

```
118        drawnode(draw,clust.left,x+ll,top+h1/2,scaling,labels)
119        drawnode(draw,clust.right,x+ll,bottom−h2/2,scaling,
              labels)
120      else:
121        draw.text((x+5,y−7),labels[clust.id],(0,0,0))
122
123  def rotatematrix(data):
124    newdata=[]
125    for i in range(len(data[0])):
126      newrow=[data[j][i] for j in range(len(data))]
127      newdata.append(newrow)
128    return newdata
129
130  import random
131
132  def kcluster(rows,distance=pearson,k=4):
133    ranges=[(min([row[i] for row in rows]),max([row[i] for row
            in rows]))
134    for i in range(len(rows[0]))]
135
136    clusters=[[random.random()*(ranges[i][1]−ranges[i][0])+
            ranges[i][0]
137    for i in range(len(rows[0]))] for j in range(k)]
138
139    lastmatches=None
140    for t in range(100):
141      print 'Iteration %d' % t
142      bestmatches=[[] for i in range(k)]
143      for j in range(len(rows)):
144        row=rows[j]
145        bestmatch=0
146        for i in range(k):
147          d=distance(clusters[i],row)
148          if d<distance(clusters[bestmatch],row): bestmatch=i
149        bestmatches[bestmatch].append(j)
150      if bestmatches==lastmatches: break
151      lastmatches=bestmatches
152
153      for i in range(k):
154        avgs=[0.0]*len(rows[0])
155        if len(bestmatches[i])>0:
156          for rowid in bestmatches[i]:
157            for m in range(len(rows[rowid])):
158              avgs[m]+=rows[rowid][m]
159          for j in range(len(avgs)):
160            avgs[j]/=len(bestmatches[i])
161          clusters[i]=avgs
162
163    return bestmatches
```

```
164
165  def tanamoto(v1,v2):
166    c1,c2,shr=0,0,0
167
168    for i in range(len(v1)):
169      if v1[i]!=0: c1+=1
170      if v2[i]!=0: c2+=1
171      if v1[i]!=0 and v2[i]!=0: shr+=1
172
173    return 1.0-(float(shr)/(c1+c2-shr))
174
175  def scaledown(data,distance=pearson,rate=0.01):
176    n=len(data)
177
178    realdist=[[distance(data[i],data[j]) for j in range(n)]
179               for i in range(0,n)]
180
181    loc=[[random.random(),random.random()] for i in range(n)]
182    fakedist=[[0.0 for j in range(n)] for i in range(n)]
183
184    lasterror=None
185    for m in range(0,1000):
186      for i in range(n):
187        for j in range(n):
188          fakedist[i][j]=sqrt(sum([pow(loc[i][x]-loc[j][x],2)
189                                    for x in range(len(loc[i]))
190                                    ]))
191      grad=[[0.0,0.0] for i in range(n)]
192
193      totalerror=0
194      for k in range(n):
195        for j in range(n):
196          if j==k: continue
197          errorterm=(fakedist[j][k]-realdist[j][k])/realdist[j][k]
198          grad[k][0]+=((loc[k][0]-loc[j][0])/fakedist[j][k])*
199                         errorterm
200          grad[k][1]+=((loc[k][1]-loc[j][1])/fakedist[j][k])*
201                         errorterm
202          totalerror+=abs(errorterm)
203      print totalerror
204
205      if lasterror and lasterror<totalerror: break
206      lasterror=totalerror
207
208      for k in range(n):
209        loc[k][0]-=rate*grad[k][0]
210        loc[k][1]-=rate*grad[k][1]
```

```
209
210      return loc
211
212  def draw2d(data, labels, jpeg='mds2d.jpg'):
213      img=Image.new('RGB',(2000,2000),(255,255,255))
214      draw=ImageDraw.Draw(img)
215      for i in range(len(data)):
216        x=(data[i][0]+0.5)*1000
217        y=(data[i][1]+0.5)*1000
218        draw.text((x,y),labels[i],(0,0,0))
219      img.save(jpeg,'JPEG')
```

**Listing 4.2.** Python code 'clusters.py' from PCI

# 5

## Extra-Credit Question-5

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.**

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

- I calculated the TFIDF values similar to assignment 3.
- I got the JPEG dendogram using the same python code as question 2 5.1

