# HW8 - Clustering
Adeniran Adeniyi
Sunday, April 18, 2021 by 11:59pm

# Q1

*Generate a list of 100 popular accounts on Twitter. The accounts must be verified, have < 10,000 followers, and have > 5000 tweets.*

*For example: https://twitter.com/weiglemc - not verified, 457 followers, 2554 tweets - don't include https://twitter.com/WNBA - verified (blue checkmark), 672,800+ followers, 77,900+ tweets - could include*
*See GET users/lookup, Tweepy's API.lookup_users(), 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.*

*How did you choose to collect the accounts?*

*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?*

## Answer

```python
# tweetparser.py
# MCW - 4/1/2021
import json
import tweepy
import pandas as pd
import numpy as np

def confirm_usermetsrequiremetn(twAuth,ids):
    #Find followers that are in this category and them to the list
    status= False
    try:
```

```
12          a = twAuth.get_user(id = ids)
13          if a.statuses_count >= 5000 and a.followers_count >= 10000 and
     a.verified == True:
14              status = True
15          else:
16              status = False
17      except:
18          print("There was an error")
19      return status
20 def setup_api(filename):
21      '''
22      filename: file where Twitter API keys are stored
23      returns Twitter API object to pass into parse()
24      '''
25
26      # load Twitter API keys from a file so they're not hard-coded
27      with open(filename, "r") as secretsFile:
28          secrets = json.load(secretsFile)
29
30      # set up Twitter API with OAuth2 procedure
31      consumer_key = secrets['consumer_key']
32      consumer_secret = secrets['consumer_secret']
33      try:
34          auth = tweepy.AppAuthHandler(consumer_key, consumer_secret)
35          api = tweepy.API(auth,wait_on_rate_limit=True,
     wait_on_rate_limit_notify=True)
36          return api
37      except tweepy.TweepError as e:
38          print ("Tweepy Error: %s" % str(e))
39
40 def parse(api, screen_name, num_tweets=200):
41      '''
42      api: Twitter API object, use setup_api() to create
43      screen_name: Twitter screen_name
44      num_tweets: Number of tweets to request (default: 200)
45      returns dict with {'screen_name': screen_name, 'tweets': [tweet1,
     tweet2, ...]}
46      '''
47
48      tweet_data = []
49      try:
50          for tweet in tweepy.Cursor(api.user_timeline, screen_name=
     screen_name, count=num_tweets, lang='en',
51                                     tweet_mode='extended', exclude_replies=
     True, include_rts=False).items():
52              tweet_data.append(tweet.full_text)
53      except tweepy.TweepError as e:
```

```
54          print ("Tweepy Error: %s" % str(e))
55
56      account_data = {'screen_name': screen_name, 'tweets': tweet_data}
57      return account_data
```

**Listing 1:** tweetparser.py

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Sat Apr 17 13:02:33 2021
4
5  @author: aadeniran
6  """
7  import tweepy
8  from tweetparser import setup_api, confirm_usermetsrequiremetn
9  import pandas as pd
10 import numpy as np
11 def  per_account(types,secret =setup_api("secrets.json") ):
12
13     #Get twitter friends screen_name of the parse account, as one
       arrayList
14     public_tweets = tweepy.Cursor(secret.friends_ids, id = types)
15     count= 1
16     resultList =[] # store the final result of screen names
17     for user in public_tweets.pages():
18         #print(user)
19         #print(user.dtypes())
20         #Transverse the list of followers ids
21         for i in user:
22             #check if the friends meet the requirement
23             #print(confirm_usermetsrequiremetn(secret,i))
24             if(confirm_usermetsrequiremetn(secret,i) ==True):
25                 if( i not in resultList):
26
27                     #instead of the user id get the user name
28                     user_sc = secret.get_user(i)
29                     print("{}:{}".format(count,user_sc.screen_name))
30                     resultList.append(user_sc.screen_name)
31                     count += 1
32             #Get maximum 30 file accounts screen_names
33             if(count >= 30):
34                 break
35     print(resultList)
36
37     #build text file and save file in Q1 directory
38     filename =  'Q1/' + types + '.txt'
39     with open(filename, 'w') as filehandle:
40         for listitem in resultList:
```

```python
41              filehandle.write('%s\n' % listitem)
42      return resultList
43 """
44 Tech account = @WIRED
45 Sport account = @WNBA
46 politics = @POTUS45
47 music = @future_of_music
48 """
49 #Get all screen_names with that fulfils 10,000 followers and have 5000
       tweets and verified
50 per_account("WIRED")
51 per_account("WNBA")
52 per_account("POTUS45")
53 per_account("future_of_music")
54
55 """
56 Bring all result from the files in to One text file
57 accounts.txt
58 """
59 #get the list of the created files
60 column_name= ["User_screen_names"]
61 final = pd.DataFrame(columns= column_name)
62 fileList = ["Q1/WIRED.txt","Q1/WNBA.txt","Q1/POTUS45.txt","Q1/
       future_of_music.txt"]
63 df = pd.DataFrame()
64
65
66 for t in fileList:
67      frame = pd.read_csv(t,header=None)
68      frame.columns = column_name
69      for ind in frame.index:
70          final.loc[len(final)] =[frame['User_screen_names'][ind]]
71
72 #drop duplicated values keep just one of them
73 # dropping duplicate values
74 final.User_screen_names.drop_duplicates(inplace=True)
75
76 #confirm that there are all unique values there are 113
77 print(final.User_screen_names.nunique())
78
79 # Number of rows to drop
80 n = 14
81
82 # Dropping last n rows using drop
83 final.drop(final.tail(n).index,
84          inplace = True)
85 #print(final.User_screen_names.nunique())
```

```
86 numpy_array = final.to_numpy()
87 #print as a text file to in the same directory for future use
88 np.savetxt(r'accounts.txt', numpy_array,fmt="%s")
```

<center>**Listing 2:** gatherId.py</center>

## Discussion

*In gathering users screen names I did took into consideration of the following:*

- I collected the screen names going through 4 major popular twitter accounts WIRED, WNBA, POTUS45, and future_of_music

