

Mesoscopic analysis of online social networks: The role of negative tiesPouya Esmailian, Seyed Ebrahim Abtahi,^{*} and Mahdi Jalili[†]*Department of Computer Engineering, Sharif University of Technology, Tehran, Iran*

(Received 24 May 2014; published 29 October 2014)

A class of networks are those with both positive and negative links. In this manuscript, we studied the interplay between positive and negative ties on mesoscopic level of these networks, i.e., their community structure. A *community* is considered as a tightly interconnected group of actors; therefore, it does not borrow any assumption from balance theory and merely uses the well-known assumption in the community detection literature. We found that if one detects the communities based on only positive relations (by ignoring the negative ones), the majority of negative relations are already placed between the communities. In other words, negative ties do not have a major role in community formation of signed networks. Moreover, regarding the internal negative ties, we proved that most unbalanced communities are maximally balanced, and hence they cannot be partitioned into k nonempty sub-clusters with higher balancedness ($k \geq 2$). Furthermore, we showed that although the mediator triad $++-$ (hostile-mediator-hostile) is underrepresented, it constitutes a considerable portion of triadic relations among communities. Hence, mediator triads should not be ignored by community detection and clustering algorithms. As a result, if one uses a clustering algorithm that operates merely based on social balance, mesoscopic structure of signed networks significantly remains hidden.

DOI: [10.1103/PhysRevE.90.042817](https://doi.org/10.1103/PhysRevE.90.042817)

PACS number(s): 89.75.Fb, 89.20.Ff, 89.75.Hc

I. INTRODUCTION

In the past two decades, there have been increasing interests toward the analysis of complex networks both empirically and theoretically [1–3]. One of the important research lines is to study networks from the structural point of view, trying to answer, *What do different types of networks look like?* This is an important issue, since it has been shown that many dynamical properties depend on the network structure [4–6]. This endeavor is constantly coevolving with the studies on theoretical models of networks trying to describe the observations and further predict new features [7,8]. Most of these works have been carried out due to abundant large-scale datasets gathered over the internet. They have attracted a lot of studies mainly to justify the long standing debates on static and/or dynamic patterns of relations [9–13].

There are a number of challenges related to signed networks. Discovering the community structure is one of these problems that has been addressed in a number of research works [14–16]. Another problem related to these networks is to predict the sign of relations [17–19].

Generally speaking, there have been two trends toward the analysis of signed networks. The first trend tries to evaluate the long-standing social balance theory and to deduce some new implications [10,20]. The social balance theory has some predictions about the grouping of people based on the analysis of network evolution toward a more balanced structure [21]. The second trend, regardless of the balance theory, tries to improve the inference tasks using the negative relations [14,17]. For example, detecting the community of densely interacting individuals is one of the issues studied in such works [16]. The notion of *community* has been introduced as a meaningful building block of networks [22]. Indeed, community structure acts as a bridge between local and global understanding of network structure [23,24]. In

signed networks, grouping the actors has been studied in both community detection and social balance literature [15,25]. In the former, the main objective is expressed as “dense positive” and “negative free” relations inside groups. In the latter, the objective is explicitly stated as minimizing the number of negative (positive) links inside (between) the groups. These two notions, despite their similarities, have fundamental differences, which are investigated in this work. The main motivation of our work is based on the recent work of Doreian and Mrvar [25]. They suggested that the $++-$ relation among groups of individuals is likely to be seen, and thus, it should not be ignored while detecting the mesoscale structure of networks.

As a connection to the above trends, our work starts with the justification of community detection in signed networks and shows that negative relations are not informative enough to improve the detection task. In other words, one can accomplish the task by considering only the positive relations. Our study also deals with the justification of the balance theory in mesoscopic level. Analogous to the local level, this theory states that no matter how (internally balanced) communities are identified, one must not see (or at least rarely see) the $++-$ triadic relation among them. We found that the observed triads are also underrepresented in mesoscopic level consistent with this theory. However, they form a considerable portion of social relations, which is far more than the corresponding local level, and cannot be simply ignored by clustering algorithms. Therefore, if the social groups are identified based on balance theory, one would miss a considerable amount of distinguishable groups by merging them into one another. Our results shed new light on mesoscale structure of signed networks.

II. PRELIMINARIES**A. Notations**

Throughout the paper, the expressions “link,” “edge,” “tie,” “relation,” and “interaction” are used interchangeably, unless

^{*}abtahi@sharif.edu[†]mjalili@sharif.edu

we explicitly make a note. A signed graph G is determined using triple (V, E, σ) . V is the set of nodes, E is the set of edges, which is defined by pair (v_i, v_j) of nodes $[(v_i, v_j) = (v_j, v_i)$ for undirected graph], and σ assigns either $+1$ or -1 to each edge. In this work, we consider only undirected signed graphs with values -1 and $+1$ for negative and positive relations, respectively. Having k *nonempty* clusters in a network, let us define the number of inconsistent or frustrated edges as follows:

$$F_k(G, C) = \sum_{C_i=C_j, i < j} A_{ij}^- + \sum_{C_i \neq C_j, i < j} A_{ij}^+, \quad (1)$$

where G is a signed graph, C determines the cluster of nodes ($C_i =$ cluster to which node i belongs), k is the number of nonempty clusters, and $A_{ij}^+ = 1$ if $\sigma_{ij} = 1$, or $A_{ij}^- = 1$ if $\sigma_{ij} = -1$, or both are zero otherwise. We denote the minimum value of the above function under all possible clusterings as

$$F_k(G) = \min_C F_k(G, C), \quad (2)$$

where the number of clusters k is a constant value. When k is tunable, one has:

$$F(G) = \min_{C, k} F_k(G, C). \quad (3)$$

In the literature, Eq. (2) is often considered as *frustration index* [26], *true frustration*, or merely *frustration* [20,27]. However, in this context, the frustration and its minimum are considered separately. For $k = 1$, frustration of a subgraph is equal to the number of negative edges, and thus $F_1(G, C)$ [or equally $F_1(G)$] is used to denote the number of negative edges inside a subgraph. We use $f_k(G, C)$ as the ratio of $F_k(G, C)$ to the edge count $m = |E|$ [similar for $f_k(G)$]. Notations $F_{k, \text{up}}(G)$ and $F_{k, \text{low}}(G)$ are used for the upper bound of $F_k(G)$ and its lower bound, respectively ($F_{k, \text{low}}(G) \leq F_k(G) \leq F_{k, \text{up}}(G)$).

Given a specific clustering C , we define balancedness of graph G as follows:

$$B_k(G) = 1 - f_k(G). \quad (4)$$

Generally, we use the term *balanced* when a given subgraph S (i.e., an extracted community) has no negative edges [$B_1(S) = 1$], and *unbalanced* when $B_1(S) < 1$. Note that a graph may have higher balancedness for $k > 1$, which is denoted explicitly throughout the paper.

B. Correlation clustering problem

In this problem, one seeks to find a clustering of nodes that minimizes inconsistent relations. This is equivalent to minimization of $F_k(G, C)$ considering k either as a constant value [28] or a tunable parameter [29]. We should mention that the maximization of consistent edges has also been considered in the above works, which has different implications from the algorithmic point of view.

III. RELATED WORKS

In this section, we introduce some of the research lines related to the clustering of signed networks. Conceptually, they could be divided into two categories, where (1) positive links between clusters are penalized, or (2) instead of this punishment, internal density of clusters is rewarded.

