

1 **COMPARING LANGUAGE USE AND NETWORK STRUCTURE** 2 **USING TWITTER***

3 CAITLIN SHENER[†], BRIAN OCEGUERA[‡], AND SEUNGJOO (SALLY) LEE[‡]

4 **Abstract.** We develop an approach for exploring questions regarding language use and con-
5 nections between people in social networks. In particular, we investigate community structure and
6 language usage in the network composed of the most followed ninety-nine users in Twitter. Our goal
7 is to measure the relation between a community of users and the words employed by those users. We
8 accomplish the investigation using k -means clustering to group users based on word choice, and we
9 use modularity maximization and InfoMap clustering to find communities based on network connec-
10 tions. Our study illustrates how to mathematically analyze and interpret language use within social
11 network structure.

12 **Key words.** social network | twitter | community detection | document clustering

13 **AMS subject classifications.** 94A15, 91D30, 90C35

14 **1. Introduction.** Network scientists often question the extent to which similar-
15 ity plays a role in forming social ties. They ask questions like, “Are similar people
16 more likely to think, act, and communicate the same way as their friends?” Or, “Are
17 we, as humans, drawn to form connections with people dissimilar from ourselves?”
18 We use the term similarity here to refer to common interests and therefore common
19 topics of conversation. The idea of “homophily” is defined as the principle that similar
20 people will be more connected to each other than dissimilar people [11]. In order to
21 easily determine if people who are more connected discuss the same topics, we develop
22 an approach to determine how homophilic word use functions in networks.

23 Our method of comparing word use to social ties involves creating groups in two
24 distinct ways, and then measuring the overlap of these groups:

- 25 1. We created **Network Communities** through modularity maximization and
26 Infomap; such an approach uses information from data on the edges in the
27 network. We explain these methods in further detail in section 3.
- 28 2. We also created **Language Clusters** through k -means clustering; such an
29 approach uses information from data on the words used in the tweets.

30 Although some have considered language use in their analysis of networks [2],
31 our approach uses distinct groupings based on the term-frequency inverse-document-
32 frequency model in section 4 of this paper. This allows researchers to avoid asking
33 the question: “How different is the language use between network communities?” af-
34 ter creating the communities. Instead, researchers can ask, “How much do language
35 clusters overlap with network communities?” The approach allows for the considera-
36 tion of language use along with network structure in community detection, which is
37 important since algorithms for community detection on networks are still contested

*Submitted to the editors November 5, 2018. Completed under the guidance of Mason A. Porter, Professor in the Department of Mathematics, University of California, Los Angeles (mason@math.ucla.edu).

[†]Department of Statistics, University of California, Los Angeles, CA, 90095 (caitlinshener13@ucla.edu).

[‡]Department of Mathematics, University of California, CA 90095 (brianoceguera@ucla.edu, sally.lee1725@gmail.com).

38 today [4]. Additionally, our approach incorporates a consideration of language use in
 39 finding clusters in the network that can be applied to many different types of net-
 40 works. The last benefit to our approach is the ability to produce comparable and
 41 interpretable values when comparing partitions on the network. The z-scores that we
 42 calculate in section 5 allow us to answer how much language use overlaps with the
 43 social ties in a given network.

44 To illustrate our approach, we use the social media site Twitter. Twitter has
 45 become a popular platform for users to connect to other users, including people who
 46 they have never met. A Twitter user can “follow” another user, which allows the
 47 former to see recent updates from the latter in the form of “tweets”. A tweet consists
 48 of a maximum of 140 characters¹ and, for example, can be a short text snippet
 49 describing a user’s current opinions or business promotions. For many users, Twitter
 50 is a stage for expression and identity creation, where language is an important currency
 51 for influencing others and spreading information.

52 We construct a network [14] of the top ninety-nine users followed on Twitter. We
 53 represent each user as a node, and create directed edges pointing from one user to
 54 the users they follow. For the Twitter follow network, one can intuitively consider
 55 the network as a social network, because users interact with others in whom they are
 56 interested and possibly know. However, Twitter differs from many other social net-
 57 working sites, because many users often follow accounts of celebrities or news sources
 58 and tweet at users who are unaffiliated with them. Because of the duality behind
 59 Twitter’s function as a social network and Twitter’s usage for the dissemination of
 60 information, a Twitter network has characteristics of both a social network and an
 61 information network [12].

62 Through our language and network analysis of the top Twitter users, our approach
 63 considers whether users who write about similar topics are more likely to be in the
 64 same network community. Essentially, we ask: “Do people in the same communities,
 65 based on network structure, discuss similar topics?”