```
49 #Get all screen_names with that fulfils 10,000 followers and have
       5000 tweets and verified
50 per_account("WIRED")
51 per_account("WNBA")
52 per_account("POTUS45")
53 per_account("future_of_music")
```

**Listing 3:** A snap shot of gatherIds.py show these accounts doing screen_name account extraction

following list, I extracted accounts the met the requirement of being verified, have > 10,000 followers, and have > 5000 tweets

- The function confirm_usermetsrequiremetn() in tweetparser.py on line 8 to 19 check to make sure the accounts fulfilled these requirement

```
8 def confirm_usermetsrequiremetn(twAuth,ids):
9     #Find followers that are in this category and them to the list
10    status= False
11    try:
12        a = twAuth.get_user(id = ids)
13        if a.statuses_count >= 5000 and a.followers_count >= 10000
      and a.verified == True:
14            status = True
15        else:
16            status = False
17    except:
18        print("There was an error")
19    return status
```

**Listing 4:** A snap shot of tweetparser.py that check screen_name account meets requirements

Imports of the file was made on line 8 of gatherIds.py

```
7 import tweepy
8 from tweetparser import setup_api, confirm_usermetsrequiremetn
9 import pandas as pd
```

```
10 import numpy as np
```

**Listing 5:** A snap shot of gatherIds.py that shows the needed imports

- Function per_account() produce a text files (names based on the arguement supplied ) that gets stored in Q1 folder

- I then merge all the produced files from each of the accounts to form the main accounts.txt file. Seen in line 55 - 85

```
55 """
56 Bring all result from the files in to One text file
57 accounts.txt
58 """
59 #get the list of the created files
60 column_name= ["User_screen_names"]
61 final = pd.DataFrame(columns= column_name)
62 fileList = ["Q1/WIRED.txt","Q1/WNBA.txt","Q1/POTUS45.txt","Q1/
      future_of_music.txt"]
63 df = pd.DataFrame()
64
65
66 for t in fileList:
67     frame = pd.read_csv(t,header=None)
68     frame.columns = column_name
69     for ind in frame.index:
70         final.loc[len(final)] =[frame['User_screen_names'][ind]]
71
72 #drop duplicated values keep just one of them
73 # dropping duplicate values
74 final.User_screen_names.drop_duplicates(inplace=True)
75
76 #confirm that there are all unique values there are 113
77 print(final.User_screen_names.nunique())
78
79 # Number of rows to drop
80 n = 14
81
82 # Dropping last n rows using drop
83 final.drop(final.tail(n).index,
84         inplace = True)
85 #print(final.User_screen_names.nunique())
```

**Listing 6:** A snap shot of gatherIds.py shows the merging process of the files generated by each account

# Q2

*Create Account-Term Matrix*

## Answer

```python
1  # generatetweetvector.py
2  # Based on generatefeedvector.py from
3  # https://github.com/arthur-e/Programming-Collective-Intelligence/blob/
      master/chapter3/generatefeedvector.py
4
5  from tweetparser import setup_api, parse
6  import re
7
8  def getwordcounts(api, screen_name):
9      """
10     api: Twitter API object
11     screen_name: Twitter screen_name
12     returns screen_name and dictionary of word counts for a Twitter
       account
13     """
14
15     # Parse the Twitter feed
16     d = parse(api, screen_name)
17     wc = {}
18
19     # Loop over all the entries
20     for tweet in d['tweets']:
21
22         # Extract a list of words
23         words = getwords(tweet)
24         for word in words:
25             wc.setdefault(word, 0)
26             wc[word] += 1
27
28     return (d['screen_name'], wc)
29
30 def getwords(tweet):
31     """
32     returns lowercase list of words after filtering
33     """
34
35     # Remove URLs
36     text = re.compile(r'(http://|https://|www\.)([^ \'\"]*)').sub('',
       tweet)
```

```python
37        "
38        # Remove other screen names (start with @)
39        text = re.compile(r'(@\w+)').sub('', text)
40
41        # Split words by all non-alpha characters
42        words = re.compile(r'[^A-Z^a-z]+').split(text)
43
44        # Filter for words between 3-15 characters, convert to lowercase,
       and return as a list
45        return [word.lower() for word in words if (len(word) >= 3 and len(
       word) <= 15)]
46
47 # MAIN CODE STARTS HERE
48
49 # set up Twitter API object
50 api = setup_api("secrets.json")
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
      top 1000)
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 """
```

```
79  Did some testing to figure out the datatype
80  and sample output for one acccount
81  Spyder ide  gives flexibility of running each variable without running
       the whole code after your ran the whole code at least once
82  """
83  #print("Done Counting words")
84  #print(type(sumcounts))
85  #print(sumcounts.keys())
86  #for elem in listofTuples :
87  #    print(elem)
88
89
90  # remove stopwords ("fake" way)
91  wordlist = []
92  for (w, ac) in apcount.items():
93      # w is the word, ac is the account count (was bc 'blog count' in
        textbook)
94      frac = float(ac) / len(accountlist)
95      if frac > 0.1 and frac < 0.5:
96          wordlist.append(w)
97
98  # Save only the 1000 most frequent terms over all accounts (see
       sumcounts) in popularlist
99  popularlist = []
100 #
101 # INSERT YOUR CODE HERE
102 #
103 # Create a list of tuples sorted by index 1 i.e. value field
104
105 listofTuples = sorted(sumcounts.items() ,  key=lambda x: x[1],reverse=
       True)
106 #extract only 1000 rows
107 listofTuples = listofTuples[:1000]
108 #store in the dictionary
109 popularlist = [i[0] for i in listofTuples]
110
111 print("Writing it to text file")
112 # write out popular word list
113 with open('popularlist.txt', 'w') as outf:
114     for word in popularlist:
115         outf.write(word + '\n')
116
117 # write out account-term matrix
118 with open('tweetdata.txt', 'w') as outf:
119     # write header row ("Account", list of words)
120     outf.write('Account')
121     for word in popularlist:
```