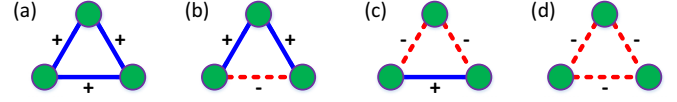


FIG. 1. (Color online) Different types of triadic signed relations between three actors. In structural balance, triads A and C are balanced, and B and D are unbalanced. In general structural balance, only triad B is unbalanced and the others are balanced.

A. Structural balance and clustering

The origin of structural balance theory is the seminal work of Heider [30], which has been further developed as a mathematical framework by Cartwright and Harary [31]. In the local level, the structural balance theory states that a triadic relation is balanced, if and only if, it has one or three positive ties.¹ As shown in Fig. 1, triads A and C are balanced, and B and D are unbalanced. In the global level, the structural balance theory states that a graph is structurally balanced (SB), if and only if it can be partitioned into two clusters with no inconsistent edges (known as *structure theorem*), or equivalently, when every cycle is positive. Inconsistent edges are negative ones inside and positive ones between the clusters. A cycle is positive (or balanced), if and only if it has an even number of negative links.

Davis [32] argued that a social network may have multiple hostile groups, implying that triad D is also balanced. In the global level, a graph is k -balanced, if and only if it can be partitioned into k -clusters with no inconsistent edge. The term *structural balance* is used for $k = 2$ and *weak- or general-structural balance* (GSB) for $k \geq 2$.

To measure the balancedness of networks, a number of research works have provided some metrics that specify the distance of a graph from GSB [33]. In this context, there are two well-known classes of metrics. The first class is based on counting all unbalanced l -cycles (cycles of length l), which can only be used for SB. The second class is based on counting the minimum number of inconsistent edges under all possible k -clusterings [$=F_k(G)$]. In this work, we base our investigations on the second class, and thus, it is briefly discussed in the following. This metric is equal to the minimum number of edges that their deletion (or sign flipping) results in a k -balanced graph, which is equivalent to distance of a graph from being k -balanced.²

The problem of finding a partition that corresponds to $F_k(G)$ is NP-hard [29], even for $k = 2$ [20]. If we set $k = 2$, the optimal solution is the best two-clustering of a graph where the number of inconsistent edges is equal to the distance of a graph from SB [$=F_2(G)$]. Iacono *et al.* proposed a graph-theoretic approach to approximate $F_2(G)$, which has been originally stated as “distance from monotonicity” for biological networks [27]; note that monotonicity has the same mathematical implication as SB. The algorithm has been further applied to social networks, validating that their distance

¹In the structural balance theory *balanced* and *unbalanced* are used only for $k = 2$.

²This equivalence holds for $k > 2$ with the same proof provided by Zaslavsky [26].

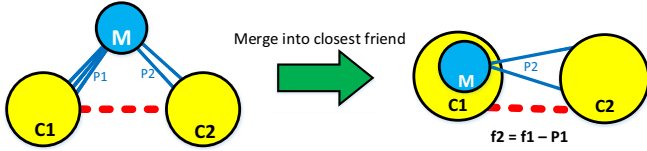


FIG. 2. (Color online) “+ + -” relation (hostile-mediator-hostile) between three clusters. The value of $F(G)$ is reduced by $P1$, if mediator cluster M is merged into the closest friend $C1$.

from SB is significantly lower than those of sign-shuffled counterparts [20,34]. Another achievement of Ref. [27] is a scalable algorithm that calculates a lower bound for $F_2(G)$, which determines, at most, how far is the proposed solution from the optimal value. For $k > 2$, Chiang *et al.* [35] proposed a scalable k -clustering algorithm by transforming an objective function similar to $F_k(G, C)$ (along with some other objectives) into weighted kernel k means. In this paper, we only use the two-clustering algorithm of Iacono *et al.*, together with a theorem that extends our results to $k > 2$.

B. Relaxed structural balance and generalized block modeling

In contrast to the implications of GSB, Doreian and Mrvar [25] argued that real-world networks are not completely balanced. Accordingly, it has been shown that in online social networks with 17%–23% negative ties, at least 7%–14% of edges are inconsistent with SB [20]. As a result, Doreian and Mrvar proposed the relaxed structural balance (RSB) theory stating that positive interactions between two clusters are also valid. This relaxation is mainly due to intermediary processes in social networks, implying that it is likely to find a mediator group with positive relations toward two hostile groups [Fig. 2 (left)].

Based on GSB, positive edges between clusters are punished. Hence, as depicted in Fig. 2 (right), a mediator cluster is merged into one of the hostile clusters with which it has more positive connection $P1$, decreasing the frustration from $F(G)$ to $F(G) - P1$. Accordingly, Doreian and Mrvar argued that, based on GSB, blocks (cluster-cluster relations) of positive ties are not allowed in off-diagonal positions of the relation matrix (as shown in Fig. 3). In Fig. 3, one can see the result of fitting a generalized block model (GB-model) [36] on hostile-mediator-hostile triad (mediator triad for short) based on GSB and RSB. However, in order to fit the relaxed model to

data, $F_k(G, C)$ is still used in Ref. [25] as the objective function. This means one should take care of each off-diagonal positive block *a priori* to refrain the optimization process, which tries to minimize $F_k(G, C)$, from merging mediator clusters into hostile parties.

C. Community detection in signed networks

In the community detection literature, mainly started after the seminal work of Girvan and Newman [22], there has been a different perspective toward the group identification. As the main assumption, a community is a group of nodes that have more connections inside than to the rest of the network. This intuition has been the basis of almost all community detection algorithms [37,38]. Regarding this, *modularity* function has been introduced that gives a better score to a cluster with denser relations than a null model [39]. The formulation of modularity allows for straightforward extension to signed networks [15]. The intuition is that the group of nodes should have more (less) positive (negative) intradensity relative to the null model. This intuition could be formulated by subtracting the modularity score of negative subgraph G^- from positive subgraph G^+ as follow:

$$Q(G, C) = \alpha Q(G^+, C) - (1 - \alpha)Q(G^-, C), \quad (5)$$

where $0 \leq \alpha \leq 1$ is the relative importance of positive ties compared to negative ones. A similar work has been carried out for Hamiltonian function of Potts model, which borrows the idea of modularity by incorporating an arbitrary null model with a resolution parameter [16]. As a summary, these methods reward (punish) the density of intra-positive (intra-negative) relations and punish (reward) their sparseness relative to the null model. Another work extends the community detection based on random walks [14], with the intuition that a random walker is more likely to be trapped inside a community. In the main step of the algorithm, the nodes are sorted according to their distance from a sink node. This step ignores the information of negative ties, which are only used as a cut criterion on the sorted list.

In all of these extensions, there is no explicit punishment strategy for positive edges between the communities, which makes them applicable to nonsigned (or sparsely signed) networks. As a connection to generalized block modeling, these algorithms work with dense diagonal blocks (relations inside a cluster) and sparse off-diagonal blocks (relations between clusters) for positive relations and the reverse for the negative ones. Figure 3 shows a toy example where signed community detection and RSB produce the same clustering that is different from the one produced by GSB.