66 The rest of the paper is organized as follows. In section 2, we discuss how we
 67 gathered and represented the data from Twitter. In section 3, we discuss the methods
 68 of network community detection that we used to explore the network. In section 4,
 69 we explain how we processed the tweet corpus to cluster the network based on tweet
 70 language. In section 5, we detail some calculations that we will need for our analysis
 71 in the following section. In section 6, we discuss how to interpret our results. In
 72 section 7, we further explore and interpret our results and draw conclusions. We also
 73 touch on possible further studies based on our initial research.

74 **2. Data Acquisition.** To gather the network information from Twitter, we
 75 made use of Twitter’s REST API [26]. The API has methods that allow one to
 76 make rate-limited queries to the company’s database for information about users.
 77 Due to time constraints and Twitter’s rate limitation (which throttled our requests
 78 to a limited selection every 15 minutes), we decided to investigate a small sample
 79 size (of one hundred Twitter users). For convenience and accessibility, we chose the
 80 one hundred most followed users on Twitter. However, we then excluded one Twitter
 81 user, @MohamedAlarefe, as he was both a language and a network outlier. As the
 82 only user in our dataset tweeting in Arabic, he did not follow or receive follows from
 83 any of the other ninety-nine users. Although the majority of users in our dataset

¹At the time of our data gathering, tweets were limited to 140 characters. Recently, Twitter increased its character limit to 280 for tweets in English and other languages. For more information, please see <https://blog.twitter.com/official/en-us/topics/product/2017/tweetingmadeeasier.html>.

84 wrote in English, we also included a smaller subgroup tweeting in Spanish.

85 **2.1. Creating a Directed Adjacency Matrix.** Because Twitter does not have
 86 publicly available datasets regarding follower networks on their website, we built an
 87 unweighted, directed adjacency matrix A_{ij} ourselves using the results from Twitter
 88 API queries. We define A_{ij} as the directed adjacency matrix where each entry a_{ij} in
 89 the matrix is 1 if j follows i and 0 otherwise. Thus, A_{ij} is a matrix of ones and zeros.
 90 In the Twitter network, A_{ij} has dimensions 99×99 . The process of creating the
 91 adjacency matrix involved taking the results from the Twitter queries (which return
 92 information as a JSON-formatted string) and parsing follower information as an edge
 93 list. In an edge list, there is an edge of weight 1 from user A to user B if B follows A .
 94 In the final representation of the network with the top ninety-nine users, we labeled
 95 nodes by Twitter user name.

96 **2.2. Network Measures.** We now define some of the concepts that we used to
 97 quantitatively measure network structure on our network of Twitter users. See [14]
 98 for further discussion.

99 **In-degree/Out-degree.** A degree of a vertex or node is the number of edges
 100 connected to that vertex. In directed networks, such as our dataset, the in-
 101 degree of a vertex represents the number of followers of a specific Twitter
 102 user. Likewise, the out-degree represents the number of users in the network
 103 that a specific Twitter user follows. The mean in-degree and out-degree over
 104 the set of all vertices in our Twitter network is 15.2755.

105 **Path.** A path is a sequence of vertices such that every consecutive pair of
 106 vertices in the sequence is connected by an edge in the network. A path
 107 is defined in both the directed and undirected case. In general, a path can
 108 intersect itself by revisiting a vertex or traversing an edge or set of edges
 109 multiple times.

110 **Shortest/geodesic path length.** A shortest path between two vertices is
 111 the shortest possible distance (in terms of the number of edges) needed to
 112 traverse from one vertex to the other. The mean geodesic path length of our
 113 network is 2.088

114 **Diameter.** A network's diameter is the length of its longest geodesic path.
 115 The diameter of our network is 5.

116 **Strongly connected component.** A strongly connected component within
 117 a directed network is a set of vertices where there exists a bidirectional path
 118 between any two pairs of vertices. The size of the largest strongly connected
 119 component in our Twitter network is 78. Because the largest connected compo-
 120 nent is close to the size of the overall network, a significant portion of the
 121 network is interconnected through mutual following.

122 **2.3. Downloading Tweet Data.** Using a combination of Python libraries that
 123 serve as wrappers for Twitter's REST API method calls [22, 25], we downloaded
 124 and cleaned the latest 200 tweets from each of the top ninety-nine users. We chose
 125 the number 200 to successfully work within the constraints of Twitter's API. Using
 126 TWEETPY [25], we queried Twitter to return the last 200 tweets from each of the top
 127 ninety-nine users. We will refer to the set of tweets from all ninety-nine users as the
 128 *tweet corpus*.

