

Detecting communities on social networks

1st Nguyen Hoang An
University of Information Technology
Ho Chi Minh City, Vietnamese
20521051@gm.uit.edu.vn

Abstract—The report discusses a multi-layer clustering algorithm for community detection on social networks, specifically Twitter. This is a community clustering algorithm based on the similarity of structural components and content components, called MCT. From an application perspective, identifying such communities can provide better recommendations, along with other benefits.

I. INTRODUCTION

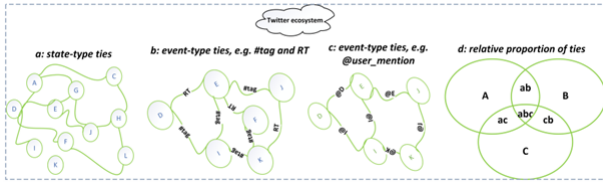


Fig. 1. Types of interaction on Twitter.

Today's social media networks like Twitter and Facebook are communities created by the diverse connections of countless users. Consequently, this ecosystem is highly interconnected and forms various types of interactions, such as interactions among friends of friends, posts with tags, direct mentions @, or retweets on Twitter in Fig. 1. These interactions are what create communities on social media at different levels.

II. RELATED WORK

A. Graph-based and hierarchical methods

Graph-based methods are algorithms that use graphs or network representations to model and analyze the structure of a system. In these methods, nodes represent entities within the system, and edges represent the relationships between those entities. Examples of graph-based methods include centrality measures, community detection algorithms, and network flow analysis.

Hierarchy-based methods are algorithms that organize data or systems into hierarchical structures, where smaller units are nested within larger units. In these methods, the data or system is recursively partitioned into smaller and smaller units based on certain criteria, such as similarity or distance. Examples of hierarchy-based methods include hierarchical clustering algorithms, decision trees, and nested sampling methods.

B. Multiview and bi-modal clustering

Multi-view clustering is an approach to clustering data that involves multiple representations or views of the same underlying dataset. In this approach, each view represents a different aspect or feature of the data, and the goal is to integrate these views to determine a consensus clustering. Multi-view clustering can be useful in cases of complex or high-dimensional data where different views can capture different aspects of the data. Examples of multi-view clustering algorithms include Co-regularized Spectral Clustering and Multi-view K-Means.

Bi-modal clustering is a type of clustering algorithm that operates on data with two distinct modes or sources of information. In this approach, the objective is to identify clusters that are consistent across both modes. Bi-modal clustering can be useful in cases where data has two separate but related sources of information, such as text and images, or gene expression and protein-protein interactions. Examples of bi-modal clustering algorithms include Bi-Clustering and Coupled Bi-Clustering.

III. DATASETS

A. Mental – Health Twitter

This data was collected using the Twitter API. The tweets have been filtered to retain only English contexts. It aims to classify users' mental health. The dataset with fields such as post_id, post_created, post_text, user_id, followers, friends, favourites, statuses, retweets, label. Refer to Fig. 2

#	post_id	post_created	post_text	user_id	followers
Index	Post ID of the post	Post created on	Uncleaned Tweet	User Identification	No of Foll.

Fig. 2. Mental – Health Twitter datasets

B. Tweets from US Capitol Riots

This data was generated using the Twitter API + Tweepy. It consists of 80,000 tweets from January 6, 2021. The tweets have been stripped of mentions, hyperlinks, emojis, and punctuation. All text has been converted to lowercase. The dataset with fields such as tweet_id, text, query, user_id,

user_name, follower_count, user_tweet_count, likes, retweets, location_name. Refer to Fig. 3

tweets.csv (19.86 MB)				
Detail Compact Column				
10 of 14 columns				
tweet_id	text	query	user_id	user_name
Unique ID assigned to each Tweet	Content of the tweet posted	What kind of query is the content related to	Unique ID assigned to the user who tweeted.	Name of the user who tweeted
	in 2017 georgias ne...	3% pandemic	7%	[null]

Fig. 3. Tweets from US Capitol Riots datasets

IV. METHODOLOGY

The algorithm is defined through 3 stages.

Stage 1: Based on Dyadic Relationships - Bidirectional relationships between two nodes. Simmelian (Transitive) Relationships - Relationships between groups of at least 3 nodes, where there exist direct or indirect relationships between pairs of nodes to calculate the compatibility between node pairs using a algorithm in Fig. 5. The output of the algorithm is an $n \times n$ adjacency matrix representing the compatibility between nodes.