```
122            outf.write('\t%s' % word)
123        outf.write('\n')
124
125        # write each row (screen_name, count for each word)
126        for (screen_name, wc) in wordcounts.items():
127            outf.write(screen_name)
128            for word in popularlist:
129                if word in wc:
130                    outf.write('\t%d' % wc[word])
131                else:
132                    outf.write('\t0')
133            outf.write('\n')
```

**Listing 7:** generatetweetvector.py

## Discussion

*The code drive started from line 47 to line 133*

- Getting the twitter api authentication at line 50, I passeding secrets.json which contained my twitter consumer_key and consummer_secret

- Read the account files containing the list of scree_names while removeing any spaces each account name

- 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. I noticed that each works and users account are store 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.

- I had done some testing to understand the result for one user in lines 78 to 87

```
78  """
79  Did some testing to figure out the datatype
80  and sample output for one acccount
81  Spyder ide  gives flexibility of running each variable without
        running the whole code after your ran the whole code at least
        once
82  """
83  #print("Done Counting words")
84  #print(type(sumcounts))
85  #print(sumcounts.keys())
86  #for elem in listofTuples :
87  #    print(elem)
```

**Listing 8:** A snapshot for generatetweetvector.py trying to test the output of sumcounts variable listing out the keys values

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

- For popularlist, From my deduced test above, I was able to sort the item using a lambda function an placed the result in a list of tuples.It goes from highest to lowest. Next I retried just 1000 row items from the tuples(by slicing) and stored it in the dictionary variable popularlist.

```
105  listofTuples = sorted(sumcounts.items() ,  key=lambda x: x[1],
         reverse=True)
106  #extract only 1000 rows
107  listofTuples = listofTuples[:1000]
108  #store in the dictionary
109  popularlist = [i[0] for i in listofTuples]
```

**Listing 9:** A snapshot for generatetweetvector.py Sort tuples Get 1000 rows Save in dictionary variable 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.

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

# Q3

*Create Dendrogram*

*Create an ASCII dendrogram and a JPEG dendrogram that uses hierarchical clustering to cluster the most similar accounts (see Module 12, slides 21, 23). Include the JPEG in your report and upload the ASCII file to GitHub (it will be too unwieldy for inclusion in the report).*

*How well did the hierarchical clustering do in grouping similar accounts together? Were there any particularly odd groupings?*

## Answer

```python
1  # -*- coding: utf-8 -*-
2  """
3  Created on Sat Apr 17 21:07:52 2021
4
5  @author: aadeniran
6  """
7
8  from PIL import Image, ImageDraw
9  from math import sqrt
10 import random
11 import csv
12 import pandas as pd
13 def readfile(filename):
14   # This is a function that I wrote from scratch    -MCW
15   data = []
16   rownames = []
17   colnames = []
18   num_rows = 0
19   with open(filename) as tsvfile:
20     reader = csv.reader(tsvfile, delimiter='\t')
21     for row in reader:
22       if num_rows > 0:
23         rownames.append(row[0])    # save the row names
24         data.append([float(x) for x in row[1:]])  # save the values as
       floats
25       else:
26         for col in row[1:]:
27           colnames.append(col)    # save the column names
28       num_rows = num_rows + 1
29   return (rownames, colnames, data)
30
31 """
32 names, word, data= readfiles('tweetdata.txt')
33 """
34
35 def pearson(v1, v2):
```

```python
36    # Simple sums
37       sum1 = sum(v1)
38       sum2 = sum(v2)
39
40    # Sums of the squares
41       sum1Sq = sum([pow(v, 2) for v in v1])
42       sum2Sq = sum([pow(v, 2) for v in v2])
43
44    # Sum of the products
45      pSum = sum([v1[i] * v2[i] for i in range(len(v1))])
46
47    # Calculate r (Pearson score)
48      num = pSum - sum1 * sum2 / len(v1)
49      den = sqrt((sum1Sq - pow(sum1, 2) / len(v1)) * (sum2Sq - pow(sum2,
      2)
50                  / len(v1)))
51      if den == 0:
52          return 0
53
54      return 1.0 - num / den
55 """
56 MD5 Scaling
57 """
58 def scaledown(data, distance=pearson, rate=0.01):
59     n = len(data)
60
61   # The real distances between every pair of items
62     realdist = [[distance(data[i], data[j]) for j in range(n)] for i in
63                 range(0, n)]
64
65   # Randomly initialize the starting points of the locations in 2D
66     loc = [[random.random(), random.random()] for i in range(n)]
67     fakedist = [[0.0 for j in range(n)] for i in range(n)]
68
69     lasterror = None
70     for m in range(0, 1000):
71     # Find projected distances
72         for i in range(n):
73             for j in range(n):
74                 fakedist[i][j] = sqrt(sum([pow(loc[i][x] - loc[j][x],
      2)
75                                           for x in range(len(loc[i]))]))
76
77     # Move points
78         grad = [[0.0, 0.0] for i in range(n)]
79
80         totalerror = 0
```

```python
 81            for k in range(n):
 82                for j in range(n):
 83                    if j == k:
 84                        continue
 85            # The error is percent difference between the distances
 86                    errorterm = (fakedist[j][k] - realdist[j][k]) /
       realdist[j][k]
 87
 88            # Each point needs to be moved away from or towards the other
 89            # point in proportion to how much error it has
 90                    grad[k][0] += (loc[k][0] - loc[j][0]) / fakedist[j][k]
       \
 91                        * errorterm
 92                    grad[k][1] += (loc[k][1] - loc[j][1]) / fakedist[j][k]
       \
 93                        * errorterm
 94
 95            # Keep track of the total error
 96                    totalerror += abs(errorterm)
 97            print (totalerror)
 98
 99        # If the answer got worse by moving the points, we are done
100            if lasterror and lasterror < totalerror:
101                break
102            lasterror = totalerror
103
104        # Move each of the points by the learning rate times the gradient
105            for k in range(n):
106                loc[k][0] -= rate * grad[k][0]
107                loc[k][1] -= rate * grad[k][1]
108
109        return loc
110
111 def draw2d(data, labels, jpeg):
112     img = Image.new('RGB', (2000, 2000), (255, 255, 255))
113     draw = ImageDraw.Draw(img)
114     for i in range(len(data)):
115         x = (data[i][0] + 0.5) * 1000
116         y = (data[i][1] + 0.5) * 1000
117         draw.text((x, y), labels[i], (0, 0, 0))
118     img.save(jpeg, 'JPEG')
119
120 def rotatematrix(data):
121     newdata = []
122     for i in range(len(data[0])):
123         newrow = [data[j][i] for j in range(len(data))]
124         newdata.append(newrow)
```

