

Detecting Geographically Dispersed Overlay Communities Using Community Networks

Madhushi Bandara, Dharshana Kasthurirathna, Danaja Maldeniya, Mahendra Piraveenan

Abstract—Community detection is an extremely useful technique in understanding the structure and function of a social network. Louvain algorithm, which is based on Newman-Girman modularity optimization technique, is extensively used as a computationally efficient method extract the communities in social networks. It has been suggested that the nodes that are in close geographical proximity have a higher tendency of forming communities. Variants of the Newman-Girman modularity measure such as dist-modularity try to normalize the effect of geographical proximity to extract geographically dispersed communities, at the expense of losing the information about the geographically proximate communities. In this work, we propose a method to extract geographically dispersed communities while preserving the information about the geographically proximate communities, by analyzing the ‘community network’, where the centroids of communities would be considered as network nodes. We suggest that the inter-community link strengths, which are normalized over the community sizes, may be used to identify and extract the ‘overlay communities’. The overlay communities would have relatively higher link strengths, despite being relatively apart in their spatial distribution. We apply this method to the Gowalla online social network, which contains the geographical signatures of its users, and identify the overlay communities within it.

Keywords—Social networks, community detection, modularity optimization, geographically dispersed communities.

I. INTRODUCTION

TOPOLOGICAL analysis of social networks have gained prominence in recent years. With the advent of network science as a separate field, modeling and characterizing self-organizing networks have been applied in a plethora of areas, ranging from biological networks, financial networks to social networks [1]-[2]. One of the most vital pieces of information that is embedded in a social network is its community structure [3]. These communities may have homogeneous features in diverse attributes. However, the communities that are extracted through topological features is becoming increasingly relevant in social network analysis and mining, as it is the most fundamental and objective form of communities that can be extracted from a social network. Identifying and extracting communities from a social network may be vital in myriad applications which involve social network analysis, such as modeling social influence, information spread and epidemic modeling and defense related

applications, among others. Moreover, extracting communities using a formal methodology may help to understand the structure and the operation of the social network in concern.

Different community detection algorithms have been proposed and have been applied in multitude of applications. Among them are the minimum-cut method, hierarchical clustering and modularity maximization [4], [5]-[6]. The most widely accepted method of community detection is the modularity maximization, where the modular behavior of a network is quantified to identify and extract communities in a network.

One of the key limitations of the modularity maximization in community detection is that it does not take into account the contribution of geographical proximity that is vital in forming communities. That is, the nodes that are in close spatial proximity may tend to form communities more than the nodes that are geographically apart. Thus, the interactions and links among the nodes that are geographically apart should carry more significance compared to the nodes that are in close proximity in extracting communities. In order to address this limitation, the dist-modularity measure has been recently proposed [7]. This particular measure attempts to normalize the strength of links formed among nodes over their geographical distance. Thus, the dist-modularity measure may be used to identify the communities that are geographically distributed. However, the dist-modularity measure has two key limitations. It requires relatively high computational time due to its computational complexity. Also, by normalizing the effect of geographical proximity of the constituent nodes in extracting communities, it actually disregards the communities consisting of geographically proximate nodes, which are equally as important as the geographically dispersed communities. Thus, it may not be used to capture the geographically proximate communities that are strongly connected with each other, while being geographically apart. Such communities can be observed in real-world networks such as migrant worker community networks and terrorist networks [8], [9], where the communities formed by geographically proximate nodes may have strong links with similar communities that may be geographically apart. Identifying and extracting such communities may provide vital information that may not be apparent in modularity based community detection algorithms.

In this work, we suggest that observing the interconnections of communities extracted through modularity maximization, in other words, analyzing the ‘community networks’ may pave way to identify and extract the geographically distributed communities. In order to do this, we suggest that the

M. Bandara and D. Maldeniya are with Lirneasias, 12 Balcombe Place, Colombo 08, Sri Lanka (e-mail: madhushi@lirneasias.net, danaja@lirneasias.net)

D. Kasthurirathna is with Sri Lanka Institute of Information technology, New Kandy Rd, Malabe 10115, Sri Lanka (e-mail: dharshana.k@slit.lk).

