Community Mining in Signed Networks: A Multiobjective Approach

Alessia Amelio National Research Council of Italy (CNR) Inst. for High Perf. Computing and Networking (ICAR) Via Pietro Bucci, 41C 87036 Rende (CS), Italy Email: {amelio}@icar.cnr.it

Abstract—Community detection in signed complex networks is a challenging research problem aiming at finding groups of entities having positive connections within the same cluster and negative relationships between different clusters. Most of the proposed approaches have been developed for networks having only positive edges. In this paper we propose a multiobjective approach to detect communities in signed networks. The method partitions a network in groups of nodes such that two objectives are contemporarily optimized. The former is that the partitioning should have dense positive intra-connections and sparse negative interconnections, the latter is that it should have as few as possible negative intra-connections and positive inter-connections. We show that the concepts of signed modularity and frustration fulfill these objectives, and that the maximization of signed modularity and the minimization of frustration allow to obtain very good solutions to the problem. An extensive set of experiments on both real-life and synthetic signed networks shows the efficacy of the approach.

Keywords-signed networks; community detection, multiobjective genetic algorithms

I. INTRODUCTION

Th field of social network analysis provides a formal way to define important social concepts and study relations among actors. Social, and more generally, complex networks are able to represent many real world systems. Community structure is an important characteristics of networks that has been receiving a lot of attention in the last few years. Many different approaches have been proposed to divide a network into groups (also called clusters) having dense intra-connections, and sparse inter-connections. However, most of the proposed approaches do not take into account additional information that could be derived by the kind of ties connecting objects, but mainly focus on link density. Recently, there has been an increasing interest in developing methods to find communities in signed networks [12], [11], [6], [1], [2]. Since the primary studies of Heider [7] on attitude and perception of social organization of individuals, it has been pointed out that relationships between individuals can be either positive or negative, such as like-dislike, friends-enemies, love-hate, trust-distrust. Signed networks are defined as extension of networks that include the additional information of positive and negative links. Delving Clara Pizzuti National Research Council of Italy (CNR) Inst. for High Perf. Computing and Networking (ICAR) Via Pietro Bucci, 41C 87036 Rende (CS), Italy Email: {pizzuti}@icar.cnr.it

community structure on these types of networks can help in deeply understand social ties and organization among the actors constituting the network. One of the fundamental and intensively investigated aspects of signed networks is the balance theory. Davis [3] defined a network k-balanced if it can be divided into k groups such that, edges within groups are positive and edges between groups are negative. In such a case the network is also said *partitionable* or *clusterable*, while the term *balanced* is generally used for 2-way balance. k-balancing is an important research topic since balancing assures stability, while imbalance generates tension inside a group. Doreian and Mrvar [4] proposed a partitioning approach to structural balance based on the optimization of a criterion function, named *frustration*, that tries to minimize the number of positive ties among different groups and the number of negative links inside the same group.

In this paper we propose to exploit both concepts of frustration and modularity [9] to detect communities in signed networks by applying multiobjective optimization. We consider, in fact, a multiobjective framework based on Genetic Algorithms, named SN-MOGA (Signed Networks with MultiObjective Genetic Algorithm), where the first objective is modularity maximization, and the second objective is frustration minimization. The multiobjective genetic algorithm evolves a population of candidate solutions by trying to obtain the best trade-off between high modularity and low frustration. Experiments on synthetic and real life networks show that the multiobjective genetic approach is capable to divide signed networks with high accuracy and low edge misclassification with respect to the true known partitioning. The paper is organized as follows. In the next section some preliminary definitions are given and the problem is defined. Section III describes the algorithm. In Section IV an extensive experimentation on both real-life networks and synthetic networks is presented. Section V, finally, concludes the paper and gives some insights on future developments.

II. PRELIMINARIES

A signed social network can be modeled as a graph G = (V, E, W), where V is the set of n nodes (vertices) and



Figure 1. A partitionable network [1] of 9 nodes divided in three communities $\{0, 1, 2\}, \{3, 4, 5\}$ and $\{6, 7, 8\}$. Red dashed lines correspond to negative links, solid lines to positive edges. Colors correspond to the true partitioning. On the right, the locus-based representation of the genotype corresponding to these three communities.