129 We downloaded the data on June 5th, 2017. We expected a total number of
 130 $99 \times 200 = 19,800$ tweets, but instead the total number of tweets in the tweet corpus
 131 is 19,579. The reason for the smaller tweet count is that there were a small number
 132 of users with an entire tweet history fewer than 200 tweets. We decided to ignore this
 133 inconsistency in our textual analysis of the tweets. For reference, the dates of the
 134 tweets range from June 3rd, 2011 to June 5th, 2017. When inspecting the dataset,
 135 we found a single user (@Adele) with infrequent tweets (mostly dating from 2011 to
 136 2014). If we chose to exclude Adele from our set of users, the oldest user’s tweet
 137 occurs on July 13th, 2014. Interestingly, the median tweet timestamp occurs on April
 138 25th, 2017, which indicates that most of these influential users tweet very often. The
 139 mean timestamp occurs on February 7th, 2017.

140 **3. Community Detection in the Network.** [4, 16] A community is a group
 141 of vertices which are densely connected by edges inside the group and sparsely (less
 142 densely) connected to vertices outside the group. The distinction between dense and
 143 sparse depends on the resolution parameter. A larger resolution parameter leads to
 144 a smaller number of total communities. There are several algorithms for determin-
 145 ing communities in a network, and different community detection algorithms may
 146 determine different communities inside a given network. As such, when referring to
 147 “the community of node i ” or c_i , we mean the community to which i belongs accord-
 148 ing to the community finding algorithm under discussion. Some formulations allow
 149 communities to overlap. In our formulation, we do not allow communities to overlap.

150 We used two methods of community detection in the network. Both are based
 151 on random walks on a network [10]. The first method we used was modularity max-
 152 imization, which seeks to partition a network into clusters such that there are more
 153 edges between nodes of the same cluster than between those of different clusters [4].
 154 In terms of random walks, the method of modularity maximization tries to maximize
 155 the length of a random walk contained in the same community [10]. We use the tra-
 156 ditional modularity function [13] with various values of the resolution parameter to
 157 achieve different numbers of communities (see Table 2). With a resolution parameter
 158 of 1, the network partitioned into six communities, but we used values above and
 159 below 1 to get partitions between four and eight communities; a larger resolution
 160 parameter leads to a smaller number of communities.

Because our network is directed, we must modify the undirected form of modu-
 larity into the directed form. Below is the undirected form of modularity,

$$Q = \frac{1}{2m} \sum_{i,j} \left[A_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j),$$

161 where m is the total number of edges in the unweighted network, A_{ij} is the ij^{th}
 162 element of the adjacency matrix, k_i and k_j represent the degree of node i and node j
 163 respectively, and $\delta(c_i, c_j)$ is the Kronecker delta function (which is 1 node i and node
 164 j are in the same community and is 0 otherwise). The null model, $P_{ij} = \frac{k_i k_j}{2m-1} \rightarrow \frac{k_i k_j}{2m}$
 165 in the limit of large m , represents the probability that node j is connected to node
 166 i according to the configuration model [9]. A graph in the configuration model is
 167 generated as follows: given a set of degrees where each degree is mapped to a node in
 168 the network (degree sequence), a node has edge “stubs” that are connected uniformly
 169 at random [14].

170 In transitioning to the directed case, P_{ij} must take into account the in-degree
 171 and out-degree of each node. As a result, P_{ij} becomes a directed version of the con-
 172 figuration model. This yields the following expression for modularity in the directed

173 case [9]:

174 (1)
$$Q = \frac{1}{m} \sum_{i,j} \left[A_{ij} - \frac{k_i^{\text{in}} k_j^{\text{out}}}{m} \right] \delta(c_i, c_j),$$

175 where k_i^{in} and k_j^{out} represent the in-degree of node i and out-degree of node j .

The algorithm that we used for calculating Q is based on minimizing the Hamiltonian \mathcal{H} . The Hamiltonian provided by [19] can be rewritten in our notation as

$$\mathcal{H} = - \sum_{i,j} (A_{ij} - \gamma P_{ij}) \delta(c_i, c_j)$$

176 where γ is the resolution parameter. Our definition is the same as in [19], where
 177 the pair use σ_i rather than c_i to denote the spin state of a node i in a graph. As
 178 shown in [19], when setting the resolution parameter $\gamma = 1$, one quickly notices that
 179 \mathcal{H} and Q are directly proportional by a constant $-\frac{1}{m}$. Consequently, modularity can
 180 be written as

181 (2)
$$Q = -\frac{1}{m} \mathcal{H}.$$

182 Since the Hamiltonian is negative, we can maximize modularity by minimizing \mathcal{H} . We
 183 used the Python library BCTPY [8] along with NETWORKX [5] to carry out the above
 184 calculations for directed modularity on the social network. Specifically, we used the
 185 method “modularity_dir” that uses a deterministic modularity maximization method
 186 and a resolution parameter of 1 to compare the clustering results to the word clustering
 187 results, which we discuss in section 4.2.

