

HW8 - Clustering

Logan Bartels
December 6, 2020

Q1

Answer

```
1 import tweepy
2 import json
3 import sys
4
5 auth = tweepy.OAuthHandler("", "")
6 auth.set_access_token("", "")
7
8 api = tweepy.API(auth, wait_on_rate_limit=True)
9
10 screenNames = open("100screen_names.txt", "r")
11
12 tweet_data = []
13
14 for screenName in screenNames:
15     try:
16         noNewlineName = screenName.replace('\n', '')
17         tweets = api.user_timeline(screen_name=noNewlineName, count
18         =200, include_rts=False, exclude_replies=True, lang="en", tweet_mode
19         ="extended")
20         for tweet in tweets:
21             tweet_dict = {}
22             tweet_dict["username"] = tweet.user.screen_name
23             tweet_dict["tweetId"] = tweet.id_str
24             tweet_dict["tweetText"] = tweet.full_text
25             tweet_data.append(tweet_dict)
26     except:
27         continue
28
29 json.dump(tweet_data, sys.stdout, indent=2)
30 screenNames.close()
```

Listing 1: Python script used to download tweets from 100 accounts listed in 100screen_names.txt

```
1 #!/bin/bash
2 input="100screen_names.txt"
3
4 while IFS= read -r line
```

```
5 do
6   count= grep -c "$line" "tweets.json"
7   echo "$line"
8   echo "$count"
9 done < "$input"
```

Listing 2: Shell script used to count usernames in tweets.json and gauge number of tweets.

Discussion

The script in listing 1 is what was used to download tweets for 100 popular Twitter accounts. Before line 10 you have your basic Twitter authentication for Python. Line 10 reads in my list of 100 popular twitter accounts stored in “100screen_Names.txt.” It then creates an empty list to store the tweets in. The for loop starting on line 14 reads “100screen_names.txt” line-by-line. Because I typed the list manually, with one screen name on each line, I have to get rid of the newline character that is read in with the screen name on line 16. Then for each username, their 200 most recent English-language tweets (excluding retweets and replies) are downloaded. Then for each individual tweet, a blank dictionary is created. The username of the account that made the tweet, the tweet id, and the full tweet text (accessed by setting tweet_mode to extended) are saved to the dictionary with an appropriately named key for each attribute. The dictionary is then appended to the empty list created on line 12. The try-except block is to prevent the script from stopping when a “page does not exist” error is encountered. The list is then dumped to stdout and directed to “tweets.json”.

The script in listing 2 performs a grep -c for each account name in “100screen_names.txt” on “tweets.json” This is to give me an idea of how many tweets were collected for each account. The results were not entirely accurate, however. The inaccuracies are due to some accounts mentioning other accounts, but the mentions were infrequent enough to be negligible. But I did have to re-select about 50 accounts because that many accounts on my first run produced less than 100 tweets that met my criteria according to the shell script.

Q2

Answer

```
1 import json
2 import re
3 import sys
4
5 tweets = open("tweets.json", "r")
6
7 data = json.load(tweets)
```

```
8 uris = re.compile('http[s]?://(?:[a-zA-Z]|[0-9]|[$-_@.&+]|[*\(\)\n\r\t]|(?:%[0-9a-fA-F][0-9a-fA-F]))+')
9 accountNames = re.compile('@+[A-Za-z0-9-_]+')
10 punctuation = re.compile('[^\w\s]')
11 smallTerms = re.compile(r'\W*\b\w{1,3}\b')
12 largeTerms = re.compile(r'\W*\b\w{15,}\b')
13
14 tweet_data = []
15
16 for tweet in data:
17     tweet_dict = {}
18     tweet_dict["username"] = tweet["username"]
19     tweet_dict["tweetId"] = tweet["tweetId"]
20
21     text = tweet["tweetText"]
22     text = re.sub(uris, "", text)
23     text = re.sub(accountNames, "", text)
24     text = text.lower()
25
26     splitText = text.split()
27
28     for word in splitText:
29         if any(ord(char) >= 128 for char in word):
30             splitText.remove(word)
31         else:
32             pass
33
34     joinedText = " ".join(splitText)
35     joinedText = re.sub(punctuation, "", joinedText)
36     joinedText = re.sub(smallTerms, "", joinedText)
37     joinedText = re.sub(largeTerms, "", joinedText)
38
39     tweet_dict["tweetText"] = joinedText
40     tweet_data.append(tweet_dict)
41
42 json.dump(tweet_data, sys.stdout, indent=2)
43 tweets.close()
```

