

Modularity maximization to find community structure in complex networks

Bilal SAOUD

dept. Electrical engineering

University of Bouira

Bouira, Algeria

bilal340@gmail.com

Abstract—Complex networks have in generally communities. These communities are very important. Network’s communities represent sets of nodes, which are very connected. In this research, we developed a new method to find the community structure in networks. Our method is based on flower pollination algorithm (FPA) witch is used in the splitting process. The splitting of networks in our method maximizes a function of quality called modularity. We provide a general framework for implementing our new method to find community structure in networks. We present the effectiveness of our method by comparison with some known methods on computer-generated and real-world networks.

Index Terms—Community detection, Networks, Flower pollination algorithm, Normalized mutual information, Modularity

I. INTRODUCTION

Many systems can be represented by network or graphs, which makes them very powerful structure. A network G is defined by two sets [1]. The first set is vertex set V (node set) and the second is edge set E . Vertexes share relationships between them. Relationships are represented by edges. In general, the number of nodes is $|V| = n$ and edges is $|E| = m$. Euler’s solution of the Seven Bridges of Königsberg problem is considered to be the first use of networks to represent systems [2]. Today networks are used to illustrate several systems. For instance, in social network, which is an interaction between entities (persons, groups of persons, organizations, web sites, ...), can be represented by a network with two sets V and E . Vertexes stand for entities (for example persons) and edges stand for relationships between entities (for examples between persons). Analyzing and understanding a network leads to understand better the system. Among features that can help to understand the structure of a network, we can find the community structure.

Community structure exists in networks and it gives more information about the network. For instance, we can understand very well the system, which is represented by a network, by finding its community structure and the relationship between communities. In addition, networks can represent many systems like social networks, electric networks, biological networks, etc. It is vital to develop new methods to find network’s communities. When we analyze networks by studying relationships between nodes we can get

extra information about networks and systems. In general, nodes in the same community have common properties or insure similar tasks in network. Basically, a network has parts that are more densely connected than other parts. In other words, the nodes in these parts share many edges between them. These parts of nodes and edges are called communities (clusters). Finally, many studies have been done around networks and how to find community structure.

Many community structure detection methods have been developed [3]. According to the type of network, we can find methods for unipartite/bipartite networks, weighted/unweighted networks and directed/undirected networks. Furthermore, methods can be classified into different classes such as hierarchical methods (merging or splitting), methods that are based on maximization of an objective function. Some methods find disjoint communities, where intersection between communities is empty. However, other methods were designed to find overlapping communities, for instance the method in [4], where the intersection between communities is not empty.

In this paper, we address the problem of finding community structure in networks. We present a new method to discover community structure in unweighted and undirected networks. Our method is based on nature-inspired metaheuristics algorithm. We have developed our method based on the pollination process of flowers [9]. Our method is an hierarchical one. It is based on the splitting of a given network $G(V, E)$, which models a system. Splitting step in our method is done by the flower pollination algorithm (FPA) [9] in order to optimize the function of quality called modularity Q . The process of splitting will be stopped when the graph G has been disconnected, which means that each node of G represents a community. Finally, our method builds a dendrogram and finds the the most optimal community structure $\pi = \{c_1, \dots, c_k\}$, such as $\bigcup_{i=1}^k c_i = V$ and $c_i \cap c_j = \emptyset$ (for $i, j = 1 : k$).

The paper is organized as follows. The concept of FPA is presented in Section II. Our approach is detailed in Section III. Experimental results and discussions are given in Section IV. Finally, Section V concludes the paper.

II. FPA PRESENTATION

Flower pollination is an interesting phenomena in nature. Based on the studying flower pollination process, a new algorithm of optimization was designed by Yang in [9]. The algorithm has been named Flower Pollination Algorithm (FPA). In nature pollination can be abiotic form or biotic form. In general, 90% of flower have biotic pollination where the pollen is transferred by animals (pollinator) like insects. Biotic pollination by bees for instance can be done at long distance.