E is the set of m edges. $W: V \times V \rightarrow \{-1, 0, 1\}$ is a function which assigns +1 to edges connecting positively a pair of nodes, -1 to edges that connect negatively a pair of nodes and 0 if an edge does not exist between the nodes. Let A denote the weighted adjacency matrix associated with G, i.e. $A_{i,j} = W(i,j)$. The matrix A can be split into two adjacency matrices corresponding to positive and negative edges by setting $A_{i,j}^+ = A_{i,j}$ if $A_{i,j} > 0$, zero otherwise, and $A^-_{i,j} = -A_{i,j}$ if $A^-_{i,j} < 0$, zero otherwise, thus $A = A^+ - A^+_{i,j}$ A^{-} . Given a node $i \in V$, a_i^+ and a_i^- are defined respectively as the positive degree and the negative degree of i. The Frustration F(C) of a network partition $C = \{C_1, \ldots, C_k\}$ of the graph G into k communities, is defined as the sum of the number of negative edges between nodes inside the same community and the number of positive edges between nodes belonging to different communities:

$$F(C) = \sum_{i,j \in V} A_{i,j}^{-} \delta(c_i, c_j) + A_{i,j}^{+} (1 - \delta(c_i, c_j))$$
(1)

where $c_i(c_j)$ is the community of the node i(j) and $\delta(c_i, c_j)$ is the Kronecker delta function which takes the value 1 if nodes i and j belong to the same community, 0 otherwise.

The concept of *modularity* has been introduced by Newman and Girvan in [9]. For signed networks the definition of modularity is modified to take into account the contribution of positive edges inside communities and negative edges between communities. The signed modularity can be defined as [6]:

$$Q_{S} = \frac{1}{2m} \sum_{i,j \in V} (A_{i,j} + \frac{a_{i}^{-}a_{j}^{-}}{2m} - \frac{a_{i}^{+}a_{j}^{+}}{2m})\delta(c_{i}, c_{j})$$
(2)

Our objective is to solve the following problem. Given a graph G = (V, E, W) modeling a signed network, find a partitioning of G in k clusters such that: 1) intra-connections are dense and most edges within clusters are positive; 2) inter-connections between clusters are sparse and most of these edges are negative.

III. Algorithm Description

In this section we give a description of the multiobjective algorithm *SN-MOGA* for signed networks, the representation

adopted for partitioning the network, and the variation operators used. The MultiObjective Genetic Algorithm (MOGA) we used is the Nondominated Sorting Genetic Algorithm (NSGA-II) proposed by Srinivas and Deb in [10] and implemented in the Genetic Algorithm and Direct Search Toolbox of MATLAB. NSGA-II builds a population of competing individuals and ranks them on the basis of nondominance. In order to employ NSGA-II, SN-MOGA has been adapted with a customized population type that suitably represents a partitioning of a network and endowed with the two complementary objectives of frustration (formula (1)) and signed modularity (formula (2)). The algorithm uses the locus-based adjacency representation employed in [5] for community discovery in dynamic unsigned networks. In this representation an individual of the population consists of ngenes g_1, \ldots, g_n and each gene can assume a value in the range $\{1, \ldots, n\}$. A value j assigned to the *i*th gene means that there is an edge (i, j) in E. A main characteristic of this representation is that the number k of clusters is automatically determined by the number of components contained in an individual. Figure 1 shows a signed network (originally reported in [1]) of 9 nodes clusterable in the three groups $\{0, 1, 2\}, \{3, 4, 5\}$ and $\{6, 7, 8\}$. Dashed lines correspond to negative links, while solid lines to positive edges. The genotype corresponding to this division is showed on the right part of the figure and it is interpreted as: node 0 is connected with node 1, node 1 with node 0, node 2 with node 1, and so on. SN-MOGA initializes a population of random individuals by assigning to each node *i* one of its neighbors. Mutation operator, analogously to initialization, randomly selects one of the neighbors of i and assigns this value to the i-th gene. The kind of crossover adopted is uniform crossover. Multiobjective optimization techniques do not return a unique solution to a problem, but a set of solutions are found through the use of *Pareto optimality theory*. In this context, since a vector of competing objectives must be simultaneously optimized, the goal is to obtain Paretooptimal solutions, i.e. nondominated solutions for which an improvement in one objective requires a degradation of another (Pareto front). Thus the Pareto front represents the best compromise solutions satisfying all the objectives as best as possible. However, a single solution, out of the Pareto front, must be selected. In our case, in order to show the differences in selecting a different solution from the Pareto front, in the experiments we show the results obtained by choosing either minimum frustration or maximum modularity.