M. Piraveenan is with the Complex Systems Research Group, Faculty of Engineering & IT, University of Sydney, New South Wales 2006, Australia (e-mail: mahendrarajah.piraveenan@sydney.edu.au).

centroids of communities may be used to form a network of communities, where the links among such communities may be used to identify geographically dispersed communities or overlay communities. Detecting such overlay communities may have interesting applications in areas such as indirect marketing, information propagation modeling and defense and counter terrorism domains. Identifying such overlay communities may be used as a computationally efficient method to extract geographically distributed communities.

The rest of this paper is organized as follows. The Background section provides an overview of the different community detection methods that are derived from modularity optimization, including the dist-modularity measure. Then, the methodology section describes the proposed method of detecting overlay communities and how it is applied to a real-world social network dataset to extract geographically distributed overlay communities. Afterwards, the results obtained from the experimental analysis is presented. Finally, the concluding remarks are presented along with a brief description on potential future work.

II. BACKGROUND

Network science emerged as a prominent field of science with the proposition of the scale-free network model [1], [4]. The study of network science has facilitated observing networks in diverse domains such as social networks, biological networks and financial networks. In the recent years, much interest have been given to extracting communities of social networks. The community information may be vital in understanding the information flows and the structure of an existing social network. While multitude of community extraction techniques such as the Hierarchical Clustering method and minimum cut method have been proposed to extract community information from social networks; modularity maximization remains the most widely used technique of extracting communities of social networks.

Louvain algorithm is a heuristic method developed by Blondel et al. [10], that partitions a social network into communities while optimizing Newman-Girvan (N-G) modularity of the partition. Louvain algorithm improves on the computational time of the modularity optimization technique, which is originally an NP-hard problem. Newman-Girvan modularity is used to measure how densely the detected communities of the partition are connected relative to connections between these communities [10], [7]. In other words, the Newman-Girvan modularity measure is the fraction of edges within communities in the observed network minus the expected value of that fraction in a null model, which serves as a reference and should characterize some features of the observed network. Equation (1) defines the Newman-Girvan modularity measure, which is used in the Louvain algorithm.

Given a network that is modeled as an undirected graph $G = (V;E)$ where V is a set of nodes and E is a set of relationships among nodes. The variables n and m represent the cardinalities of V and E respectively. Each edge $(v_i; v_j)$

is assumed to have an associated weight w_{ij} . For a given node $v_i \in V$, $\eta_i = \{v_j | (v_i, v_j) \in E \vee (v_j, v_i) \in E\}$ and $k_i = |\eta_i|$. Accordingly, the modularity $M(C)$ of a given partition C is given as;

$$M(C) = \frac{1}{2m} \sum_{c \in C} \sum_{i,j \in c} w_{ij} - P_{ij} \quad (1)$$

$$P_{ij} = \frac{k_i \cdot k_j}{2m} \quad (2)$$

Here, P_{ik} refers to the null model that is used as a reference model, where the edges of the network are rewired randomly while preserving the degree distribution.

Multitude of social networks, including online social networks incorporate location information of the nodes in the network, in addition to the nodes and relationships among them, which may be utilized to extract the geographically specific information of the individual nodes. One important aspect in geographically distributed social networks is that the nodes in close proximity have an inherent nature of connecting with each other [7]. Thus, the community detection algorithms should ideally take into account this feature and normalize the effect of proximity to identify the actual communities in a social network.

As a result, a subsequent modularity measure called dist-modularity [7] has been proposed to normalize the effect of geographical proximity. This measure tries to identify the geographically distributed communities with a distance decaying function, under the assumption that the nodes that are in close geographical proximity have a higher tendency of forming community structures. This is an important assumption that we too employ indirectly, in formulating the idea of overlay communities that are geographically distributed.