FPA has three steps [9] described in the following:

- In the first step, the algorithm initializes its parameters and generate the initial population. The best solution is found also in the first step.
- The second step, flowers in population start doing pollination in d-dimensional search (solution space). Flowers can choose a local or global pollination at every iteration in the search space. The algorithm switch between local pollination and global pollination based on probability $p \in [0, 1]$. Flowers location represent the vector of solutions vector and the value of objective function for every solutions is estimated. According to the value of objective function the new solution is evaluated and updated at every iteration and the best solution will may be improved.
- In the final step, the algorithm stops after some iterations and the best solution will be selected.

FPA can converge very fast and can escape the problem of local minima because it makes the long distances movement based on levy flight [10]. FPA can be used to solve diffrent problems like in [8].

III. A NEW METHOD TO FIND COMMUNITIES IN NETWORKS

In this section, we present our hierarchical method to discover community structure in networks. Hierarchical methods can be divisive or agglomerative. Our method is hierarchical divisive method. Network is divided by our method based on the maximization of the function of quality called modularity [6]. Our method is designed to find community structure in networks with only a single type of vertex and undirected, unweighted edge.

We can measure the strength of a community structure by the function of quality called modularity [6]. Modularity function Q is based on the observed edges fraction $e(c_i)$ within communities and the expected edges fraction $a(c_i)$ within the same communities, $Q = \sum_{c_i} e(c_i) - a(c_i)^2$. Modularity can be estimated for undirected and unweighted graph $G(V, E)$ as:

$$Q = \frac{1}{2m} \sum_i \sum_j (A[i, j] - P[i, j]) \delta(c_i, c_j) \quad (1)$$

where n is the number of nodes in G ($n = |V|$), m is the number of edges in G ($m = |E|$) and the community structure is $\pi = \{c_1, \dots, c_k\}$. $A_{n,n}$ represents the adjacency matrix of $G(V, E)$. For any vertex $i \in V$, d_i is the degree of node i

and c_i its community. The matrix A takes two values 1 or 0 if there is an edge between node i and j then $A[i, j] = 1$ or $A[i, j] = 0$ if there is not a connection between i and j . $P_{n,n}$ represents the adjacency matrix that corresponds null model. In the null model the probability of an existing edge between vertexes i and j is $P_{i,j} = \frac{d_i \times d_j}{2m}$. Finally, δ function is given as follows:

$$\delta(c_i, c_j) = \begin{cases} 1 & \text{if } c_i = c_j \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

Values of Q are between 0 and 1. Q closer to 1 indicates stronger community structures. According to Clauset et al. [6], a value above about 0.3 is a good indicator of significant community structure in a network.

Let $G(V, E)$ be an undirected and unweighted network, where $V = (v_1, \dots, v_n)$ is the set of vertexes, $E = (e_1, e_2, \dots, e_m)$ is the set of edges. The goal of our community detection method is to partition the network G into k communities (groups): $\pi = \{c_1, c_2, \dots, c_k\}$, where $c_i \neq \emptyset$, $c_i \cap c_j = \emptyset$, ($i = 1 : k, j = 1 : k$) and $V = \bigcup_{i=1}^k c_i$. In addition, our method finds the community structure π of the network G with the greatest value of modularity Q . To reach this goal, we used a FPA. Our method splits $G(V, E)$ into two new networks G_1 and G_2 . Nodes of each new network represent a community. Nodes of G_1 represent a community c_1 and nodes of G_2 represent a community c_2 . The splitting is based on FPA in order to maximize the value of modularity function Q . Then, G_1 and G_2 will be split until the network G has been disconnected. At the end of our method each node in $G(V, E)$ represents a community. Finally, we get a dendrogram for our method and the community structure will be chosen based on value of modularity Q or the number of communities.

