Influence maximization on correlated networks through community identification

Didier Augusto Vega-Oliveros Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo, São Carlos, SP, Brazil.

Luciano da Fontoura Costa Instituto de Física de São Carlos, Universidade de São Paulo, São Carlos, SP, Brazil.

Francisco Aparecido Rodrigues*

Departamento de Matemtica Aplicada e Estatstica,
Instituto de Ciências Matemáticas e de Computação,
Universidade de São Paulo, São Carlos, SP, Brazil.

The identification of the minimal set of nodes that maximizes the propagation of information is one of the most important problems in network science. In this paper, we introduce a new method to find the set of initial spreaders to maximize the information propagation in complex networks. We evaluate this method in assortative networks and verify that degree-degree correlation plays a fundamental role on the spreading dynamics. Simulation results show that our algorithm is statistically similar, in terms of the average size of outbreaks, to the greedy approach. However, our method is much less time consuming than the greedy algorithm.

I. INTRODUCTION

With the popularization of Internet access by mobile devices, online social networks have emerged as a significant medium for information transmission [1–3]. News, rumors and advertisements propagate fast in these networks due to the low average degree of separation between users [3]. Information is also exchanged in communication networks, where users share files related to multiple contents, including images, audio, and video. Communication and social networks are also characterized by a very heterogeneous structure, in which most of the users are low connected, whereas a very small set of them have many connections [3]. Moreover, in some social networks, high degree vertices tend to connect to low degree vertices, defining a disassortative wiring pattern. This complex structure of networks affects the information propagation, defining a hierarchy among the nodes [1, 2, 4]. This means that networks present a special set of nodes that are the most efficient spreaders of information [2, 4–6], i.e., nodes that maximize the average size of outbreaks.

The identification of these influential nodes is important to understand and control the spreading process on social networks [6]. Particularly, the influence maximization problem (IMP) is faced with the selection of a set of η spreaders that trigger the largest cascade of new adopters according to a spreading dynamics [1]. The problem of finding this set of initial spreaders is NP-hard for most of the spreading models [2], which makes the IMP as a challenge for network scientists. Thus, since it is not possible to obtain the optimal results for most of the networks,

the IMP is addressed by optimization and heuristic algorithms. For instance, one of the most studied methods is a hill climbing greedy approach [1, 2], which covers about 63% of optimal for several classes of influence models. Morone and Makse [6] mapped the IMP onto optimal percolation in random networks to identify the nodes that should be removed to minimize the average size of outbreaks. They verified that this set is given by the nodes whose removal break down the network into many disconnected subgraphs. However, this set of nodes does not correspond necessarily to optimal spreaders, as verified by Radicchi and Castellano [7]. Although all these works advanced the study of influence maximization, they disregard patterns of connections, such as degree-degree correlation and community structure, which have a fundamental impact on spreading dynamics [3].

Degree-degree correlations (or assortativity) is a network property in which nodes with similar features, such as degree, tend to be connected. Previous works verified that epidemics spreads faster in assortative networks, but the reach is larger on disassortative structures [8]. Assortativity also influences the spreading threshold [9] and the diffusion time [10]. Although degree-degree correlation influence the spreading dynamics, the role of this network property on the influence maximization problem has not been addressed yet (see for instance [2, 3]). Here analyze how degree-degree correlation affects the average size of outbreaks in rumor dynamics.

We also propose a method for identification of the most influential spreaders based on community organization. Communities are groups of nodes densely connected among them, but with few connections with other groups [11]. Some authors verified that to improve the spreading efficiency, a good strategy is to distribute the seeds among the communities [12–14]. If the community structure is not considered, then only suboptimal solu-

^{*} francisco@icmc.usp.br

tions can be obtained [12]. This happens because vertices belonging to the same community are likely to be more similar to each other and share the same set of neighbors. Although communities influence the information diffusion, only a few studies have considered the community organization to study the influence maximization problem [5, 12–16]. Indeed, most of these works try to reduce the number of candidate vertices according to some evaluation method and the community structure. For instance, Galstyan et al. [12] employed the greedy approach for selecting the seeds in the smallest community and verified that this might cause a global activation cascade even for a small number of seeds. However, the results are restricted to random networks made up of two communities. Wang et al. [13] introduced a community-based greedy algorithm to find the η most influential nodes. The idea is to divide the network into communities and then, by a dynamic programming algorithm, incrementally select the community from which the next influential node is taken. The method involves high computational cost, although it is an order of magnitude faster than the greedy algorithm. In a similar approach, Cao et al., [5] transformed the influence maximization problem into an optimal resource allocation problem in the network communities. Initially, the method assumes that the communities are disconnected. Then, the method selects η candidates from each community according to the degree centrality and a dynamic programming algorithm identifies the final target nodes.