```
125      return newdata
126
127  """
128  Hierarchical Clustering
129  class bicluster - data structure to hold the clustering information
130  hcluster(rows, distance=pearson) - does the hierarchical clustering,
         default distance function is pearson()
131  printclust(clust, labels=None, n=0) - traverses the cluster and prints
         an ASCII text representation
132  """
133  class bicluster:
134
135      def __init__(self, vec, left=None, right=None, distance=0.0, id=
         None,):
136          self.left = left
137          self.right = right
138          self.vec = vec
139          self.id = id
140          self.distance = distance
141
142  def hcluster(rows, distance=pearson):
143      distances = {}
144      currentclustid = -1
145
146    # Clusters are initially just the rows
147      clust = [bicluster(rows[i], id=i) for i in range(len(rows))]
148
149      while len(clust) > 1:
150          lowestpair = (0, 1)
151          closest = distance(clust[0].vec, clust[1].vec)
152
153      # loop through every pair looking for the smallest distance
154          for i in range(len(clust)):
155              for j in range(i + 1, len(clust)):
156          # distances is the cache of distance calculations
157                  if (clust[i].id, clust[j].id) not in distances:
158                      distances[(clust[i].id, clust[j].id)] = \
159                          distance(clust[i].vec, clust[j].vec)
160
161                  d = distances[(clust[i].id, clust[j].id)]
162
163                  if d < closest:
164                      closest = d
165                      lowestpair = (i, j)
166
167      # calculate the average of the two clusters
```

```python
168          mergevec = [(clust[lowestpair[0]].vec[i] + clust[lowestpair
     [1]].vec[i])
169                      / 2.0 for i in range(len(clust[0].vec))]
170
171      # create the new cluster
172          newcluster = bicluster(mergevec, left=clust[lowestpair[0]],
173                              right=clust[lowestpair[1]], distance=
     closest,
174                              id=currentclustid)
175
176      # cluster ids that weren't in the original set are negative
177          currentclustid -= 1
178          del clust[lowestpair[1]]
179          del clust[lowestpair[0]]
180          clust.append(newcluster)
181
182      return clust[0]
183
184
185 def printclust(clust, labels=None, n=0):
186   # indent to make a hierarchy layout
187      for i in range(n):
188          print (' ', end =" ")
189      if clust.id < 0:
190      # negative id means that this is branch
191          print ('-')
192      else:
193      # positive id means that this is an endpoint
194          if labels == None:
195              print (clust.id)
196          else:
197              print (labels[clust.id])
198
199   # now print the right and left branches
200      if clust.left != None:
201          printclust(clust.left, labels=labels, n=n + 1)
202      if clust.right != None:
203          printclust(clust.right, labels=labels, n=n + 1)
204
205 """
206 Dendrogram
207 """
208 def getheight(clust):
209   # Is this an endpoint? Then the height is just 1
210      if clust.left == None and clust.right == None:
211          return 1
212
```

```python
213    # Otherwise the height is the same of the heights of
214    # each branch
215      return getheight(clust.left) + getheight(clust.right)
216
217
218 def getdepth(clust):
219    # The distance of an endpoint is 0.0
220      if clust.left == None and clust.right == None:
221          return 0
222
223    # The distance of a branch is the greater of its two sides
224    # plus its own distance
225      return max(getdepth(clust.left), getdepth(clust.right)) + clust.
     distance
226
227 def drawdendrogram(clust, labels, jpeg='clusters.jpg'):
228    # height and width
229      h = getheight(clust) * 20
230      w = 1200
231      depth = getdepth(clust)
232
233    # width is fixed, so scale distances accordingly
234      scaling = float(w - 150) / depth
235
236    # Create a new image with a white background
237      img = Image.new('RGB', (w, h), (255, 255, 255))
238      draw = ImageDraw.Draw(img)
239
240      draw.line((0, h / 2, 10, h / 2), fill=(255, 0, 0))
241
242    # Draw the first node
243      drawnode(
244          draw,
245          clust,
246          10,
247          h / 2,
248          scaling,
249          labels,
250          )
251      img.save(jpeg, 'JPEG')
252
253 def drawnode(
254      draw,
255      clust,
256      x,
257      y,
258      scaling,
```

```
259      labels,
260      ):
261      if clust.id < 0:
262          h1 = getheight(clust.left) * 20
263          h2 = getheight(clust.right) * 20
264          top = y - (h1 + h2) / 2
265          bottom = y + (h1 + h2) / 2
266      # Line length
267          ll = clust.distance * scaling
268      # Vertical line from this cluster to children
269          draw.line((x, top + h1 / 2, x, bottom - h2 / 2), fill=(255, 0,
     0))
270
271      # Horizontal line to left item
272          draw.line((x, top + h1 / 2, x + ll, top + h1 / 2), fill=(255,
     0, 0))
273
274      # Horizontal line to right item
275          draw.line((x, bottom - h2 / 2, x + ll, bottom - h2 / 2), fill
     =(255, 0,
276                    0))
277
278      # Call the function to draw the left and right nodes
279          drawnode(
280              draw,
281              clust.left,
282              x + ll,
283              top + h1 / 2,
284              scaling,
285              labels,
286              )
287          drawnode(
288              draw,
289              clust.right,
290              x + ll,
291              bottom - h2 / 2,
292              scaling,
293              labels,
294              )
295      else:
296      # If this is an endpoint, draw the item label
297          draw.text((x + 5, y - 7), labels[clust.id], (0, 0, 0))
298 """
299 K-Means Clustering
300 """
301 def kcluster(rows,k,distance=pearson ):
302   # Determine the minimum and maximum values for each point
```

