Assignment 9
CS 595: Introduction to Web Science Fall 2013 Shawn M. Jones Finished on December 4, 2013

Question

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.

Answer

The script fetchFeeds.py shown in Listing 1 outputs a list of URIs that dereference to Atom feeds.

This script uses a canned URI on line 54 to acquire the next blog URI via the call on line 82. If it failed to acquire this URI, it sleeps for 5 seconds and tries again. If successful, it then dereferences the URI using the function defined on line 12. This function returns the representation of the resource, as well as its URI after all redirects have been followed.

If that is successful, it uses the function defined on line 21 to extract the Atom URI from the blog's HTML. If that is successful, it then tries to dereference the Atom feed URI, again acquiring the representation and the redirect-free version of the atom URI. Then it tests the URI using the function defined on line 33. This function checks to ensure that the blog has at least 25 entries. Attempts to ensure that the blog was English did not work well, as all Blogger pages use a character encoding of UTF-8, and detecting 'English' (maybe using a dictionary?) appeared to be outside the realm of this assignment.

Once the URI gets through all of those wickets, it is saved to a list, with the query string "?max-results=1000" appended. This ensures that we get at most 1000 entries from each blog, providing the maximum amount of data

for the following questions. Because all blogs come from Blogger, which has a consistent API, this query string works every time.

```
1
   #!/usr/local/bin/python
2
3
   import sys
   import os
4
   import os.path
   import feedparser
   import urllib2
   import time
   import chardet
10
   from bs4 import BeautifulSoup
11
   def dereferenceUri(uri):
12
13
14
        pagehandle = urllib2.urlopen(uri)
15
        pagedata = pagehandle.read()
        derefurl = pagehandle.geturl()
16
17
        pagehandle.close()
18
19
        return pagedata, derefurl
20
   def getAtomFeedUri(html):
21
22
23
        soup = BeautifulSoup(html)
24
        atomLinks = soup.find_all('link',
attrs = { 'rel' : 'alternate', 'type' : 'application/
25
26
                atom+xml' })
27
        # we assume there is only one atom link
28
29
        atomURI = atomLinks[0].attrs['href']
30
31
        return atomURI
32
   def meetsCriteria(feedText):
33
34
35
        parsedData = feedparser.parse(feedText)
36
37
        # assume we're good to go by default (fail optimistic?)
38
        goodToGo = True
39
40
        sys.stderr.write("blog has " + str(len(parsedData.entries))
           + " entries \n")
41
42
        if (len(parsedData.entries) < 25):
43
            goodToGo = False
44
```

```
#if (chardet.detect(feedText)['encoding'] != 'ascii'):
45
             sys.stderr.write("blog charset is " + chardet.detect(
46
           feedText)['encoding'] +
                 ", likely won't parse well for feed vector\n")
47
48
             goodToGo = False
49
50
       return goodToGo
51
52
   def getNextUri():
53
        uri = "http://www.blogger.com/next-blog?navBar=true&blogID
54
           =3471633091411211117"
55
        pagehandle = urllib2.urlopen(uri)
56
57
        nexturi = pagehandle.geturl()
58
        pagehandle.close()
59
60
        return nexturi
61
62
63
   if __name__ == "__main__":
64
65
        feedlist = []
        feedlist.append("http://f-measure.blogspot.com/feeds/posts/
66
            default?max-results=200")
67
        feedlist.append("http://ws-dl.blogspot.com/feeds/posts/
           default?max-results=200")
68
        while (len(feedlist) \le 100):
69
70
71
            try:
72
                # Tried these steps, but realized I would have to
73
                # parse JavaScript
                \# * look for the iframe containing the next button
74
                        <iframe name='navbar-iframe'...</pre>
75
                #
                # iframeURI = getIframeUri(html)
76
77
                \# * dereference the link from that iframe
78
                # iframeText = getIframeText(iframeURI)
79
                \# * extract the uri form
80
                # <a class="b-link" href="...">Next Blog</a>
81
                # uri = getNextUri(iframeText)
82
                uri = getNextUri()
83
            except urllib2.HTTPError as e:
                sys.stderr.write("failed to acquire next uri,
84
                    delaying 5 seconds\n")
85
                time.sleep(5)
86
            else:
87
                sys.stderr.write("working on URI" + uri + "\n")
88
```

```
89
                 try:
90
                      # dereference the uri and get text
91
                      html, derefuri = dereferenceUri(uri)
                 except urllib2.HTTPError as e:
92
                      sys.stderr.write("failed to dereference" + uri
93
                          ", delaying 5 seconds\n")
94
95
                      time.sleep(5)
96
                 else:
97
                      try:
                          # fetch the atom feed URI
98
99
                          feedURI = getAtomFeedUri(html)
100
                          sys.stderr.write("acquired feed URI" +
                              feedURI + "\n")
                      except IndexError as e:
101
102
                          sys.stderr.write(
                              "failed to acquire Atom feed from HTML,
103
                                  delaying 5 seconds\n")
104
                      else:
105
                          \operatorname{try}:
106
                              # get the atom feed text
                              feedText , feedURI = dereferenceUri(
107
                                  feedURI)
108
                          except urllib2.HTTPError as e:
109
                               sys.stderr.write("failed to dereference
                                  " + feedURI +
                                  ", delaying 5 seconds\n")
110
111
                               time.sleep(5)
112
                          else:
113
                              # if it meets the criteria, save the
                                  file
114
                               if meetsCriteria (feedText):
115
                                   sys.stderr.write("Saving blog feed "
                                        + feedURI + "?\max-results=1000\n
116
                                   feedlist.append(feedURI + "?max-
                                       results = 1000")
117
118
                              # be nice to the site
119
                              time.sleep(1)
120
121
         for feed in feedlist:
122
             print feed
```