IV. EXPERIMENTAL RESULTS

In this section we evaluate the capability of our approach in obtaining meaningful partitions of signed networks. As regards parameters needed by the genetic approach, we set crossover rate to 0.8, mutation rate to 0.2, elite reproduction 10% of the population size, roulette selection function,

Table I

Error obtained by *SN-MOGA* on different networks when the solution having minimum frustration is chosen from the Pareto front, with the corresponding modularity and NMI values, and the maximum modularity (max Mod.) obtained from the Pareto front with the corresponding error and NMI values. In parenthesis the standard deviation is reported. For each network the number of nodes and the number of positive (E^+) and negative (E^-) edges are also reported. For Wikipedia the error has been computed as in [2].

Name	nodes	E+	Е-	Error	Modularity	NMI	max Mod.	Error	NMI
Network1	9	9	6	0 (0)	0.5333 (0)	1 (0)	0.5333 (0)	0 (0)	1 (0)
Network2	28	30	12	0 (0)	0.5612 (0)	1 (0)	0.5612 (0)	0 (0)	1 (0)
Network3	28	30	19	0 (0)	0.5257 (0)	1 (0)	0.5257 (0.0058)	0 (0)	1 (0)
Gahuku-Gama Subtribes	16	29	29	0.0345 (0)	0.4483 (0)	0.7528 (0)	0.4483 (0)	0.0345 (0)	0.7528 (0)
Karate	34	68	10	0 (0)	0.4997 (0)	1 (0)	0.5127 (0.0046)	0.0462 (0.0162)	0.8430 (0.0552)
Football	115	394	219	0.0571 (0.0157)	0.5516 (0.0180)	0.8744 (0.0302)	0.5520 (0.0177)	0.0573 (0.0159)	0.8762 (0.0294)
Dolphins	62	153	6	0 (0)	0.4112 (0)	1 (0)	0.5412 (0.0064)	0.1673 (0.0168)	0.6358 (0.0153)
Krebs	105	371	69	0.0677 (0.0136)	0.4456 (0.0126)	0.6955 (0.0483)	0.4469 (0.0133)	0.0711 (0.0155)	0.6946 (0.0377)
Wikipedia	7118	83953	23118	0.001007158 (0.000063)	0.0105 (0.004132)	-	0.07738 (0.011084)	0.001823565 (0.000081)	-

population size was 100, number of generations 200. The algorithm has been executed 10 times and the average values of error rate and NMI have been computed together with standard deviation. It is worth to note that there is still no a general measure to validate and compare methods for signed networks. Yang et al. [12] employed the frustration concept to define the error rate of a signed network partitioning C as

$$error(C) = \frac{F(C)}{\sum_{i} \sum_{j} |A_{i,j}|} \times 100\%$$
(3)

However, as pointed out by the authors, this error function considers only the sign of the links, and completely disregards the edge density, thus we also used *Normalized Mutual Information (NMI)*, a well known entropy measure in information theory. The normalized mutual information NMI(A, B) of two divisions A and B of a network is defined as follows. Let C be the confusion matrix whose element C_{ij} is the number of nodes of community *i* of the partition A that are also in the community *j* of the partition B.