```
303      ranges = [(min([row[i] for row in rows]), max([row[i] for row in
     rows]))
304              for i in range(len(rows[0]))]
305
306   # Create k randomly placed centroids
307     clusters = [[random.random() * (ranges[i][1] - ranges[i][0]) +
     ranges[i][0]
308              for i in range(len(rows[0]))] for j in range(k)]
309
310     lastmatches = None
311     for t in range(100):
312         print ('Iteration %d' % t)
313         bestmatches = [[] for i in range(k)]
314
315     # Find which centroid is the closest for each row
316         for j in range(len(rows)):
317             row = rows[j]
318             bestmatch = 0
319             for i in range(k):
320                 d = distance(clusters[i], row)
321                 if d < distance(clusters[bestmatch], row):
322                     bestmatch = i
323             bestmatches[bestmatch].append(j)
324
325     # If the results are the same as last time, this is complete
326         if bestmatches == lastmatches:
327             break
328         lastmatches = bestmatches
329
330     # Move the centroids to the average of their members
331         for i in range(k):
332             avgs = [0.0] * len(rows[0])
333             if len(bestmatches[i]) > 0:
334                 for rowid in bestmatches[i]:
335                     for m in range(len(rows[rowid])):
336                         avgs[m] += rows[rowid][m]
337                 for j in range(len(avgs)):
338                     avgs[j] /= len(bestmatches[i])
339                 clusters[i] = avgs
340
341     return bestmatches
342
343 """
344 Question 3
345 """
346 #running the code
347 tweetdata, word,data =readfile("tweetdata.txt")
```
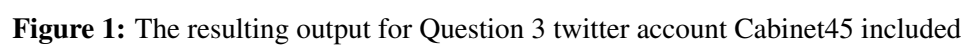
```
348 #clust = hcluster(data)
349 #print(clust.vec)
350
351 """
352 Question 3  ASCIII
353 To view  cluster
354
355 printclust(clust,labels=tweetdata)
356 """
357 """
358 Question 3 Dendrogram
359
360
361 drawdendrogram(clust,tweetdata,jpeg="Q3/tweetdata.jpeg")
362
363 """
364 """
365 Question 4
366 """
367 """
368 For 5  kcluster
369 number of iteration:
370     Iteration 0
371     Iteration 1
372     Iteration 2
373     Iteration 3
374     Iteration 4
375     Iteration 5
376     Iteration 6
377     Iteration 7
378     Iteration 8
379     Iteration 9
380     Iteration 10
381     Iteration 11
382     Iteration 12
383 Cluster summary:
384     cluster  1 :  0
385     cluster  2 :  42
386     cluster  3 :  0
387     cluster  4 :  0
388     cluster  5 :  56
389
390
391 """
392
393 kclust5 = kcluster(data,5)
394 #print(len(kclust5))
```