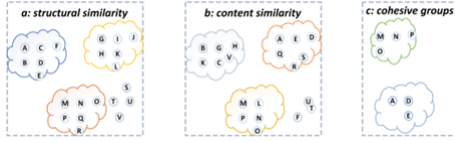


Fig. 4. Node clustering in Detecting communities according to stages

$$p(R_{v_i, v_j}) = \frac{1}{1 + \exp \varphi}$$

Trong đó: $\varphi = -\log(\epsilon_{v_i, v_j} + J(v_i, v_j)) \times (\epsilon_{v_i, v_j} + J(v_i, v_j))$

Fig. 5. Algorithm for calculating compatibility between Node pairs

Stage 2: We will convert the text components into individual vectors and apply the Latent Dirichlet Allocation (LDA) algorithm to identify a set of topics within the text corpus. The distribution of topics for each tweet stream from related nodes is represented as a distribution of topics in the LDA model. Next, we utilize the Jensen-Shannon distance (JSD) to measure the similarity between two probability distributions. This distance is symmetric and always has a finite value. The result returned is a similarity matrix of size $m \times m$ regarding content.

Stage 3: Apply clustering methods on both structural and textual aspects, then perform a function to combine the structural similarity matrix and the content similarity matrix to form an aggregate matrix. Next, transform it into a distance matrix and cluster on this matrix until the stopping criteria are met. The result of the algorithm is the clustered community

V. EXPERIMENTS

Before proceeding with the main algorithm, we need to evaluate clustering algorithms for structural and content-based clustering.

A. Structural

The execution of the algorithm will produce two sets: Sr for structurally similar nodes and Su for structurally dissimilar nodes. The results are as depicted in the Fig. 6. The structurally similar parts will be saved in the file