The general algorithm of our method to find community structure is as follows:

Algorithm 1: The algorithm of our method

Data: $G(V, E)$

Result: dendrogram

- 1 $\pi = FPA()$, find a partition π based on FPA;
 - 2 Divide G based π , $G = G_1 + G_2$;
 - 3 Update the matrix of merge M for a final dendrogram;
 - 4 Go to *Steps 1* for each graph G_1 and G_2 ;
 - 5 Return the final dendrogram;
-

Fig. 1 shows the dendrogram that was built by our method on Zachary's Karate Club network [11]. Our method gave a community structure with two communities, which were separated by vertical lines on dendrogram. Labels of dendrogram are nodes (members of Zachary's Karate Club). The community structure that was found by our method on the same network is also represented in Fig. 2. In this figure, communities' nodes have different colors and shapes.

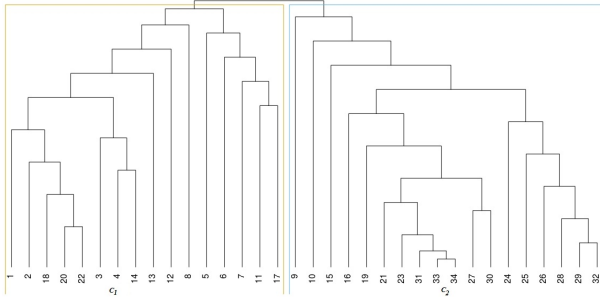


Fig. 1. The dendrogram of Zachary's Karate Club network created by our method.

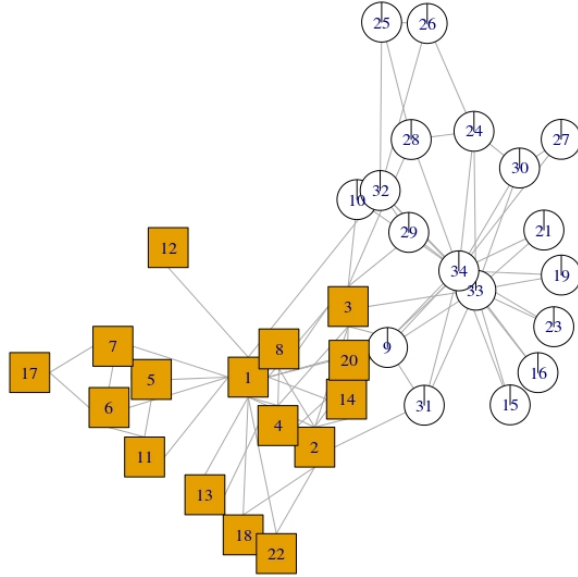


Fig. 2. Zachary's karate club network community structure is detected by our method.

IV. EXPERIMENTS AND RESULTS

To evaluate our method to find community structures in networks, we have tested it on computer-generated and several real networks (Zachary's Karate Club [11], American College Football [5], Dolphins [16], Books about US Politics [17], Jazz musicians [18], Word adjacencies [19] and Les Miserables [20]). We have compared our method with some well-known methods: fast greedy method [6], label propagation method [7], and infomap method [12].

A. Normalized Mutual Information

The comparison of our method with other methods is done based on the normalized mutual information (NMI) function [13]. The NMI is a powerful function to compare a community structure that was founded by methods with the real community structure. The value of NMI is based on defining a confusion matrix N , where the rows represent the real communities, and the columns represent the found communities. N_{ij} is the number of nodes in the real community that appears in the found community j . For two partitions A

and B , the partition A represents the real partition with c_A communities and B represents the found partition with c_B communities, The normalized mutual information (NMI) is estimated as follows:

$$NMI(A, B) = \frac{-2 \sum_{i=1}^{c_A} \sum_{j=1}^{c_B} N_{ij} \log\left(\frac{N_{ij}N}{N_i N_j}\right)}{\sum_{i=1}^{c_A} N_i \log\left(\frac{N_i}{N}\right) + \sum_{j=1}^{c_B} N_j \log\left(\frac{N_j}{N}\right)} \quad (3)$$

NMI values are in the range $[0, 1]$. Partitions A and B are identical if $NMI(A, B) = 1$.