```python
395 for i in range(len(kclust5)):
396   print ("cluster ", i+1, ": ", len(kclust5[i]))
397   for r in kclust5[i]:
398       filen= 'Q4/kclust5cluster'+ str(i+1)+'.csv'
399       with open(filen, 'a+') as myfile:
400           myfile.write(tweetdata[r])
401           myfile.write("\n")
402
403 """
404 For 10  kcluster
405
406 number of iteration:
407     Iteration 0
408     Iteration 1
409     Iteration 2
410     Iteration 3
411     Iteration 4
412 Cluster summary:
413     cluster  1 :   0
414     cluster  2 :   0
415     cluster  3 :   0
416     cluster  4 :   5
417     cluster  5 :   3
418     cluster  6 :   0
419     cluster  7 :   1
420     cluster  8 :   87
421     cluster  9 :   0
422     cluster  10 :   2
423 """
424
425 kclust10 = kcluster(data,10)
426 for i in range(len(kclust10)):
427   print ("cluster ", i+1, ": ", len(kclust10[i]))
428   for r in kclust10[i]:
429       filen= 'Q4/kclust10cluster'+ str(i+1)+'.csv'
430       with open(filen, 'a+') as myfile:
431           myfile.write(tweetdata[r])
432           myfile.write("\n")
433
434 """
435 For 20  kcluster
436
437 number of iteration:
438     Iteration 0
439     Iteration 1
440     Iteration 2
441     Iteration 3
```
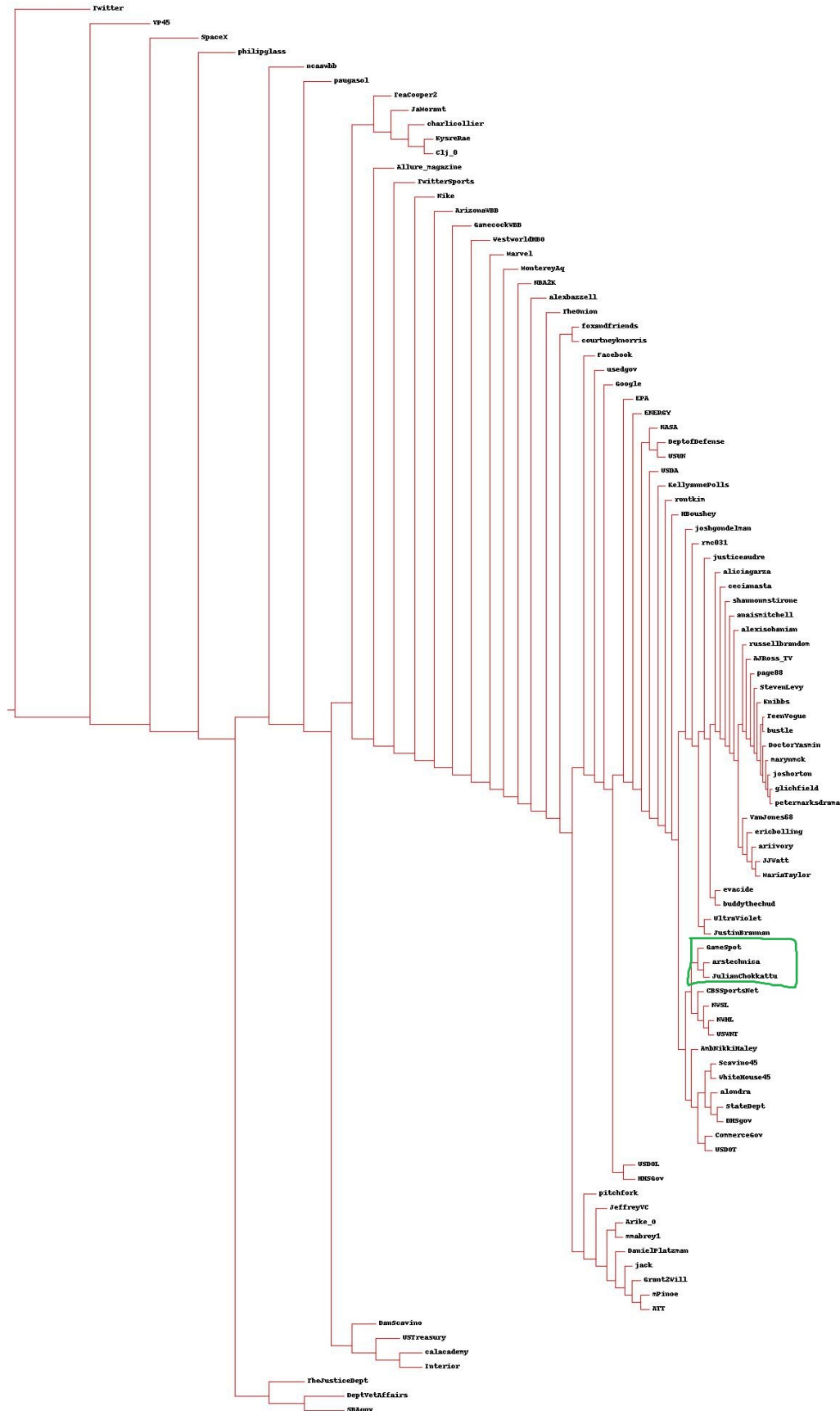
```
442      Iteration 4
443      Iteration 5
444 clusters summary:
445      cluster  1 :   0
446      cluster  2 :   54
447      cluster  3 :   0
448      cluster  4 :   0
449      cluster  5 :   7
450      cluster  6 :   0
451      cluster  7 :   0
452      cluster  8 :   0
453      cluster  9 :   0
454      cluster  10 :   0
455      cluster  11 :   0
456      cluster  12 :   0
457      cluster  13 :   0
458      cluster  14 :   1
459      cluster  15 :   0
460      cluster  16 :   0
461      cluster  17 :   6
462      cluster  18 :   26
463      cluster  19 :   0
464      cluster  20 :   4
465 """
466
467
468 kclust20 = kcluster(data,20)
469 for i in range(len(kclust20)):
470   print ("cluster ", i+1, ": ", len(kclust20[i]))
471   filen= 'Q4/kclust20cluster'+ str(i+1)+'.csv'
472   for r in kclust20[i]:
473       with open(filen, 'a+') as myfile:
474           myfile.write(tweetdata[r])
475           myfile.write("\n")
476
477 """
478 Question 5 ='mds2d.jpg'
479
480 23380.029615352396
481 140722.48343299245
482
483 98 lenght of coord variable
484 """
485 coords = scaledown(data)
486 print(len(coords))
487 draw2d(coords, tweetdata, jpeg='Q5/tweets2d.jpg')
```

**Listing 10:** question3.py for questions 3 4 and 5

**Figure 1:** The resulting output for Question 3 twitter account Cabinet45 included

**Figure 2:** The resulting output for Question 3 after removing tweet account Cabinet45

## Discussion

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

- Line 349 to 358 is the drive for the code.

```
340
341     return bestmatches
342
343 """
344 Question 3
345 """
346 #running the code
347 tweetdata, word,data =readfile("tweetdata.txt")
348 #clust = hcluster(data)
349 #print(clust.vec)
350
351 """
352 Question 3  ASCIII
353 To view  cluster
354
355 printclust(clust,labels=tweetdata)
356 """
357 """
358 Question 3 Dendrogram
```

**Listing 11:** A snapshot question3.py driver for Q3

- Used the readfile() function to read in the textdata.txt, it gets all data and returns the usernames as rownames, the words as colnames, and the value as data(data stored as a 2D array of float values).

- Next I parsed in data in hcluster(data) aka 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.

- printcluster() 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. I basically copied and pasted in a text file.

- drawdendrogram() does the similar objective as the printcluster but it produces a jpeg format of the recursive output of the cluster.

- Before i removed the Cabinet45.

I noticed the cluster algorithm was not doing a great job grouping the scree name properly mainly because the account Cabinet45 came out with all zero result during account-term matrix was generation.

This also caused odd grouping with GameSpot (tech)
arstechnica (tech)
ulianChokkattu(he is a senior associate editor for a tech company WIRED - I chose this company at the beginning ) GameSpot(tech) NWSL (sport) NWHL (sport) USWNT(sport)

After removing the account the clustering was much logical. It broke account even further as seen in figure 1 and figure 2.

For figure 2 the circled part we see that it only constitute tech screen names as mentioned above.The hierarchical clustering did a good job to me at the end of the day especially just a single variable could make it more accurate.

# Q4

*Cluster using k-Means*

*Cluster the accounts using k-Means, using k=5,10,20 (see Module 12, slide 37). For each value of k, create a file that lists the accounts in each cluster and upload to your GitHub repo.*

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

*How many iterations were required for each value of k?*

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

## Answer

```
351  """
352  Question 3   ASCIII
353  To view   cluster
354
355  printclust(clust,labels=tweetdata)
356  """
357  """
358  Question 3 Dendrogram
359
360
361  drawdendrogram(clust,tweetdata,jpeg="Q3/tweetdata.jpeg")
362
363  """
364  """
365  Question 4
366  """
367  """
368  For 5   kcluster
369  number of iteration:
370      Iteration 0
371      Iteration 1
372      Iteration 2
373      Iteration 3
374      Iteration 4
375      Iteration 5
376      Iteration 6
377      Iteration 7
378      Iteration 8
379      Iteration 9
380      Iteration 10
381      Iteration 11
382      Iteration 12
383  Cluster summary:
384      cluster  1 :   0
385      cluster  2 :   42
386      cluster  3 :   0
387      cluster  4 :   0
388      cluster  5 :   56
389
390
391  """
392
393  kclust5 = kcluster(data,5)
394  #print(len(kclust5))
```