10011277, 10011278, 10011279, 10011280, 10011281, 10011282, 10011283, 10011284, 10011285, 10011286, 10011287, 10011288, 10011289, 10011290, 10011291, 10011292, 10011293, 10011294, 10011295, 10011296, 10011297, 10011298, 10011299, 10011300, 10011301, 10011302, 10011303, 10011304, 10011305, 10011306, 10011307, 10011308, 10011309, 10011310, 10011311, 10011312, 10011313, 10011314, 10011315, 10011316, 10011317, 10011318, 10011319, 10011320, 10011321, 10011322, 10011323, 10011324, 10011325, 10011326, 10011327, 10011328, 10011329, 10011330, 10011331, 10011332, 10011333, 10011334, 10011335, 10011336, 10011337, 10011338, 10011339, 10011340, 10011341, 10011342, 10011343, 10011344, 10011345, 10011346, 10011347, 10011348, 10011349, 10011350, 10011351, 10011352, 10011353, 10011354, 10011355, 10011356, 10011357, 10011358, 10011359, 10011360, 10011361, 10011362, 10011363, 10011364, 10011365, 10011366, 10011367, 10011368, 10011369, 10011370, 10011371, 10011372, 10011373, 10011374, 10011375, 10011376, 10011377, 10011378, 10011379, 10011380, 10011381, 10011382, 10011383, 10011384, 10011385, 10011386, 10011387, 10011388, 10011389, 10011390, 10011391, 10011392, 10011393, 10011394, 10011395, 10011396, 10011397, 10011398, 10011399, 10011400, 10011401, 10011402, 10011403, 10011404, 10011405, 10011406, 10011407, 10011408, 10011409, 10011410, 10011411, 10011412, 10011413, 10011414, 10011415, 10011416, 10011417, 10011418, 10011419, 10011420, 10011421, 10011422, 10011423, 10011424, 10011425, 10011426, 10011427, 10011428, 10011429, 10011430, 10011431, 10011432, 10011433, 10011434, 10011435, 10011436, 10011437, 10011438, 10011439, 10011440, 10011441, 10011442, 10011443, 10011444, 10011445, 10011446, 10011447, 10011448, 10011449, 10011450, 10011451, 10011452, 10011453, 10011454, 10011455, 10011456, 10011457, 10011458, 10011459, 10011460, 10011461, 10011462, 10011463, 10011464, 10011465, 10011466, 10011467, 10011468, 10011469, 10011470, 10011471, 10011472, 10011473, 10011474, 10011475, 10011476, 10011477, 10011478, 10011479, 10011480, 10011481, 10011482, 10011483, 10011484, 10011485, 10011486, 10011487, 10011488, 10011489, 10011490, 10011491, 10011492, 10011493, 10011494, 10011495, 10011496, 10011497, 10011498, 10011499, 10011500, 10011501, 10011502, 10011503, 10011504, 10011505, 10011506, 10011507, 10011508, 10011509, 10011510, 10011511, 10011512, 10011513, 10011514, 10011515, 10011516, 10011517, 10011518, 10011519, 10011520, 10011521, 10011522, 10011523, 10011524, 10011525, 10011526, 10011527, 10011528, 10011529, 10011530, 10011531, 10011532, 10011533, 10011534, 10011535, 10011536, 10011537, 10011538, 10011539, 10011540, 10011541, 10011542, 10011543, 10011544, 10011545, 10011546, 10011547, 10011548, 10011549, 10011550, 10011551, 10011552, 10011553, 10011554, 10011555, 10011556, 10011557, 10011558, 10011559, 10011560, 10011561, 10011562, 10011563, 10011564, 10011565, 10011566, 10011567, 10011568, 10011569, 10011570, 10011571, 10011572, 10011573, 10011574, 10011575, 10011576, 10011577, 10011578, 10011579, 10011580, 10011581, 10011582, 10011583, 10011584, 10011585, 10011586, 10011587, 10011588, 10011589, 10011590, 10011591, 10011592, 10011593, 10011594, 10011595, 10011596, 10011597, 10011598, 10011599, 10011600, 10011601, 10011602, 10011603, 10011604, 10011605, 10011606, 10011607, 10011608, 10011609, 10011610, 10011611, 10011612, 10011613, 10011614, 10011615, 10011616, 10011617, 10011618, 10011619, 10011620, 10011621, 10011622, 10011623, 10011624, 10011625, 10011626, 10011627, 10011628, 10011629, 10011630, 10011631, 10011632, 10011633, 10011634, 10011635, 10011636, 10011637, 10011638, 10011639, 10011640, 10011641, 10011642, 10011643, 10011644, 10011645, 10011646, 10011647, 10011648, 10011649, 10011650, 10011651, 10011652, 10011653, 10011654, 10011655, 10011656, 10011657, 10011658, 10011659, 10011660, 10011661, 10011662, 10011663, 10011664, 10011665, 10011666, 10011667, 10011668, 10011669, 10011670, 10011671, 10011672, 10011673, 10011674, 10011675, 10011676, 10011677, 10011678, 10011679, 10011680, 10011681, 10011682, 10011683, 10011684, 10011685, 10011686, 10011687, 10011688, 10011689, 10011690, 10011691, 10011692, 10011693, 10011694, 10011695, 10011696, 10011697, 10011698, 10011699, 10011700, 10011701, 10011702, 10011703, 10011704, 10011705, 10011706, 10011707, 10011708, 10011709, 10011710, 10011711, 10011712, 10011713, 10011714, 10011715, 10011716, 10011717, 10011718, 10011719, 10011720, 10011721, 10011722, 10011723, 10011724, 10011725, 10011726, 10011727, 10011728, 10011729, 10011730, 10011731, 10011732, 10011733, 10011734, 10011735, 10011736, 10011737, 10011738, 10011739, 10011740, 10011741, 10011742, 10011743, 10011744, 10011745, 10011746, 10011747, 10011748, 10011749, 10011750, 10011751, 10011752, 10011753, 10011754, 10011755, 10011756, 10011757, 10011758, 10011759, 10011760, 10011761, 10011762, 10011763, 10011764, 10011765, 10011766, 10011767, 10011768, 10011769, 10011770, 10011771, 10011772, 10011773, 10011774, 10011775, 10011776, 10011777, 10011778, 10011779, 10011780, 10011781, 10011782, 10011783, 10011784, 10011785, 10011786, 10011787, 10011788, 10011789, 10011790, 10011791, 10011792, 10011793, 10011794, 10011795, 10011796, 10011797, 10011798, 10011799, 10011800, 10011801, 10011802, 10011803, 10011804, 10011805, 10011806, 10011807, 10011808, 10011809, 10011810, 10011811, 10011812, 10011813, 