$$NMI(A,B) = \frac{-2\sum_{i=1}^{c_A}\sum_{j=1}^{c_B}C_{ij}log(C_{ij}N/C_{i.}C_{.j})}{\sum_{i=1}^{c_A}C_{i.}log(C_{i.}/N) + \sum_{j=1}^{c_B}C_{.j}log(C_{.j}/N)}$$
(4)

where $c_A(c_B)$ is the number of groups in the partition A(B), $C_i(C_j)$ is the sum of the elements of C in row i (column j), and N is the number of nodes. If A = B, NMI(A, B) = 1. If A and B are completely different, NMI(A, B) = 0.

A. Evaluation on clusterable networks

We first consider the toy example reported in [1] and shown in Figure 1(a), and the two artificial signed networks considered by Yang et al. [12]. Network *Network2* in Figure 2(a) is partitionable and can be divided into the three groups $\{4, 5, 6, 7, 22, 23, 24, 25, 13, 14, 15, 16\}$, $\{8, 9, 26, 27, 17, 18\}$, and $\{20, 21, 10, 11, 12, 1, 2, 3, 19, 28\}$. *Network3* (Figure 2(b)) is also partitionable in the same three groups of *Network2*. The main difference between these two networks is that *Network2* is also balanced, since it has a two-way partitioning constituted by the first group and the union of the other two groups, while *Network3* is not balanced. Table I reports, for each network, the



Figure 2. Synthetic networks reported in [12]. Network2 (a) is partitionable and can be divided into the three groups $\{4, 5, 6, 7, 22, 23, 24, 25, 13, 14, 15, 16\}$, $\{8, 9, 26, 27, 17, 18\}$, and $\{20, 21, 10, 11, 12, 1, 2, 3, 19, 28\}$. Network3 (b) is also partitionable in the same three groups of network (a). For each network, colors correspond to the true partitioning, red dotted lines to the negative edges and black solid lines to the positive edges.

number of nodes, the number of positive and negative edges, the error rate obtained when the solution having minimum frustration is chosen from the Pareto front, with the corresponding modularity and NMI values, and the maximum modularity obtained from the Pareto front with the corresponding error rate and NMI values. For these three networks *SN-MOGA* finds a unique solution having both zero frustration and maximum modularity. Anchuri and Magdon-Ismail' approach [1], because of the used parameter setting, needs an improvement step to correctly assign node 8 in Network1. Note that, for *Network2 SN-MOGA* could not find a 2-way partition, since there are no connections between the second and the third groups.

B. Evaluation on real-life networks

Next we consider five real-life networks well known in the literature. The *Gahuku-Gama Subtribes* social network has been studied by Yang et al. [12]. The other four networks are very popular networks used to compare community detection methods: *Zackary's Karate Club network, The American College Football network, Bottlenose Dolphins,* and *Krebs' books on American politics.* Since these networks are unsigned, we transformed them in signed networks by assigning a positive sign to edges between nodes in the same community, and negative sign to edges between communities. From Table I we can observe that *SN-MOGA* finds a unique solution for the *Gahuku-Gama Subtribes*. The error obtained by *SN-MOGA* for this network is the

same of that reported by Yang et al. [12]. As regards Karate and Dolphins networks, when the Pareto front solution we select is that having minimum frustration, the solutions found actually correspond to the ground truth division in two groups of these networks. However, when we choose the solution having maximum modularity the error increases to 0.0462 while NMI diminishes to 0.8430 for Karate, and the error increases to 0.1673, while NMI diminishes to 0.6358 for Dolphins. These solutions are also significant since, as regards Karate, the subgroup constituted by nodes $\{5, 6, 7, 11, 17\}$ is separated from one of the two ground truth clusters, while for Dolphins SN-MOGA divides the bigger ground truth community in three smaller communities. The values obtained on the other two networks, Football and Krebs are also very good and show the capability of SN-MOGA in finding meaningful partitioning of signed networks.



