# HW 8 - Clustering
Swathi Venkatesh
12/05/2021

# Q1

Generate a list of 100 popular accounts on Twitter. The accounts must be verified, have greater than 10,000 followers, and have greater than 5000 tweets. For example:

weiglemc - not verified, 509 followers, 2813 tweets - don't include wnba - verified (blue checkmark), 739,600+ followers, 84,200+ tweets - could include See Twarc API user_lookup, GET users/lookup, and User object for details on obtaining this information for a set of accounts.

You may also generate this information manually by visiting individual account pages.

Because we're trying to cluster the accounts based on the text in their tweets, you should choose several sets of accounts that are similar (political, tech, sports, etc.) to see if they'll get clustered together later.

Save the list of accounts (screen_names), one per line, in a text file named accounts.txt and upload to your GitHub repo.

## Answer

```python
1  #!/usr/local/bin/python3
2  from twarc import Twarc2, expansions
3  from configparser import ConfigParser
4
5  def userstat(twAuth,ids):
6      #Find followers that are in this category and them to the list
7      stat= False
8      try:
9          user = twAuth.get_user(id = ids)
10         if user.verified and user.statuses_count >= 5000 and user.
     followers_count >= 10000 == True:
11             stat = True
12         else:
13             stat = False
14     except:
15         print("Error")
16     return stat
17
18  def setup_api(filename):
```

```python
19      '''
20      filename: file where Twitter API keys are stored
21      returns Twitter API object to pass into parse()
22      '''
23
24      # read Twitter API keys from twarc config file, setup twarc2 object
25      config = ConfigParser(interpolation=None)
26      with open(filename) as twarc_config:
27          config.read_string("[TWARC]\n" + twarc_config.read())
28      bearer_token = config['TWARC']['bearer_token'].strip('\'')
29      t = Twarc2(bearer_token=bearer_token)
30      return t
31
32  def parse(api, screen_name, num_tweets=100):
33      '''
34      api: Twitter API object, use setup_api() to create
35      screen_name: Twitter screen_name
36      num_tweets: Number of tweets to request (default: 100)
37      returns dict with {'screen_name': screen_name, 'tweets': [tweet1,
        tweet2, ...]}
38      '''
39
40      tweet_data = []
41      try:
42          timeline = api.timeline(screen_name, max_results=num_tweets,
        exclude_replies=True, exclude_retweets=True)
43          for page in timeline:
44              result = expansions.flatten(page)
45              for tweet in result:
46                  tweet_data.append(tweet["text"])
47                  if len(tweet_data) == num_tweets:
48                      # must include this to stop after a certain # of
        tweets
49                      timeline.close()
50      except Exception as e:
51          print ("Twarc Error: %s" % str(e))
52
53      account_data = {'screen_name': screen_name, 'tweets': tweet_data}
54      return account_data
```

**Listing 1:** tweetparser.py

```python
1  #!/usr/local/bin/python3
2  from twarc import Twarc2, expansions
3  from tweetparser import setup_api, user_stat
4  import pandas as pd
5  import numpy as np
6  def  twitter_account(types,api= setup_api("/config") ):
```

```
 7
 8      #Get twitter friends screen_name of the parse account, as one
     arrayList
 9     public_tweets = twarc.user_lookup(api.friends_ids, id = types)
10     count= 1
11     resultList =[] # store the final result of screen names
12     for user in public_tweets.pages():
13          #print(user)
14          #print(user.dtypes())
15          #Transverse the list of followers ids
16          for i in user:
17              #check if requirement is met
18              if(user_stat(api,i) ==True):
19                  if( i not in resultList):
20                      #instead of the user id get the user name
21                      user_sc = api.get_user(i)
22                      print("{}:{}".format(count,user_sc.screen_name))
23                      resultList.append(user_sc.screen_name)
24                      count += 1
25              #Get maximum 30 file accounts screen_names
26              if(count >= 30):
27                  break
28      print(resultList)
29
30      #build text file and save file in q1 directory
31      filename =  'one/' + types + '.txt'
32      with open(filename, 'w') as filehandle:
33          for listitem in resultList:
34              filehandle.write('%s\n' % listitem)
35      return resultList
36 """
37 Tech= @WIRED
38 Sport= @WNBA
39 politics = @POTUS45
40 music = @future_of_music
41
42
43
44 """
45 #Get all screen_names with that fulfils 10,000 followers and have 5000
    tweets and verified
46 twitter_account("WIRED")
47 twitter_account("WNBA")
48 twitter_account("POTUS45")
49 twitter_account("future_of_music")
50
51
```

```python
52 """
53 Bring all result from the files in to One text file
54 accounts.txt
55 """
56 #get the list of the created files
57 column_name= ["User_screen_names"]
58 final = pd.DataFrame(columns= column_name)
59 fileList = ["one/WIRED.txt","one/WNBA.txt","one/POTUS45.txt","one/
       future_of_music.txt"]
60 df = pd.DataFrame()
61
62
63 for t in fileList:
64     frame = pd.read_csv(t,header=None)
65     frame.columns = column_name
66     for ind in frame.index:
67         final.loc[len(final)] =[frame['User_screen_names'][ind]]
68
69
70 # dropping duplicate values
71 final.User_screen_names.drop_duplicates(inplace=True)
72
73 #confirm that there are all unique values
74 print(final.User_screen_names.nunique())
75
76 # Number of rows to drop
77 n = 14
78
79 # Dropping last n rows using drop
80 final.drop(final.tail(n).index,
81         inplace = True)
82 #print(final.User_screen_names.nunique())
83 numpy_array = final.to_numpy()
84 #print as a text file
85 np.savetxt(r'accounts.txt', numpy_array,fmt="%s")
```