Listing 3: Python script to filter the text of all tweets. Output directed to filteredTweets.json.

```
1 import json
2 import sys
3
4 tweets = open("filteredTweets.json", "r")
5
6 data = json.load(tweets)
7 accounts = {}
8
```

```
9 for tweet in data:
10     accountName = tweet["username"]
11     accounts[accountName] = []
12 '''
13 for terms in accounts.values():
14     terms.append("test")
15
16 print(accounts)
17
18 '''
19 for account in accounts:
20     for tweet in data:
21         if account == tweet["username"]:
22             text = tweet["tweetText"]
23             splitText = text.split()
24             for word in splitText:
25                 accounts[account].append(word)
26
27
28 json.dump(accounts, sys.stdout, indent=2)
29 tweets.close()
```

Listing 4: Python script to split the text of each account's tweets into a list.

```
1 import json
2 import sys
3
4 tweetTerms = open("accountsAndTerms.json", "r")
5
6 data = json.load(tweetTerms)
7 wordCount = {}
8 wordList = {}
9
10
11 for terms in data.values():
12     for term in terms:
13         wordCount[term] = ""
14
15 for term in wordCount:
16     count = 0
17     for termList in data.values():
18         if term in termList:
19             count += 1
20         else:
21             continue
22     wordCount[term] = count
23
24 for word, accountCount in wordCount.items():
```

```
25     frac=float(accountCount)/100
26     if frac>0.1 and frac<0.5:
27         wordList[word] = ""
28
29 for word in wordList:
30     appearanceCount = 0
31     for listOfWords in data.values():
32         if word in listOfWords:
33             appearanceCount = appearanceCount + listOfWords.count(word)
34         else:
35             continue
36     wordList[word] = appearanceCount
37
38 popularWords = sorted(wordList.items(), key=lambda word: word[1],
39                        reverse=True)
40
41 json.dump(popularWords[0:1000], sys.stdout, indent=2)
42 #json.dump(wordCount, sys.stdout, indent=2)
43 tweetTerms.close()
```

Listing 5: Python script that gets the 1000 most popular words across all accounts.

```
1 import json
2 import sys
3
4 accountsAllTerms = open("accountsAndTerms.json", "r")
5 popularTerms = open("1000popularWords.json", "r")
6
7 accounts = json.load(accountsAllTerms)
8 terms = json.load(popularTerms)
9
10 accountsPopularWords = {}
11
12 def getWordCount(accountName, term):
13     wordAndCount = (term, accounts[accountName].count(term))
14     return wordAndCount
15
16 for account in accounts:
17     accountsPopularWords[account] = []
18
19
20 for account in accountsPopularWords:
21     for term in terms:
22         wordWithCount = getWordCount(account, term[0])
23         accountsPopularWords[account].append(wordWithCount)
24
25
26 json.dump(accountsPopularWords, sys.stdout, indent=2)
```

```
27 accountsAllTerms.close()
28 popularTerms.close()
```

Listing 6: Python script that checks the 1000 most popular words against each account's list of terms and produces a count of how many times each of the 1000 popular words appeared their tweets.

```
1 import json
2 import sys
3 import csv
4
5 accountsAndWords = open("accountsAndPopularTerms.json", "r")
6 words = open("1000popularWords.json", "r")
7
8 accountsPopularWords = json.load(accountsAndWords)
9 popularWords = json.load(words)
10 headers = ["Blog"]
11
12 for word in popularWords:
13     headers.append(word[0])
14 '''
15 for account, terms in accountsPopularWords.items():
16     row = []
17     row.append(account)
18     for term in terms:
19         row.append(term[1])
20 '''
21 with open("accountTermMatrix.csv", "w") as file:
22     writer = csv.writer(file)
23     writer.writerow(headers)
24     for account, terms in accountsPopularWords.items():
25         row = []
26         row.append(account)
27         for term in terms:
28             row.append(term[1])
29         writer.writerow(row)
30
31 #json.dump(accountsPopularWords, sys.stdout, indent=2)
32 accountsAndWords.close()
33 words.close()
```

Listing 7: Python script to construct the account-term matrix.