188 The second algorithm we used for partitioning is InfoMap, a method that uses
 189 random walks to determine community structure ² to represent information flow in a
 190 network [21]. A typical random walk across a network can be represented by a string of
 191 code words, where each word in the string corresponds to a node that the random walk
 192 visited in sequential order. InfoMap uses communities to shorten this representation
 193 by using a two-layered coding system where each community has a separate code
 194 word associated with the random walk entering and leaving the community, and
 195 each community also has a code system for the nodes within the community that
 196 are specific to that community [10]. Consequently, nodes in different communities
 197 can have the same code word, but because of the coded entry and exit words of
 198 the communities, these nodes can be differentiated from each other in the random
 199 walk. One finds an optimal partition of a network by solving the problem of how to
 200 most concisely represent random walks across the network by changing the encoding
 201 method of these walks based on the partitions. With the network partitioned using
 202 this optimal partition, a random walker should remain in the same community for a
 203 long time before exiting to a different community [10]. We used the implementation
 204 of the InfoMap algorithm from the IGRAPH software library [3], in R [17], using a
 205 method called “cluster_infomap.”

206 **4. Clustering By Language.**

²Community structure refers to the state of having grouped nodes in a graph according to their respective community. To “determine community structure,” means the process of determining all of the communities in a graph (including which node belongs to which community) by using a community detection algorithm.

207 **4.1. Preprocessing.** While one can, in principle, download the entire tweet
 208 history of a set of Twitter users and compare language use across individual tweets,
 209 we instead seek to acquire a general sense of what each user discusses over many
 210 tweets rather than in an individual tweet. We seek to find a way to summarize a
 211 Twitter user’s language usage by extracting their most-used words while filtering out
 212 irrelevant words (such as articles) that do not add useful information for study. After
 213 doing so, we can compare users through language usage. Therefore, we grouped the
 214 last 200 tweets of each Twitter user into one large text: the tweet corpus. We parsed
 215 the tweet corpus into a dictionary of key-value pairs, where each key is a user and each
 216 value is a string of that user’s 200 tweets put together. We then utilized Python’s
 217 NLTK library [1] to process the words used by each user. Before applying text-based
 218 clustering, the words must be filtered so that only relevant words remain. In following
 219 paragraphs, we discuss the text filtration process.

220 First, we use NLTK’s built-in list of English *stop words*. Stop words are words
 221 that do not have any real meaning associated with them. In particular, stop words
 222 consist of articles, conjunctions, and prepositions that only serve to connect other,
 223 more important words together. Although NLTK’s list of stop words only include
 224 English words, we decided not to include other languages because an overwhelming
 225 majority of the tweets were in English and we suspect that the use of other languages
 226 will play a role in connections.

227 After removing the stop words, we used NLTK’s built-in word stemmers to stem
 228 each word [1]. The process of stemming words involves taking a word and condensing
 229 it into its word stem. For example, the words “walking” and “walks” both become the
 230 word stem “walk”. To account for words occurring too frequently or too infrequently
 231 potentially skewing our results, we remove any words that appear in under twenty
 232 percent or over eighty percent of the users’ tweets.

233 After preprocessing the list of tweets for each user, we used the NLTK library
 234 to build a term frequency-inverse document frequency matrix (tf-idf). The tf-idf
 235 matrix is a product of two statistical measures: term frequency and inverse document
 236 frequency. The term frequency statistic measures the number of times a word appears
 237 in a text, such as the collection of a user’s tweets. The dimension of the term frequency
 238 matrix is $99 \times p$, where p is the number of stemmed words across the entirety of the
 239 tweet corpus. The inverse document frequency statistic weighs words by how often
 240 they occur across a set of documents, such as the entire tweet corpus from all ninety-
 241 nine users in our dataset. The dimension of the idf matrix is $p \times 1$. The idf statistic
 242 assigns less weight to very common words and more weight to unique words in the
 243 entire tweet corpus.

244 For our dataset, the term frequency matrix contains each user’s tweet collection
 245 as a row header and each stemmed word as a column header. To fill in the elements of
 246 the matrix, we calculate the term frequency of each word with respect to each user’s
 247 tweet collection. As a result, many entries of the matrix are zeros because it is rare
 248 when the same term is used by many users. We then multiply the term frequency by
 249 the idf matrix. The idf matrix is computed as follows:

$$250 \quad (3) \quad \text{idf}(t) = \ln \left(\frac{n_d}{f_d(t)} \right),$$

251 where t is a given term from our tweet corpus and $n_d = 99$ is the total number
 252 of documents or tweet collections (one per user). Also, $f_d(t)$ is the frequency of
 253 documents or tweet collections that include the word t . Note $f_d(t) \geq 1$, because any

254 given word t appears in at least one document.