Listing 1: Python script for fetching valid Atom feeds from Blogger

The script is run like so:

```
./fetchFeeds.py > feedlist.txt
```

Once feedlist.txt exists, it can be used by Toby Segaran's generatefeedvector.py[1]. I only modified it on line 59 so that it could handle UTF-8 encodings for some of the blogs. It is captured in Listing 6 on page 16.

The script ${\tt eliminateWords.py}$ shown in Listing 2 takes in the ${\tt blogdata1.txt}$ file produced by ${\tt generatefeedvector.py}$ and removes the top N terms from it.

It does this by generating scores for each word by summing all of its frequencies across all blogs together. Once it has those scores, it gets the top n words. It then determines the index of each of those words in the list. Using these indices, it regenerates the format of blogdata.txt, only printing out each column if it is an "approved" index.

```
#!/usr/local/bin/python
1
2
3
   import sys
4
5
   sys.path.insert(0, '../libs')
7
   import clusters
8
   import operator
9
   import pprint
10
11
12
    def getWordscores (words, data):
13
        wordscores = \{\}
14
15
        for i in range(len(words)):
16
            sys.stderr.write('examining' + words[i] + '\n')
17
18
            for j in range(len(data)):
19
20
21
                 if words[i] in wordscores:
22
                     wordscores [words [i]] += data [j][i]
23
                 else:
24
                     wordscores [words [i]] = data [j] [i]
25
26
        return wordscores
27
28
    def getTopNWords(wordscores, n):
29
        topNWords = []
30
31
32
        # thanks Stack Overflow:
33
        # http://stackoverflow.com/questions/613183/python-sort-a-
            dictionary-by-value
34
        reversedWordscores = sorted(
```

```
35
            wordscores.iteritems(), key=operator.itemgetter(1),
                reverse=True
36
            )
37
        for i in range(n):
38
39
            sys.stderr.write(
                "adding " + reversedWordscores[i][0] + " with a
40
                    score of "
41
                    + str(reversedWordscores[i][1]) + '\n'
42
            topNWords.append(reversedWordscores[i][0])
43
44
45
        return topNWords
46
47
48
   if -name_{-} = '-main_{-}':
49
        n = int(sys.argv[1])
50
51
52
        blognames, words, data = clusters.readfile('../q1/blogdata1.
            txt')
53
54
        wordscores = getWordscores(words, data)
55
56
        topNWords = getTopNWords(wordscores, n)
57
58
        indexlist = []
59
        for word in topNWords:
60
            indexlist.append(words.index(word))
61
62
        lines = []
63
64
65
        line = []
66
        line.append('Blog')
67
68
        for i in range(len(words)):
69
70
            if i in indexlist:
71
                line.append(words[i])
72
73
        lines.append(line)
74
75
        for i in range(len(blognames)):
76
            line = []
            line.append(blognames[i])
77
78
79
            for j in range(len(words)):
                if j in indexlist:
80
```

Listing 2: Python script for eliminating words from blog data

Question

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).