**Listing 2:** twittergatherid.py

## Discussion

The following consideration is made for gathering the user screen names. The screen names were collected through popular twitter accounts WNBA, POTUS45, future of music and WIRED. The lines 5-16 handles account that check screen_name account meets requirements.

```python
5 def userstat(twAuth,ids):
6     #Find followers that are in this category and them to the list
```

```
7      stat= False
8      try:
9          user = twAuth.get_user(id = ids)
10         if user.verified and user.statuses_count >= 5000 and user.
   followers_count >= 10000 == True:
11             stat = True
12         else:
13             stat = False
14     except:
15         print("Error")
16     return stat
```

**Listing 3:** A snap shot of tweetparser.py

*Q: How did you choose to collect the accounts?*

Function twitter_account() produce a text files (names based on the arguement supplied) that gets stored in /one folder.

```
52 """
53 Bring all result from the files in to One text file
54 accounts.txt
55 """
56 #get the list of the created files
57 column_name= ["User_screen_names"]
58 final = pd.DataFrame(columns= column_name)
59 fileList = ["one/WIRED.txt","one/WNBA.txt","one/POTUS45.txt","one/
   future_of_music.txt"]
60 df = pd.DataFrame()
61
62
63 for t in fileList:
64     frame = pd.read_csv(t,header=None)
65     frame.columns = column_name
66     for ind in frame.index:
67         final.loc[len(final)] =[frame['User_screen_names'][ind]]
68
69
70 # dropping duplicate values
71 final.User_screen_names.drop_duplicates(inplace=True)
72
73 #confirm that there are all unique values
74 print(final.User_screen_names.nunique())
75
76 # Number of rows to drop
77 n = 14
78
79 # Dropping last n rows using drop
80 final.drop(final.tail(n).index,
```

```
81          inplace = True)
82 #print(final.User_screen_names.nunique())
83 numpy_array = final.to_numpy()
84 #print as a text file
85 np.savetxt(r'accounts.txt', numpy_array,fmt="%s")
```

**Listing 4:** A snap shot of twittergatherid.py

*Q: What topics/categories do the accounts belong to? You don't need to specify a grouping for each account, but what general topics/categories will you expect to be revealed by the clustering?*

The topics belong to technology, sports, politics and music categories.

# Q2

## Answer

```
1 #!/usr/local/bin/python3
2 from tweetparser import setup_api, parse
3 import re
4
5 def getwordcounts(api, screen_name):
6     """
7      api: Twitter API object
8      screen_name: Twitter screen_name
9      returns screen_name and dictionary of word counts for a Twitter
    account
10     """
11
12     # Parse the Twitter feed
13     d = parse(api, screen_name)
14     wc = {}
15
16     # Loop over all the entries
17     for tweet in d['tweets']:
18
19         # Extract a list of words
20         words = getwords(tweet)
21         for word in words:
22             wc.setdefault(word, 0)
23             wc[word] += 1
24
25     return (d['screen_name'], wc)
26
27 def getwords(tweet):
```

```
28        """
29        returns lowercase list of words after filtering
30        """
31
32        # Remove URLs
33        text = re.compile(r'(http://|https://|www\.)([^ \'\"]*)').sub('',
      tweet)
34        "
35        # Remove other screen names (start with @)
36        text = re.compile(r'(@\w+)').sub('', text)
37
38        # Split words by all non-alpha characters
39        words = re.compile(r'[^A-Z^a-z]+').split(text)
40
41        # Filter for words between 3-15 characters, convert to lowercase,
      and return as a list
42        return [word.lower() for word in words if (len(word) >= 3 and len(
      word) <= 15)]
43
44 #####
45 # MAIN CODE STARTS HERE
46 #####
47
48
49 # set up Twitter API object
50 api = setup_api("/config")
51
52 apcount = {}       # number of accounts each word appears in
53 wordcounts = {}    # words and frequency in each account
54 sumcounts = {}     # words and frequency over all accounts (to determine
       most popular)
55
56 # list of screen names should be in 'accounts.txt', one per line
57 accountlist = [line.strip() for line in open('accounts.txt')]
58 #print(accountlist)
59 #print(len(accountlist))
60
61 for screen_name in accountlist:
62     try:
63         # get tweets, filter and count words
64         (user, wc) = getwordcounts(api, screen_name)
65         wordcounts[user] = wc
66
67         # count number of accounts each term appears in
68         for (word, count) in wc.items():
69             apcount.setdefault(word, 0)
70             sumcounts.setdefault(word, 0)
```

```
71              if count > 1:
72                  apcount[word] += 1         # counting accounts with the
    word
73                  sumcounts[word] += count  # summing total counts for
    the word
74       except:
75           print ('Failed to parse account %s' % screen_name)
76
77
78 #print("Counting words done")
79 #print(sumcounts.keys())
80
81 # remove stopwords ("fake" way)
82 wordlist = []
83 for (w, ac) in apcount.items():
84     # w is the word, ac is the account count (was bc 'blog count' in
    textbook)
85     frac = float(ac) / len(accountlist)
86     if frac > 0.1 and frac < 0.5:
87         wordlist.append(w)
88
89 popularlist = []
90
91 ####
92 # BEGIN YOUR CODE BLOCK
93 ####
94
95 #tuple list sorted by index 1 i.e. value field
96 l = sorted(sumcounts.items() , key=lambda x: x[1],reverse=True)
97 #extract 500 rows
98 l = l[:500]
99 #store in the dictionary
100 popularlist = [i[0] for i in l]
101
102 # write out popular word list
103 with open('popularlist.txt', 'w') as outf:
104     for word in popularlist:
105         outf.write(word + '\n')
106
107 # write out account-term matrix
108 with open('tweetdata.txt', 'w') as outf:
109     # write header row ("Account", list of words)
110     outf.write('Account')
111     for word in popularlist:
112         outf.write('\t%s' % word)
113     outf.write('\n')
114
```

```
115      # write each row (screen_name, count for each word)
116      for (screen_name, wc) in wordcounts.items():
117          outf.write(screen_name)
118          for word in popularlist:
119              if word in wc:
120                  outf.write('\t%d' % wc[word])
121              else:
122                  outf.write('\t0')
123          outf.write('\n')
```