Although these works provided important results on the influence maximization problem, none of them addressed the impact of the assortativity on the propagation dynamics. Moreover, these methods are computationally expensive and consider a relatively low number of initial spreaders, i.e., up to $\eta = 50$ spreaders. Moreover, classical rumor models are not addressed by these studies, although the models by Daley and Kendall [17] and Maki and Thompson [18] are often used to study information dynamics in networks [19–22]. Thus, in the present work, we provide an analysis of the impact of degree-degree correlation on the influence maximization problem, where the information spreading is modeled by the Maki-Thompson algorithm. A simple approach to maximize the information diffusion considering the community structure of the network is introduced. We perform exhaustive simulations in eight real and six artificial complex networks and verify that assortativity plays a significant role on the influence maximization problem. For instance, increasing the number of initial spreaders may not increase the size of the outbreak. Moreover, the selection of influential spreaders through communities is statistically similar to the greedy optimization algorithm. However, our method requires much lower computational cost and, therefore, is more suitable in practice.

II. CONCEPTS AND METHODS

A social network can be represented as a graph G =(V, E) made up of a set V of vertices (nodes) and a set E of edges that connect pairs of vertices. The cardinalities of V and E are denoted by |V| = n and |E| = m, respectively. Here, we consider only undirected and static networks. The degree k_i of a vertex i corresponds to the number of edges attached to i. The degree distribution of a network P(k) gives the probability that a given randomly selected vertex has degree k. Social networks are characterized by highly heterogeneous degree distribution, presenting a scale-free organization [23], where most of the nodes are low connected, but a small set of nodes have high degree. The connection pattern of vertices can also be analyzed by the degree-degree correlation. In assortative, or positively correlated, networks nodes of similar degree tend to be connected. In disassortative, or negatively correlated, networks low-degree nodes tend to connect with strongly connected vertices. If the tendency of connection is independent of the node degree, then the network is called non-assortative. The level of assortativity can be quantified by the Pearson correlation coefficient, ρ , of the degrees of nodes at either end of an edge [24]. According to this measure, a network can be classified as (i) assortative ($\rho > 0$), (ii) dissassortative ($\rho < 0$), or (iii) non-assortative ($\rho \approx 0$). Degree-degree correlation plays a fundalmental role in the analysis of several dynamical processes in networks, like synchronization [25] and epidemic spreading [24].

A. Influence models

The spreading of rumors or information can be approached as a psychological contagion where an idea "contaminates" the mind of a population [26]. In the general rumor approach [17, 18] ignorant vertices are those who are unaware of the information, spreaders are informed individuals that can transmit the rumor, and stiflers are individuals who have heard the rumor, but do not spread the information anymore [3, 22, 26, 27]. Thus, each subject can be in one of the three states. i.e., ignorant, spreader or stifler, at each time step. Notice that stiflers act as recovered individuals in a disease spreading model, as they do not participate in the spreading process anymore [3]. Rumor models are different from the traditional susceptible-infected-recovered (SIR) spreading model, in the sense that the transition between states occurs only through contacts, whereas the transition from infected to recovery in the SIR model occurs spontaneously, independent of the connections.

In the rumor model proposed by Maki and Thompson [18], a node that knows the rumor tries to pass the information to each of its neighbors according to a probability β . When the contact is performed between two informed individuals, the active spreader becomes a stifler according to probability μ . Here, we consider this

model for information propagation because the contagion happens only through contact between individuals, which implies that the network structure impacts the dynamics. In addition, this model simulates real social dynamics more accurate than sociology-based models, as people may have multiple opinions, i.e., positive, hesitating or negative [28]. In rumor propagation, spreaders behave like individuals with positive views and stiflers with negative ones.

B. Influence maximization

The influence maximization problem (IMP) seeks the set \mathbf{S} of vertices which contain $|\mathbf{S}| = \eta$ initial seeds that maximize the reach of information. The propagation impact $(\sigma(\mathbf{S}))$ for the set \mathbf{S} of seeds corresponds to the expected fraction of vertices that were informed during the spreading process.

Let us consider a discrete diffusion scenario in which each vertex i can be in only one state at each time step. The initial conditions for the influence maximization problem is defined as $S(0) = V \setminus \mathbf{S}$, $I(0) = \mathbf{S}$ and $R(0) = \{\emptyset\}$, where S(0) represents the set of ignorants in time t = 0 and I(0) the set of initial spreaders or initial seeds S. At each time step, all spreaders uniformly try to infect their neighbors with probability β , or stop the diffusion with probability μ according to the truncated dynamics [22]. More specifically, an spreader tries to inform each of its neighbor until meets another informed node and becomes a stifler. The process ends when $I(\infty) = \{\emptyset\}$ and then we can calculate the final fraction of informed individuals $(\widehat{\sigma}_{\mathbf{S}})$, i.e. $\widehat{\sigma}_{\mathbf{S}} = |R(\infty)|/|V|$ or $\widehat{\sigma}_{\mathbf{S}} = 1 - |S(\infty)|/|V|$. However, $\widehat{\sigma}_{\mathbf{S}}$ is a function with stochastic fluctuations. Thus, the influence function, $\sigma(\mathbf{S})$, is estimated by performing a sufficient number of calculations of the final fraction of informed individuals $\widehat{\sigma}_{\mathbf{S}}$, with the set of initial spreaders **S**:

$$\sigma(\mathbf{S}) = \frac{1}{K} \sum_{\kappa=1}^{K} \widehat{\sigma}_{\mathbf{S}}^{\kappa} \tag{1}$$

where $\widehat{\sigma}_{\mathbf{S}}^{\kappa}$ represents the final fraction of informed individuals for a particular run κ , and K is the total number of simulations in order to obtain a good estimate of the mean value of $\sigma(\mathbf{S})$. We also calculate the fraction of spreaders over time $\widehat{\phi}_{\mathbf{S}}(t) = |I(t)|/|V|$. In this way, the expected fraction of spreaders in each timestep is defined by:

$$\phi_{\mathbf{S}}(t) = \frac{1}{K} \sum_{\kappa=1}^{K} \widehat{\phi}_{\mathbf{S}}^{\kappa}(t), \tag{2}$$

which is a concave downward function. With this formulation it is possible to calculate the expected timestep (t_p) in which the maximum number of spreaders occurs, i.e., $\{t_p \geq 0 \mid \phi_{\mathbf{S}}(t_p) - \phi_{\mathbf{S}}(t) \geq 0, \forall t \geq 0\}$, and $\phi_{\mathbf{S}}(t_p)$ is

the maximum expected fraction of spreaders during the propagation.

The velocity of the propagation (\bar{V}) given a set of initial spreaders **S** is defined as the variation of the number influenced individuals in given time interval,

$$\bar{V}_p(\mathbf{S}) = \frac{\phi_{\mathbf{S}}(t_p) - \phi_{\mathbf{S}}(0)}{t_p} \quad \text{for} \quad t_p > 0$$
 (3)

where $\phi_{\mathbf{S}}(0) = |\mathbf{S}|/|V|$ is the fraction of initial spreaders. We propose this new measure for evaluating how the methods for selecting the most influential spreaders and the assortativity impact in the propagation velocity.

To study the influence maximization problem, we can consider two different approaches:

- 1. Heuristic methods that analyze the topology of the network, assuming that there exists a strong relation between the structure and the propagation process, but without guarantee of optimal results.
- 2. Optimization methods that analyze the influence models and their properties, assuming that, by some optimization strategies, it is possible to provide an approximation that guarantees at least 63% of the optimal value, but neglecting the network structure at the expense of considerable computational cost.

Figure 1 shows the three methods considered here for solving the influence maximization problem: (i) by selecting the η vertices with the highest value of a given centrality measure; (ii) by detecting the η communities on the network and selecting the most central nodes inside each community; and (iii) by a greedy approach, that is a hill-climbing optimization that returns the η most influential nodes. These methods receive as input the network G and the number of η initial spreaders and return the set \mathbf{S} $(|\mathbf{S}| = \eta)$ of influential spreaders that maximize the information propagation.

Methods based on network centrality assumes that the most central nodes convince the largest number of individuals on the network [29, 30]. However, the problem in this case is how to select the most suitable centrality measure to identify the most influential spreaders [4, 27, 29], since centrality can be defined in terms of distance, flow and random walks, for instance [4]. After defining the centrality measure, the influence maximization problem select the η most central vertices as the initial spreaders (see Figure 1). As in previous works [4, 27, 29–31], here we consider degree (DG), betweenness centrality (BE) and PageRank (PR) to measure the centrality of each node.

Another heuristical approach considered here is based on network community structure. A community is a group of nodes that has more connections between them than with nodes in other groups. In social networks, communities represent people that share affinities, defining the phenomenon of homophily [14]. This condition is the reason that information or sentiments propagate

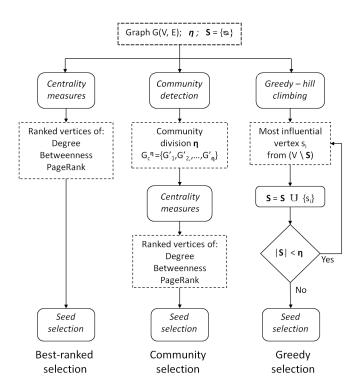


FIG. 1. Methods considered here to address the influence maximization problem (IMP).

better in a community, making people more open to the information shared by their peers [14]. Few works have considered the community structure in the influence maximization problem [5, 12–16, 32]. Here, we propose a new approach to select the initial spreaders, i.e., we consider the most central vertices within each community. Several methods have been developed to detect the communities on networks [11, 13]. Here we employ the fastgreedy algorithm, which is a fast and accurate method for community identification [11, 33]. As depicted in Figure 1, the main η communities of the network are detected and for each G'_i community, we calculate a centrality measure and select the most central vertex within the community. Notice that we fix the number of influential spreaders η and the communities obtained may not correspond to the best division of the network into communities, which vields the maximum modularity. Therefore, we obtain η seeds in which the influence overlap is minimized but the influence within communities is maximized.