Answer

The script makeDendogram.py, shown in Listing 3 uses Toby Segaran's *clusters.py* [1] shown on page 18.

```
#!/usr/local/bin/python
1
   # all code here stolen shamelessly from
   # "Programming Collective Intelligence, Chapter 3"
5
6
   import sys
7
8
   sys.path.insert(0, '../libs')
10
   import clusters
11
   blognames\,,words\,,data = clusters\,.\,readfile\,(\ '.../\,q1/\,blogdata500\,.\,txt\ ')
12
13
   clust = clusters.hcluster(data)
14
   # print ASCII dendrogram
15
   clusters.printclust(clust, labels=blognames)
16
17
   # save JPEG dendrogram
18
   clusters.drawdendrogram(clust, blognames, jpeg='blogclust.jpg')
```

Listing 3: Python script for generating dendograms from the blog data captured in question 1

It is run like so:

```
./makeDendogram.py > ascii-dendogram.txt
```

The printclust function on line 16 prints out the dendogram, so shell redirection is used to save it. The drawdendogram function on line 19 saves a JPEG of the dendogram, which can be seen in Figure 1.

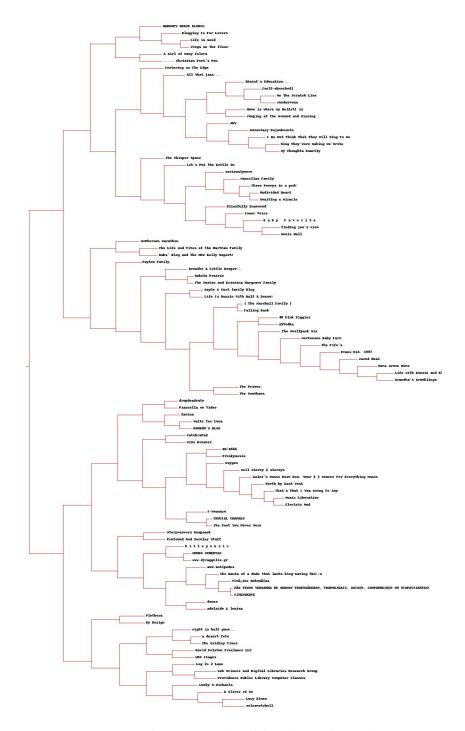


Figure 1: Dendogram produced by the makeDendogram.py script

Unfortunately, it is difficult to see, but this dendogram shows that the blogs calculated to be most like *F-Measure* are *CRUCIAL CHANGES* and *The Poet You Never Were*. The blog calculated to be most like *Web Science* and Digital Libraries Research Group is Providence Public Library Computer Classes.

Question

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

Answer

The blog clustering is performed by the script shown in Listing 4, which makes use of Toby Segaran's *clusters.py* [1] on page 18., using the function *kcluster* on lines 15, 19, and 23.

```
#!/usr/local/bin/python
1
2
3
   # all code here stolen shamelessly from
   # "Programming Collective Intelligence, Chapter 3"
6
   import sys
7
   sys.path.insert(0, '../libs')
8
9
10
   import clusters
11
12
   blognames, words, data=clusters.readfile('../q1/blogdata500.txt')
13
   print "For k=5"
14
   kclust=clusters.kcluster(data, k=5)
15
16
   print
17
   print "For k=10"
18
   kclust=clusters.kcluster(data, k=10)
19
20
   print
21
22
   print "For k=20"
23 | kclust=clusters.kcluster(data, k=20)
```

Listing 4: Python script for clustering the blogs using K-means, using k=5, 10, and 20

This script is run like so:

```
./makeClusters.py
```

From its output, we see how many iterations each value of k produces.

```
For k=5
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
For k=10
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
Iteration 5
Iteration 6
Iteration 7
For k=20
Iteration 0
Iteration 1
Iteration 2
Iteration 3
```

Thus, for k=5 we get 5 iterations, for k=10 we get 8 iterations, and for k=20 we get 4 iterations.