```
395 for i in range(len(kclust5)):
396   print ("cluster ", i+1, ": ", len(kclust5[i]))
397   for r in kclust5[i]:
398       filen= 'Q4/kclust5cluster'+ str(i+1)+'.csv'
399       with open(filen, 'a+') as myfile:
400           myfile.write(tweetdata[r])
401           myfile.write("\n")
402
403 """
404 For 10  kcluster
405
406 number of iteration:
407     Iteration 0
408     Iteration 1
409     Iteration 2
410     Iteration 3
411     Iteration 4
412 Cluster summary:
413     cluster  1 :   0
414     cluster  2 :   0
415     cluster  3 :   0
416     cluster  4 :   5
417     cluster  5 :   3
418     cluster  6 :   0
419     cluster  7 :   1
420     cluster  8 :   87
421     cluster  9 :   0
422     cluster  10 :   2
423 """
424
425 kclust10 = kcluster(data,10)
426 for i in range(len(kclust10)):
427   print ("cluster ", i+1, ": ", len(kclust10[i]))
428   for r in kclust10[i]:
429       filen= 'Q4/kclust10cluster'+ str(i+1)+'.csv'
430       with open(filen, 'a+') as myfile:
431           myfile.write(tweetdata[r])
432           myfile.write("\n")
433
434 """
435 For 20  kcluster
436
437 number of iteration:
438     Iteration 0
439     Iteration 1
440     Iteration 2
441     Iteration 3
```

```
442      Iteration 4
443      Iteration 5
444 clusters summary:
445      cluster   1 :    0
446      cluster   2 :    54
447      cluster   3 :    0
448      cluster   4 :    0
449      cluster   5 :    7
450      cluster   6 :    0
451      cluster   7 :    0
452      cluster   8 :    0
453      cluster   9 :    0
454      cluster   10 :    0
455      cluster   11 :    0
456      cluster   12 :    0
457      cluster   13 :    0
458      cluster   14 :    1
459      cluster   15 :    0
460      cluster   16 :    0
461      cluster   17 :    6
462      cluster   18 :    26
463      cluster   19 :    0
464      cluster   20 :    4
465 """
466
467
468 kclust20 = kcluster(data,20)
469 for i in range(len(kclust20)):
470   print ("cluster ", i+1, ": ", len(kclust20[i]))
471   filen= 'Q4/kclust20cluster'+ str(i+1)+'.csv'
472   for r in kclust20[i]:
473       with open(filen, 'a+') as myfile:
474           myfile.write(tweetdata[r])
475           myfile.write("\n")
476
477 """
478 Question 5 ='mds2d.jpg'
479
480 23380.029615352396
481 140722.48343299245
482
483 98 lenght of coord variable
484 """
485 coords = scaledown(data)
486 print(len(coords))
487 draw2d(coords, tweetdata, jpeg='Q5/tweets2d.jpg')
```

**Listing 12:** A snapshot question3.py solution driver for Q4

## Discussion

*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).

- This coordinates then is used to determine best clusters closest to each row of data. The close the value is to zero the closer it is to the cluster or centroid.

- 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 its 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 cluster of 5 starting from 1

```
368 For 5   kcluster
369 number of iteration:
370      Iteration 0
371      Iteration 1
372      Iteration 2
373      Iteration 3
374      Iteration 4
375      Iteration 5
376      Iteration 6
377      Iteration 7
378      Iteration 8
379      Iteration 9
380      Iteration 10
381      Iteration 11
382      Iteration 12
383 Cluster summary:
384      cluster  1 :  0
385      cluster  2 :  42
386      cluster  3 :  0
387      cluster  4 :  0
388      cluster  5 :  56
```

**Listing 13:** A snapshot question3.py cluster 5 showing interations and cluster summary

- kcluster = 10 hand 4 iteration starting from zero and produced a total cluster of 5 starting from 1

```
403 """
404 For 10  kcluster
405
406 number of iteration:
407     Iteration 0
408     Iteration 1
409     Iteration 2
410     Iteration 3
411     Iteration 4
412 Cluster summary:
413     cluster  1 :   0
414     cluster  2 :   0
415     cluster  3 :   0
416     cluster  4 :   5
417     cluster  5 :   3
418     cluster  6 :   0
419     cluster  7 :   1
420     cluster  8 :   87
421     cluster  9 :   0
422     cluster  10 :   2
```

**Listing 14:** A snapshot question3.py cluster 10 showing interations and cluster summary

- kcluster = 20 had 5 iteration starting from zero and produced a total cluster of 20 starting from 1

```
434 """
435 For 20  kcluster
436
437 number of iteration:
438     Iteration 0
439     Iteration 1
440     Iteration 2
441     Iteration 3
442     Iteration 4
443     Iteration 5
444 clusters summary:
445     cluster  1 :   0
446     cluster  2 :   54
447     cluster  3 :   0
448     cluster  4 :   0
449     cluster  5 :   7
450     cluster  6 :   0
451     cluster  7 :   0
452     cluster  8 :   0
453     cluster  9 :   0
454     cluster  10 :   0
455     cluster  11 :   0
456     cluster  12 :   0
```

```
457      cluster  13 :  0
458      cluster  14 :  1
459      cluster  15 :  0
460      cluster  16 :  0
461      cluster  17 :  6
462      cluster  18 :  26
463      cluster  19 :  0
464      cluster  20 :  4
```

**Listing 15:** A snapshot question3.py cluster 20 showing interations and cluster summary

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

  In its cluster 4, it had a odd grouping of Twitter which does not fit any of my previous group(tech, politics, sport and music), the rest of the members where all politics.

  For cluster 5, it had all political screen_name groupings in it which is good.

  In cluster 7 it has just only one sport grouping. In cluster 8,had majorly tech scree name grouping, followed my political grouping then sports, this is why classification I said the clustering did do a good job. The finally cluster 10 was just about tech.

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