Equation (3) gives the formal definition of the dist-modularity function.

$$M_{dist}(C) = \frac{1}{2m} \sum_{c \in C} \sum_{i,j \in c} w_{ij} - P_{ij} \quad (3)$$

$$P_{ij} = \frac{\widehat{P}_{ij} + \widehat{P}_{ji}}{2} \quad (4)$$

$$\widehat{P}_{ij} = \frac{k_i k_j f(d(v_i, v_j))}{\sum_{v_q \in V} k_q f(d(v_q, v_i))}; f: R^+ \rightarrow (0, 1] \quad (5)$$

Here, f is the distance-decay function. The basic assumption in dist-modularity optimization is that each node exerts a field on the surrounding nodes, which is inversely proportional to the distance from it. Thus, the null model used in the dist-modularity calculation assumes that nodes which are closer based on the distance function are more likely to be connected. This is the same assumption that we'd be utilizing to propose the idea of overlay communities where the Newman-Girvan modularity is considered to be likely to extract communities of members within the same geographical proximity.

The distance decaying function used in dist-modularity measure may further be extended to a gravity model where

the inherent node properties are taken into account to capture the heterogeneity of the nodes [11].

$$P_{ij} = N_i N_j f(d_{ij} | d_{ij} = d) \quad (6)$$

where N_i captures the importance of the node. Thus, the distance-decaying function may be modified to capture the node heterogeneity as:

$$f(d) = \frac{\sum_{i,j|d_{i,j}=d} A_{ij}}{\sum_{i,j|d_{i,j}=d} N_i N_j} \quad (7)$$

which is the weighted average of the probability for a link to exist at distance d .

While the dist-modularity measure and its variants attempt to normalize of effect of geographical proximity in extracting communities, another branch of modularity maximization techniques attempt to harness the spatial information to extract the communities based on their geographical closeness. Spatially-near Modularity [12], which correlates with the spatial proximity of nodes, is an example of this particular approach of modularity maximization. Another interesting application of the modularity optimization is where it is been applied to extract multilevel communities based on a 'similarity attribute'. This particular application works under the assumption that nodes with similar features have a higher probability of being connected to each other. While this particular measure is useful in extracting communities that share the same geographical space, it is not much useful in extracting communities that are geographically distributed [13].

Based on the existing literature, two main approaches can be observed in extracting communities with geographical constraints. One is to extract geographically dispersed communities by normalizing the effect of geographical proximity. Dist-modularity and its variations are used for this purpose. The other approach is to harness or exploit the geographical proximity of communities and purposely consider the spatial nearness in extracting the communities. While these two approaches seem to contradict each other, the communities in social networks may encompass both geographically proximate as well as geographically dispersed communities. Thus, we attempt to propose a method to extract both the geographically proximate as well as geographically dispersed communities in a complimentary fashion. In other words, we propose a method to extract the geographically distributed communities, based on the interconnections of geographically proximate communities.

III. METHODOLOGY

In order to resolve this apparent dilemma where the geographically distributed communities have to be extracted without losing the information on geographically proximate communities, we propose the concept of 'overlay-communities', quite similar to the idea of 'overlay-networks' in peer-to-peer computing [14]. The idea is to extract the communities using the Louvain algorithm and then connect the extracted communities with

inter-community links assuming that the nodes that are in close proximity have a higher probability of being in the same community [7].

When connecting the communities, we consider the centroids of each community as the 'node' of the community network, in order to assign a geographical location to each community. Afterwards, the connections among the members in communities are aggregated into 'links'. This way, we can easily quantify the geographical alignment of each community along with their inter-community link strengths. The link strength of each link are then divided by the multiplication of the sizes of the communities that it connects, in order to normalize the effect of the heterogeneity of community sizes. Normalizing over community sizes would help to identify the communities that are geographically distributed and yet strongly connected with each other, irrespective of the sizes of the underlying geographically proximate communities.

The communities that are strongly connected over the community network are termed the 'overlay communities', within the context of this work. Based on the assumption that the nodes that are in close geographical proximity tend to form communities, we may argue that the communities that are in close geographical proximity may tend to form strong connections with each other. Thus, the most interesting overlay communities would be those which are strongly connected yet whose centroids are further apart. Extracting such overlay communities may reveal information about the geographically distributed communities social networks that are not apparent and that cannot be identified using the existing community detection algorithms. The algorithm 1 explains the proposed technique in detail.