**Listing 5:** generatetweetvector.py

## Discussion

- *Q:Explain the general operation of generatetweetvector.py and how the tweets are converted to the account-term matrix.*

- The code drive started from line 56 to line 123.

- From lines 61 to 75, uses a major function called getwordcounts(api,screen name). In this function the parse(api,screen name) is called. The parse function returns users tweets full text that are not retweets nor replies. When that full tweet text is gotten for a particular user, getwords(tweet) function removes unwanted text from the tweets gotten such as URLS and Mentions. It then extracts word that has at least 3 to 15 length size in each sentence in the tweets and converts them to lower cases.

- The result is stored for each words in a dictionary where by if the word repeats itself again the word(which is the key of the dictionary) increments the values by one.

- getwordcounts returns the screen name with a dictionary or words with its frequency count as well.

- Moving we count the number of accounts each term appear. It is noticed that each works and users account are stored in a variable outside the for loop; so basically result gotten gets added on each screen names.

- It appears to keep track of the overall frequency of the word too for every user combined.

- Then the removal of stop words (this is calculated, not gotten from a list of unwanted words) is done. It basically gets the word count divided by the total number of account names gotten from the accounts.txt, once it satisfies a particular fraction between 0.1 and 0.5 the word is added to the group of wordlist.

- *Q: Explain in detail the code that you added to filter for the 500 most frequent non-stopword terms.*

- For the popular list the items were sorted using the lambda function and the results were placed in the variable l. This variable l is a list of tuple. It goes from highest to lowest. Only 500 row items were considered by splicing the tuple and they were stored in a dictionary variable called popularlist.

- Finally, results of popularlist variable words are saved in a text file while tweetdata.txt as a header of the popularlist or words and the screen name with count of each word on a single row.

- *Q: Do the 500 most frequent terms make sense based on the accounts that you chose?*

- The word i viewed in popularlist.txt makes a lot of sense because it is as a form of connection to sports, politics, music or tech considered.