We also consider an approximation method based on a greedy hill-climbing algorithm [1]. Many approaches have been derived from this general greedy method [2], such that most of them try to reduce the computational complexity to some polynomial order [2]. Here, we consider only the general greedy method [2]. The algorithm determines among all vertices $\mathbf{s}_i \in V \setminus \mathbf{S}$, i.e. $\{\mathbf{s}_i \in V \mid \mathbf{s}_i \notin \mathbf{S}\}$, the node that maximizes the function $\sigma(\mathbf{S} \cup \{\mathbf{s}_i\})$, recalling that \mathbf{S} is initially empty. Afterwards, the vertex \mathbf{s}_i is added to the set of seeds $\mathbf{S} = \mathbf{S} \cup \{\mathbf{s}_i\}$ and the procedure runs until the target

TABLE I. Topological measures of the networks considered here. ρ is the assortativity coefficient, |V| the network size, $\langle k \rangle$ the average degree, $\langle g \rangle$ is the average shortest path length and $\langle C_c \rangle$ is the average clustering coefficient. The community-related parameters are the modularity Q and the number of communities N_c .

Network	ρ	V	$\langle k \rangle$	$\langle g \rangle$	$\langle C_c \rangle$	FastGreedy	
						Q	Nc
BA	-0.43	1000	11.9	2.96	0.017	0.26	8
BA	-0.31	1000	11.9	2.87	0.028	0.25	9
BA	-0.21	1000	11.9	2.86	0.031	0.25	9
BA	0.02	1000	11.9	2.91	0.035	0.25	10
BA	0.11	1000	11.9	2.94	0.034	0.25	11
BA	0.34	1000	11.9	3.11	0.026	0.27	8
Google+	-0.39	23613	3.32	4.03	0.174	0.74	33
Internet	-0.20	22963	4.22	3.84	0.231	0.63	57
Caida	-0.20	26475	4.03	3.87	0.208	0.64	43
Advogato	-0.09	5054	15.6	3.27	0.253	0.34	49
email	0.01	1133	9.62	3.60	0.220	0.49	16
Hamsterster	0.02	2000	16.1	3.58	0.539	0.46	57
PGP	0.23	10680	4.55	7.48	0.266	0.85	179
Astrophysics	0.23	14845	16.1	4.79	0.638	0.63	1172

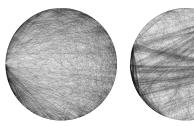
set achieves the size $|\mathbf{S}| = \eta$.

III. DATABASES

We perform extensive numerical simulations in several artificial and real-world networks, evaluating the impact of the degree correlation in the influence maximization problem. The structural properties of the networks are summarized in Table I, with the respective assortativity ρ , number of vertices |V|, average degree $\langle k \rangle$, average shortest path length $\langle g \rangle$ and the average clustering coefficient $\langle C_c \rangle$. Also, the highest modularity Q value and number of communities N_c identified by the fastgreedy algorithm are reported. We can see that real-world networks are more modular than artificial networks. However, most networks present similar values of average shortest path length $(\langle g \rangle)$.

A. Artificial networks

We employ the algorithm proposed by Xulvi-Brunet-Sokolov [34] for controlling the degree-degree correlation in Barabási-Albert (BA) networks. This algorithm performs rewirings in order to increase or decrease the degree-degree correlation, i.e., by favoring the connection between nodes with similar degrees or by hubs and low degree nodes. We adopt this particular model as social networks also present scale-free degree distribution and degree-degree correlation. Figure 2(b) presents examples of assortative and disassortative networks.



(a)Disassortative network

(b)Assortative network

FIG. 2. Circular graph representation of a BA network with the degree correlation modified by the Xulvi-Brunet-Sokolov algorithm [34]: (a) disassortative BA with $\rho = -0.43$ and (b) assortativy BA with $\rho = 0.34$.