Algorithm 1: Extracting overlay communities using community networks

- 1 Extract the community set C using the Louvain method of N-G modularity optimization;
 - 2 **for each** community c in the set of communities C **do**
 - 3 Identify the centroid of each community based on geographical location of each node in the community ;
 - 4 Assign the centroid as the node representing that particular community in the community network ;
 - 5 **for each** community pair p in the set of communities C **do**
 - 6 Compute the strength of the link connecting the community pair p by aggregating the connections among the nodes in community pair p ;
 - 7 Normalize the link strengths by the community sizes by dividing the link strengths by the multiplication of community sizes of the community pair p ;
 - 8 Identify the communities that are relatively further apart geographically yet have relatively higher link strengths as the 'overlay communities' ;
-

In order to test the effectiveness of the proposed algorithm, it was applied to a real-world social networks with geographical information. We used a dataset from the Gowlla online social network [15] which has the geographical signatures of the users included in it for this purpose. By applying the above algorithm to the Gowlla network, we

could extract the geographically distributed communities by identifying the communities that are strongly connected while being geographically distributed.

The Gowalla network data set has 196,591 connected users as nodes. Check-in details of only 107,092 users were available for analysis, from which the home locations could be derived. There were 1,900,654 edges connecting the 196,591 users indicating the friendship between them. There were 6,442,892 total check-in records. It was observed that a relatively large portion of users did not have check-in records, resulting in unknown home locations for them. Hence, only the users with known location information were considered for the location analysis. We made the assumption that the other members were in similar vicinity. We ignored communities where locations of all members were unknown.

To derive approximate home location of each user, most frequent location for check-in was calculated. Based on that, all the check-in locations greater than 95% of the distance from the most frequent location were filtered out, assuming they represent anomaly trips of the user. Then, for each user, the weighted center of gravity of his check-in locations was calculated. That was considered the home location for them.

To detect basic communities we considered the results of Louvain modularity optimization algorithm as it is computationally efficient and yield better results in comparison to many existing techniques. When the Gowalla social network was processed into communities by Louvain multi-level community detection algorithm, 5 non-overlapping community levels were detected. At level 1 network was broken into 19,396 communities, 2875 at level 2, 1025 at level 3, 839 and 820 at level 4 and 5 respectively. Level 5 produced the maximum modularity value for the network.

To study how the resulting communities are dispersed geographically, we used home locations of community members to calculate the standard distance deviation which calculates the centroid of the community (with respect to the dispersion of community members) as well as the radius and the area of dispersion. The results indicated that higher the modularity results, the radius and area of the community become small. This further supports the assumption that the nodes that are in close geographical proximity have a higher tendency of forming communities.

The link strengths of the 820 communities extracted were measured. There were 4138 links connecting the 820 communities. The link density was observed to be relatively low. We then normalized the link strengths over the community sizes of the community pairs connected by each link in order to remove the effect of the heterogeneity of the community sizes, which could otherwise invariably affect the strengths of the inter-community links. The next section presents some of the results obtained using the analysis performed on the extracted community network.

IV. RESULTS

Fig. 1 depicts the distribution of the communities over the community size, in the community network formed using the above methodology. Based on the figure, it is evident that

the community network contains relatively few communities with large number of members while relatively high number of communities with relatively smaller number of members. This is characteristic of the scale-free model.