Discussion

Starting with script in listing 3, this script reads in “tweets.json” and filters each tweet’s text to match the requirements for question 2. To break it down further, the list in “tweets.json” is stored in the “data” variable on line 7. Lines 8-12 have variables that store the regular expressions to be run against each term, with each variable being appropriately named to match its regular

expression. Before I go further, I should state that the intent of this script is to create a 1:1 replica of “tweets.json,” just with the text filtered to be rid of uris, account names, punctuation, etc. So similar to listing 1, an empty list is created to store each tweet. Then for each tweet dictionary in “data,” an empty dictionary is created to store each tweet’s attributes as in listing 1. Only this time, the username and tweet id of each tweet is immediately stored since I do not have to parse them. Then the tweet text is stored in “text.” After that, the “re.sub” function is used to get rid of any terms that match the regular expression passed in as the first parameter on lines 22 and 23. The second parameter (what to substitute each match with) is the empty string. The third parameter is the text that is passed in to be run against the regular expression in the first parameter. The text is entirely converted into lowercase on line 24. Then, to prepare the text to have any words with non-ASCII characters removed, the text is split into a list, with each element being a term in the text. For each term in the split text, if it contains a non-ASCII character, remove the term. I should note that this method does not remove all terms with non-ASCII characters because not all non-ASCII characters presented themselves as something like “-backslash-u2019.” For example: after filtering, I noticed a “youll” in one of the texts. Looking at the unfiltered version, it was listed as “you’ll” instead of “you-backslash-u2019ll” (I couldn’t get latex to print the backslash). After the non-ASCII terms are removed, the text is rejoined using the “join” function. From there, just like lines 22 and 23, “re.sub” is used to remove matches of the regular expressions passed in. After filtering is complete, the text is added as a value to the dictionary created on line 17. The dictionary is then appended to the “tweet_data” list. The “tweet_data” list is then dumped to stdout using “json.dump.” The output was directed to “filteredTweets.json.”

Moving on to listing 4. This script’s purpose is to create a dictionary, with each account being the key and a list of every term in the accounts tweets being the value. The script reads in “filteredTweets.json” and stores that list of dictionaries in “data.” An empty dictionary, “accounts” is created. Then the username of every tweet in “data” is added as the key in the “accounts” dictionary. The nature of dictionaries ensures that there are no duplicate account names. The value for each account name in the “accounts” dictionary is initially made an empty list. Lines 13-16 are just a test to confirm I could properly add terms to the values of “accounts.” For each account in “accounts,” and for each tweet in “data,” if the account matches the username in the tweet, then get the text and split it into a list of words. Then for each word in the split text, append the word to the list of the account’s value in “accounts.” The result is a dictionary with the account names as keys, and a list of all the terms from that user’s tweets. The “accounts” dictionary is dumped using “json.dump” to stdout, with the output being directed to “accountsAndTerms.json.”

Next is listing 5. This script is probably the most complex in this report. It ends up getting the 1000 most popular terms across all accounts. Let’s break down how it gets there. It starts by reading in “accountsAndTerms.json,” and being stored in “data” once again. Two empty dictionaries are created: “wordCount” and “wordList.” The “wordCount” dictionary ends up representing the number of accounts each term appears in, with the terms being the keys, and the number of accounts the term appears in being the values. The script builds “wordCount” by iterating through “data.values” starting on line 11. Recall that data.values is a list of terms, and that “terms” is a single list of terms. The script also iterates through each term in “terms.” Then each term is added as a key in “wordCount,” with its value being declared an empty string. The next loop starting on line 15 iterates through each term in “wordCount.” A counting variable named “count” is