IV. COMMUNITY VERSUS CLUSTER

In this section, we try to pinpoint some implications of the algorithms that try to optimize $F_k(G, C)$ against those that are frequently used for community detection. The main differences are illustrated in Fig. 4. In all cases, detected clusters result in $F_2(G) = 0$ as the optimal solution. In Fig. 4(a), the output of clustering algorithms is consistent with the notion of *community*, which is also produced by relaxed GB-modeling. However, in Fig. 4(b), members of each

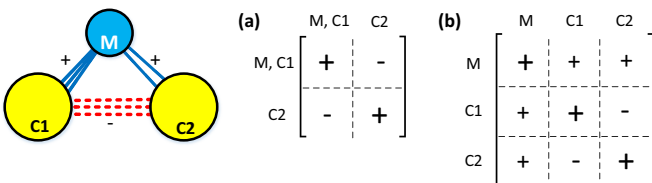


FIG. 3. (Color online) (a) Result of fitting a GB model based on general structural balance, which merges $C1$ and M to get a lower frustration. (b) Result of fitting a GB model based on relaxed structural balance, which allows off-diagonal positive blocks similar to the output of signed community detection methods, if one restricts the clusters to be internally cohesive.

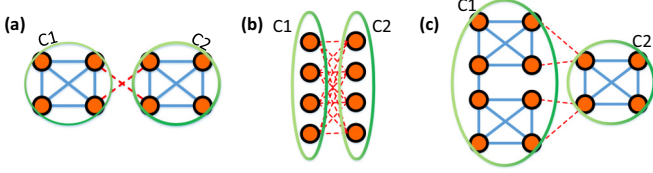


FIG. 4. (Color online) All three two-clusterings result in $F_2(G, C) = 0$. (a) *Cluster* and *community* are consistent with each other. (b) Clusters are not *communities*. (c) Two distinguishable communities get clustered into C_1 .

cluster are disconnected, and thus, despite their similar role in the network, they cannot be considered as a social group of interacting individuals. Also, in Fig. 4(c) two distinguishable communities are grouped together missing an obvious pattern of relations. These cases [Figs. 4(a)–4(c)], as well as the mediator triads, are the main shortcomings of clustering algorithms in social networks. Trying to relate these two notions, a community is a cluster of nodes that is internally well-connected. One of the aims of this work is to investigate the frequency of such cases in real networks. We show that the case as shown in Fig. 4(c), as well as mediator triads, are frequent enough that cannot be ignored when one deals with large-scale social networks.

V. METHODS

A. Extracting communities

We want to extract groups of densely interconnected nodes that are consistent with the notion of *community*. To this end, we use InfoMap [38,40], which is proved to be one of the most accurate community detectors [41]. It confidently extracts the communities from large-scale networks that have heterogeneous group sizes [6,42]. We used the open source code provided in Ref. [43] utilizing the hierarchical mode that refines a few big communities into smaller ones, and leaves other communities intact. As studied by Lancichinetti and Fortunato [42], if a group of nodes is well-separated from the environment, it could be accurately detected by InfoMap. However, if the density among some groups passes a threshold, InfoMap mistakenly considers them as a single community. Indeed, this problem happens for all methods that merely consider the structure of a network. In the case of InfoMap, we are confident about the internal density of detected communities relative to the environment [44], and as the only problem, there might be more than one group in a single community. Nevertheless, as we illustrate in the results, this problem does not significantly affect the outcome, and the conclusion drawn from the results remains valid.

B. Computing the distance from structural balance

As we mentioned in Sec. III A, for graph G , the distance from SB is $F_2(G)$, which is equal to the minimum inconsistent edges under all two-clusterings. Although the computation of this value is NP-hard, the scalable algorithm of Iacono *et al.* [27] outputs a two-clustering, which is an upper bound for $F_2(G)$, as well as a lower bound for $F_2(G)$. Thus, we always know, at most, how far is the suboptimal solution from an

optimal one. Same as Ref. [27], quantity $F_{2,\text{low}}(G)/F_{2,\text{up}}(G)$ is used to measure the precision of a solution. Considering the following inequality:

$$\frac{F_{2,\text{low}}(G)}{F_{2,\text{up}}(G)} \leq \frac{F_2(G)}{F_{2,\text{up}}(G)} \leq 1, \quad (6)$$

if $F_{2,\text{up}}(G) = F_{2,\text{low}}(G)$, an optimal solution is found. We propose a theorem that generalizes our results to k -clustering for $k > 2$:

Theorem 1. If $F_1(G) \leq F_2(G)$, then $F_1(G) \leq F_k(G)$ for every $k > 2$; where every cluster is nonempty.

Proof. The proof is through induction. Suppose the theorem holds for $k = 2, \dots, k-1$ and there exists a k -clustering that results in $F_k(G) < F_1(G)$. Consider $A_{C,C'}^+$ ($A_{C,C'}^-$) as the number of positive (negative) links between clusters C and C' . In such k -clustering, links from each C toward every other C' must satisfy $A_{C,C'}^- \geq A_{C,C'}^+$. Otherwise, by merging C into such C' , inequality $F_{k-1}(G) < F_k(G)$ is reached, implying the $F_k(G) < F_1(G) < F_k(G)$ contradiction. With this restriction, if there is no C_1 satisfying $A_{C_1,C'}^- > A_{C_1,C'}^+$ for some C' , the $F_k(G) = F_1(G)$ contradiction is reached via merging all clusters into one cluster. Otherwise, we select such C_1 and merge all other clusters into C_2 . Consequently, we find a two-clustering that satisfies $A_{C_1,C_2}^- > A_{C_1,C_2}^+$, and therefore, results in the $F_2(G) < F_1(G) \leq F_2(G)$ contradiction. The proof is complete with this. ■

Reminding that $F_1(G)$ is the number of negative edges in graph G , theorem 1 states that if an optimal two-clustering has worse frustration than a one-clustering, then every k -clustering is also worse than one-clustering, and thus, it is maximally balanced. As a result, if we get $F_1(G) \leq F_{2,\text{low}}(G)$ from Iacono algorithm, which signifies $F_1(G) \leq F_2(G)$, we conclude that inconsistent edges in G cannot be reduced (or equally, balancedness cannot be increased) via k -clustering for $k \geq 2$. Thus, G is optimally clustered into one cluster.

VI. DATASETS

We used two widely studied online signed networks known as *Slashdot* and *Epinions* [17], which have been frequently used as benchmarks for studying signed social relations.³ These datasets have special characteristics that make them suitable for the analysis of social relations. For example, all the links have been explicitly established by the users, either positive (for friendship or trust) or negative (for enmity or distrust). Hence, the links neither have been inferred indirectly nor been asked from a person, which may introduce biasedness into data.

Data preprocessing

We performed some preprocessings on the datasets preparing them for our purpose:

(1) In order to get an undirected network, reciprocal links with inconsistent signs were omitted, and the remaining links were considered as undirected [inconsistent relations were 0.7% (0.4%) of relations in *Epinions* (*Slashdot*)].

³All datasets are publicly available at <http://snap.stanford.edu>. For more detailed statistics refer to <http://konect.uni-koblenz.de/>

TABLE I. Basic statistics of datasets after preprocessing. Average members is the mean of community sizes.

	Node	Edge	Negative edge	Community	Average members
Slashdot	68 409	327 490	69 682 (21.27%)	4598	14.88
Epinions	76 653	220 932	13 921 (6.30%)	6032	12.71

(2) Only the largest connected component of each network was considered (90% of nodes in Epinions and nearly 100% of nodes in Slashdot).

(3) Nodes incident to zero positive edges were removed as they, trivially, belong to an isolated cluster.

(4) After detecting communities, we kept only those of size 3 to 2000 with all connections between them; the reason is provided in the following.