# Q3

## Answer

```python
#!/usr/local/bin/python3
from PIL import Image, ImageDraw
from math import import sqrt
import random
import csv
import pandas as pd
def readfile(filename):
  data = []
  rownames = []
  colnames = []
  num_rows = 0
  with open(filename) as tsvfile:
    reader = csv.reader(tsvfile, delimiter='\t')
    for row in reader:
      if num_rows > 0:
        rownames.append(row[0])    # save the row names
        data.append([float(x) for x in row[1:]])  # save the values as
    floats
      else:
        for col in row[1:]:
          colnames.append(col)    # save the column names
      num_rows = num_rows + 1
  return (rownames, colnames, data)

def pearson(v1, v2):
  # Simple sums
    sum1 = sum(v1)
```

```python
27      sum2 = sum(v2)
28
29    # Sums of the squares
30      sum1Sq = sum([pow(v, 2) for v in v1])
31      sum2Sq = sum([pow(v, 2) for v in v2])
32
33    # Sum of the products
34      pSum = sum([v1[i] * v2[i] for i in range(len(v1))])
35
36    # Calculate r (Pearson score)
37      num = pSum - sum1 * sum2 / len(v1)
38      den = sqrt((sum1Sq - pow(sum1, 2) / len(v1)) * (sum2Sq - pow(sum2,
    2)
39                 / len(v1)))
40      if den == 0:
41          return 0
42
43      return 1.0 - num / den
44 """
45 MD5 Scaling
46 """
47 def scaledown(data, distance=pearson, rate=0.01):
48      n = len(data)
49
50    # The real distances between every pair of items
51      realdist = [[distance(data[i], data[j]) for j in range(n)] for i in
52                  range(0, n)]
53
54    # Randomly initialize the starting points of the locations in 2D
55      loc = [[random.random(), random.random()] for i in range(n)]
56      fakedist = [[0.0 for j in range(n)] for i in range(n)]
57
58      lasterror = None
59      for m in range(0, 1000):
60      # Find projected distances
61          for i in range(n):
62              for j in range(n):
63                  fakedist[i][j] = sqrt(sum([pow(loc[i][x] - loc[j][x],
    2)
64                                        for x in range(len(loc[i]))]))
65
66      # Move points
67          grad = [[0.0, 0.0] for i in range(n)]
68
69          totalerror = 0
70          for k in range(n):
71              for j in range(n):
```

```python
72                if j == k:
73                    continue
74          # The error is percent difference between the distances
75                errorterm = (fakedist[j][k] - realdist[j][k]) /
     realdist[j][k]
76
77          # Each point needs to be moved away from or towards the other
78          # point in proportion to how much error it has
79                grad[k][0] += (loc[k][0] - loc[j][0]) / fakedist[j][k]
     \
80                    * errorterm
81                grad[k][1] += (loc[k][1] - loc[j][1]) / fakedist[j][k]
     \
82                    * errorterm
83
84          # Keep track of the total error
85                totalerror += abs(errorterm)
86          print (totalerror)
87
88      # If the answer got worse by moving the points, we are done
89          if lasterror and lasterror < totalerror:
90              break
91          lasterror = totalerror
92
93      # Move each of the points by the learning rate times the gradient
94          for k in range(n):
95              loc[k][0] -= rate * grad[k][0]
96              loc[k][1] -= rate * grad[k][1]
97
98      return loc
99
100 def draw2d(data, labels, jpeg):
101     img = Image.new('RGB', (2000, 2000), (255, 255, 255))
102     draw = ImageDraw.Draw(img)
103     for i in range(len(data)):
104         x = (data[i][0] + 0.5) * 1000
105         y = (data[i][1] + 0.5) * 1000
106         draw.text((x, y), labels[i], (0, 0, 0))
107     img.save(jpeg, 'JPEG')
108
109 def rotatematrix(data):
110     newdata = []
111     for i in range(len(data[0])):
112         newrow = [data[j][i] for j in range(len(data))]
113         newdata.append(newrow)
114     return newdata
115
```

```python
116  """
117  Hierarchical Clustering
118  class bicluster - data structure to hold the clustering information
119  hcluster(rows, distance=pearson) - does the hierarchical clustering,
         default distance function is pearson()
120  printclust(clust, labels=None, n=0) - traverses the cluster and prints
         an ASCII text representation
121  """
122  class bicluster:
123
124      def __init__(self, vec, left=None, right=None, distance=0.0, id=
         None,):
125          self.left = left
126          self.right = right
127          self.vec = vec
128          self.id = id
129          self.distance = distance
130
131  def hcluster(rows, distance=pearson):
132      distances = {}
133      currentclustid = -1
134
135    # Clusters are initially just the rows
136      clust = [bicluster(rows[i], id=i) for i in range(len(rows))]
137
138      while len(clust) > 1:
139          lowestpair = (0, 1)
140          closest = distance(clust[0].vec, clust[1].vec)
141
142      # loop through every pair looking for the smallest distance
143          for i in range(len(clust)):
144              for j in range(i + 1, len(clust)):
145          # distances is the cache of distance calculations
146                  if (clust[i].id, clust[j].id) not in distances:
147                      distances[(clust[i].id, clust[j].id)] = \
148                          distance(clust[i].vec, clust[j].vec)
149
150                  d = distances[(clust[i].id, clust[j].id)]
151
152                  if d < closest:
153                      closest = d
154                      lowestpair = (i, j)
155
156      # calculate the average of the two clusters
157          mergevec = [(clust[lowestpair[0]].vec[i] + clust[lowestpair
         [1]].vec[i])
158                      / 2.0 for i in range(len(clust[0].vec))]
```

```
159
160      # create the new cluster
161          newcluster = bicluster(mergevec, left=clust[lowestpair[0]],
162                              right=clust[lowestpair[1]], distance=
     closest,
163                              id=currentclustid)
164
165      # cluster ids that weren't in the original set are negative
166          currentclustid -= 1
167          del clust[lowestpair[1]]
168          del clust[lowestpair[0]]
169          clust.append(newcluster)
170
171      return clust[0]
172
173
174 def printclust(clust, labels=None, n=0):
175   # indent to make a hierarchy layout
176      for i in range(n):
177          print (' ', end =" ")
178      if clust.id < 0:
179      # negative id means that this is branch
180          print ('-')
181      else:
182      # positive id means that this is an endpoint
183          if labels == None:
184              print (clust.id)
185          else:
186              print (labels[clust.id])
187
188   # now print the right and left branches
189      if clust.left != None:
190          printclust(clust.left, labels=labels, n=n + 1)
191      if clust.right != None:
192          printclust(clust.right, labels=labels, n=n + 1)
193
194 """
195 Dendrogram
196 """
197 def getheight(clust):
198   # Is this an endpoint? Then the height is just 1
199      if clust.left == None and clust.right == None:
200          return 1
201
202   # Otherwise the height is the same of the heights of
203   # each branch
204      return getheight(clust.left) + getheight(clust.right)
```

```
205
206
207 def getdepth(clust):
208   # The distance of an endpoint is 0.0
209     if clust.left == None and clust.right == None:
210         return 0
211
212   # The distance of a branch is the greater of its two sides
213   # plus its own distance
214     return max(getdepth(clust.left), getdepth(clust.right)) + clust.
     distance
215
216 def drawdendrogram(clust, labels, jpeg='clusters.jpg'):
217   # height and width
218     h = getheight(clust) * 20
219     w = 1200
220     depth = getdepth(clust)
221
222   # width is fixed, so scale distances accordingly
223     scaling = float(w - 150) / depth
224
225   # Create a new image with a white background
226     img = Image.new('RGB', (w, h), (255, 255, 255))
227     draw = ImageDraw.Draw(img)
228
229     draw.line((0, h / 2, 10, h / 2), fill=(255, 0, 0))
230
231   # Draw the first node
232     drawnode(
233         draw,
234         clust,
235         10,
236         h / 2,
237         scaling,
238         labels,
239         )
240     img.save(jpeg, 'JPEG')
241
242 def drawnode(
243     draw,
244     clust,
245     x,
246     y,
247     scaling,
248     labels,
249     ):
250     if clust.id < 0:
```

```
251          h1 = getheight(clust.left) * 20
252          h2 = getheight(clust.right) * 20
253          top = y - (h1 + h2) / 2
254          bottom = y + (h1 + h2) / 2
255      # Line length
256          ll = clust.distance * scaling
257      # Vertical line from this cluster to children
258          draw.line((x, top + h1 / 2, x, bottom - h2 / 2), fill=(255, 0,
     0))
259
260      # Horizontal line to left item
261          draw.line((x, top + h1 / 2, x + ll, top + h1 / 2), fill=(255,
     0, 0))
262
263      # Horizontal line to right item
264          draw.line((x, bottom - h2 / 2, x + ll, bottom - h2 / 2), fill
     =(255, 0,
265                  0))
266
267      # Call the function to draw the left and right nodes
268          drawnode(
269              draw,
270              clust.left,
271              x + ll,
272              top + h1 / 2,
273              scaling,
274              labels,
275              )
276          drawnode(
277              draw,
278              clust.right,
279              x + ll,
280              bottom - h2 / 2,
281              scaling,
282              labels,
283              )
284      else:
285      # If this is an endpoint, draw the item label
286          draw.text((x + 5, y - 7), labels[clust.id], (0, 0, 0))
287 """
288 K-Means Clustering
289 """
290 def kcluster(rows,k,distance=pearson ):
291   # Determine the minimum and maximum values for each point
292     ranges = [(min([row[i] for row in rows]), max([row[i] for row in
     rows]))
293                for i in range(len(rows[0]))]
```

```
294
295    # Create k randomly placed centroids
296     clusters = [[random.random() * (ranges[i][1] - ranges[i][0]) +
     ranges[i][0]
297              for i in range(len(rows[0]))] for j in range(k)]
298
299     lastmatches = None
300     for t in range(100):
301         print ('Iteration %d' % t)
302         bestmatches = [[] for i in range(k)]
303
304     # Find which centroid is the closest for each row
305         for j in range(len(rows)):
306             row = rows[j]
307             bestmatch = 0
308             for i in range(k):
309                 d = distance(clusters[i], row)
310                 if d < distance(clusters[bestmatch], row):
311                     bestmatch = i
312             bestmatches[bestmatch].append(j)
313
314     # If the results are the same as last time, this is complete
315         if bestmatches == lastmatches:
316             break
317         lastmatches = bestmatches
318
319     # Move the centroids to the average of their members
320         for i in range(k):
321             avgs = [0.0] * len(rows[0])
322             if len(bestmatches[i]) > 0:
323                 for rowid in bestmatches[i]:
324                     for m in range(len(rows[rowid])):
325                         avgs[m] += rows[rowid][m]
326                 for j in range(len(avgs)):
327                     avgs[j] /= len(bestmatches[i])
328                 clusters[i] = avgs
329
330     return bestmatches
331
332 """
333 Q3
334 """
335 tweetdata, word,data =readfile("tweetdata.txt")
336 clust = hcluster(data)
337 #print(clust.vec)
338
339 """
```

```python
340 Q3  ASCIII
341 To view  cluster
342 """
343 printclust(clust,labels=tweetdata)
344
345
346 """
347 Q3 Dendrogram
348 """
349 drawdendrogram(clust,tweetdata,jpeg="three/tweetdata.jpeg")
350
351 """
352 Q4
353 """
354 """
355 For 5  kcluster
356 number of iteration:
357     Iteration 0
358     Iteration 1
359     Iteration 2
360     Iteration 3
361     Iteration 4
362     Iteration 5
363     Iteration 6
364     Iteration 7
365     Iteration 8
366     Iteration 9
367     Iteration 10
368     Iteration 11
369     Iteration 12
370 Cluster summary:
371     cluster  1 :   1
372     cluster  2 :   39
373     cluster  3 :   36
374     cluster  4 :   20
375     cluster  5 :   2
376 """
377
378
379 clust5 = kcluster(data,5)
380 for i in range(len(clust5)):
381   print ("cluster ", i+1, ": ", len(clust5[i]))
382   for r in clust5[i]:
383       f= 'four/clust5cluster'+ str(i+1)+'.csv'
384       with open(f, 'a+') as file:
385           file.write(tweetdata[r])
386           file.write("\n")
```

```
387
388 """
389 For 10  kcluster
390 number of iteration:
391     Iteration 0
392     Iteration 1
393     Iteration 2
394     Iteration 3
395 Cluster summary:
396     cluster  1 :   0
397     cluster  2 :   0
398     cluster  3 :   88
399     cluster  4 :   4
400     cluster  5 :   0
401     cluster  6 :   0
402     cluster  7 :   6
403     cluster  8 :   0
404     cluster  9 :   0
405     cluster  10 :   0
406 """
407
408
409 clust10 = kcluster(data,10)
410 for i in range(len(clust10)):
411   print ("cluster ", i+1, ": ", len(clust10[i]))
412   for r in clust10[i]:
413       f= 'four/clust10cluster'+ str(i+1)+'.csv'
414       with open(f, 'a+') as file:
415           file.write(tweetdata[r])
416           file.write("\n")
417
418 """
419 For 20  kcluster
420 number of iteration:
421     Iteration 0
422     Iteration 1
423     Iteration 2
424     Iteration 3
425     Iteration 4
426     Iteration 5
427     Iteration 6
428     Iteration 7
429     Iteration 8
430     Iteration 9
431     Iteration 10
432     Iteration 11
433     Iteration 12
```

```
434     Iteration 13
435     Iteration 14
436 clusters summary:
437     cluster  1 :   1
438     cluster  2 :   1
439     cluster  3 :   0
440     cluster  4 :   0
441     cluster  5 :   0
442     cluster  6 :   0
443     cluster  7 :   0
444     cluster  8 :   0
445     cluster  9 :   39
446     cluster  10 :   2
447     cluster  11 :   0
448     cluster  12 :   0
449     cluster  13 :   6
450     cluster  14 :   0
451     cluster  15 :   0
452     cluster  16 :   5
453     cluster  17 :   0
454     cluster  18 :   0
455     cluster  19 :   44
456     cluster  20 :   0
457 """
458
459
460 clust20 = kcluster(data,20)
461 for i in range(len(clust20)):
462   print ("cluster ", i+1, ": ", len(clust20[i]))
463   f= 'four/clust20cluster'+ str(i+1)+'.csv'
464   for r in clust20[i]:
465       with open(f, 'a+') as file:
466           file.write(tweetdata[r])
467           file.write("\n")
468
469 """
470 Q5='mds2d.jpg'
471 25602.68055067695
472 175394.90708243262
473 98
474 """
475
476
477 coords = scaledown(data)
478 print(len(coords))
479 draw2d(coords, tweetdata, jpeg='five/mds2d.jpg')
```