declared and initialized to 0 at the start of every iteration. Then it iterates through each term list in “data.values” and checks if the term is in the term list. If it is, increment the counter, if not, skip to the next term list. After every term list has been parsed, pass in the final value of the “count” variable as the term’s value. The resulting “wordCount” dictionary can be viewed in “termCounts.json.” That file was created by dumping “wordCount” to stdout, and directing the output to it on line 41. This was done before the rest of the script was written. The next two loops deal with building the “wordList” dictionary. The first loop, starting on line 24, iterates through the now-built “wordCount.” The loop iterates through both “wordCount’s” keys and values: “word” and “accountCount,” respectively. All this loop does is fake the removal of stopwords by only adding the word as a key to “wordList” if it appears in more than ten percent and less than fifty percent of all 100 accounts. The key’s value is initialized to an empty string. The next loop (starting on line 29) starts by iterating through the keys (words) of “wordList.” Similar to line 16, a counting variable, “appearanceCount,” is initialized to 0 at the start of each iteration. Then for each list of words in “data.values,” if the word is in the list of words, then increment “appearanceCount” by adding its current value to the count of that word in the list of words. If the word isn’t in the list of words, skip to the next list of words. This method ensures that the total number of appearances of the term across all 100 accounts is counted correctly. After all the lists of words have been parsed, assign the final value of “appearanceCount” to the value of the word. The “wordList” dictionary can be viewed in “wordList.json.” That file was created by dumping “wordList” to stdout and directing the output to that file on line 41 (just replace “wordCount” with “wordList”). This was done before line 38 was written. If you compare “wordList.json” and “termCounts.json,” you will notice that the counts for each word are higher in “wordList.json” than “termCounts.json.” This is supposed to be the case, as “termCounts” only counts the number of accounts each term appears in, versus “wordList” that counts the word’s appearances over all accounts. The 1000 most popular words are gotten by sorting “wordList” by value in descending order on line 38. The sorted dictionary, now a nested list, is stored in “popularWords.” Each item in the list contains a list with 2 values: a term, and that term’s count. Then the first 1000 items are dumped to stdout on line 40, with the output being directed to “1000popularWords.json”

Now we have to track how many times each popular word appears in each user’s tweets. This is done using the script in listing 6. This time 2 json files are opened: “accountsAndTerms.json” and “1000popularWords.json.” They are stored in “accounts” and “terms,” respectively. An empty dictionary, “accountsPopularWords” is created. Then there is a function that takes an account name and a term. It takes that term and finds the count of how many times that term appears in an account’s term list, and stores the term and its count in a tuple. That tuple is stored in “wordAndCount,” and returned. The first loop iterates through the accounts, and assigning each account as the key in the dictionary created on line 10. The key’s value is initialized to an empty list. The next loop iterates through the new dictionary and each item in “terms.” The function on line 12 is called, and the account and term (not the term’s count, note the index of the passed in parameter) are passed in. The returned value is stored in “wordWithCount.” Then that variable is appended to the value of the “account” key in the “accountsPopularWords” dictionary. This dictionary is then dumped to stdout, with the output being directed to “accountsAndPopularTerms.json.” This file contains a dictionary, with the accounts being the keys, and a nested list of terms and their number of appearances in that account’s tweets as the values.

Now it's time to build the account-term matrix. This is done using the script in listing 7. Two files are opened: "accountsAndPopularTerms.json" and "1000popularWords.json." The files are stored in "accountsPopularWords" and "popularWords," respectively. Since this script will write to a csv file, a "headers" list is created with only one item: "Blog." This list will serve as the first row in the account-term matrix, which is modeled after blogdata.txt (minus the tab delimiter). To build the "headers" list, the script iterates through the "popularWords" nested list, and appends each word to the "headers" list. Lines 15-19 were a test to build subsequent rows of the csv file. Those rows consist of an account, and its counts for each term. Now the script is ready to write to "accountTermMatrix.csv." We start by writing the headers to the first row. Then the loop starting on line 24 iterates through each account and its list of terms. At the beginning of every iteration an empty list called "row" is created, and the account name is the first thing added to it. Then for each term in the list of terms, append the term's count (once again, note the index). Once all the counts are added, write the row to the csv file. One thing of note about the csv file is I wrote it with a comma delimiter instead of a tab delimiter like blogdata.txt. Line 31 can be disregarded. It was a method of testing to make sure I was correctly reading in files.