- The output JPEG that clusters the most similar blogs based on TFIDF values is illustrated in the Figure 5.1
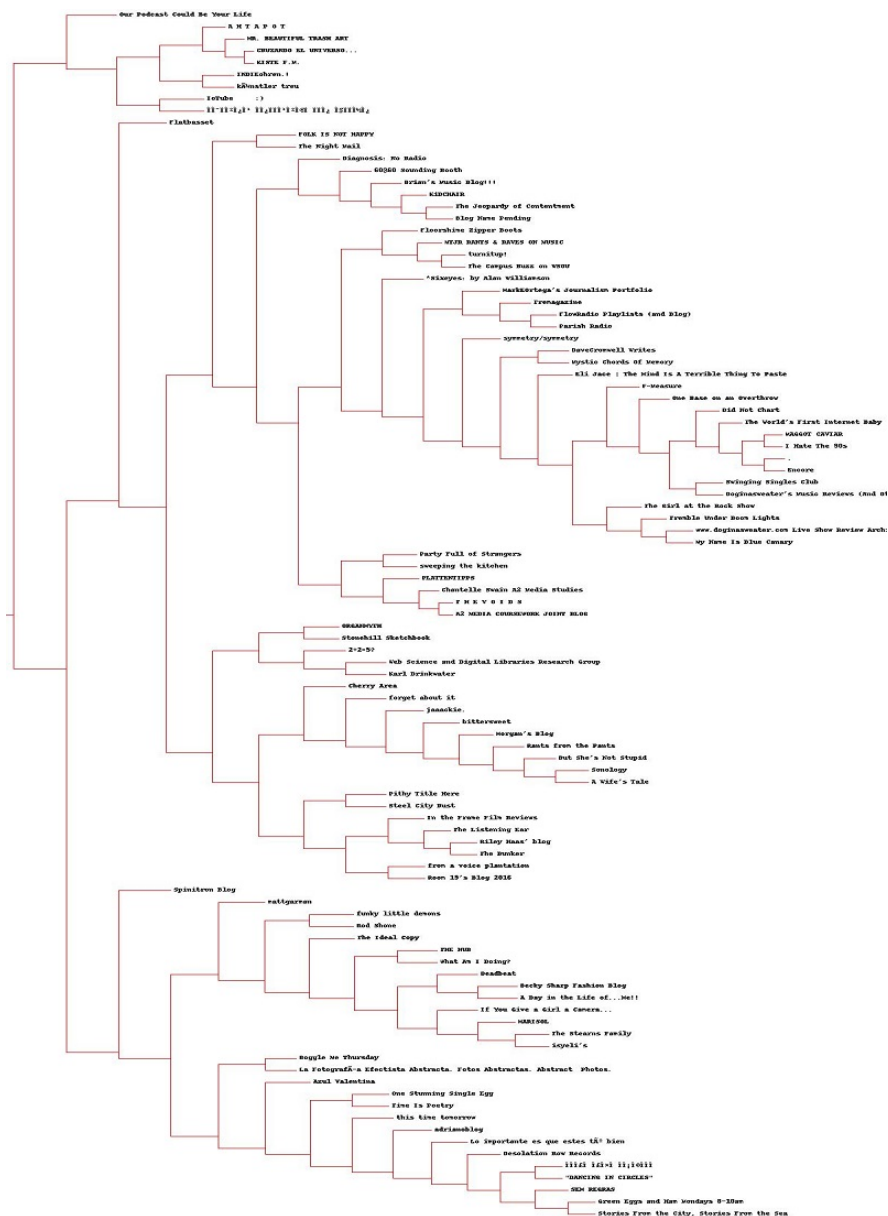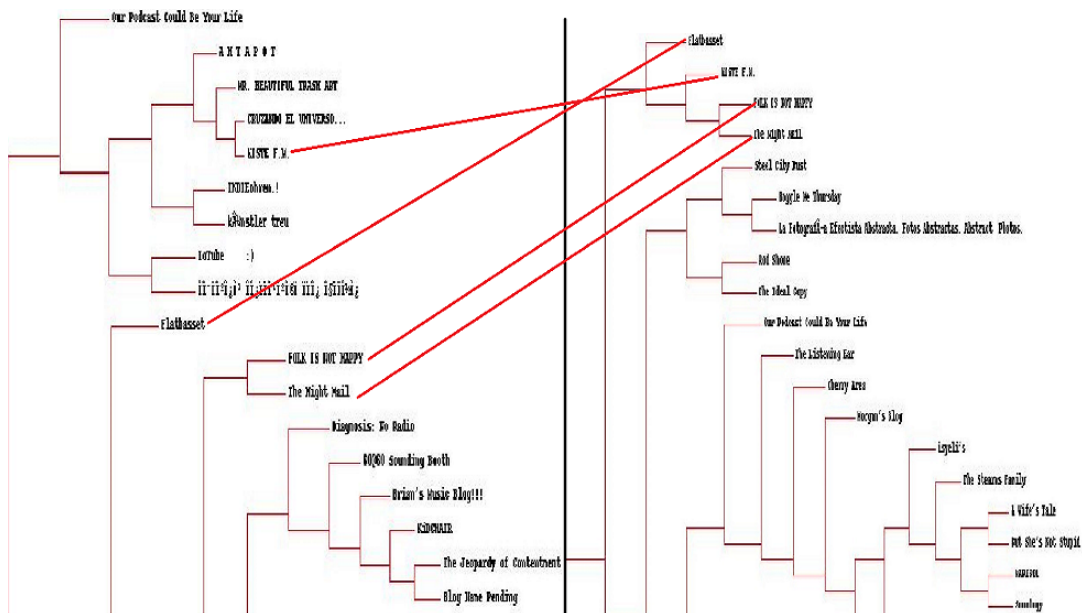