B. Real world networks

We consider eight real world representing connections in social and communication networks. The disassortative networks are (i) Google + [35], which is an user-user social network; (ii) Internet [36], which represents a fraction of the symmetrical snapshot of the Internet structure at the level of autonomous systems, reconstructed from BGP tables published on ROUTEVIEWS.ORG project; (iii) Caida [37], which is an undirected network whose nodes are autonomous system on the Internet, collected in 2007 from the CAIDA project; and (iv) advogato network [38], which is an online platform for free software community launched in 1999 that considers trust relationship between developers. The assortative networks are (i) the email network [39], which is a network of emails exchanged between members of the Rovira i Virgili University; (ii) hamsterster [40], which is a network based on the friend and family relationship among users of the HAMSTERSTER.COM website; (iii) PGP [41], which is the largest component of the network of users of the Pretty-Good-Privacy algorithm for secure information interchange; and (iv) astrophysics [42], which is a collaborative network between scientists on previous studies of astrophysics on arXiv. We assume that all these networks are undirected and unweighted. Only the largest network component is considered in our analysis.

IV. RESULTS AND DISCUSSION

A. Impact of assortativity on artificial networks

We calculate the final fraction of influenced individuals according to Equation (1). The number of initial seeds varies from two nodes to 90% of the total number of nodes. Here, we define the number of simulations K = 600. The impact of degree correlation on the influence maximization problem is shown in Figure 3. We observe that an unexpected phenomenon occurs when networks are disassortative (see Fig. 3(a),(b) and (c)) —

the curve of influenced nodes $(\sigma(S))$ has a peak when the number of seeds corresponds to about 10% of the network and then starts to decline when the number of seeds is increased. This peak is due to the low interaction between the initial spreaders in disassortative networks, i.e., central nodes (e.g. hubs) are connected through low degree nodes, increasing the distance between them. If the number of seeds is higher than 10% of the total number of nodes, then the overlap of influence between spreaders occurs and they become stiflers more frequently than for a smaller number of seeds. The number of informed nodes increases again when the number of seeds is large enough to compensate this effect. For assortative networks, this fraction corresponds to about 30% of the network nodes. Thus, the increase in the number of seeds does not always improve the reach of a rumor in disassortative networks.

For non-correlated and assortative networks (see Fig. 3(d),(e) and (f)), selecting the central nodes by considering the communities or the whole network lead to the lowest $\sigma(\mathbf{S})$ and the propagation influence is even worse than the uniform selection of seeds (see the curves for the random case). This effect is also due to the interaction between the initial spreaders.

We also observe in Fig. 3 that the selection according to communities yields similar results for all centrality measures considered. The same happens to the selection according to global centrality (best-ranked approach). Thus, the selection of different centrality measures does not affect the prediction of the fraction of influenced nodes significantly. This can be explained by the high correlation between degree, PageRank and betweenness centrality in BA networks.

Comparing all the cases and the whole interval of seed selection, we observe that the propagation of information is enhanced in networks with higher degree-degree correlations. However, when the number of seeds corresponds to less than 10% of the network, the highest number of informed nodes is obtained in disassortative networks. Thus, the reach of a rumor depends on the level of network assortativity, the method for selecting the initial spreaders and the number of seeds. This result was not expected, since the number of infected nodes should always increase with the number of seeds. This is approximately observed for the uniform selection of initial spreaders, without considering the network topology.

Figure 4 shows the time evolution of the fraction of informed nodes calculated according to Equation (2). As the number of seed nodes is increased, the peak of informed nodes occurs earlier. The highest peak is obtained for the choice according to the large scale centrality measures in disassortative networks, overcoming the greedy approach. The velocity to reach the peak of informed nodes is shown in Figure 5. Only in networks whose seed nodes are selected according to large-scale centrality, are affected by the level of degree-degree correlation. For the other two methods, i.e., according to communities or the greedy approach (Fig. 5(b) and (c)), the degree correlation of the network is not apparently affecting this

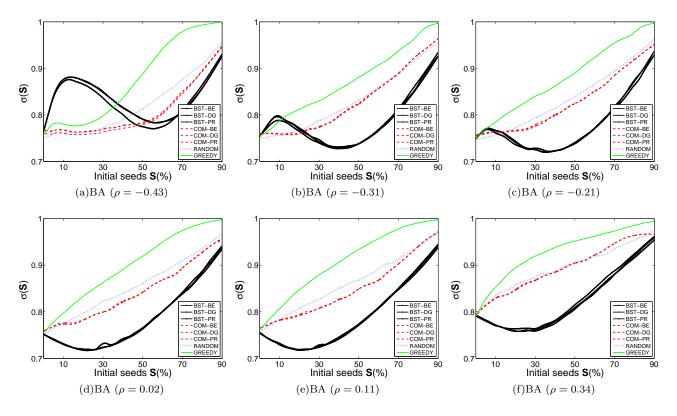


FIG. 3. Impact of degree-degree correlation on the influence maximization problem with $|\mathbf{S}|$ in the full range from 2 to 90% of vertices. For each artificial scale-free network, we calculate the set of initial seeds according to: (BST) the best-ranked vertices of the network; (COM) the most central vertices from communities; (RANDOM) randomly selecting the initial seeds; and (GREEDY) the greedy optimization method. The adopted centrality measures are betweenness centrality (BE), degree (DG) and PageRank (PR).

velocity. In these cases, the velocity is increased until a maximum value and then start to decrease, mainly in the case of seeds selected according to communities. This is also an effect of the increase in the interaction between initial spreaders when their fraction is around 10% of the network size.