Note

Questions 3-5 were completed using a copy of the week 12 Google Colab notebook. The modifications made to it include, and *are* limited to the following:

- Eliminated the tab delimiter in the readfile function.
- Changed the file in the call to readfile (blogdata.txt) to accountTermMatrix.csv.
- Additional cells to output blognames in kclusters were added under the "kmeans" section.

The notebook can be accessed from the repository and here: <https://colab.research.google.com/drive/118JsEQwzDPXRJAde7Tavzk9dDvXgUBIf?usp=sharing>

Q3

Answer

Figure 1 is on page 19.

Discussion

The dendrogram shown in figure 1 is accurate for the most part in its clustering. Many gaming accounts are clustered together such as PlayStation, pcgamer, gamesradar, Xbox, endgadgetgaming, gameinformer, IGN, GameSpot, and Polygon. Although there are a couple of outliers in Halo and GameStop being clustered with celebrities. Many political accounts are also clustered together such as JoeBiden, KamalaHarris, RealJamesWoods(he's very political these days), scrowder, and realDonaldTrump. There are also many news outlets clustered at the bottom of the dendrogram.

Q4

Answer

```
1 k=5
2 5 iterations
3 k=10
4 5 iterations
5 k=20
6 4 iterations
```

Listing 8: Text file with the number of iterations for each value of k

```
1 Cluster 1:
2 ['KevinHart4real',
3  'kourtneykardash',
4  'MLB',
5  'WIRED',
6  'TheSun',
7  'CBS',
8  'EW',
9  'GMA',
10 'enews',
11 'people',
12 'TMZ',
13 'MTV',
14 'usweekly',
15 'THR',
16 'Variety',
17 'HBO',
18 'Marvel',
19 'RollingStone',
20 'billboard',
21 'DEADLINE',
22 'voguemagazine',
23 'glamourmag',
```

```
24 'InStyle',
25 'wwd']
26 Cluster 2:
27 ['SportsCenter',
28 'KylieJenner',
29 'NBA',
30 'Xbox',
31 'PlayStation',
32 'IGN',
33 'GameSpot',
34 'pcgamer',
35 'Halo',
36 'Polygon',
37 'NFL',
38 'starwars',
39 'GameStop',
40 'GamesRadar',
41 'gameinformer',
42 'engadgetgaming']
43 Cluster 3:
44 ['BreitbartNews',
45 'cnnbrk',
46 'CNN',
47 'BBCBreaking',
48 'BBCWorld',
49 'TheEconomist',
50 'Reuters',
51 'FoxNews',
52 'AP',
53 'washingtonpost',
54 'guardian',
55 'Forbes',
56 'businessinsider',
57 'USATODAY',
58 'BreakingNews',
59 'CNBC',
60 'business',
61 'WSJ',
62 'CBSNews',
63 'NPR',
64 'MSNBC',
65 'Independent',
66 'NewYorker',
67 'HuffPost',
68 'SkyNews',
69 'TIME',
70 'FortuneMagazine',
```

```
71 'TODAYshow',
72 'latimes',
73 'TheView',
74 'VanityFair',
75 'NYMag',
76 'newsmax',
77 'DailyCaller',
78 'NBCNews',
79 'seanhannity']
80 Cluster 4:
81 ['BarackObama',
82 'katyperry',
83 'rihanna',
84 'realDonaldTrump',
85 'TheEllenShow',
86 'KingJames',
87 'espn',
88 'Harry_Styles',
89 'LiamPayne',
90 'NICKIMINAJ',
91 'JoeBiden',
92 'KamalaHarris',
93 'scrowder',
94 'RealJamesWoods',
95 'WNBA',
96 'Entrepreneur',
97 'Inc',
98 'HarvardBiz',
99 'nbc',
100 'FOXTV',
101 'LindseyGrahamSC']
102 Cluster 5:
103 ['NASA', 'NatGeo', 'DailyMirror']
```