**Listing 6:** three.py

## Discussion

In getting the ASCII and Dendrogram the following had to be put in place:

- Line 339 to 350 is the drive for the code.

```
339 """
340 Q3  ASCIII
341 To view  cluster
342 """
343 printclust(clust,labels=tweetdata)
344
345
346 """
347 Q3 Dendrogram
348 """
349 drawdendrogram(clust,tweetdata,jpeg="three/tweetdata.jpeg")
```

**Listing 7:** Driver code for Q3

- Used the readfile() function to read in the text, it gets all data and returns the user-names as row names, the words as col names, and the value as data(data stored as a 2D array of float values).

- Then the data is passed in hcluster(data) also known as hierarchical clustering, this function used the pearson() function when called.

- The pearson function takes two vector arrays as arguments and then returns the pearson correlation between these values.

- hcluster() function also uses a bi-cluster class which is an helper for comparing two clusters.

- hcluster does the hierarchical clustering by agglomerative clustering which involves merging the best matches cluster with a new cluster. It repeats this cycle untill there is just one cluster left.

- printclust() function prints recursively the final end product of the hclusters returned function.The printcluster has an argument called labels, the label when removed prints the cluster based on the position number of accounts names but when the label is supplied it prints the actual name of the labels. Print produces the ASCII text of the output. It is saved in as ascii.txt in /three.

- drawdendrogram() does the similar objective as the printcluster but it produces a jpeg format of the recursive output of the cluster.
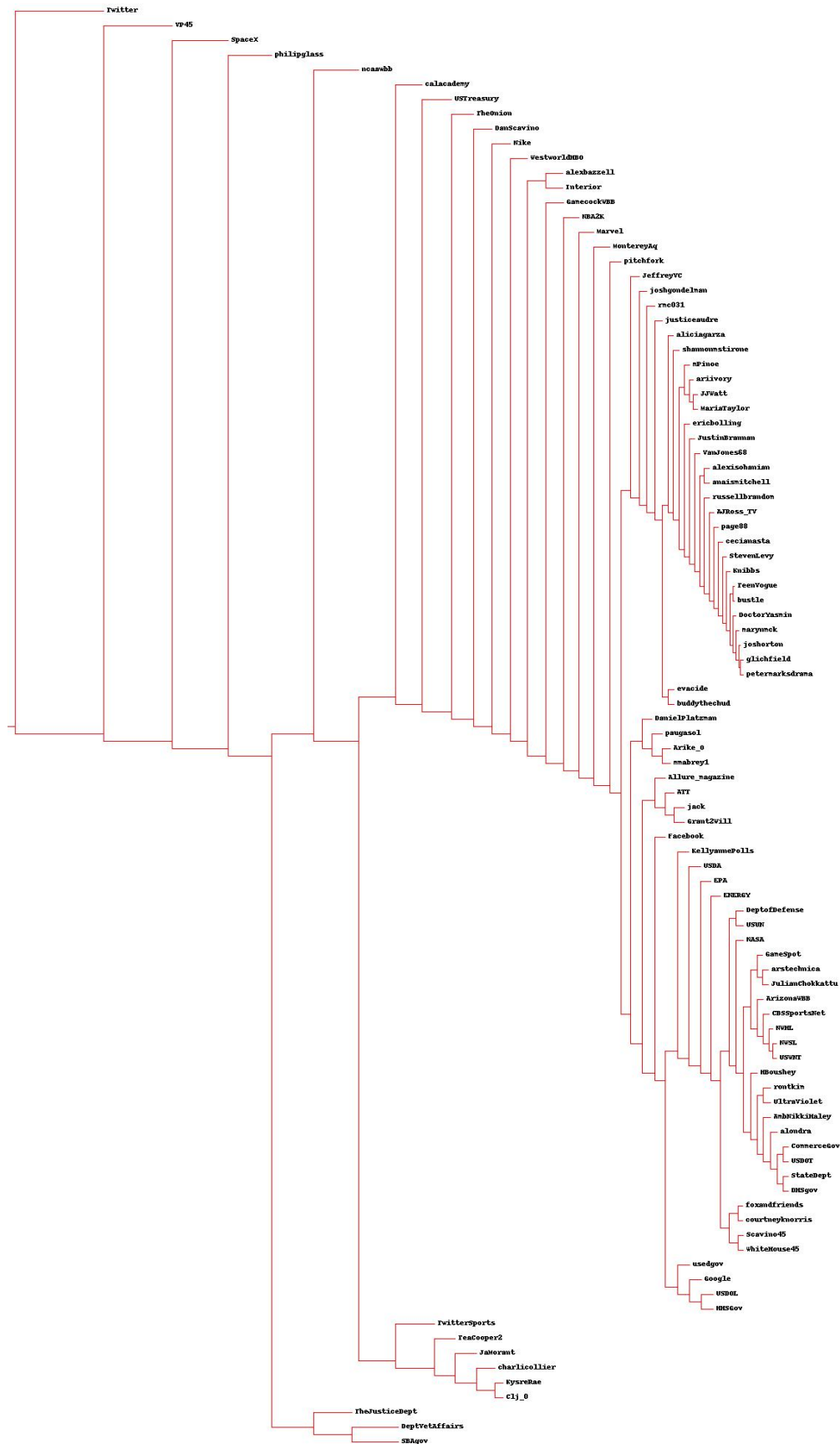
**Figure 1:** Question 3 output