Figure 3. NMI corresponding to the maximum modularity values obtained from *SN-MOGA* for all the possible p+ and p- values at different values of the μ parameter.



Figure 4. NMI corresponding to the minimum frustration values obtained from our algorithm for all the possible p+ and p- at different values of the μ parameter.

C. Evaluation on synthetic networks

In this section a more deep study on synthetic networks generated with control parameters that determine the structure of communities, is performed. In particular, we modified the benchmark proposed by Lancichinetti et al. [8], which is an extension of the classical benchmark of Girvan and Newman, by assigning a controlled sign to edges. The networks consist of 128 nodes divided into four communities of 32 nodes each. Every node has an average degree of 16 and shares a fraction μ of edges with the other nodes of the network, and $1 - \mu$ of links with the nodes of its community. μ is called the *mixing parameter*. When $\mu < 0.5$ the neighbors of a node inside its group are more than the neighbors belonging to the other three groups, thus a good algorithm should discover them. We generated 10 different networks for values of μ ranging from 0.1 to 0.5. In order to make the networks signed, analogously to Yang et al. [12], we used two parameters p_{-} , denoting the probability of negative links appearing within communities, and p_+ , denoting the probability of positive links appearing between communities. Thus, for all the combinations of p_{-} and p_{+} values ranging in the interval $[0, 0.1, \ldots, 1]$, we randomly assigned a negative sign to edges inside a community with probability p_{-} , and a positive sign to edges between two different communities with probability p_+ . Figure 3 depicts the NMI values obtained by running SN-MOGA for all the combinations of parameters $\mu = [0.1, \dots, 0.5], p_{-} = [0, 0.1, \dots, 1],$ and $p_{+} = [0, 0.1, \dots, 1]$ and selecting from the Pareto front the community structure having the highest modularity value. The figure points out that, as the network structure becomes more noisy, i.e. μ increases, the corresponding NMI value decreases, as expected. However, fixed a μ value, the method obtains slightly decreasing values of normalized mutual information until $p_+ \leq 0.6$ and $p_- \leq 0.3$, thus SN-MOGA is less sensitive to the number of positive edges between communities, but it is negatively biased by the augmentation of negative links within a community. A different behavior can be observed in Figure 4, where the solutions having the minimum frustration are now selected from the Pareto front. In this case the NMI values obtained are lower with respect to the previous case, however SN-MOGA is insensitive to the variation of both positive and negative edges for $p_{-} \leq 1$ and $p_+ \leq 0.4$. This behavior is very interesting because it means that even if the structure is highly unbalanced, the method is able to unveil the underlying community structure by searching for dense, but with low frustration, groups of nodes. Finally Figure 5 shows the error rate obtained by SN-MOGA, when minimum frustration solutions are chosen, for mixing parameter $\mu = 0.2$, and the combination values of p_{-} and p_{+} from 0 to 0.5. The figure points out that, fixed a p_{-} value, the error rate shows a very slight increase for increasing values of p_+ , i. e. the augmentation of positive links between different communities does not provoke abrupt changes in the frustration value.

D. Comparison with other approaches

Comparing *SN-MOGA* with the other state-of-the-art methods is rather difficult since there do not still exist neither synthetic benchmarks to use like those defined for unsigned



Figure 5. Error rate values obtained from *SN-MOGA* for p+ and p- varying in the interval [0, 0.5] and $\mu = 0.2$.