B. Impact of assortativity on real world networks

We also analyze the information maximization problem in eight real-world networks (see Table I). These networks present different levels of degree-degree correlation and community organization. The final fraction of stiflers, i.e., informed individuals, $\sigma(\mathbf{S})$ is calculated by considering the number of initial spreaders (seeds) in the interval $[2, \eta_{\text{max}}]$, where η_{max} corresponds to 10% of nodes. The set of initial spreaders are selected uniformly at random or according to centrality measures calculated from the whole network or inside communities. Fig. 6 shows the final fraction of stiflers obtained according to different number of seeds. We can see that a peak in the curve $\sigma(\mathbf{S})$ occurs in disassortative networks. These results are similar to those obtained in artificial networks (see Figure 3).

Fig. 6 also shows that the selection of initial spreaders according to communities provides more informed individuals than the greedy approach in the Google+network. For the remainder networks, the greedy and community-based methods provide similar fraction of informed nodes. Thus, since the greedy method is computationally expensive, the selection of seeds according to communities revealed to be more suitable. Moreover, the choice of seeds according to global centrality measures provides lesser informed nodes than the random selection of seeds, mainly in assortative networks. This is due to the interaction of central spreaders in the early steps of the process, as central nodes tend to be connected in these networks.

Figure 7 shows the velocity for achieving the peak of spreaders, $\bar{V}_p(\mathbf{S})$, for the three methods considered here and different sizes of \mathbf{S} , grouped by the networks: Google+ (dissassortative) (Fig. 7(a)), email (noncorrelated) (Fig. 7(b)); and astrophysics (assortative) (Fig. 7(c)). We notice that the method based on communities presents propagation velocities higher than the greedy approach, for sizes of initial seeds lower than 2% of network size. For non-correlated and assortative networks, all the methods present close behavior. In the case of assortative networks, the selection by communi-

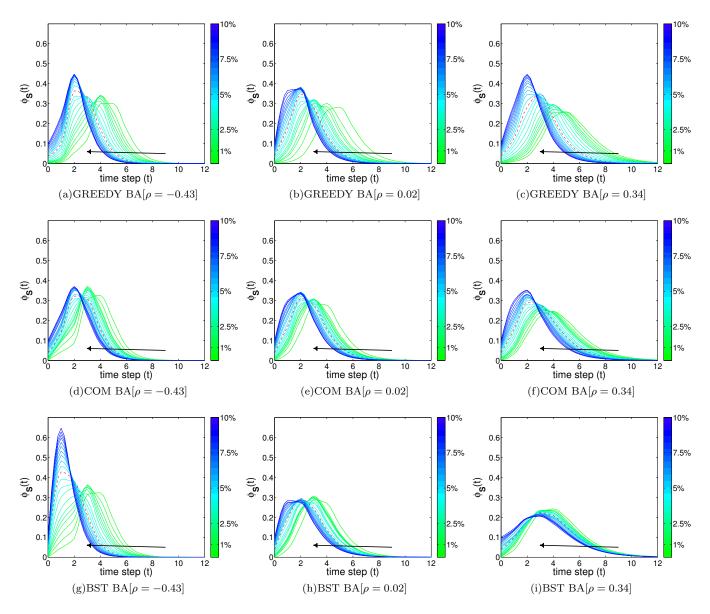


FIG. 4. Impact of degree-degree correlation on the peak of infected nodes. The number of seeds (η) is in the range from 2 to 10% of network size, indicated by the color bars. Arrows indicate increasing η . Initial spreaders are selected by the greedy-optimization (GREEDY), communities (COM) and best-ranked (BST) methods.

ties also provides the higher velocities. Thus, the information propagate faster when the initial seeds are selected by communities.

C. Distribution of initial spreaders

We expect that rumors is better spread when the seeds are distributed equally among communities. This should occur because each seed tries to infect its own community and the interaction between pairs of spreaders is minimized. This hypothesis is verified here by inspecting the distribution of seeds among communities. We consider the normalized variation of information (NVI), which is

an information-theoretic metric that obeys the triangle inequality and is normalized in a stochastic sense [43]. This measure is built upon fundamental concepts from information theory [44] defined as follow: Given two sets of discrete variables or sets \mathbf{X} and \mathbf{Y} , their joint information entropy (\mathcal{H}) and mutual information (\mathcal{I}) are expressed respectively in terms of the marginal and joint distributions of \mathbf{X} and \mathbf{Y} as:

$$\mathcal{H}(\mathbf{X}, \mathbf{Y}) = -\sum_{\boldsymbol{x} \in \mathbf{X}} \sum_{\boldsymbol{y} \in \mathbf{Y}} p(\boldsymbol{x}, \boldsymbol{y}) \log p(\boldsymbol{x}, \boldsymbol{y}), \qquad (4)$$

$$\mathcal{I}(\mathbf{X}, \mathbf{Y}) = \sum_{\boldsymbol{x} \in \mathbf{X}} \sum_{\boldsymbol{y} \in \mathbf{Y}} p(\boldsymbol{x}, \boldsymbol{y}) \log \left(\frac{p(\boldsymbol{x}, \boldsymbol{y})}{p(\boldsymbol{x})p(\boldsymbol{y})} \right), \quad (5)$$