# Q4

## Answer

```
351 """
352 Q4
353 """
354 """
355 For 5  kcluster
356 number of iteration:
357     Iteration 0
358     Iteration 1
359     Iteration 2
360     Iteration 3
361     Iteration 4
362     Iteration 5
363     Iteration 6
364     Iteration 7
365     Iteration 8
366     Iteration 9
367     Iteration 10
368     Iteration 11
369     Iteration 12
370 Cluster summary:
371     cluster  1 :   1
372     cluster  2 :   39
373     cluster  3 :   36
374     cluster  4 :   20
375     cluster  5 :   2
376 """
377
378
379 clust5 = kcluster(data,5)
380 for i in range(len(clust5)):
381   print ("cluster ", i+1, ": ", len(clust5[i]))
382   for r in clust5[i]:
383       f= 'four/clust5cluster'+ str(i+1)+'.csv'
384       with open(f, 'a+') as file:
385           file.write(tweetdata[r])
386           file.write("\n")
387
388 """
389 For 10  kcluster
390 number of iteration:
391     Iteration 0
392     Iteration 1
```

```
393     Iteration 2
394     Iteration 3
395 Cluster summary:
396     cluster  1 :   0
397     cluster  2 :   0
398     cluster  3 :   88
399     cluster  4 :   4
400     cluster  5 :   0
401     cluster  6 :   0
402     cluster  7 :   6
403     cluster  8 :   0
404     cluster  9 :   0
405     cluster  10 :   0
406 """
407
408
409 clust10 = kcluster(data,10)
410 for i in range(len(clust10)):
411   print ("cluster ", i+1, ": ", len(clust10[i]))
412   for r in clust10[i]:
413       f= 'four/clust10cluster'+ str(i+1)+'.csv'
414       with open(f, 'a+') as file:
415           file.write(tweetdata[r])
416           file.write("\n")
417
418 """
419 For 20  kcluster
420 number of iteration:
421     Iteration 0
422     Iteration 1
423     Iteration 2
424     Iteration 3
425     Iteration 4
426     Iteration 5
427     Iteration 6
428     Iteration 7
429     Iteration 8
430     Iteration 9
431     Iteration 10
432     Iteration 11
433     Iteration 12
434     Iteration 13
435     Iteration 14
436 clusters summary:
437     cluster  1 :   1
438     cluster  2 :   1
439     cluster  3 :   0
```

```
440      cluster   4 :    0
441      cluster   5 :    0
442      cluster   6 :    0
443      cluster   7 :    0
444      cluster   8 :    0
445      cluster   9 :   39
446      cluster  10 :    2
447      cluster  11 :    0
448      cluster  12 :    0
449      cluster  13 :    6
450      cluster  14 :    0
451      cluster  15 :    0
452      cluster  16 :    5
453      cluster  17 :    0
454      cluster  18 :    0
455      cluster  19 :   44
456      cluster  20 :    0
457 """
458
459
460 clust20 = kcluster(data,20)
461 for i in range(len(clust20)):
462   print ("cluster ", i+1, ": ", len(clust20[i]))
463   f= 'four/clust20cluster'+ str(i+1)+'.csv'
464   for r in clust20[i]:
465       with open(f, 'a+') as file:
466           file.write(tweetdata[r])
467           file.write("\n")
```

**Listing 8:** Driver code for Q4

## Discussion

- Q: Give a brief explanation of how the k-Means algorithm operates on this data. What features is the algorithm considering?