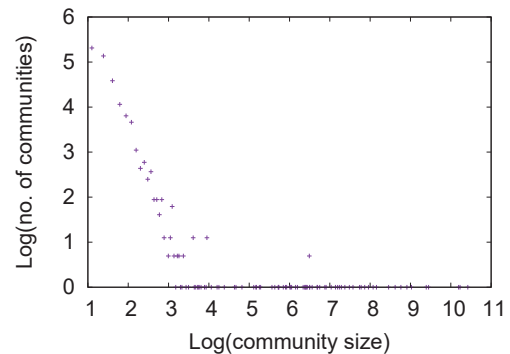


Fig. 1 Distribution of communities over community size in logarithmic scale

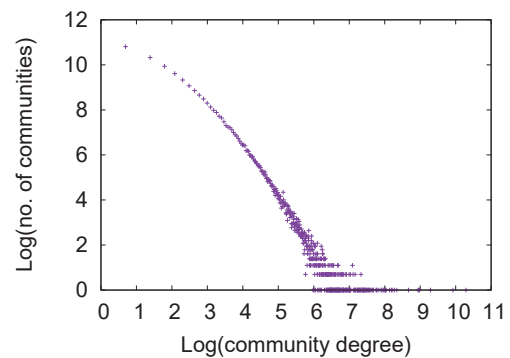


Fig. 2 The degree distribution of the community network in logarithmic scale

We further observe the degree distribution of communities in Fig. 2. According to the figure, the degree distribution fits well into a power-law degree distribution. The scale-free correlation and the scale-free exponent of the network were measured to be 0.74 and 0.67, respectively, further indicating that the community network fits into the scale-free model.

Fig. 3 depicts a graphical representation of the community network obtained, where the link strengths were normalized over community sizes of the communities connected by each link. As the figure depicts, the community sizes and link strengths are heterogeneous and non-correlated in nature, suggesting that certain communities may be strongly connected, despite being geographically apart.

Fig. 4 (a) depicts that strength of each link against the Euclidean distance between the centroids of the communities that it connects. The link strengths are not normalized over community sizes. According to the figure, it is evident that there are certain links that show relatively higher strengths, while they connect communities that are geographically dispersed.

Fig. 4 (b) shows the strength of each link, where the link strength of each link is normalized over the community sizes

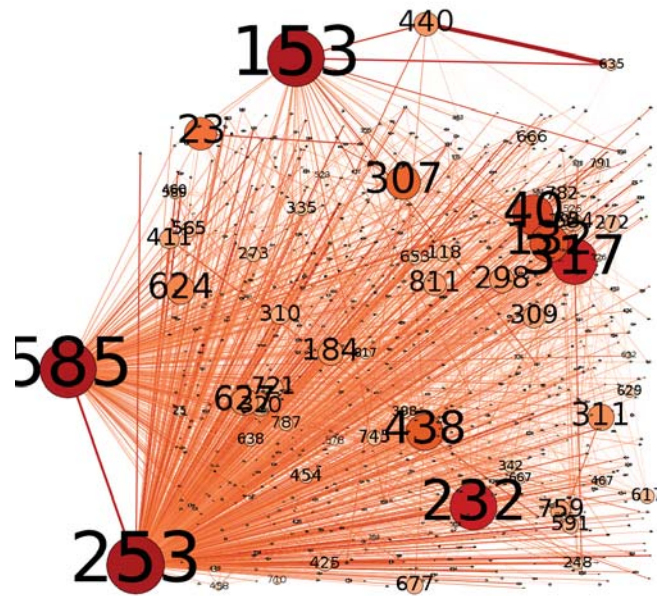


Fig. 3 Community network with heterogeneous node sizes and normalized link strengths