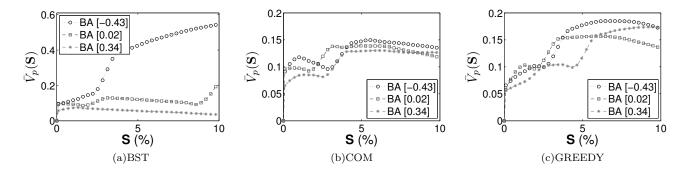


FIG. 5. Velocity of the peak of spreaders in the IMP according to the degree correlation. The number of seeds is in the range from 2 to 10% of network size, selected by: (a) best-ranked (BST), (b) communities (COM) and (c) the greedy-optimization methods (GREEDY).

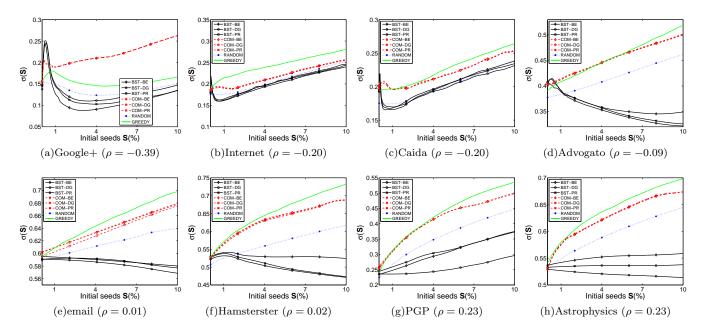


FIG. 6. Impact of degree correlation in the real-world networks. The number of seeds varies from two nodes to 10% of network size, selected according to communities (COM), ranking of the most central nodes (best-ranked, BST), random (RANDOM) and by the greedy-optimization method (GREEDY). Betweenness centrality (BE), degree (DG) and PageRank (PR) are the centrality measures used.

where the joint probability p(x, y) denotes the probability that a data item belongs to **X** and **Y**. The mutual information measures the overlap between the two sets, however it is not a metric nor is it normalized. Thus, for this study we employ the NVI,

$$NVI(\mathbf{X}, \mathbf{Y}) = 1 - \frac{\mathcal{I}(\mathbf{X}, \mathbf{Y})}{\mathcal{H}(\mathbf{X}, \mathbf{Y})},$$
 (6)

which is normalized in [0,1]. It takes a value of 0 when the two sets are identical in the information of the item distribution and 1 when we have a complete dissimilarity between the partitions, i.e., they do not share information.

In terms of the seed and community distribution, we identify η main communities, where each com-

munity is defined by the subgraph $G_c(V_c, E_c)$ and $\sum_{c=1}^{\eta} |G_c(V_c, E_c)| = |V|$. Since these communities may have different sizes, we define \mathbf{X} as the set that describe the size of the communities, in which the probability of selecting a vertex belonging to some community \mathbf{x} is $p(\mathbf{x}) = |V_x|/N$. On the other hand, we say that \mathbf{Y} is the set that denotes the number of seeds in the communities, where $p(\mathbf{y}) = |S^{\mathbf{y}}|/\eta$ is the probability that a seed belongs to some community \mathbf{y} .

We calculate the NVI measure for the seeds selected according to the greedy approach and degree centrality, as shown in Table II. We notice that NVI is lower for the greedy algorithm than for the case in which seeds are selected according to degree. Thus, the greedy approach tends to select seeds homogeneously distributed among

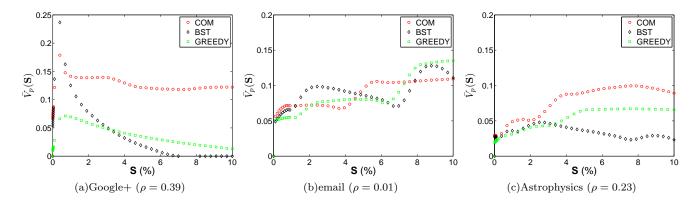


FIG. 7. Velocity of propagation to reach the peak of spreaders. The number of seeds varies from two nodes to 10% of vertices, selected according to the greedy-optimization (GREEDY), communities (COM) and best-ranked (BST) methods.