Table I summarizes the properties of networks after the above operations. The community size is lower bounded to 3, which is the minimum trivial group size. We did not consider megascale communities of size larger than 2000 (4 out of 10 000 communities that have 3000, 5000, 8000, and 10 000 nodes), because either they are a composition of many highly interconnected sub-communities, or they have no community structure at all. As a result, they cannot be counted as reliable social communities. Indeed, the size of these communities is significantly far from 150, which is the expected upper-limit for human community [5]. In addition, the significantly high $f_1(G)$ of the largest community in each network fortifies this conjecture. Nonetheless, we further analyze them along with the other unbalanced communities in Sec. VIII, and found similar results for the role of negative edges that lie inside them.

VII. INTERPLAY BETWEEN DENSE POSITIVE AND NEGATIVE TIES

We discussed that the community detection problem in signed networks is to find groups of densely connected positive ties that are as balanced as possible. First, one needs to get an image of interplay between dense positive ties and those with negative sign. To this end, we first detect the communities from positive subgraph of preprocessed networks using InfoMap. In other words, we exclude the negative subgraph and ignore the information given by negative ties. Next, we bring the negative ties back to the network, noting that the communities have been detected beforehand. Considering only the communities of size 3 to 2000 and the connections between them, we find that more than 98% of communities are completely balanced, meaning they contain no negative relations (lower bounding the size to 10 also gives similar results). Knowing that unbalanced

communities are mostly the bigger ones, we also find that more than 98% of negative ties lie between communities (see Table II for more detailed statistics). These results are interesting, since we based our community detection merely on positive ties and ignored the negative ones. One immediate conclusion is that negative ties naturally lie between densely connected positive ties, and thus, both objectives “densely connected positive ties inside cluster” and “negative ties between clusters” could be reasonably satisfied without considering the latter. In other words, positive ties have the major role in detecting the community structure in signed networks, whereas negative ties have a minor effect. These results, somehow, legitimize the idea behind FEC algorithm [14], which scores the nodes regardless of negative ties; however, this may not be the case for other types of networks. This observation is consistent with the findings of Leskovec *et al.* that are based on the analysis of triads [10]. In particular, they concluded that negative ties tend to act like bridges in signed social networks. Nevertheless, due to relatively low amount of negative ties (around 21% in Slashdot and 6% in Epinions after preprocessing), it may not be a significant observation and could be highly probable in random counterparts of observed networks; this issue is investigated in Sec. VII A. Moreover, unbalanced communities, which are mostly the big ones, are analyzed separately in Sec. VIII to investigate the role of negative ties inside them.

A. How significant are the observed statistics?

In order to show the significance of observed statistics in signed networks, first we should define a proper null model to estimate the probability of desired statistics being as extreme as the observed ones. If the estimated probability is small enough, one can conclude that the observed statistics cannot be due to the chance and depend on the characteristics that have been randomized in the null model. We want to show that this significance is due to the particular position of negative edges between dense positive regions. In order to achieve this goal, we proposed null model $M_r(G)$ that is sampled by perturbing r percent of negative links on graph G while keeping the structure of the network fixed. In particular, for a given graph G , we select r percent of negative edges uniformly at random and flip their sign to positive, then randomly select the same amount from positive edges and flip their sign to negative. In this setting, the relative number of negative edges and the structure of networks generated from $M_r(G)$ resemble the observed one, and only, the position of r percent of negative edges is randomly shuffled. This null model has been previously used in Refs. [20] and [10] for $r = 100$. The complete procedure of acquiring a sample statistic from $M_r(G)$ is as follows:

- (1) Perturb r percent of negative ties.
- (2) Apply InfoMap on the positive subgraph.

TABLE II. Community statistics of studied online social networks. Solved negative ties are links that lie between communities. Average frustration is the mean of $f_1(G)$ over communities.

	Slashdot		Epinions	
	Count	percentage	Count	percentage
Balanced communities	4543	98.80%	5952	98.67%
Solved negative ties	68 794	98.73%	13 737	98.67%
Average frustration	0.08%		0.05%	

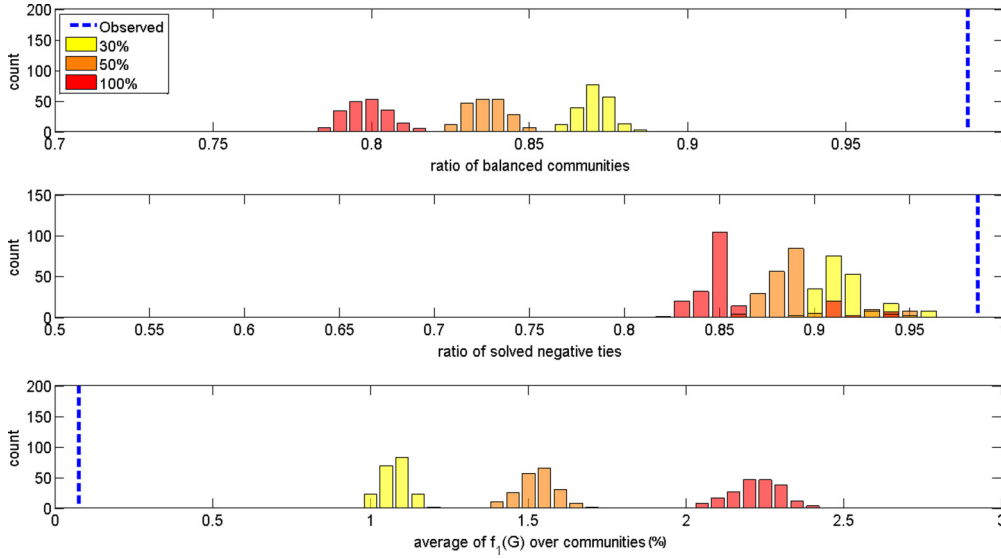


FIG. 5. (Color online) Observed statistics (dashed blue line) in Slashdot dataset as compared to those of 200 realizations from $M_r(G)$ for different values of r . In all cases, p -value is less than 0.001.

- (3) Bring the negative ties back.
- (4) Measure the desired statistic.

In Figs. 5 and 6, the observed statistics are plotted in dashed blue line along with 200 realizations of $M_r(G)$ for $r = 30, 50, \text{ and } 100$ on Slashdot and Epinions networks, respectively. Due to the sufficient number of samples (far more than 30), z test could be confidently used for computing the p -value. This probability is very small ($p \ll 0.001$) for all statistics in both networks compared to the null models with $r = 30\%, 50\%, \text{ and } 100\%$. Consequently, it could be concluded that the ratio of balanced communities and solved negative ties are significantly higher, and average frustration [average of $f_1(G)$ over communities] is significantly lower than being created by chance. Furthermore, since we merely flipped the sign of negative ties, observed phenomenon is due to the topological

position of negative ties implying that *negative ties almost entirely lie between dense positive ties*.

B. Are considerable parts of the networks isolated from negative ties?

One of the plausible causes for the observed statistics would be the isolation of major parts of the networks from negative ties, which can lead to numerous balanced communities. However, over 91% of balanced communities in both networks are incident to at least one negative edge, and the average percentage of external negative ties is around 28% for balanced communities. Although, this is 5–10% lower than that of unbalanced communities, it is sufficient enough to reject the *major parts of networks are free from negative ties* hypothesis.