  Q: How many iterations were required for each value of k?

  Q: Which k value created the most reasonable clusters? For that grouping, characterize the accounts that were clustered into each group.

- K-mean algorithm work in the manner.

- in the kcluster function, it first try to organize the data in ranges of minimum and maximum values for each row so that the clusters can be represents as a coordinate in 2d form.

- The last part is update the clusters average(mean) to their new members, this changes the location of the clusters and groups them together and closer based on the average of all the group members.

- Final check is to make sure that the best is the matches is the list of rows in each clusters. This ensure the group members are more close together on average.

- kcluster = 5 had 12 iteration starting from zero and produced a total of 5 clusters.

```
354 """
355 For 5  kcluster
356 number of iteration:
357     Iteration 0
358     Iteration 1
359     Iteration 2
360     Iteration 3
361     Iteration 4
362     Iteration 5
363     Iteration 6
364     Iteration 7
365     Iteration 8
366     Iteration 9
367     Iteration 10
368     Iteration 11
369     Iteration 12
370 Cluster summary:
371     cluster  1 :  1
372     cluster  2 :  39
373     cluster  3 :  36
374     cluster  4 :  20
375     cluster  5 :  2
376 """
```

Listing 9: kcluster 5

- kcluster = 10 had 3 iteration starting from zero and produced a total of 10 clusters.

```
388 """
389 For 10  kcluster
390 number of iteration:
391     Iteration 0
392     Iteration 1
393     Iteration 2
394     Iteration 3
395 Cluster summary:
396     cluster  1 :  0
397     cluster  2 :  0
398     cluster  3 :  88
399     cluster  4 :  4
400     cluster  5 :  0
401     cluster  6 :  0
402     cluster  7 :  6
403     cluster  8 :  0
```

```
404     cluster  9 :  0
405     cluster  10 :  0
406 """
```

**Listing 10:** kcluster 10

- kcluster = 20 had 14 iteration starting from zero and produced a total of 20 clusters.

```
418 """
419 For 20  kcluster
420 number of iteration:
421     Iteration 0
422     Iteration 1
423     Iteration 2
424     Iteration 3
425     Iteration 4
426     Iteration 5
427     Iteration 6
428     Iteration 7
429     Iteration 8
430     Iteration 9
431     Iteration 10
432     Iteration 11
433     Iteration 12
434     Iteration 13
435     Iteration 14
436 clusters summary:
437     cluster  1 :  1
438     cluster  2 :  1
439     cluster  3 :  0
440     cluster  4 :  0
441     cluster  5 :  0
442     cluster  6 :  0
443     cluster  7 :  0
444     cluster  8 :  0
445     cluster  9 :  39
446     cluster  10 :  2
447     cluster  11 :  0
448     cluster  12 :  0
449     cluster  13 :  6
450     cluster  14 :  0
451     cluster  15 :  0
452     cluster  16 :  5
453     cluster  17 :  0
454     cluster  18 :  0
455     cluster  19 :  44
456     cluster  20 :  0
457 """
```

**Listing 11:** kcluster 20

- From all the list of clusters made I would say none of them satisfied a uniform grouping but if I would have to choose k equals 10 is better.

- The code for each k values of 5,10 and 20 writes the twitter screen names row by row for each cluster based on the cluster summary. Each cluster gets a particular file name.

```
379 clust5 = kcluster(data,5)
380 for i in range(len(clust5)):
381   print ("cluster ", i+1, ": ", len(clust5[i]))
382   for r in clust5[i]:
383       f= 'four/clust5cluster'+ str(i+1)+'.csv'
384       with open(f, 'a+') as file:
385            file.write(tweetdata[r])
386            file.write("\n")
```

**Listing 12:** kcluster 20

# Q5

## Answer

The mds2d.jpg is added in the next page.

## Discussion

- Q: How many iterations were required?

  Q: How well did the MDS do in grouping similar accounts together? Were there any particularly odd groupings?

- This is the resulting output.

- There were 98 iteration in total by checking the length of coord in line 478.

```
477 coords = scaledown(data)
478 print(len(coords))
479 draw2d(coords, tweetdata, jpeg='five/mds2d.jpg')
```

**Listing 13:** Driver for Q5

- MDS did fairly good in grouping similar accounts together. Considering philipglass and SBAgov which is grouped together, actually is a odd category.

- For the code to work, we would have to read in the tweetdata.txt and parse in the data variable into scaledown(data) function. The functions are from the lecture notes. The function scaledown then creates the file mds2d.jpg Figure 2.

**Figure 2:** Question 5 output

# References

– https://www.geeksforgeeks.org/python-pandas-dataframe-drop_
  duplicates/

– https://www.kite.com/python/answers/how-to-write-contents-
  of-a-dataframe-into-a-text-file-in-python

– https://stackoverflow.com/questions/48230230/typeerror-mismatch-
  between-array-dtype-object-and-format-specifier-18e/48231106

– https://www.geeksforgeeks.org/remove-last-n-rows-of-a-pandas-
  dataframe/

– https://stackoverflow.com/questions/10897339/python-fetch-
  first-10-results-from-a-list

– https://www.geeksforgeeks.org/how-to-add-one-row-in-an-existing-
  pandas-dataframe/

– https://stackoverflow.com/questions/22412258/get-the-first-
  element-of-each-tuple-in-a-list-in-python