255 In other words, the dimension of the tf-idf matrix is 99×1 . We then normalize
256 the tf-idf vector v by the Euclidean norm:

$$257 \quad (4) \quad \hat{v} = v_{\text{norm}} = \frac{v}{\|v\|_2} = \frac{v}{\sqrt{v_1^2 + v_2^2 + \dots + v_n^2}}.$$

258 **4.2. K-means Clustering.** Given the tf-idf matrix, we seek to partition the
259 data into distinct groups based on word usage. To obtain a partition, we use k -means
260 clustering [7], which attempts to minimize:

$$261 \quad (5) \quad \min_{C_1, \dots, C_K} \left\{ \sum_{k=1}^K W(C_k) \right\},$$

262 where

$$263 \quad (6) \quad W(C_k) = \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2,$$

264 and C_k denotes the set of observations in the k th cluster. Equation (6) is the sum
265 of all the pairwise squared Euclidean distances between the observations in the k th
266 cluster, divided by the total number of observations in the k th cluster. Combining
267 equations (5) and (6), the clustering procedure becomes an optimization problem:

$$268 \quad (7) \quad \min_{C_1, \dots, C_K} \left\{ \sum_{k=1}^K \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 \right\}.$$

269 To perform k -means clustering, we use the SCIKIT-LEARN library in Python [15].
270 One important thing to note about k -means clustering is that it requires specifying the
271 desired number, k , of clusters. We decided to match k to the number of communities
272 generated by the community-detection algorithms from section 3. We computed k -
273 means with $k = 4, 5, 6, 7, 8$.

274 We present an example of 5-means clustering in Table 1. As one can see from
275 Table 1, the algorithm was able to group users based on common tweet topics. For
276 example, Cluster 3 is centered around people who recently talked about the NBA
277 Finals. Cluster 2 is centered around current politics and the then-recent attack on
278 the London Bridge. Table 1 shows that the use of k -means clustering yields clusters
279 with words that have topical similarity, as opposed to clusters based on arbitrary word
280 usage. We suspect that the temporal nature of conversation topics within Twitter play
281 a role in our clustering. Also, recall that our tweet corpus consists of the last 200
282 tweets from each user. In future work, one could collect the entire tweet history of
283 these users to see how tweets over a longer time period affect the clustering results.

284 Through modularity-based community detection with a resolution of 0.8
285 we found 5 different communities. In Figure 1 we label these communities with letters
286 A through E. The size of each community from A to E is 23, 29, 1, 44, and 2. As
287 seen in Figure 1, the clusters found by k -means clustering overlap slightly with the
288 communities found through modularity-based community detection with a resolution
289 of 0.8. For example, Community A is the only community with the pink cluster,
290 and Community B is the only community with the orange cluster. However, both
291 shades of green and purple are in multiple communities. We provide a detailed way
292 to determine the overlap of clusters and communities in the next section.

TABLE 1
K-Means Clusters for $K = 5$

K-means	Number of Users	Most-Used Terms Per Cluster
Cluster 1	49	rt, new, love, thank, tonight, day
Cluster 2	6	attack, london, police, trump, terror, say
Cluster 3	5	nbafinals, game, espn, valverde, sportscenter, warriors
Cluster 4	31	love, u, thank, rt, tonight, happy
Cluster 5	7	en, el, que, la, para, por

Modularity Parameter 0.8 Algorithmic Community Detection and 5-Means Language Clustering

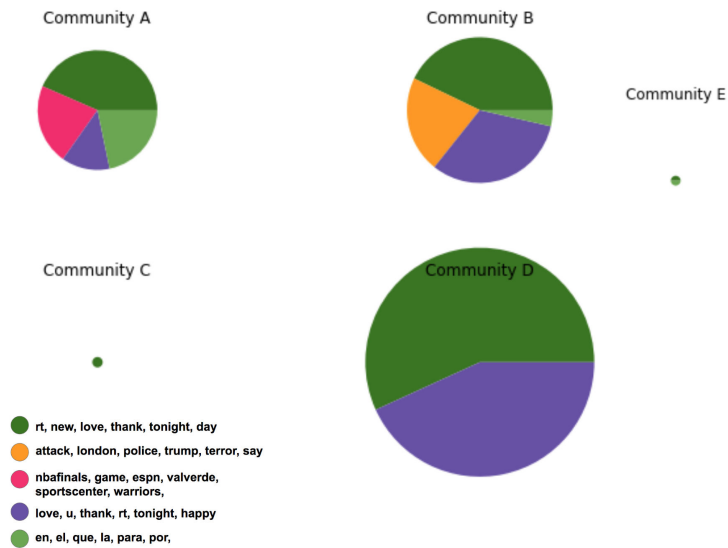


FIG. 1. Each community, which we indicate with a letter, is the output of detecting communities by maximizing modularity as described in section 3 with a resolution parameter of .8. The size of the pies reflects the number of users in an algorithmic cluster. The color represents text-based clusters outlined in 1. This figure serves as a visual demonstration of the extent to which clusters determined by language usage and communities determined by network structure overlap.