**Fig. 5.1.** JPEG dendogram

- Therefore I generated the output file 'blogdataWithTFIDF.txt' which has a blog matrix with blog title as identifier for each blog. This text file is uploaded to github at `https://github.com/majetisiri/cs532-s16/blob/master/a8/q5-blogdataWithTFIDF.txt`
- By comparing and contrasting we can say that the blog dendogram based on TFIDF has more clusters than the blog dendogram with frequency count in question 1. By looking at the dendograms it seems there is a much difference, only the hierarchy of clusters are changing but the clusters are around the same blog with different hierarchy. This is illustrated in Figure 5.2.



**Fig. 5.2.** comparing and contrasting dendograms from question 1 and question 5

**Code Listing**

```
1   import feedparser
2   import re
3   import sys
4   import math
5
6   def getwordcounts(url):
7     d=feedparser.parse(url)
8     wc={}
9
10    for e in d.entries:
11      if 'summary' in e:
12        summary=e.summary
13      else:
14        summary=e.description
15
16      words=getwords(e.title+' '+summary)
17      for word in words:
18        wc.setdefault(word,0)
19        wc[word]+=1
20    print d.feed.title
21    return d.feed.title,wc
22
23  def getwords(html):
24    txt=re.compile(r'<[^>]+>').sub('',html)
25    words=re.compile(r'[^A-Z^a-z]+').split(txt)
26    return [word.lower() for word in words if word!='']
27
28  def generateFeedVector():
29    apcount={}
30    wordcounts={}
31    iteration = 1
32    feedlist=[line for line in file('getAtomsFor100Urls')]
33    for feedurl in feedlist:
34      try:
35        title,wc=getwordcounts(feedurl)
36        wordcounts[title]=wc
37        for word,count in wc.items():
38          apcount.setdefault(word,0)
39          if count>1:
40            apcount[word]+=1
41      except:
42        print 'Failed to parse feed %s' % feedurl
43      iteration+=1
44
45    wordlist=[]
46    countFrequentWords=[]
47    for w,bc in apcount.items():
```

```
48        frac=float(bc)/len(feedlist)
49        if frac >0.1 and frac <0.5:
50            countFrequentWords.append((w,bc))
51
52    countFrequentWords=sorted(countFrequentWords,key=lambda x:
          x[1], reverse = True)
53
54    for value in countFrequentWords:
55        value1 = value[0]
56        value2 = value[1]
57        length = len(wordlist)
58        if(length < 500):
59            wordlist.append(value1)
60        else:
61            break
62
63    out=file('blogdata.txt','w')
64    out.write('Blog')
65    for word in wordlist:
66        word1 = word.encode('UTF-8')
67        out.write('\t%s' % word1)
68    out.write('\n')
69    for blog,wc in wordcounts.items():
70        blogName = blog.encode('UTF-8')
71        print blog
72        print blogName
73        out.write(blogName)
74        for word in wordlist:
75            if word in wc:
76            # EDITED CODE
77                termFrequency = wc[word]/float(len(wc))
78                inverseDocumentFrequency = logBase2(iteration/float(
                    apcount[word]))
79                tfIdf = termFrequency*inverseDocumentFrequency
80                out.write('\t%f' % tfIdf)
81            else:
82                out.write('\t0')
83        out.write('\n')
84
85  def logBase2(number):
86    return math.log(number) / math.log(2)
87
88  generateFeedVector()
```

**Listing 5.1.** Python code for calculating TFIDF

# References

1. get requests for URI: `http://docs.python-requests.org/en/master/user/quickstart/`, 2016, A Kenneth Reitz Project
2. Python code for generating feed vector from PCI book. Igraph Tutorial. `https://github.com/cataska/programming-collective-intelligence-code/blob/master/chapter3/generatefeedvector.py`, 2007, Toby Segaran
3. Python code for clusters from PCI book. Download GraphML for Karate Club. `https://github.com/cataska/programming-collective-intelligence-code/blob/master/chapter3/clusters.py`, 2007, Toby Segaran