Listing 9: List of accounts in clusters for k equal to 5.

```
1 Cluster 1:
2 ['SportsCenter', 'espn', 'NBA', 'PlayStation', 'CBS']
3 Cluster 2:
4 ['BBCBreaking',
5 'BreakingNews',
6 'NPR',
7 'MSNBC',
8 'SkyNews',
9 'DailyMirror',
10 'Entrepreneur',
11 'Inc',
12 'NBCNews']
```

```
13 Cluster 3:
14 ['Harry_Styles', 'MLB', 'LindseyGrahamSC']
15 Cluster 4:
16 ['kourtneykardash',
17  'Halo',
18  'TheSun',
19  'EW',
20  'TODAYshow',
21  'GMA',
22  'enews',
23  'people',
24  'TMZ',
25  'usweekly',
26  'Variety',
27  'billboard',
28  'DEADLINE',
29  'voguemagazine',
30  'glamourmag',
31  'InStyle',
32  'wwd']
33 Cluster 5:
34 ['BarackObama',
35  'katyperry',
36  'rihanna',
37  'TheEllenShow',
38  'KingJames',
39  'KevinHart4real',
40  'KylieJenner',
41  'NICKIMINAJ',
42  'NFL',
43  'FOXTV']
44 Cluster 6:
45 ['Reuters', 'guardian', 'nbc', 'HBO', 'Marvel', 'RollingStone']
46 Cluster 7:
47 ['scrowder', 'WIRED']
48 Cluster 8:
49 ['Xbox',
50  'IGN',
51  'GameSpot',
52  'pcgamer',
53  'Polygon',
54  'THR',
55  'GameStop',
56  'GamesRadar',
57  'gameinformer',
58  'engadgetgaming']
59 Cluster 9:
```

```
60 ['BreitbartNews',
61  'realDonaldTrump',
62  'cnnbrk',
63  'CNN',
64  'BBCWorld',
65  'TheEconomist',
66  'NatGeo',
67  'JoeBiden',
68  'KamalaHarris',
69  'FoxNews',
70  'RealJamesWoods',
71  'AP',
72  'washingtonpost',
73  'Forbes',
74  'businessinsider',
75  'USATODAY',
76  'CNBC',
77  'business',
78  'WSJ',
79  'CBSNews',
80  'Independent',
81  'NewYorker',
82  'HuffPost',
83  'TIME',
84  'FortuneMagazine',
85  'HarvardBiz',
86  'latimes',
87  'TheView',
88  'VanityFair',
89  'NYMag',
90  'newsmax',
91  'DailyCaller',
92  'seanhannity']
93 Cluster 10:
94 ['NASA', 'LiamPayne', 'WNBA', 'MTV', 'starwars']
```

Listing 10: List of accounts in clusters for k equal to 10.

```
1 Cluster 1:
2 ['washingtonpost', 'MSNBC', 'HuffPost', 'newsmax', 'NBCNews']
3 Cluster 2:
4 ['Entrepreneur', 'Inc']
5 Cluster 3:
6 ['BreitbartNews',
7  'CNN',
8  'USATODAY',
9  'BreakingNews',
10 'WSJ',
```

```
11 'SkyNews',
12 'latimes']
13 Cluster 4:
14 ['cnnbrk',
15 'FoxNews',
16 'AP',
17 'business',
18 'CBSNews',
19 'Independent',
20 'TheView',
21 'DailyCaller']
22 Cluster 5:
23 ['NASA',
24 'NatGeo',
25 'WIRED',
26 'EW',
27 'GMA',
28 'enews',
29 'people',
30 'TMZ',
31 'usweekly',
32 'Variety',
33 'starwars',
34 'DEADLINE']
35 Cluster 6:
36 ['TheEconomist', 'Reuters', 'CNBC', 'NPR']
37 Cluster 7:
38 ['KylieJenner', 'kourtneykardash', 'Marvel', 'wwd']
39 Cluster 8:
40 ['KevinHart4real', 'LiamPayne', 'NICKIMINAJ']
41 Cluster 9:
42 ['SportsCenter', 'espn', 'NFL']
43 Cluster 10:
44 ['realDonaldTrump',
45 'RealJamesWoods',
46 'Forbes',
47 'businessinsider',
48 'DailyMirror',
49 'TIME',
50 'FortuneMagazine',
51 'HarvardBiz',
52 'VanityFair']
53 Cluster 11:
54 ['NBA', 'nbc', 'MTV', 'HBO']
55 Cluster 12:
56 ['billboard']
57 Cluster 13:
```

```
58 ['pcgamer', 'seanhannity']
59 Cluster 14:
60 ['Halo', 'MLB', 'RollingStone']
61 Cluster 15:
62 ['NewYorker', 'CBS', 'TODAYshow', 'NYMag']
63 Cluster 16:
64 ['BBCBreaking',
65  'BBCWorld',
66  'guardian',
67  'TheSun',
68  'THR',
69  'voguemagazine',
70  'glamourmag',
71  'InStyle']
72 Cluster 17:
73 ['Xbox',
74  'PlayStation',
75  'IGN',
76  'GameSpot',
77  'Polygon',
78  'GameStop',
79  'GamesRadar',
80  'gameinformer',
81  'engadgetgaming']
82 Cluster 18:
83 ['WNBA']
84 Cluster 19:
85 ['rihanna', 'KingJames', 'Harry_Styles', 'FOXTV']
86 Cluster 20:
87 ['BarackObama',
88  'katyperry',
89  'TheEllenShow',
90  'JoeBiden',
91  'KamalaHarris',
92  'scrowder',
93  'LindseyGrahamSC']
```