293 **5. Z-Score Calculations.** To answer our original question of whether connec-
 294 tions between users in the network reflected similar content in their tweets, we looked
 295 at how the users grouped based on the structure of the network (so-called “network
 296 communities”), as opposed to the language that they used (so-called “language clus-
 297 ters”). To compare the language clusters and network communities, we calculated the
 298 Rand coefficients [18] of the combination of one network grouping and one language
 299 grouping, and then determined the z-scores of these coefficients [6, 23]. The Rand
 300 coefficient is a measure of the similarity of two partitions on the same data. It ranges
 301 from 0 to 1, where 0 indicates no similarity and 1 indicates the two partitions are iden-
 302 tical. Although the distribution of the Rand coefficients is asymptotically Gaussian,
 303 and the z-scores and associated p-values of the coefficients cannot be interpreted as
 304 exact measures of significance, the z-scores and p-values are still good approximations

305 of the probabilities of seeing the given outcomes.

306 In comparing the network communities against our language clusters, finding
 307 pairs of people is of great importance. Each node in our network is a person, and if
 308 node pairs are assigned to the same group in the network community and also the
 309 language cluster then our clustering methods have some similarity. In adopting the
 310 Rand coefficient to compare partitions, we define M to be the total number of pairs of
 311 nodes in our network, M_1 to be the number of pairs in the same network community,
 312 M_2 to be the number of pairs in the same language cluster, and w to be the number
 313 of pairs that are in the same network community and in the same language cluster.
 314 The Rand coefficient is the $S = [w + (M - M_1 - M_2 + w)]/M$ [18] and the z-score is

315 (8)
$$z = \frac{1}{\sigma_w} \left(w - \frac{M_1 M_2}{M} \right),$$

316 where

317
$$\sigma_w^2 = \frac{M}{16} - \frac{(4M_1 - 2M)^2(4M_2 - 2M)^2}{256M^2} + \frac{C_1 C_2}{16n(n-1)(n-2)}$$

318
$$+ \frac{[(4M_1 - 2M)^2 - 4C_1 - 4M][(4M_2 - 2M)^2 - 4C_2 - 4M]}{64n(n-1)(n-2)(n-3)},$$

319

320 $n = 99$ is the total number of nodes in the network,

$$C_1 = n(n^2 - 3n - 2) - 8(n + 1)M_1 + 4 \sum_i n_i^3,$$

and

$$C_2 = n(n^2 - 3n - 2) - 8(n + 1)M_2 + 4 \sum_j n_j^3.$$

321 The summation terms in the calculation are based on a contingency table of
 322 network communities and language clusters, where the ij^{th} element of the table is the
 323 number of nodes in the i^{th} cluster based on network structure and the j^{th} cluster based
 324 on language use [23]. Here, $n_i = \sum_j n_{ij}$ is the row sum of the contingency table,
 325 signifying the number of nodes in the i^{th} network cluster. Similarly, $n_j = \sum_i n_{ij}$
 326 is the column sum of the contingency table, signifying the number of nodes in the j^{th}
 327 language cluster [23].

328 **6. Results.** As mentioned in [23], there are issues with using z-scores that are
 329 important to consider when interpreting results. Because the distribution of Rand
 330 coefficients is not Gaussian and often has heavy tails, it can be common to observe
 331 extreme z-score values. Nevertheless, even using an approximation of a Gaussian
 332 distribution is sufficient to claim significance for large enough z-scores. Although we
 333 ultimately do not claim significance, calculating z-scores is one viable way to compare
 334 our language clusters and structural communities.