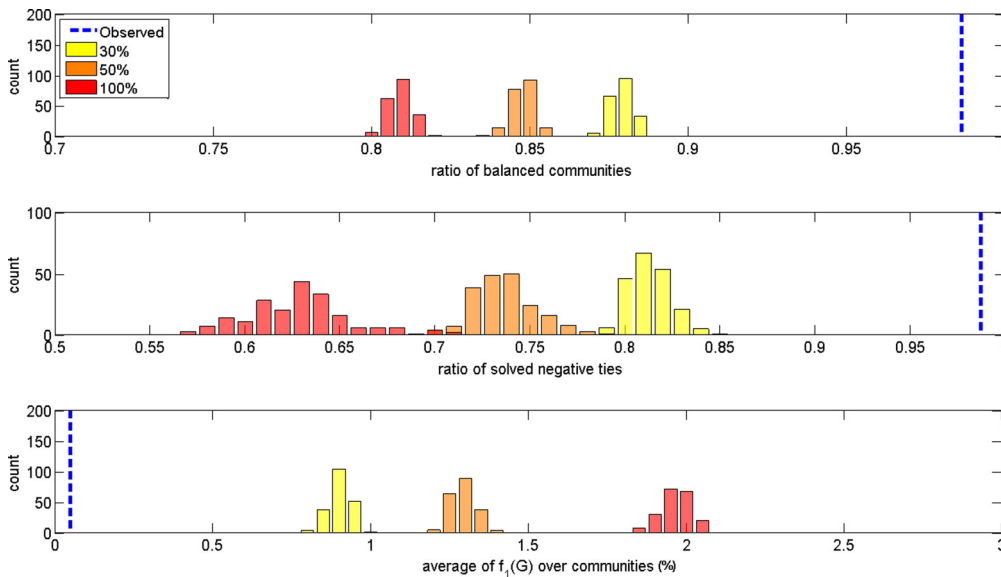


FIG. 6. (Color online) Observed statistics (dashed blue line) in Epinions dataset as compared to those of 200 realizations from $M_r(G)$ for different values of r . In all cases, p -value is less than 0.001.

C. Effect of InfoMap on the observed statistics

As we discussed in Sec. V A, InfoMap may merge highly interconnected communities into each other. First, 98% of communities are balanced and more than 85% of nodes in unbalanced communities are incident to only positive edges, thus by further partitioning these communities, the 98% statistics, if not increasing, would not considerably decrease. Second, by separating mistakenly glued communities, the number of inter-negative (or -positive) ties does not decrease, which is a trivial case. Finally, the average frustration also follows the first case, and thus, it does not considerably increase. Consequently, the three main reported statistics are valid and do not considerably depend on InfoMap algorithm.

VIII. ANALYSIS OF UNBALANCED COMMUNITIES

So far, we have shown that over 98% of communities are balanced, which means they have no internal negative ties. On the other hand, unbalanced communities, which are less than 2% of total communities, are mostly the bigger ones in terms of the number of nodes and ties, and thus, they should be analyzed to find the role of negative ties lying inside them.

Table IV shows the major unbalanced communities including those of size larger than 2000. Considering only communities smaller than 2000 nodes, they have far less negative ties compared to positive ones [$f_1(G)$]. However, this does not mean negative ties are useless from signed community detection or k -clustering point of view. In order to measure the usefulness of these negative edges for extraction of balanced clusters, we propose a simple information-theoretic measure that is based on structural balance theory.

How informative are negative ties?

From the structural balance perspective, at one extreme, negative edges inside community S are in the most informative position, if there exists an optimal two-clustering for S that results in two equally sized clusters. From another point of view, one can argue that by ignoring negative ties, two maximally balanced, equally sized clusters are mistakenly considered as a single community. At the other extreme, negative ties are in the less-informative position, if the

TABLE III. Result of applying two-clustering algorithm of Iacono on unbalanced communities of InfoMap. Percentages are based on total number of unbalanced communities.

	Unbalanced communities	Optimal two-clustering	Zero information
Slashdot	56	55 (98%)	48 (86%)
Epinions	83	81 (98%)	77 (93%)

minimum number of inconsistent ties is achieved by putting S into one cluster. In other words, the same community is achieved with or without considering the negative edges. With this intuition in mind, we define the information of negative ties as follows:

$$I(G) = -\log_2 (\text{ratio of nodes in larger partition}), \quad (7)$$

where *larger partition* is obtained from an optimal two-clustering, which has $F_2(G)$ inconsistent ties. Whenever $F_1(G) \leq F_2(G)$ (no improvement upon one-clustering), the two-clustering is set to one-clustering. This measure is illustrated in Fig. 7 for some toy graphs.

As shown in Table III and detailed in Table IV, for the five largest unbalanced communities, by applying Iacono algorithm we find an optimal two-clustering for 136 out of 139 unbalanced communities. The information of 125 unbalanced communities is zero, for which an optimal solution has been found. Using theorem 1, this means unbalanced communities cannot reach a higher balancedness by being further partitioned into k nonempty clusters, and thus, they are inseparable from GSB point of view. Moreover, for those with $I(G) > 0$ (11 of 139), separated clusters are relatively very small, and also they are, internally, highly disconnected. This indicates there is no significant subcommunity that can be separated from the original one.

In this section, we showed that negative ties inside unbalanced communities are not effectively informative from community detection or the k -clustering point of view. However, the established results are from two widely studied social networks and should be further investigated on other large-scale ones.

IX. RELATION BETWEEN INFOMAP AND SIGNED MODULARITY

We argued in Sec. III C that the objective of signed modularity is in line with the community detection literature. Therefore, in the absence of negative ties, the goal is the same for both InfoMap and modularity. However, there still remain some major problems. First, as our experiments show, nonsigned modularity ($\alpha = 1$) is incapable of distinguishing communities effectively. In particular, the output is mostly made of a few megascale communities of size 2000 to 20 000, which cannot be reliably considered as *single* community (especially for Epinions). Second, modularity suffers from the well-known resolution limit, stating that it is expected to have trivially distinct communities being grouped even in medium-scale networks [45].

In agreement with our results, by sliding the parameter α from 1 to 0.5 (increasing the effect of negative ties with

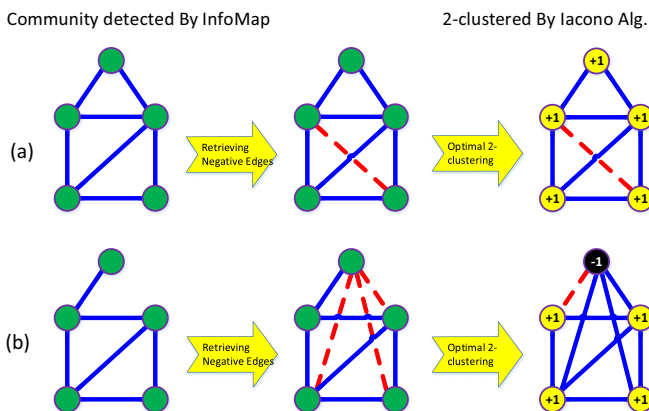


FIG. 7. (Color online) (a) Information of negative edge is equal to $-\log_2(1) = 0$, which also means $F_1(G) \leq F_2(G)$. (b) Information of negative edges is equal to $-\log_2(\frac{4}{5}) = 0.322$.

TABLE IV. Detailed statistics of the five largest unbalanced communities in Slashdot and Epinions networks. Bold communities have sizes larger than 2 000, which were excluded from preprocessed networks (C1 in Slashdot, C1-2-4 in Epinions). Optimal two-clustering is achieved for those communities that have $\frac{F_{2,\text{low}}(G)}{F_{2,\text{up}}(G)} = 1$.