Listing 11: List of accounts in clusters for k equal to 20.

Discussion

For the clusters in listing 9, I think I can characterize the accounts in each cluster. Cluster 1 has celebrities and general entertainment. Cluster 2 has mostly gaming accounts, with some outliers in SportsCenter, KylieJenner, NBA, and NFL. Cluster 3 has mostly political figures and news outlets. Cluster 4 is a combination of celebrities and political figures. Cluster 5 is mostly science, minus DailyMirror.

For listing 10, Cluster 1 is mostly sports with 2 outliers in PlayStation and CBS. Cluster 2 is news outlets. I cannot characterize Cluster 3. Cluster 4 is celebrities and general entertainment. Cluster 5 is mostly celebrities. Cluster 6 has publishers, 2 of which (HBO and Marvel) are film publishers. I cannot characterize cluster 7. Cluster 8 is mostly gaming. Cluster 9 is political figures and news outlets. I'm not sure what to make of cluster 10.

For listing 11, I'm only going to talk about the clusters I can categorize. If I don't mention one, then either I don't know what to make of it, or it only has 1 account. Cluster 1 is news outlets. Cluster 2 is business. Cluster 3 is more news, as is Cluster 4. Cluster 5 is entertainment. Cluster 6 is more news. Cluster 8 is celebrities. Cluster 9 is sports. Cluster 10 is mostly magazines. Cluster 15 is more news, as is 16. Cluster 17 is gaming. Cluster 19 is mostly celebrities. Cluster 20 is a mixture of celebrities and political figures.

I would say that the k value that created the most reasonable clusters is k equals 5. I feel like those clusters are the easiest to label with a small amount of outliers.

Q5

Answer

Figure 2 is on page 20.

Discussion

I would say the graph shown in figure 2 is very accurate in its clustering. All the gaming accounts are clustered. All of the news outlets are clustered together, as well as celebrities and political figures. I think this plot is more accurate than the dendrogram.

References

- Parse Twitter Hashtags, Usernames & URLs with JS, <https://www.benmarshall.me/parse-twitter-hashtags/>
- Best Practice to Extract and Remove URLs from Python String – Python Tutorial, <https://www.tutorialexample.com/best-practice-to-extract-and-remove-urls-from-python-string-python-tutorial/>
- The Perfect URL Regular Expression, <https://urlregex.com/>

- Read JSON file using Python, <https://www.geeksforgeeks.org/read-json-file-using-python/>
- Python — Remove punctuation from string, <https://www.geeksforgeeks.org/python-remove-punctuation-from-string/>
- Python: Remove words from a string of length between 1 and a given number, <https://www.w3resource.com/python-exercises/re/python-re-exercise-49.php>
- Lowercase in Python, <https://www.educba.com/lowercase-in-python/>
- Python String encode() Method, https://www.w3schools.com/python/ref_string_encode.asp
- Python 3 - String decode() Method, https://www.tutorialspoint.com/python3/string_decode.htm
- Python String split() Method, https://www.w3schools.com/python/ref_string_split.asp
- Python String join() Method, https://www.w3schools.com/python/ref_string_join.asp
- Week 12 Colab Notebook https://colab.research.google.com/github/cs432-websci-fall20/assignments/blob/master/432_PCI_Ch03.ipynb
- Week 12 Slides, <https://docs.google.com/presentation/d/1Sz5tSqXBjMCL0Iq7oIEq7aVE/edit?usp=sharing>