335 Table 2 gives the z-scores between different comparisons of community divisions
 336 based on different methods. We compared the language clusters using various values
 337 of k to network communities derived using different methods (modularity at different
 338 resolutions and InfoMap clustering). As seen in Table 2, our methodology allows us
 339 to use a variety of k values and community-detection algorithms. From Table 2, we
 340 can also see that if we assume a Gaussian distribution, fourteen of the twenty-five

TABLE 2
Z-Scores for Comparison of Partitions

K-means	Mod (res=0.75)	Mod (res=0.80)	Mod (res=1.0)	Mod (res=1.055)	InfoMap
K = 4	3.8115	4.2733	1.3161	0.8090	1.4234
K = 5	3.9427	3.7199	1.5309	0.9298	1.9257
K = 6	3.9427	3.7199	1.5309	0.9298	1.9257
K = 7	7.7907	6.8025	4.7976	3.0046	0.0911
K = 8	12.4150	9.4211	5.7226	4.4967	1.6520

“K-means” represents the language clusters based on k-means clustering using the given value of k . “Mod” indicates partitioning using modularity maximization using the associated resolution. “InfoMap” indicates partitioning using InfoMap. The bold numbers across the diagonal signify comparisons of language-based partitions with network-based partitions with the same number of clusters.

z-scores would be statistically significant at a 1% level. Consequently, for the two specific methods of clustering that we consider, it is unlikely that the partitions of the network based on structure compared to partitions based on language would agree to the extent that we observed based on chance alone. Because our results did not exceed what previous papers have used as a threshold for significance [24] and the distribution of z-scores is not Gaussian, we only claim that our partitions are slightly similar, but not at a significant level. Our z-scores, however, reveal that comparing network structure to language can be compressed to a few numbers for easy analysis. For example, from the z-scores we can infer that while language use may play a role in how Twitter users create connections (i.e., follow others), language use is not a dominant factor in driving Twitter connections, and there are other factors that seem to be influencing social connections.

7. Conclusions and Discussion. Through our study we developed a new approach for answering questions about how people communicate and form connections with others.

In our example, we found that there is a slight relationship between language clusters and network communities in how similar they are among the top ninety-nine users on Twitter. Because we only used a small sample of users and only looked at their word usage in their tweets at a specific time, our conclusions only apply to this specific group of users within the time frame we examined them. Our results are influenced by our small sample size and our selective choice of subjects because many of the users in our dataset come from similar industries (specifically, the entertainment industry). Our dataset reflects a snapshot of popular culture at that time; many of the clustered words center around current events. To build on our work, we suggest that others examine larger datasets, also using other social networking sites with similar friendship structure, to see if there also exists similarities between network communities and clusters based on other characteristics, like language.

Our analysis of the top ninety-nine Twitter network illustrates that the words that users post can play a role in their connections on Twitter. As seen in Table 1, words that users of a certain network community use can center around certain topics, such as sports or news. Accordingly, what people talk about is fairly related to the people with whom they are connected. However, it is necessary to be cautious when interpreting our results. The z-values are not as large as in other papers [24]. While a user’s language may play some role in community structure, it is not the only

375 determining factor for connectivity between users.

376 Future studies also focusing on Twitter could adapt our approach to look at
 377 other network groups, either larger or smaller, and perhaps selected based on other
 378 criteria. Determining if the results are consistent across varying sizes and structures of
 379 networks might lead to more insight on what underlying factors contribute to network
 380 structure.

381 Additional studies could also explore larger and more diverse social and infor-
 382 mation networks (in addition to Twitter). Another example of a network that uses
 383 language would be the user network of Instagram users and the use of language in
 384 hashtags. Future studies could apply our approach to this network to see if similar-
 385 ities in language clusters and network communities appear in this network as well.
 386 In addition, the communities in a social network change over time, and the content
 387 depends heavily on news, especially groundbreaking pieces. Capturing the changes in
 388 communities both through language and network structure over time, and especially
 389 as they change with current events, would also be an interesting extension of our
 390 study.

391 **Acknowledgments.** We would like to thank the UCLA Mathematics Depart-
 392 ment as well as Professor Mason Porter for his guidance throughout our project. We
 393 would also like to acknowledge brandonrose.org [20] for their post on “Document
 394 Clustering with Python,” which provided guidance for the text clustering part of our
 395 work.