Slashdot								
	Edge	Node	$f_1(C_i)$	$f_{2,\text{up}}(C_i)$	$f_{2,\text{low}}(C_i)$	$\frac{F_{2,\text{low}}(C_i)}{F_{2,\text{up}}(C_i)}$	Larger Partition	$I(G)$
C1	33 552	5378	8.57%	7.97%	7.94%	0.9955	99.00%	0.015
C2	16 156	1705	1.46%	1.40%	1.40%	1	99.65%	0.005
C3	6485	1204	5.97%	4.61%	4.61%	1	98.50%	0.022
C4	4700	1320	2.55%	2.49%	2.49%	1	99.92%	0.001
C5	3896	1095	0.49%	0.49%	0.49%	1	100%	0
Epinions								
	Edge	Node	$f_1(C_i)$	$f_{2,\text{up}}(C_i)$	$f_{2,\text{low}}(C_i)$	$\frac{F_{2,\text{low}}(C_i)}{F_{2,\text{up}}(C_i)}$	Larger Partition	$I(G)$
C1	114 989	10 568	11.40%	8.54%	8.50%	0.9949	97.56%	0.036
C2	66 314	8312	5.93%	5.15%	5.14%	0.9988	99.01%	0.014
C3	36 931	1033	0.05%	0.05%	0.05%	1	100%	0
C4	12 366	2981	0.31%	0.31%	0.31%	1	100%	0
C5	6346	1043	0.77%	0.77%	0.77%	1	100%	0

respect to the density of positive ties), the percentage of solved negative ties (those placed between communities) remains almost constant around 88% in Slashdot. That is, signed modularity is incapable of effectively reducing the internal negative ties. However, in the case of Epinions, this percentage goes from 71% to 82%, which means 11% of negative ties manage to break the communities apart and get placed between them. Nonetheless, we argue that these negative ties are not informative in a general sense. The motivation is that the less an algorithm is capable of distinguishing the communities, the more information it gets from negative ties. In other words, the usefulness of negative ties depends on the performance of a detector. This is intuitively correct, since one can build detector D as

$$D = \begin{cases} \text{One-clustering} & \frac{m^-}{m} < x \\ \text{Signed modularity} & \text{o.w.} \end{cases}, \quad (8)$$

which puts all the nodes in one community until a certain ratio of negative ties is reached (x), and uses the signed modularity afterward. In this case, even if negative ties are truly placed between dense positive ones, they are still useful for detector D , since by exceeding x , the percentage of solved negative ties increases. This example suggests that the presence of megascale communities along with the resolution limit of modularity refrains us from saying that *negative ties are informative* for Epinions. However, if one can find a detector P (i.e., InfoMap), which is more powerful than detector D (i.e., modularity), and the output of P places almost all negative ties between communities, one can confidently state that the negative ties are not informative for the detection task (“more powerful” qualitatively refers to a detector that finds more cohesive groups, and wrongly clusters distinct communities due to having more interconnections). Furthermore, if one can find a detector M that is more powerful than P , the statement is still valid, since detector M further splits the communities of P , rather than clustering them together.

Knowing that InfoMap is more powerful than modularity in nonsigned mode (see Refs. [42,44,46,47]), we consider the

objective function of detector P as follows:

$$L(G, C) = \alpha L(G^+, C) - (1 - \alpha)L(G^-, C), \quad (9)$$

where $L(G, C)$ is the generalization of InfoMap’s objective function. Note that there is still no exact formulation for $\alpha < 1$, nonetheless, we suppose it will be devised in the future, and will outperform modularity for $\alpha < 1$.

According to the results, for $\alpha = 1$, detector P places almost all of negative ties between communities, and thus they have no contribution to $L(G^-, C)$, as it only punishes internal negative ties. In addition, we showed in Sec. VIII A that the remaining internal negative ties are incapable of ripping the unbalanced communities apart, mainly because they are supported by a large number of positive ties. Therefore, even if a general objective $L(G, C)$ is proposed, it cannot considerably improve upon $L(G^+, C)$ for $\alpha < 1$. Moreover, as we previously argued, this statement also holds for even more powerful detectors than P .

It could be concluded that, in the case of Epinions, information of negative ties is helpful for signed modularity, which is also the case for weaker methods like detector D . However, by the use of nonsigned InfoMap, which performs at least as well as modularity, together with the information analysis of internal negative ties, one can conclude that negative ties are not considerably informative for community detection in Slashdot and Epinions.

X. COMMUNITY-COMMUNITY INTERACTIONS

As discussed in Sec. IV, one should let positively related communities be separated from each other. This is the main goal of all community detection methods for nonsigned networks, which have been extended to be fit for signed networks. However, from the GSB perspective, this discrimination is not allowed and has been questioned by Doreian and Mrvar [25]. As a result, they proposed RSB that allows positive relationships between two clusters. In particular, RSB was successfully applied on some small-scale networks that a

TABLE V. Statistics of community networks. Each node represents a balanced community and each link is positive (negative) if $F(i, j) = 1$ [$E(i, j) = 1$]. For each community, friendship (enmity) is the percentage of positive (negative) degree to total degree. A community is friendly if its friendship is larger than 80%. A community is aggressive if its enmity is larger than 50% (similar thresholds do not considerably change the results).

	Slashdot	Epinions
Friendly relations	66%	79%
Enmity relations	29%	19%
Average no. neighbors	58	28
Average friendship	69%	77%
Average enmity	28%	22%
Friendly communities	31%	58%
Aggressive communities	12%	12%

complete scenario of their relations was known. However, this assertion, to the best of our knowledge, has not yet been examined on large-scale networks. Indeed, in this work we try to answer the question, “Are mediator triads frequent enough in signed networks to be considered in clustering algorithms?”

First, we should provide a connection between the output of InfoMap and GB modeling. The GB model does not impose any restrictions on built-in structure of each cluster and leaves it to the algorithm to find a suitable type or to the researcher to prespecify it based on his or her knowledge [36]. However, for large-scale datasets, like those analyzed in this work, this cannot be efficiently done mainly due to the lack of data about the history of individuals. This leaves us with only one option, to use the general assumption that social clusters are likely to be densely connected [9,23,24]. This assumption is a special case in GB modeling known as *complete block*, which is used for detection of cohesive subgroups [36]. With this restriction, we are allowed to investigate the mediator clusters in social networks based on the output of InfoMap and further probe the community-community relations.

Let us define some quantities to investigate community-community interactions quantitatively:

$$F(i, j) = \begin{cases} 1 & \text{all interedges are positive} \\ 0 & \text{o.w.} \end{cases}, \quad (10)$$

$$E(i, j) = \begin{cases} 1 & \text{all interedges are negative} \\ 0 & \text{o.w.} \end{cases}, \quad (11)$$

where (i, j) is a pair of communities. Less than 6% (2%) of community-community relations are ignored due to partial negative-positive ties in Slashdot (Epinions).

As depicted in Table V, more than 66% (in Slashdot) and 79% (In Epinions) of community-community relations are friendship, and on average, a balanced community is friends with around 69% (in Slashdot) and 77% (In Epinions) of its neighbor communities. On the other hand, less than 12% of communities are mostly enemies with their neighbors.