B. Dataset based on computer-generated networks

Our method is tested on computer-generated networks benchmark proposed by Lancichinetti et al. [14]. The benchmark parameters are the number of nodes N , the exponents γ and β of the degree and community size distribution respectively (both distributions are power laws), the number of average degree $\langle k \rangle$, number of communities N_c , and the mixing parameter μ . Each node shares a fraction $(1 - \mu)$ of its links with other nodes of its community and a fraction μ with the other nodes of the network.

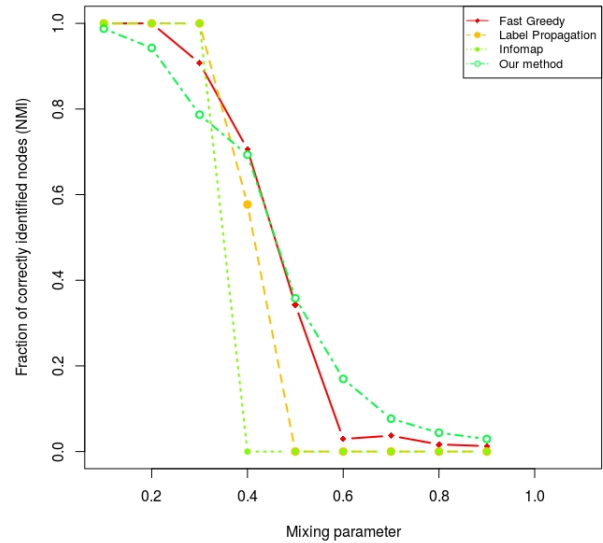


Fig. 3. NMI vs. mixing parameter μ .

Fig. 3 shows the variation of the NMI obtained by our method, fast greedy method, label propagation method and infomap method on the benchmarks networks, with the parameters: mixing parameter μ between 0.1 and 0.9, $\langle k \rangle = 16$, $\gamma = 3$, $\beta = 2$, $N = 128$ and $N_c = 4$. The value of NMI obtained by our method is high when μ changes from 0 to 0.5 and the same thing with other methods. At this range, nodes share many edges with nodes of its community that makes the community structure very clear and easy to find. Methods could group the most nodes in the correct communities when the mixing parameter μ is in $[0, 0.5]$. When μ is in $[0.5 - 0.9]$ range, it is difficult for all methods to find the true community structure. At this range

nodes share few edges with nodes of its community and many edges with nodes from other communities, which makes the community structure unclear and difficult to find. However, our method is still more accurate than the other methods. Our method evaluates the community structure at each step of splitting process and at the end our method selects the best community structure based on modularity value. From Fig. 3, we see that our method can discover community structure better than fast greedy, label propagation method and infomap method when μ is greater than 0.5.

Fig. 4 illustrates the result of our method on network generated by computer with mixing parameter $\mu = 0.8$. Fig. 4 shows the different communities that were found by our method. On this network with a mixing parameter $\mu = 0.8$, our method found a community structure (π) with eight communities ($\pi = \{c_1, c_2, \dots, c_8\}$). Dendrogram labels stand for nodes. In this example, we have a network with 128 nodes. We mention that the community structure can be found by breaking the dendrogram (Fig. 4) at different levels [15]. In our case, we have chosen to break the dendrogram at the level which maximizes the modularity function.

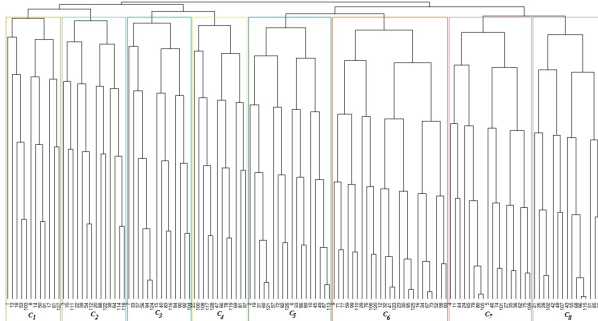


Fig. 4. The dendrogram and community structure by our method on computer generated network with mixing parameter $\mu = 0.8$.