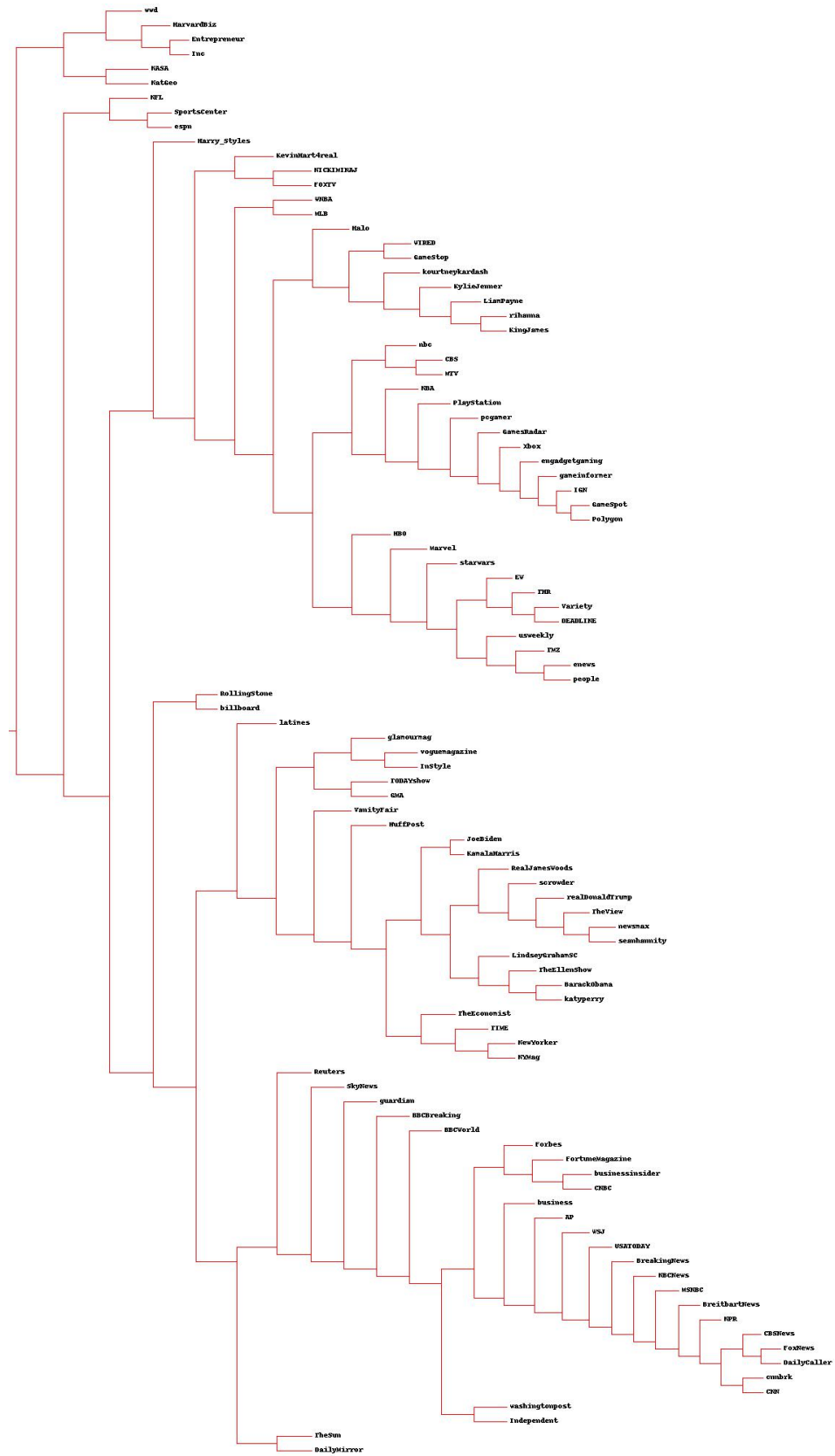


Figure 1: Hierarchical clustering dendrogram.