Question

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

Answer

The blog space is generated using multidimensional scaling from the script makeMDS.py shown in Listing 5, which makes use of Toby Segaran's *clusters.py* [1] on page 18, using the functions *scaledown* on line 14 and *draw2d* on line 16.

```
#!/usr/local/bin/python
   # all code here stolen shamelessly from
   # "Programming Collective Intelligence, Chapter 3"
5
6
   import sys
7
   sys.path.insert(0, '../libs')
8
9
10
   import clusters
11
   blognames, words, data=clusters.readfile('../q1/blogdata500.txt')
12
13
   coords = clusters.scaledown(data)
14
15
   clusters.draw2d(coords, blognames, jpeg='blogs2d.jpg')
16
```

Listing 5: Python script for generating a MDS from the blog data

Again, unfortunately the blog space produced does not fit well on a letter-sized page shown in Figure 2. According to this two-dimensional representation, the blog closest to F-Measure is b if l e y m u s i c. The blogs closest to Web Science and Digital Libraries Research Group are MUNDO OCHENTAS and Providence Public Library Computer Classes.



Figure 2: Dendogram produced by the makeMDS.py script

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.

Answer

Not attempted.

A Segaran's generatefeedvector.py

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

```
47
   wordlist = []
48
   for w, bc in apcount.items():
      frac=float(bc)/len(feedlist)
49
50
      if frac > 0.1 and frac < 0.5:
51
        wordlist.append(w)
52
53
   out=file ('blogdata1.txt', 'w')
   out.write('Blog')
54
   for word in wordlist: out.write('\t%s' % word)
55
   out.write('\n')
56
57
   for blog, wc in wordcounts.items():
58
      print blog
     blog = blog.encode('UTF-8')
59
60
     out.write(blog)
61
      for word in wordlist:
62
        if word in wc: out.write('\t%d' % wc[word])
63
        else: out.write('\t0')
64
     out.write('\n')
```

Listing 6: Segaran's generatefeedvector.py, with small modification added on line 59 to handle UTF-8 encodings