C. Dataset based on real networks

In this section, we give the simulation results of our method, fast greedy, label propagation and infomap on real networks. We considered some real networks drawn from disparate fields (Zachary [11], Dolphins [16], Football [5] and Books about US politics [17]), where the community structure is known, which made them suitable to evaluate community detection methods.

- 1) Zachary's club network [11] is a real network that corresponds to a social network of friendships between 34 members of a karate club at a university in the United States in the 1970 ($n = 34$ and $m = 78$). The network has two clusters.
- 2) Dolphins Network [16] is an undirected social network of frequent associations between 62 dolphins in a community living off Doubtful Sound, New Zealand. This network ($n = 62$ and $m = 159$) has two communities.

- 3) College football network [5] represents the schedule of Division I Games for the year 2000 season. This network is made of 115 teams (nodes) and 613 edges. It is divided into 12 groups.
- 4) Books about US politics Network [17] is a network of books about US politics published around the time of the 2004 presidential election and sold by the online bookseller Amazon.com. Edges between books represent frequent purchasing of books by the same buyers. Compiled by Valdis Krebs. Books network has three communities.

Table I gives obtained results on networks. In this table, for each network we have estimated the value of modularity function according to equation 1, NMI values (according to equation 3) and we have mentioned the number of communities. As can be seen from Tables 1, methods find community structure with different number of communities. According to NMI values, our method can regroup the most nodes in the correct communities on Zachary's karate club, dolphin social network, American college football and books about US politics. On Zachary's karate club our method finds the same real community structure. On dolphin social network our method grouped more than 80% of nodes in the correct communities. Our method grouped 78% and 52% of nodes in correct communities on American college football and books about US politics respectively. The value of modularity by our method on these networks was above 0.3.

Fig. 5 and Fig. 6 show the community structure that was found by our method on Dolphins network and Books about US politics network. Each label represents a node and edges stand for the relationship between nodes. The community structure that was found by our method was represented by different shapes and colors. Nodes of the same community are represented by the same color and shape. From these Fig. 5 and Fig. 6, we can see that nodes in the same community are more connected between them and have a few connection with nodes from other communities.

We evaluated the performance of our method with other different real networks without a known community structure. A brief description of these networks is given below.

- Jazz network is a collaborative network [18], which represents the association between jazz musicians. The jazz musicians are represented by nodes and edge existing between nodes just if two musicians played together. The network has $n = 198$ nodes and $m = 2742$ edges.
- Word adjacencies network represents the adjacency network of common adjectives and nouns in the novel *David Copperfield* by *Charles Dickens* [19]. It has $n = 112$ nodes and $m = 425$ edges.
- Les Miserables network is co-appearance network of characters in the novel *Les Miserables* [20]. The network has $n = 77$ nodes and $m = 254$ edges.

Table II gives results of our method, fast greedy, label

TABLE I
PERFORMANCE RESULTS ON REAL NETWORKS WITH KNOWN COMMUNITY STRUCTURE.

Methods	Karate			Dolphins			Football			Books		
	$ c $	NMI	Q	$ c $	NMI	Q	$ c $	NMI	Q	$ c $	NMI	Q
<i>Fast greedy</i>	3	0.69	0.38	4	0.55	0.49	6	0.70	0.54	4	0.53	0.50
<i>Label propagation</i>	4	0.70	0.41	3	0.76	0.48	11	0.85	0.58	3	0.50	0.47
<i>Infomap</i>	3	0.50	0.40	5	0.53	0.52	12	0.91	0.60	6	0.49	0.52
<i>Our method</i>	2	1	0.37	2	0.81	0.38	10	0.78	0.51	2	0.52	0.43