networks, nor a standard and recognized measure to evaluate the results. In the existing approaches, frustration is used to compute the percentage of misclassified edges. However, while Yang et al. [12] define the error by dividing frustration by the total number of edges, the error defined by Chiang et al. [2] divides by the square of the number of nodes, while that employed by Anchuri and Magdon-Ismail [1] divides by the number of negative edges. Since we did not have at disposal the synthetic networks employed by Yang et al. [12], we used the generator of Lancichinetti et al. [8] with the same parameters adopted by Yang et al. to generate similar, though not equal, networks and compute the fraction of nodes correctly clustered, i.e. 1 - error (see formula (3)). For these networks, the accuracy values are above 80% for $p^{+} \leq 1$ and $p_{-} \leq 0.25$, while when $0.25 < p_{-} \leq 1$, the accuracy is never below 60%, independently the p_+ value. Yang et al. obtained a percentage of correctly clustered edges near 100% for $p_{-} \leq 0.35$ and $p_{+} \leq 1$. However the clustering accuracy of their method is not less than 50% for 0.4 $< p_{-} \leq$ 0.6. Clustering accuracy of SN-MOGA, instead, as outlined above, is never less than 60%. We emphasize that this comparison is between two kind of networks with similar characteristics, but not exactly the same. Chiang et al. [2] in their paper computed the error as formula (3), where denominator is substituted by n^2 , i.e. the square of the number of nodes. They applied their algorithm, among the others, to English Wikipedia network for admin elections, downloadable from http://konect.unikoblenz.de/networks/elec. The empirical error rate they reported is 0.2186, for values of number of communities kranging from 3 to 30. The authors observe that, for each k, the errors are very close. We executed SN-MOGA on this network and, as reported in Table I, obtained an error rate of 0.001007158, which is much lower than that obtained by Chiang et al. The number of clusters found by SN-MOGA has been, on average, about 100. This means that the range of values used by Chiang et al. was insufficient to obtain a reasonable partitioning of the Wikipedia network. This result confirms the advantage of applying SN-MOGA, which is capable of finding meaningful k-way divisions with small frustration values, without any knowledge on the network structure. Modularity values obtained are however very low, the maximum value being 0.07738.

V. CONCLUSIONS

The paper proposed a multiobjective approach to detect communities in signed networks. The method obtains network partitioning by minimizing the number of negative edges inside communities and positive edges between communities, while maximizing cluster modularity. The optimization of these two objectives allows to find network divisions such that intra-connections are dense and most edges within clusters are positive, and inter-connections between clusters are sparse and most of these edges are negative. An experimental evaluation on both real-life and randomly generated networks for which the true partitioning is known proved the ability of the method to find solutions having low frustration and high NMI values. Future work aims at extending the method to dynamic and signed networks.

REFERENCES

- P. Anchuri and M. Magdon-Ismail. Communities and balance in signed networks: A spectral approach. In ASONAM'12, pages 235–242, 2012.
- [2] K. Chiang, J. Jiyoung Whang, and I. S. Dhillon. Scalable clustering of signed networks using balance normalized cut. In *CIKM'12*, pages 615–624, 2012.
- [3] J.A. Davis. Clustering and structural balance in graphs. *Human Relations*, 20:181–187, 1967.
- [4] P. Doreian and A. Mrvar. A partitioning approach to structural balance. *Social Networks*, 18:149–168, 1996.
- [5] F. Folino and C. Pizzuti. A multi-objective genetic algorithm for community detection in networks. In ASONAM'10, pages 256–263, 2010.
- [6] S. Gómez, P. Jensen, and A. Arenas. Analysis of community structure in networks of correlated data. *Physical Review*, E80:016114, 2009.
- [7] F. Heider. Attitudes and cognitive organization. J. Psycology, 21:107–112, 1946.
- [8] A. Lancichinetti, S. Fortunato, and F. Radicchi. Benchmark graphs for testing community detection algorithms. *Physical Review E*, 78(046110), 2008.
- [9] M. E. J. Newman and M. Girvan. Finding and evaluating community structure in networks. *Physical Review*, E69:026113, 2004.
- [10] N. Srinivas and K. Deb. Multiobjective optimization using nondominated sorting in genetic algorithms. *Evolutionary Computation*, 2(3):221–248, 1994.
- [11] V.A. Traag and Jeroen Bruggeman. Community detection in networks with positive and negative links. *Physical Review E*, 80(3):036115, 2009.
- [12] B. Yang, W. K. Cheung, and J. Liu. Community mining from signed social networks. *IEEE Transactions on Knowledge and Data Engineering*, 19(10):1333–1348, 2007.