In Fig. 8, the Slashdot network of balanced communities is visualized using Gephi.⁴ It is worth mentioning that the work of

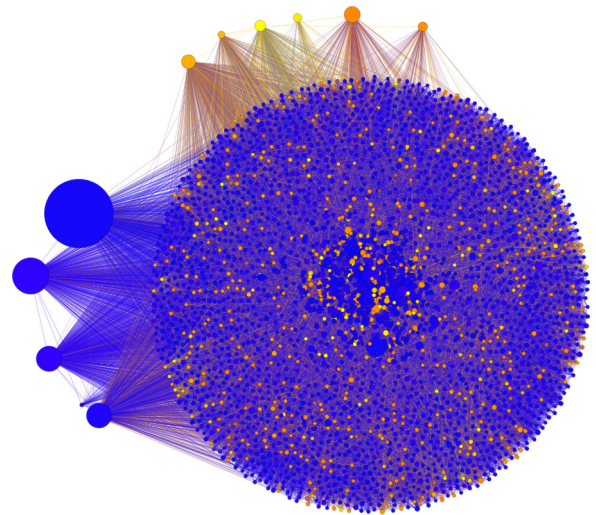


FIG. 8. (Color online) Constructed community network of Slashdot. All communities are internally balanced. The size of each community is proportional to the number of its members. Communities with higher ratio of negative relations are closer to yellow color (lighter gray). Some of the larger aggressive communities, which are the enemies of the majority of their neighbors, have been placed on the top of the network. These communities indicate a meaningful group of people that are allies, and troll the neighbor communities.

Kunegis *et al.* [48] investigated the Slashdot trolls individually (users that are the enemy of most of their neighbors). In the mesoscopic level, our result in Fig. 8 shows some of the major trolling communities unearthing a novel view of Slashdot. The k -clustering algorithms on social networks merely distinguish the aggressive groups. However, low amounts of these clusters indicates that clustering algorithms, even in optimal case, could miss detecting a considerable amount of distinguishable and positively related groups by merging them into one another, similar to Fig. 4(c). Therefore, clustering based on social balance hides a major part of the mesoscopic structure.

Mediator triads should not be overlooked

The relaxation of Doreian and Mrvar could be justified only if the mediator triad appears frequently enough in the mesoscopic level of social relations. In particular, RSB does not claim that the frequency of mediator triads, relative to a null model, is either underrepresented or overrepresented. Nonetheless, it legitimizes the absolute presence of such relations. Therefore, if there is a considerable amount of such relations, even underrepresented, the RSB should be utilized instead of GSB.

TABLE VI. Total number of triads and percentage of negative links in the original Slashdot and Epinions networks and the constructed community network. M and K stand for 10^6 and 10^3 , respectively.

	Local		Mesoscopic	
	No. Triads	Negative ties	No. Triads	Negative ties
Slashdot	0.6M	23.60%	284K	30.40%
Epinions	4.8M	16.80%	35K	20.20%

⁴Gephi is an open source software for visualizing large-scale graphs: <https://gephi.org/>

TABLE VII. Percentage and z scores of each triad type in the local and mesoscopic levels compared to $M_{100}(G)$. Statistics for the $++-$ triad, which is the only unbalanced triad in GSB (excluding $k = 2$), are shown in bold.

	Slashdot						Epinions					
	Local			Mesoscopic			Local			Mesoscopic		
	Real	Random	z score	Real	Random	z score	Real	Random	z score	Real	Random	z score
+++	73	44.6	90	38.2	33.7	16	82.6	57.6	149	66.3	50.7	24
++-	11.2	41.3	-232	33.7	44.2	-105	8.3	34.9	-261	20.8	38.6	-45
+--	13.6	12.7	5	21.9	19.3	13	7.9	7	14	11	9.8	4
---	2.1	1.3	23	6.2	2.8	58	1.1	0.5	93	1.9	0.8	17

In the local level, as shown in Table VII, the $++-$ triad is the only one that is considerably underrepresented, which is 11% (in Slashdot) and 8% (in Epinions) of triads in real networks against 41% (in Slashdot) and 35% (in Epinions) of triads in randomized counterparts. This observation is consistent with the previous findings in favor of GSB (excluding $k = 2$) [10,49].

In the mesoscopic level, we conduct a similar experiment on the network of balanced communities. To this end, we consider each community as a node, and connect a pair of communities with $+1(-1)$ edge, if they were friend (enemy). As summarized in Table VI, the community network has similar percentage of negative edges compared to the local level network. This makes the comparison of two levels more meaningful.

As depicted in Table VII, the ratio of mediator triads is 34% (in Slashdot) and 21% (in Epinions), which is by far higher than 11% (in Slashdot) and 8% (in Epinions) in local level. As listed in the z -score column, mediator triads are also underrepresented according to randomized networks. This means even if this type of triad is a less desirable relation between groups of individuals, nonetheless, this is a notable pattern among them. Hence, as suggested by Doreian and Mrvar, ignoring the intermediary processes leads to merging or even splitting a considerable amount of mediators into hostile parties. In other words, if one has low amount of mediator triads close to that of a k -balanced network, ignoring the mediator triads does not conceal the true mesoscopic structure. However, observed frequencies suggest that although mesoscopic structures are driven away from mediator triads (according to corresponding z scores), the considerable amount of these relations refrain us from simply ignoring them. Moreover, explicitly established relations between users leave no room for the assumption that the mediator triads are mainly due to noise.

In conclusion, although the mediator triad is underrepresented in social networks, it should not be overlooked. In particular, the implication of GSB remains valid in the sense that social dynamics drive the relations away from the mediator triad. Nonetheless, the relaxation of Doreian and Mrvar is still necessary to account for mediator triads, which are still surviving the social dynamics, as a remarkable aspect of social relations.

XI. CONCLUSION

In this work, we investigated the mesoscopic level of online signed social networks. First, we observed that communities

(extracted based on merely positive edges) in signed social networks are highly balanced. This indicates that negative edges mostly lie between dense positive clusters. Also, when negative edges lie inside the communities, they have either no or weak divisive power. In other words, negative edges do not have a significant effect on the community structure of signed networks, and it is mainly determined by positive relations. Furthermore, we showed that this salient characteristic is almost impossible to be created by randomly placed negative edges. This assertion is consistent with the previous studies both on the local level, where it was shown that the clustering coefficient of positive subgraph is much higher than that of negative subgraph [10], and the global level, where it was demonstrated that social networks are highly balanced compared to sign-shuffled ones [20]. This role of negative ties partially explains why sign prediction models that are based on machine learning techniques can perform highly accurately, despite the fact that they utilize the information of merely adjacent nodes [17,18]. Our second observation was that the $++-$ mediator triad between communities is underrepresented consistent with GSB; however, it is highly frequent compared to the triad of the same type between users. Hence, mediator triads cannot be simply ignored as they still survived the social dynamics and form a considerable portion of social relations. As a result, if one only tries to minimize $F_k(G, C)$ regardless of the mediator triads, many intermediary clusters are lost by merging or splitting them into hostile parties, and hence, major parts of the mesoscopic structure remain hidden. Consequently, the routes of RSB-based GB modeling and signed community detection seem to be more consistent with the structure of networks similar to Slashdot and Epinions.

XII. FUTURE WORKS

There are some interesting issues that can be investigated in future works, including:

(i) In this work, we measured the informativeness of negative edges for each community separately. It is fruitful to have a procedure that measures this information in a network as a whole. Although signed modularity can do this work, its major shortcomings make it an unreliable measure for real-world networks [45,50,51]. Nonetheless, along with $F_k(G, C)$, it can be a baseline for future measures.

(ii) An improvement in accuracy of the link prediction is likely to be achieved by augmenting the (nontrivial) statistics of InfoMap communities into machine learning methods. Noting that a successful work has been carried out by

extracting clusters [detected via minimizing $F_k(G, C)$] and further applying collaborative filtering methods [19].