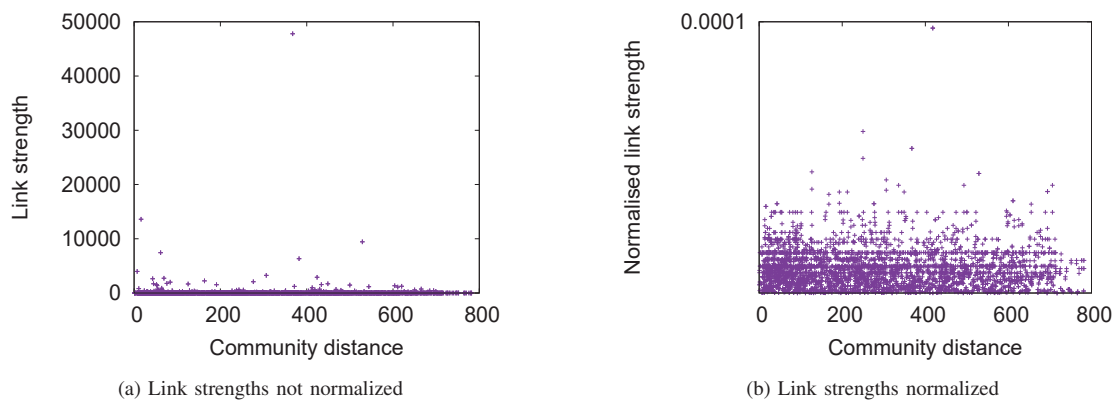


Fig. 4 The link strength among communities against distance between the community centroids. The link strength is determined based on the number of interconnections among communities. The distance between communities is measured in geographical coordinate based distance

connecting them. Even after normalizing the link strengths, there are certain communities that are tightly connected despite being geographically apart. In the given figure, there is a link that particular stands out in its link strength, although it connects two communities that are spatially apart. If we consider the link with the highest strength in the given distribution, it is the link connecting the community id 253 (population of size 26702) with the community with the id 585 (population of size 33560). The relative distance between these two communities is 368 units. When we normalize the tie strength by population, the link strength of these communities is higher than 82% of community pairs observed. Hence we cannot ignore this strong connection between the two communities consider it to be a result of the community size of one or both communities. It is important to note that these two communities are not overlapping and geographically apart significantly, when extracted using the Louvain algorithm,

even though they have a strong connection between them. Thus, we may identify this particular community pair as a geographically dispersed single overlay community.

V. CONCLUSION AND FUTURE WORK

In this work, we propose a computationally efficient method to extract geographically dispersed communities, while preserving the information about the geographically proximate communities. We suggest that the Louvain algorithm may be used to identify the geographically proximate communities, which may then be connected by aggregating the connections among the nodes between each community. This is based on the assumption that the nodes that are in close geographical proximity have a higher tendency to form communities with each other. The centroids of each community would represent the nodes in the community network. These inter-community links could be used to form a 'community network'. The link

strengths are then normalized over the sizes of the community pairs. In the resulting network, the community pairs that have relatively higher link strengths, while connecting relatively further communities are identified as geographically dispersed 'overlay communities'. These overlay communities may be used to identify geographically dispersed communities in applications such as migrant community detection and terrorist network detection, since the community networks in such scenarios have geographically proximate communities interacting with similar communities that are geographically apart.

Though we measure the Euclidean distance between two community centroids as the distance between communities it may not be fair for the cases where communities are very large, making their centers far apart. Yet they can be overlapping at the periphery. Thus, non-overlapping relatively smaller communities would be more appropriate to be considered for overlay community detection in our approach.

To our knowledge, there has not been an attempt to analyze a network of communities based on its topological properties. It may be possible to analyze the community networks to extract more information about numerous network properties and behavior such as network resilience [16], assortativity [17], growth [1] and evolution.

Though we only consider community pairs in this work, the overlay communities could be in the form of sub-networks. Further, the overlay networks may be extracted at multiple levels of hierarchy. Thus, extracting these sub-networks and the hierarchical overlay networks could be the potential extensions of this work. Further, while we consider the number of interactions within each community to denote a link and to measure link strengths, different network attributes may be used to form links and assign link strengths, resulting in networks of varying dimensions. Such community networks may be analyzed to extract information about the network that may not be visible with existing network analysis techniques.

ACKNOWLEDGMENT

The authors would like to thank Mr. Sriganesh Lokanathan, Prof. Rohan Samarajiva and Mr. Isuru Jayasooriya of *Lirneasia* [18] for their support and contribution.

REFERENCES

- [1] R. Albert and A.-L. Barabási, "Statistical mechanics of complex networks," *Reviews of Modern Physics*, vol. 74, pp. 47–97, 2002.
- [2] D. Kasthurirathna, A. Dong, M. Piraveenan, and I. Y. Tumer, "The failure tolerance of mechatronic software systems to random and targeted attacks," in *ASME 2013 International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*. American Society of Mechanical Engineers, 2013, pp. V005T06A036–V005T06A036.
- [3] M. Girvan and M. E. Newman, "Community structure in social and biological networks," *Proceedings of the national academy of sciences*, vol. 99, no. 12, pp. 7821–7826, 2002.
- [4] M. E. Newman, "The structure and function of complex networks," *SIAM review*, vol. 45, no. 2, pp. 167–256, 2003.
- [5] J. Leskovec, K. J. Lang, A. Dasgupta, and M. W. Mahoney, "Statistical properties of community structure in large social and information networks," in *Proceedings of the 17th international conference on World Wide Web*. ACM, 2008, pp. 695–704.