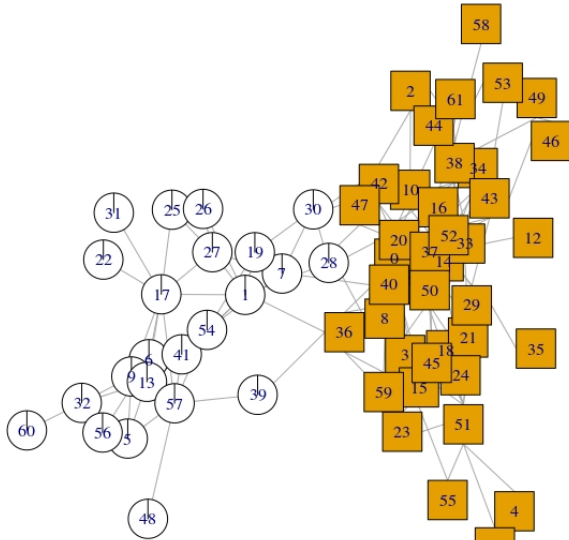


Fig. 5. Community structure of Dolphins network detected by our method and represented by different colors and shapes.

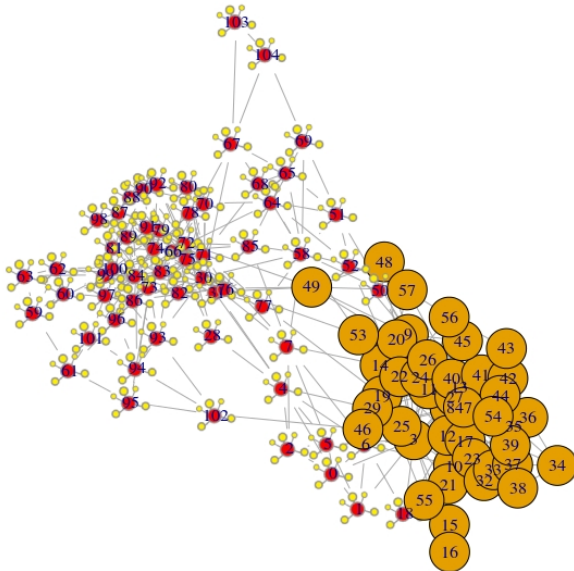


Fig. 6. Community structure of Books about US politics network detected by our method and represented by different colors and shapes.

propagation and Infomap. The number of communities and the estimated value of modularity were mentioned in Table II. From Table II, we can see that our method finds community structures with a high value of modularity. It is difficult to compare methods between them because we do not have a reference (a known community structure).

Fig. 7 shows the dendrogram, that was built by our method, and the community structure for Jazz network. Community structure that was found by our method has three communities. Labels of dendrogram represent nodes.

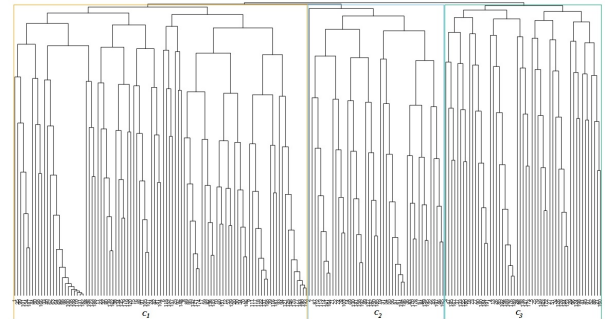


Fig. 7. The dendrogram of Jazz network and the community structure by our method.

V. CONCLUSION AND FUTURE WORK

A new hierarchical method to discover the community structure for unweighted and undirected networks was presented in this paper. Our new method was developed based on maximization of function of modularity by FPA. Results obtained on computer-generated networks and real benchmark networks prove the efficiency of our method in terms of finding community structures with high values of modularity and accuracy.