(iii) We showed when negative ties lie between dense positive ties, their informativeness vanishes for the task of community detection. On the contrary, as demonstrated in Ref. [17], they are really useful for inferring hidden links

due to this apparent pattern. Roughly, the less the usefulness of negative ties for community detection, the more their usefulness for link prediction. Thus, an interesting task would be a quantitative analysis of interplay between the information of negative ties in the local level (for link inference) and that of the mesoscopic level (for community detection).

-
- [1] R. Albert and A.-L. Barabási, *Rev. Mod. Phys.* **74**, 47 (2002).
- [2] S. Boccaletti, V. Latora, Y. Moreno, M. Chavez, and D.-U. Hwang, *Phys. Rep.* **424**, 175 (2006).
- [3] S. N. Dorogovtsev and J. F. Mendes, *Adv. Phys.* **51**, 1079 (2002).
- [4] A.-L. Barabási and R. Albert, *Science* **286**, 509 (1999).
- [5] J. Leskovec, K. J. Lang, A. Dasgupta, and M. W. Mahoney, in *Proceedings of the 17th International Conference on World Wide Web* (ACM, Beijing, 2008), pp. 695–704.
- [6] A. Lancichinetti, M. Kivela, J. Saramaki, and S. Fortunato, *PLoS ONE* **5**, e11976 (2010).
- [7] S. H. Strogatz, *Nature* **410**, 268 (2001).
- [8] M. E. J. Newman, S. H. Strogatz, and D. J. Watts, *Phys. Rev. E* **64**, 026118 (2001).
- [9] L. Backstrom, D. Huttenlocher, J. Kleinberg, and X. Lan, in *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (ACM, Philadelphia, PA, 2006), pp. 44–54.
- [10] J. Leskovec, D. Huttenlocher, and J. Kleinberg, in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (ACM, Atlanta, Georgia, 2010), pp. 1361–1370.
- [11] P. A. Grabowicz, L. M. Aiello, V. M. Eguiluz, and A. Jaimes, in *Proceedings of the Sixth ACM International Conference on Web Search and Data Mining* (ACM, Rome, 2013), pp. 627–636.
- [12] P. A. Grabowicz, J. J. Ramasco, E. Moro, J. M. Pujol, and V. M. Eguiluz, *PLoS ONE* **7**, e29358 (2012).
- [13] J.-P. Onnela, S. Arbesman, M. C. González, A.-L. Barabási, and N. A. Christakis, *PLoS ONE* **6**, e16939 (2011).
- [14] B. Yang, W. K. Cheung, and J. Liu, *IEEE Trans. Knowl. Data Eng.* **19**, 1333 (2007).
- [15] S. Gómez, P. Jensen, and A. Arenas, *Phys. Rev. E* **80**, 016114 (2009).
- [16] V. A. Traag and J. Bruggeman, *Phys. Rev. E* **80**, 036115 (2009).
- [17] J. Leskovec, D. Huttenlocher, and J. Kleinberg, in *Proceedings of the 19th International Conference on World Wide Web* (ACM, Raleigh, NC, 2010), pp. 641–650.
- [18] J. Tang, T. Lou, and J. Kleinberg, in *Proceedings of the Fifth ACM International Conference on Web Search and Data Mining* (ACM, Seattle, Washington, 2012), pp. 743–752.
- [19] A. Javari and M. Jalili, *ACM Trans. Intell. Syst. Technol.* **5**, 24 (2014).
- [20] G. Facchetti, G. Iacono, and C. Altafini, *Proc. Natl. Acad. Sci. USA* **108**, 20953 (2011).
- [21] X. Zheng, D. Zeng, and F.-Y. Wang, *Inf. Syst. Frontiers*, doi:10.1007/s10796-014-9483-8 (2014).
- [22] M. Girvan and M. E. Newman, *Proc. Natl. Acad. Sci. USA* **99**, 7821 (2002).
- [23] M. E. J. Newman and J. Park, *Phys. Rev. E* **68**, 036122 (2003).
- [24] M. E. J. Newman, *SIAM Rev.* **45**, 167 (2003).
- [25] P. Doreian and A. Mrvar, *Social Networks* **31**, 1 (2009).
- [26] T. Zaslavsky and C. R. Rao, *Balance and Clustering in Signed Graphs* (Binghamton University, New York, 2010).
- [27] G. Iacono, F. Ramezani, N. Soranzo, and C. Altafini, *IET Syst. Biol.* **4**, 223 (2010).
- [28] I. Giotis and V. Guruswami, in *Proceedings of the 17th Annual ACM-SIAM Symposium on Discrete Algorithm* (ACM, Philadelphia, PA, 2006), pp. 1167–1176.
- [29] N. Bansal, A. Blum, and S. Chawla, *Mach. Learn.* **56**, 89 (2004).
- [30] F. Heider, *J. Psychol.* **21**, 107 (1946).
- [31] D. Cartwright and F. Harary, *Psychol. Rev.* **63**, 277 (1956).
- [32] J. A. Davis, *Human Relations* **20**, 181 (1967).
- [33] P. Doreian and A. Mrvar, *Social Networks* **18**, 149 (1996).
- [34] G. Facchetti, G. Iacono, and C. Altafini, *Phys. Rev. E* **86**, 036116 (2012).
- [35] K.-Y. Chiang, J. J. Whang, and I. S. Dhillon, in *Proceedings of the 21st ACM International Conference on Information and Knowledge Management* (ACM, Maui, HI, 2012), pp. 615–624.
- [36] P. Doreian, V. Batagelj, and A. Ferligoj, *Generalized Blockmodeling* (Cambridge University Press, Cambridge, 2005), Vol. 25.
- [37] S. Fortunato, *Phys. Rep.* **486**, 75 (2010).
- [38] M. Rosvall and C. T. Bergstrom, *Proc. Natl. Acad. Sci. USA* **105**, 1118 (2008).
- [39] M. E. J. Newman and M. Girvan, *Phys. Rev. E* **69**, 026113 (2004).
- [40] M. Rosvall and C. T. Bergstrom, *PLoS ONE* **6**, e18209 (2011).
- [41] Y.-Y. Ahn, J. P. Bagrow, and S. Lehmann, *Nature* **466**, 761 (2010).
- [42] A. Lancichinetti and S. Fortunato, *Phys. Rev. E* **80**, 056117 (2009).
- [43] D. Edler and M. Rosvall, Mapequation software package, available online at <http://www.mapequation.org> (2013).
- [44] T. Kawamoto and M. Rosvall, *arXiv:1402.4385* (2014).
- [45] S. Fortunato and M. Barthelemy, *Proc. Natl. Acad. Sci. USA* **104**, 36 (2007).
- [46] R. Aldecoa and I. Marín, *Sci. Rep.* **3**, 2216 (2013).
- [47] G. K. Orman, V. Labatut, and H. Cherifi, *J. Stat. Mech.: Theory Exp.* (2012) P08001.
- [48] J. Kunegis, A. Lommatzsch, and C. Bauckhage, in *Proceedings of the 18th International Conference on World Wide Web* (ACM, Madrid, 2009), pp. 741–750.
- [49] M. Szell, R. Lambiotte, and S. Thurner, *Proc. Natl. Acad. Sci. USA* **107**, 13636 (2010).
- [50] B. H. Good, Y.-A. de Montjoye, and A. Clauset, *Phys. Rev. E* **81**, 046106 (2010).
- [51] A. Lancichinetti and S. Fortunato, *Phys. Rev. E* **84**, 066122 (2011).