```
395 for i in range(len(kclust5)):
396   print ("cluster ", i+1, ": ", len(kclust5[i]))
397   for r in kclust5[i]:
398       filen= 'Q4/kclust5cluster'+ str(i+1)+'.csv'
399       with open(filen, 'a+') as myfile:
400            myfile.write(tweetdata[r])
401            myfile.write("\n")
```

**Listing 16:** A snapshot question3.py sample of how the code writes to each file based on k value
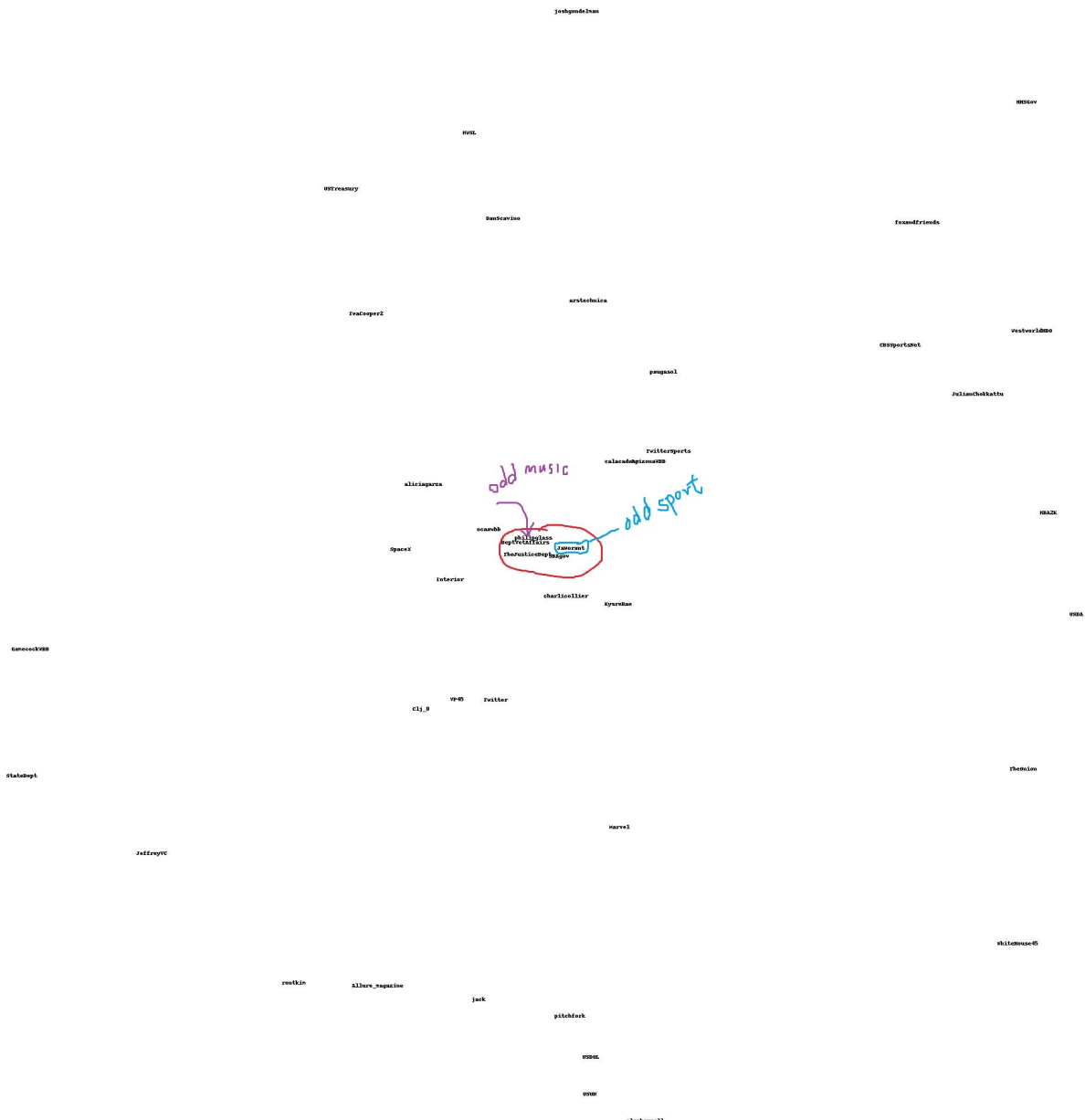
# Q5

*Create MDS Image*

*Use MDS to create a JPEG of the accounts (see Module 12, slide 50). Include the JPEG in your report.*

*How many iterations were required?*

*How well did the MDS do in grouping similar accounts together?*

*Were there any particularly odd groupings?*

# Answer



**Figure 3:** MDS Image Question 5

## Discussion

*I had the following result:*

- There were 98 iteration in total by checking the length of coord in line 486

```
477  """
478  Question 5 ='mds2d.jpg'
479
480  23380.029615352396
481  140722.4343299245
482
483  98 lenght of coord variable
484  """
485  coords = scaledown(data)
486  print(len(coords))
487  draw2d(coords, tweetdata, jpeg='Q5/tweets2d.jpg')
```

**Listing 17:** A snapshot question3.py cluster 20 showing iterations and cluster summary

- MDS did fairly good but there was a group containing two particular odd group members Jamorant (Sport category) and philipglass(Music category), the two were odd because the group was filled with mostly political screen names.

- 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 tweets2d.jpg Figure 3.

- I remember running into and error when I parse in data in the scaledown() function.

```
1        errorterm = (fakedist[j][k] – realdist[j][k]) / realdist[j
    ][k]
2        ZeroDivisionError: float division by zero
3
```

Apparently when the code was calculating teh percentage difference between the distance of two pairs, result in a zero value which should not be the case. This meant that one of my screen_names() the Cabinet45 in account-term matric had an output of all zeros values. (View tweetdataWithCabinet45.txt in github line 80 )

## References

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

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

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

- https://www.geeksforgeeks.org/python-pandas-dataframe-drop_duplicates/

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

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

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