- [6] M. E. Newman, "Modularity and community structure in networks," *Proceedings of the national academy of sciences*, vol. 103, no. 23, pp. 8577–8582, 2006.
- [7] P. Shakarian, P. Roos, D. Callahan, and C. Kirk, "Mining for geographically disperse communities in social networks by leveraging distance modularity," in *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2013, pp. 1402–1409.
- [8] M. G. Herander and L. A. Saavedra, "Exports and the structure of immigrant-based networks: the role of geographic proximity," *Review of Economics and Statistics*, vol. 87, no. 2, pp. 323–335, 2005.
- [9] R. Medina and G. Hepner, *Geospatial analysis of dynamic terrorist networks*. Springer, 2008.
- [10] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *Journal of statistical mechanics: theory and experiment*, vol. 2008, no. 10, p. P10008, 2008.
- [11] P. Expert, T. S. Evans, V. D. Blondel, and R. Lambiotte, "Uncovering space-independent communities in spatial networks," *Proceedings of the National Academy of Sciences*, vol. 108, no. 19, pp. 7663–7668, 2011.
- [12] J. Hannigan, G. Hernandez, R. M. Medina, P. Roos, and P. Shakarian, "Mining for spatially-near communities in geo-located social networks," *arXiv preprint arXiv:1309.2900*, 2013.
- [13] X. Liu, T. Murata, and K. Wakita, "Extracting the multilevel communities based on network structural and nonstructural information," in *Proceedings of the 22nd international conference on World Wide Web companion*. International World Wide Web Conferences Steering Committee, 2013, pp. 191–192.
- [14] K. Aberer, L. O. Alima, A. Ghodsi, S. Girdzijauskas, S. Haridi, and M. Hauswirth, "The essence of p2p: a reference architecture for overlay networks," in *Peer-to-Peer Computing, 2005. P2P 2005. Fifth IEEE International Conference on*. IEEE, 2005, pp. 11–20.
- [15] J. Leskovec and A. Krevl, "{SNAP Datasets}://{Stanford} large network dataset collection," 2014.
- [16] R. Albert, H. Jeong, and A.-L. Barabási, "Error and attack tolerance of complex networks," *nature*, vol. 406, no. 6794, pp. 378–382, 2000.
- [17] M. E. J. Newman, "Mixing patterns in networks," *Physical Review E*, vol. 67, no. 2, p. 026126, 2003.
- [18] Lirneasia. [Online]. Available: <http://www.lirneasia.net>

Dharshana Kasthurirathna Dharshana Kasthurirathna holds a PhD in Complex Systems from the University of Sydney. Currently he is working as a Senior Lecturer at the Faculty of Computing of Sri Lanka Institute of Information Technology.

Madhushi Bandara Madhushi Bandara graduated holds a Bachelor of Science in Computer Science & Engineering from the University of Moratuwa. Currently she's working as a Junior Researcher at Lirneasia and as a lecturer at the Department of Computer Science & Engineering of University of Moratuwa.

Danaja Maldeniya Danaja Maldeniya holds a Bachelor of Science in Computer Science & Engineering from the University of Moratuwa. Currently he's working as a Senior Researcher at Lirneasia.

Mahendra Piraveenan Dr. Mahendra Piraveenan is a lecturer affiliated to the complex systems research group, within the faculty of engineering and information technologies at University of Sydney. He holds a Bachelor of Engineering in computer systems engineering (first class) from University of Adelaide and a Ph.D from University of Sydney. His research interests include complex systems, networked game theory, social networks and infectious disease dynamics.