10011814, 10011815, 10011816, 10011817, 10011818, 10011819, 10011820, 10011821, 10011822, 10011823, 10011824, 10011825, 10011826, 10011827, 10011828, 10011829, 10011830, 10011831, 10011832, 10011833, 10011834, 10011835, 10011836, 10011837, 10011838, 10011839, 10011840, 10011841, 10011842, 10011843, 10011844, 10011845, 10011846, 10011847, 10011848, 10011849, 10011850, 10011851, 10011852, 10011853, 10011854, 10011855, 10011856, 10011857, 10011858, 10011859, 10011860, 10011861, 10011862, 10011863, 10011864, 10011865, 10011866, 10011867, 10011868, 10011869, 10011870, 10011871, 10011872, 10011873, 10011874, 10011875, 10011876, 10011877, 10011878, 10011879, 10011880, 10011881, 10011882, 10011883, 10011884, 10011885, 10011886, 10011887, 10011888, 10011889, 10011890, 10011891, 10011892, 10011893, 10011894, 10011895, 10011896, 10011897, 10011898, 10011899, 10011900, 10011901, 10011902, 10011903, 10011904, 10011905, 10011906, 10011907, 10011908, 10011909, 10011910, 10011911, 10011912, 10011913, 10011914, 10011915, 10011916, 10011917, 10011918, 10011919, 10011920, 10011921, 10011922, 10011923, 10011924, 10011925, 10011926, 10011927, 10011928, 10011929, 10011930, 10011931, 10011932, 10011933, 10011934, 10011935, 10011936, 10011937, 10011938, 10011939, 10011940, 10011941, 10011942, 10011943, 10011944, 10011945, 10011946, 10011947, 10011948, 10011949, 10011950, 10011951, 10011952, 10011953, 10011954, 10011955, 10011956, 10011957, 10011958, 10011959, 10011960, 10011961, 10011962, 10011963, 10011964, 10011965, 10011966, 10011967, 10011968, 10011969, 10011970, 10011971, 10011972, 10011973, 10011974, 10011975, 10011976, 10011977, 10011978, 10011979, 10011980, 10011981, 10011982, 10011983, 10011984, 10011985, 10011986, 10011987, 10011988, 10011989, 10011990, 10011991, 10011992, 10011993, 10011994, 10011995, 10011996, 10011997, 10011998, 10011999, 10012000, 10012001, 10012002, 10012003, 10012004, 10012005, 10012006, 10012007, 10012008, 10012009, 10012010, 10012011, 10012012, 10012013, 10012014, 10012015, 10012016, 10012017, 10012018, 10012019, 10012020, 10012021, 10012022, 10012023, 10012024, 10012025, 10012026, 10012027, 10012028, 10012029, 10012030, 10012031, 10012032, 10012033, 10012034, 10012035, 10012036, 10012037, 10012038, 10012039, 10012040, 10012041, 10012042, 10012043, 10012044, 10012045, 10012046, 10012047, 10012048, 10012049, 10012050, 10012051, 10012052, 10012053, 10012054, 10012055, 10012056, 10012057, 10012058, 10012059, 10012060, 10012061, 10012062, 10012063, 10012064, 10012065, 10012066, 10012067, 10012068, 10012069, 10012070, 10012071, 10012072, 10012073, 10012074, 10012075, 10012076, 10012077, 10012078, 10012079, 10012080, 10012081, 10012082, 10012083, 10012084, 10012085, 10012086, 10012087, 10012088, 10012089, 10012090, 10012091, 10012092, 10012093, 10012094, 10012095, 10012096, 10012097, 10012098, 10012099, 10012100, 10012101, 10012102, 10012103, 10012104, 10012105, 10012106, 10012107, 10012108, 10012109, 10012110, 10012111, 10012112, 10012113, 10012114, 10012115, 10012116, 10012117, 10012118, 10012119, 10012120, 10012121, 10012122, 10012123, 10012124, 10012125, 10012126, 10012127, 10012128, 10012129, 10012130, 10012131, 10012132, 10012133, 10012134, 10012135, 10012136, 10012137, 10012138, 10012139, 10012140, 10012141, 10012142, 10012143, 10012144, 10012145, 10012146, 10012147, 10012148, 10012149, 10012150, 10012151, 10012152, 10012153, 10012154, 10012155, 10012156, 10012157, 10012158, 10012159, 10012160, 10012161, 10012162, 10012163, 10012164, 10012165, 10012166, 10012167, 10012168, 10012169, 10012170, 10012171, 10012172, 10012173, 10012174, 10012175, 10012176, 10012177, 10012178, 10012179, 10012180, 10012181, 10012182, 10012183, 10012184, 10012185, 10012186, 10012187, 10012188, 10012189, 10012190, 10012191, 10012192, 10012193, 10012194, 10012195, 10012196, 10012197, 10012198, 10012199, 10012200, 10012201, 10012202, 10012203, 10012204, 10012205, 10012206, 10012207, 10012208, 10012209, 10012210, 10012211, 10012212, 10012213, 10012214, 10012215, 10012216, 10012217, 10012218, 10012219, 10012220, 10012221, 10012222, 10012223, 10012224, 10012225, 10012226, 10012227, 10012228, 10012229, 10012230, 10012231, 10012232, 10012233, 10012234, 10012235, 10012236, 10012237, 10012238, 10012239, 10012240, 10012241, 10012242, 10012243, 10012244, 10012245, 10012246, 10012247, 10012248, 10012249, 10012250, 10012251, 10012252, 10012253, 10012254, 10012255, 10012256, 10012257, 10012258, 10012259, 10012260, 10012261, 10012262, 10012263, 10012264, 10012265, 10012266, 10012267, 10012268, 10012269, 10012270, 10012271, 10012272, 10012273, 10012274, 10012275, 10012276, 10012277, 10012278, 10012279, 10012280, 10012281, 10012282, 10012283, 10012284, 10012285, 10012286, 10012287, 10012288, 10012289, 10012290, 10012291, 10012292, 10012293, 10012294, 10012295, 10012296, 10012297, 10012298, 10012299, 10012300, 10012301, 10012302, 10012303, 10012304, 10012305, 10012306, 10012307, 10012308, 10012309, 10012310, 10012311, 10012312, 10012313, 10012314, 10012315, 10012316, 10012317, 10012318, 10012319, 10012320, 10012321, 10012322, 10012323, 10012324, 10012325, 10012326, 10012327, 10012328, 10012329, 10012330, 10012331, 10012332, 10012333, 10012334, 10012335, 10012336, 10012337, 10012338, 10012339, 10012340, 10012341, 10012342, 10012343, 10012344, 10012345, 10012346, 10012347, 10012348, 10012349, 100123
--