B Segaran's clusters.py

```
1
          from PIL import Image, ImageDraw
  2
  3
          def readfile (filename):
                lines = [line for line in file (filename)]
  4
  5
                # First line is the column titles
  6
  7
                colnames=lines [0].strip().split('\t')[1:]
  8
                rownames = []
                data = []
  9
10
                for line in lines [1:]:
11
                      p=line.strip().split('\t')
12
                      # First column in each row is the rowname
13
                      rownames.append(p[0])
                      # The data for this row is the remainder of the row
14
15
                      data.append([float(x) for x in p[1:]])
16
                return rownames, colnames, data
17
18
19
          from math import sqrt
20
21
           def pearson (v1, v2):
22
               # Simple sums
23
                sum1=sum(v1)
24
                sum2=sum(v2)
25
               # Sums of the squares
26
27
                sum1Sq=sum([pow(v,2) for v in v1])
                sum2Sq=sum([pow(v,2) for v in v2])
28
29
                # Sum of the products
30
31
                pSum=sum([v1[i]*v2[i] for i in range(len(v1))])
32
33
               # Calculate r (Pearson score)
34
                num = pSum - (sum1*sum2/len(v1))
35
                den = sqrt((sum1Sq-pow(sum1,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1)*(sum2Sq-pow(sum2,2)/len(v1)*(sum2Sq-pow(sum2,2)/len(v1)*(sum2Sq-pow(sum2,2)/len(v1)*(sum2Sq-pow(sum2,2)/len(v1)*(sum2Sq-
                           v1)))
36
                 if den==0: return 0
37
38
                return 1.0-num/den
39
40
           class bicluster:
41
                def __init__ (self, vec, left=None, right=None, distance=0.0,id=
                           None):
42
                       self.left = left
43
                       self.right=right
44
                       self.vec=vec
```

```
self.id=id
45
46
        self.distance=distance
47
    def hcluster (rows, distance=pearson):
48
49
      distances = \{\}
50
      currentclustid=-1
51
52
     # Clusters are initially just the rows
      clust = [bicluster(rows[i],id=i) for i in range(len(rows))]
53
54
      while len(clust) > 1:
55
56
        lowestpair = (0,1)
57
        closest=distance(clust[0].vec,clust[1].vec)
58
59
        # loop through every pair looking for the smallest distance
        for i in range(len(clust)):
60
61
          for j in range(i+1,len(clust)):
62
            # distances is the cache of distance calculations
63
             if (clust [i].id, clust [j].id) not in distances:
               distances [(clust[i].id,clust[j].id)] = distance(clust[i].id)] = distance(clust[i].id)
64
                   ].vec, clust[j].vec)
65
            d=distances [(clust[i].id,clust[j].id)]
66
67
68
             if d<closest:
69
               c losest=d
70
               lowestpair=(i,j)
71
72
        # calculate the average of the two clusters
73
        mergevec=[
        (clust [lowestpair [0]]. vec [i]+clust [lowestpair [1]]. vec [i])
74
75
        for i in range(len(clust[0].vec))]
76
77
        # create the new cluster
        newcluster=bicluster (mergevec, left=clust [lowestpair [0]],
78
79
                                right=clust [lowestpair [1]],
80
                                distance=closest, id=currentclustid)
81
82
        # cluster ids that weren't in the original set are negative
83
        currentclustid -=1
84
        del clust [lowestpair [1]]
85
        del clust [lowestpair [0]]
86
        clust.append(newcluster)
87
88
      return clust [0]
89
90
   def printclust (clust, labels=None, n=0):
91
     # indent to make a hierarchy layout
```

```
92
      for i in range(n): print '',
93
      if clust.id < 0:
        # negative id means that this is branch
94
        print '-'
95
96
       else:
        # positive id means that this is an endpoint
97
        if labels=None: print clust.id
98
99
         else: print labels [clust.id]
100
101
      # now print the right and left branches
      if clust.left!=None: printclust(clust.left,labels=labels,n=n
102
103
      if clust.right!=None: printclust(clust.right, labels=labels, n=n
          +1)
104
105
    def getheight (clust):
      # Is this an endpoint? Then the height is just 1
106
107
      if clust.left=None and clust.right=None: return 1
108
      # Otherwise the height is the same of the heights of
109
110
      # each branch
111
      return getheight (clust.left)+getheight (clust.right)
112
    def getdepth(clust):
113
114
      # The distance of an endpoint is 0.0
115
      if clust.left—None and clust.right—None: return 0
116
117
      # The distance of a branch is the greater of its two sides
118
      # plus its own distance
      return max(getdepth(clust.left),getdepth(clust.right))+clust.
119
          distance
120
121
122
    def drawdendrogram(clust, labels, jpeg='clusters.jpg'):
123
      # height and width
124
      h=getheight(clust)*20
125
      w = 1200
126
      depth=getdepth(clust)
127
128
      # width is fixed, so scale distances accordingly
129
      scaling=float(w-150)/depth
130
      # Create a new image with a white background
131
      img=Image.new('RGB',(w,h),(255,255,255))
132
133
      draw=ImageDraw.Draw(img)
134
135
      draw.line((0,h/2,10,h/2),fill=(255,0,0))
136
137
      # Draw the first node
```

```
138
       drawnode (draw, clust, 10, (h/2), scaling, labels)
139
       img.save(jpeg,'JPEG')
140
    def drawnode (draw, clust, x, y, scaling, labels):
141
142
       if clust.id < 0:
143
         h1=getheight (clust.left) *20
144
         h2=getheight (clust.right) *20
145
         top=y-(h1+h2)/2
146
         bottom=y+(h1+h2)/2
         # Line length
147
148
         ll=clust.distance*scaling