396

REFERENCES

- 397 [1] S. BIRD, E. KLEIN, AND E. LOPER, *Natural Language Processing with Python: Analyzing Text*
 398 *with the Natural Language Toolkit*, O’Reilly Media, Inc., 2009.
- 399 [2] J. BRYDEN, S. FUNK, AND V. A. JANSEN, *Word Usage Mirrors Community Structure in the*
 400 *Online Social Network Twitter*, EPJ Data Science, 2 (2013), p. 3.
- 401 [3] G. CSARDI AND T. NEPUSZ, *The IGRAPH Software Package for Complex Network Research*,
 402 *InterJournal, Complex Systems* (2006), p. 1695, <http://igraph.org>.
- 403 [4] S. FORTUNATO AND D. HRIC, *Community Detection in Networks: A User Guide*, Physics
 404 Reports, 659 (2016), pp. 1–44.
- 405 [5] A. A. HAGBERG, D. A. SCHULT, AND P. J. SWART, *Exploring Network Structure, Dynamics,*
 406 *and Function Using NetworkX*, in Proceedings of the 7th Python in Science Conference
 407 (SciPy2008), Pasadena, CA USA, Aug 2008, pp. 11–15.
- 408 [6] L. HUBERT, *Nominal Scale Response Agreement as a Generalized Correlation*, British Journal
 409 of Mathematical and Statistical Psychology, 30 (1977), pp. 98–103.
- 410 [7] G. JAMES, D. WITTEN, T. HASTIE, ET AL., *An Introduction to Statistical Learning*, vol. 6,
 411 Springer, 2013.
- 412 [8] R. LAPLANTE, *Brain Connectivity Toolbox for Python (version 0.5.0)*, 2016, <https://github.com/aestrivex/bctpy> (accessed 2017-06-16).
- 413 [9] E. A. LEICHT AND M. E. J. NEWMAN, *Community Structure in Directed Networks*, Physical
 414 Review Letters, 100 (2008).
- 415 [10] N. MASUDA, M. A. PORTER, AND R. LAMBIOTTE, *Random Walks and Diffusion on Networks*,
 416 Physics Reports, (2017).
- 417 [11] M. MCPHERSON, L. SMITH-LOVIN, AND J. M. COOK, *Birds of a Feather: Homophily in Social*
 418 *Networks*, Annual review of sociology, 27 (2001), pp. 415–444.
- 419 [12] S. A. MYERS, A. SHARMA, P. GUPTA, AND J. LIN, *Information Network or Social Network?:*
 420 *The Structure of the Twitter Follow Graph*, in Proceedings of the 23rd International Con-
 421 ference on World Wide Web, ACM, 2014, pp. 493–498.
- 422 [13] M. E. J. NEWMAN, *Finding Community Structure in Networks Using the Eigenvectors of*
 423 *Matrices*, Phys. Rev. E, 74 (2006).
- 424 [14] M. E. J. NEWMAN, *Networks: An Introduction*, Oxford University Press, 2010.
- 425 [15] F. PEDREGOSA, G. VAROQUAUX, A. GRAMFORT, V. MICHEL, ET AL., *Scikit-learn: Machine*
 426 *Learning in Python*, Journal of Machine Learning Research, 12 (2011), pp. 2825–2830.
- 427 [16] M. A. PORTER, J.-P. ONNELA, AND P. J. MUCHA, *Communities in Networks*, Notices of the
 428

- 429 AMS, 56 (2009), pp. 1082–1097.
430 [17] R DEVELOPMENT CORE TEAM, *R: A Language and Environment for Statistical Computing*, R
431 Foundation for Statistical Computing, Vienna, Austria, 2011, <http://www.R-project.org>.
432 ISBN 3-900051-07-0.
433 [18] W. M. RAND, *Objective Criteria for the Evaluation of Clustering Methods*, Journal of the
434 American Statistical Association, 66 (1971), pp. 846–850.
435 [19] J. REICHARDT AND S. BORNHOLDT, *Statistical Mechanics of Community Detection*, Physical
436 Review E, 74 (2006).
437 [20] B. ROSE, *Document Clustering with Python*, 2015, <http://brandonrose.org/clustering> (accessed
438 2017-06-15).
439 [21] M. ROSVALL AND C. T. BERGSTROM, *Maps of Random Walks on Complex Networks Re-*
440 *veal Community Structure*, Proceedings of the National Academy of Sciences, 105 (2008),
441 pp. 1118–1123.
442 [22] THE PYTHON-TWITTER DEVELOPERS, *A Python Wrapper Around the Twitter API*, 2017, [http:](http://python-twitter.readthedocs.io/en/latest/)
443 [//python-twitter.readthedocs.io/en/latest/](http://python-twitter.readthedocs.io/en/latest/) (accessed 2017-06-16).
444 [23] A. L. TRAUD, E. D. KELSIC, P. J. MUCHA, AND M. A. PORTER, *Comparing Community Struc-*
445 *ture to Characteristics in Online Collegiate Social Networks*, SIAM Review, 53 (2011),
446 pp. 526–543.
447 [24] A. L. TRAUD, P. J. MUCHA, AND M. A. PORTER, *Social Structure of Facebook Networks*,
448 Physica A: Statistical Mechanics and its Applications, 391 (2012), pp. 4165–4180.
449 [25] TWEETPY, *Tweepy v3.5.0*, 2017, <http://docs.tweepy.org/en/v3.5.0/api.html> (accessed 2017-06-
450 16).
451 [26] TWITTER, *REST APIs v1.1 - Twitter Developers*, 2017, <https://dev.twitter.com/rest/public>
452 (accessed 2017-06-16).