Perform clustering on the network data 'D' with a standard number of communities 't'.

Initialize functions like 'compute_reciprocity_prob', 'compute_text_similarity' , 'f_structural', 'f_textual' , 'compute_pairwise_similarities', 'detecting_communities' to mine communities based on sample data. Due to large data and not powerful personal computers, it leads to very long waiting times when taking large standard community numbers. For that reason, I do it with $t = 10$.

VI. RESULT ANALYSIS

The clustering method has demonstrated its ability to detect and identify highly cohesive communities within data. In doing so, it not only helps address challenges stemming from incomplete data but also provides a comprehensive view of the community network. Additionally, analyzing communities at multiple levels aids in a deeper understanding of the relationships and structures within each community. This not only enhances the scalability of the method but also increases its flexibility when applied to various contexts and domains. This has proven to be quite beneficial for intelligent recommendation systems.

VII. CONCLUSION AND FUTURE WORK

Regarding the complexity of the proposed method, when the network size is large, the challenge lies in accurately determining when the algorithm for community detection will converge. Therefore, the research team only uses one iteration to analyze the complexity of the algorithm, providing information on future performance. Assuming the execution complexity of a basic algorithm parameterized as $f(C)$, the complexity of the algorithm will be approximately $O(f(C) \times s \times r \times m)$, where s is the number of comparisons to decide the next cluster, r is the size, and m is the number of runs. In terms of execution, the complexity of the algorithm is relatively low; however, it tends to increase as the amount of data increases, thus requiring continued improvement in future work.

The complexity depends on the number of comparisons to decide the next cluster and the number of iterations to determine the optimal number of clusters. So there is a need to research a method that can identify more intrinsic components accurately to minimize the number of comparisons as well as the number of iterations.

REFERENCES

- [1] Isa Inuwa-Dutse, Mark Liptrott, Ioannis Korkontzelos. *A multilevel clustering technique for community detection*. 2021.
- [2] Researchgate. *Simmelian Ties on Twitter: Empirical Analysis and Prediction*. 2019.
- [3] Researchgate. *Analysis and Prediction of Dyads in Twitter*. 2019.
- [4] Infamouscoder, Kaggle. *Mental – Health Twitter datasets*. 2022.