149
         # Vertical line from this cluster to children
150
         draw. line ((x, top+h1/2, x, bottom-h2/2), fill = (255, 0, 0))
151
152
         # Horizontal line to left item
153
         draw.line((x, top+h1/2, x+l1, top+h1/2), fill = (255, 0, 0))
154
155
         # Horizontal line to right item
156
         draw. line ((x, bottom-h2/2, x+l1, bottom-h2/2), fill = (255, 0, 0))
157
158
         # Call the function to draw the left and right nodes
         drawnode (draw, clust.left, x+ll, top+h1/2, scaling, labels)
159
160
         drawnode (draw, clust.right, x+ll, bottom-h2/2, scaling, labels)
161
       else:
162
         # If this is an endpoint, draw the item label
163
         draw.text((x+5,y-7), labels[clust.id],(0,0,0))
164
165
    def rotatematrix (data):
166
       newdata=[]
       for i in range (len (data [0])):
167
168
         newrow=[data[j][i] for j in range(len(data))]
169
         newdata.append(newrow)
170
       return newdata
171
172
    import random
173
174
    def kcluster (rows, distance=pearson, k=4):
175
      # Determine the minimum and maximum values for each point
176
       ranges = [(min([row[i] for row in rows]), max([row[i] for row in
           rows]))
177
       for i in range(len(rows[0]))]
178
179
      # Create k randomly placed centroids
       clusters = [[random.random()*(ranges[i][1] - ranges[i][0]) + ranges[
180
           i ] [0]
181
       for i in range(len(rows[0])) for j in range(k)]
182
183
       lastmatches=None
184
       for t in range (100):
```

```
185
         print 'Iteration %d' % t
         bestmatches = [[] for i in range(k)]
186
187
         # Find which centroid is the closest for each row
188
189
         for j in range(len(rows)):
190
           row=rows[j]
191
           bestmatch=0
192
           for i in range(k):
             d=distance(clusters[i],row)
193
             if d<distance(clusters[bestmatch],row): bestmatch=i
194
195
           bestmatches [bestmatch].append(j)
196
197
         # If the results are the same as last time, this is complete
198
         if bestmatches=lastmatches: break
199
         lastmatches=bestmatches
200
         # Move the centroids to the average of their members
201
202
         for i in range(k):
203
           avgs = [0.0] * len(rows [0])
204
           if len(bestmatches[i])>0:
205
             for rowid in bestmatches[i]:
206
               for m in range (len (rows [rowid])):
                  avgs [m]+=rows [rowid][m]
207
208
             for j in range(len(avgs)):
209
               avgs [j]/=len (bestmatches [i])
210
             clusters [i]=avgs
211
212
       return bestmatches
213
    def tanamoto(v1, v2):
214
215
       c1, c2, shr = 0, 0, 0
216
217
       for i in range(len(v1)):
218
         if v1[i]!=0: c1+=1 # in v1
         if v2[i]!=0: c2+=1 # in v2
219
220
         if v1[i]!=0 and v2[i]!=0: shr+=1 \# in both
221
222
       return 1.0 - (float(shr)/(c1+c2-shr))
223
224
    def scaledown (data, distance=pearson, rate=0.01):
225
      n=len (data)
226
227
      # The real distances between every pair of items
228
       realdist = [[distance(data[i], data[j]) for j in range(n)]
229
                   for i in range (0,n)
230
231
      # Randomly initialize the starting points of the locations in
232
       loc = [[random.random(), random.random()]] for i in range(n)]
```

```
fakedist = [[0.0 \text{ for } j \text{ in } range(n)] \text{ for } i \text{ in } range(n)]
233
234
235
       lasterror=None
       for m in range (0,1000):
236
237
         # Find projected distances
238
         for i in range(n):
239
            for j in range(n):
240
              fakedist[i][j]=sqrt(sum([pow(loc[i][x]-loc[j][x],2)
241
                                            for x in range(len(loc[i]))])
242
243
         # Move points
244
         \operatorname{grad} = [[0.0, 0.0] \text{ for i in range}(n)]
245
246
          totalerror=0
247
          for k in range(n):
248
            for j in range(n):
249
              if j=k: continue
              # The error is percent difference between the distances
250
251
              errorterm=(fakedist [j][k]-realdist [j][k])/realdist [j][k]
252
253
              # Each point needs to be moved away from or towards the
                  other
254
              # point in proportion to how much error it has
255
              \operatorname{grad}[k][0] + = ((\operatorname{loc}[k][0] - \operatorname{loc}[j][0]) / \operatorname{fakedist}[j][k]) *
                  errorterm
256
              grad[k][1]+=((loc[k][1]-loc[j][1])/fakedist[j][k])*
                  errorterm
257
258
              # Keep track of the total error
259
              totalerror+=abs(errorterm)
260
          print totalerror
261
262
         # If the answer got worse by moving the points, we are done
263
         if lasterror and lasterror < totalerror: break
264
         lasterror=totalerror
265
266
         # Move each of the points by the learning rate times the
              gradient
267
          for k in range(n):
268
            loc[k][0] -= rate * grad[k][0]
269
            loc[k][1] -= rate*grad[k][1]
270
271
       return loc
272
273
     def draw2d(data, labels, jpeg='mds2d.jpg'):
274
       img=Image.new('RGB',(2000,2000),(255,255,255))
275
       draw=ImageDraw.Draw(img)
276
       for i in range(len(data)):
277
         x = (data[i][0] + 0.5) *1000
```

```
278 | y=(data[i][1]+0.5)*1000

279 | draw.text((x,y),labels[i],(0,0,0))

280 | img.save(jpeg,'JPEG')
```

Listing 7: Segaran's clusters.py

References

- [1] SEGARAN, T. *Programming Collective Intelligence*, first ed. O'Reilly, O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472, August 2007.
- [2] Segaran, T. *Programming Collective Intelligence*, first ed. O'Reilly, 2007.