Our method can be tested on large scale networks. We can develop it to find community structure in weighted or directed network. It can be extended to find overlapping communities.

REFERENCES

- [1] MEJ. Newman, "Networks: an introduction," Oxford university press, 2010.
- [2] B. Hopkins, RJ. Wilson, "The truth about Königsberg," The College Mathematics Journal, 35(3), pp. 198-207, 2004.

TABLE II
PERFORMANCE RESULTS ON REAL NETWORKS WITH UNKNOWN COMMUNITY STRUCTURE.

Methods	Jazz		Word adjacencies		Miserables	
	$ c $	Q	$ c $	Q	$ c $	Q
<i>Fast greedy</i>	4	0.438	7	0.294	5	0.500
<i>Label propagation</i>	2	0.281	1	0	4	0.475
<i>Infomap</i>	7	0.280	2	0.009	9	0.546
<i>Our method</i>	3	0.346	6	0.264	7	0.505

- [3] S. Fortunato, "Community detection in graphs," Phys Rep. 486, pp. 75-174, 2010.
- [4] J. Chen, M. Liu, X. Liu, "Research on of overlapping community detection algorithm based on tag influence," Cluster Computing, 22(3), pp. 6669-6679, 2019.
- [5] M. Girvan, MEJ. Newman, "Community structure in social and biological networks," Proc Natl Acad Sci USA, 99, pp. 7821-7826, 2002.
- [6] A. Clauset, MEJ. Newman, C. Moore, "Finding community structure in very large networks," Phys Rev, 70, 66111, 2004.
- [7] U. Raghavan, R. Albert, S. Kumara, "Near linear time algorithm to detect community structures in large-scale networks," Phys. Rev. E 76, 036106, 2007.
- [8] G. Zhou, R. Wang, Y. Zhou, "Flower pollination algorithm with runway balance strategy for the aircraft landing scheduling problem," Cluster Computing, 21(3), pp. 1543-1560, 2018.
- [9] XS. Yang, "Flower pollination algorithm for global optimization," in: Unconventional Computation and Natural Computation, Lecture Notes in Computer Science, 7445, pp. 240-249, 2012.
- [10] E. Emary, HM.Zawbaa, M. Sharawi, "Impact of Lvy flight on modern meta-heuristic optimizers," Appl Soft Comput, 75, pp. 775-789, 2019.
- [11] WW. Zachary : "An information flow model for conflict and fission in small groups," J Anthropol Res. 33, pp. 452-473, 1977.
- [12] M. Rosvall, C. T. Bergstrom, "Maps of random walks on complex networks reveal community structure," Proceedings of the National Academy of Sciences of the United States of America, 105, pp. 1118-1123, 2008.
- [13] L. Danon, A. Diaz-guilera, J. Duch, A. Arenas, "Comparing community structure identification," J Stat Mech, P09008, 2005.
- [14] A. Lancichinetti, S. Fortunato, F. Radicchi, "Benchmark graphs for testing community detection algorithms," Phys. Rev. E 78, 046110, 2008.
- [15] J. Abonyi, B. Feil, "Cluster analysis for data mining and system identification," Springer Science and Business Media, 2007.
- [16] D. Lusseau, K. Schneider, O.J. Boisseau, P. Haase, E. Slooten, S.M. Dawson, "The bottlenose dolphin community of doubtful sound features a large proportion of long-lasting associations," Behav Ecol Sociobiol, 54, pp. 396-405, 2003.
- [17] V. Krebs, unpublished, <http://www.orgnet.com>, 2019.
- [18] P. M. Gleiser, L. Danon, "Community structure in jazz," Advances in Complex Systems, 6(4), pp. 565-573, 2003.
- [19] MEJ. Newman, "Finding community structure in networks using the eigenvectors of matrices," Phys. Rev. E 74, 036104, 2006.
- [20] D. E. Knuth, "The Stanford GraphBase: A Platform for Combinatorial Computing," Addison-Wesley, Reading, MA, 1993.