TABLE II. The normalized variation of information (NVI) measure calculated for the selection of seeds according to the greedy (GREEDY) approach or degree centrality (BST-DG).

	/ 11		0	0 (
	Network	$ \mathbf{S} $	GREEDY	BST-DG
		(%)	NVI	NVI
	BA[-0.43]	0.5	0.409	0.409
	BA[-0.43]	1	0.556	0.535
DN	BA[-0.43]	10	0.706	0.731
	google+	0.5	0.311	0.408
	google+	1	0.196	0.334
	google+	10	0.464	0.725
	BA[0.02]	0.5	0.409	0.409
	BA[0.02]	1	0.469	0.651
NA	BA[0.02]	10	0.708	0.843
	email	0.5	0.368	0.306
	email	1	0.524	0.695
	email	10	0.603	0.809
	BA[0.34]	0.5	0.344	0.689
	BA[0.34]	1	0.492	0.858
AN	BA[0.34]	10	0.693	0.861
	astrophysics	0.5	0.523	0.940
	astrophysics	1	0.589	0.948
	astrophysics	10	0.558	0.933

communities. Thus, this result supports the hypothesis that the seeds of the greedy method are distributed according to the communities of the network. This also indicates that the influence maximization problem can be addressed by the identification of the most central nodes inside communities, instead of considering the greedy approach, which is computationally more expensive.

1. Statistical analysis

Since the greedy approach selects seeds uniformly among communities, it is expected that the community-based and greedy methods provide similar number of informed nodes. Thus, we perform a statistical test to compare the performance of the four methods consid-

TABLE III. Final fraction of informed nodes ($\sigma(\mathbf{S})$) according to different methods.

Network	GREEDY		COM-DG		BST-DG	RANDOM
	$\sigma(1\%)$	$\sigma(10\%)$	$\sigma(1\%)$	$\sigma(10\%)$	$\sigma(1\%)$	$\sigma(1\%)$
Google+	0.1686	0.1663	0.1897	0.2627	0.1549	0.1501
internet	0.2130	0.2810	0.1906	0.2570	0.1639	0.1755
caida	0.1959	0.2640	0.1966	0.2531	0.1696	0.1743
advogato	0.4071	0.5179	0.4138	0.5002	0.3933	0.3821
email	0.6086	0.6976	0.6089	0.6707	0.5942	0.5941
ham ster ster	0.5738	0.7318	0.5693	0.6871	0.5408	0.5236
PGP	0.3126	0.5356	0.3122	0.4985	0.2593	0.2654
astrophysics	0.5785	0.6979	0.5735	0.6745	0.5418	0.5456

ered here. Initially, we analyze the influence maximization considering 1% of the network as initial spreaders (see Table III). We consider a statistically significance test employing the Friedman and Nemenyi approach [45]. The Friedman test is a non-parametric counterpart of the well-known ANOVA (analysis of variance), with the corresponding Nemenyi post-hoc test for comparing the average ranks of the algorithms. If the null hypothesis of similar performance is rejected by the Friedman test, we proceed with the Nemenyi post-hoc test for pairwise comparisons, verifying whether the differences in rank values are statistically significant.

The critical diagram representation suggested by Demšar [45] provides a visual method to compare the results. In the diagram, a horizontal line represents the axis with the average rank values of the methods. In this axis, the lowest (highest) ranked methods are on the left (right) side. Algorithms that are not significantly different from each other are connected through a bold horizontal line. The performance between methods is significantly different if their corresponding average ranks differ by at least the critical difference CD. The value of CD given by the Nemenyi test is presented on the top of the diagram.

According to the result of $\sigma(1\%)$ in the Table III, the chi-square statistics for the methods is 19.20, and the critical value of the chi-square statistics with 3 degrees of

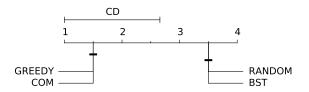


FIG. 8. The critical difference (CD), according to the Nemenyi test, for comparing the mean-ranking of two different methods at 95 percentile is 1.66. Mean-ranking differences above this value are significant and unconnected.

freedom at 95 percentile is 7.81. Thus, for the Freidman test using the chi-square statistics, the null-hypothesis that all methods behave similarly should be rejected. Moreover, we calculate the F-statistics of the methods, obtaining the value 28.00. With 3 and 21 degrees of freedom and at 95 percentile, the critical value of the F-statistics is 3.07, indicating that the null-hypothesis should be rejected again. Therefore, the method do not provide statistically similar results.

Since the methods do not provide the same fraction of informed nodes, we apply the post-hoc Nemenyi test in order to find which method achieves the maximum influence. The critical diagram of the Nemenyi test is shown in Figure 8. The CD for comparing the mean ranking between two methods at 95 percentile is 1.66. Mean-ranking differences above this value are statistically significant. Thus, we conclude that there is no statistically significant differences in the influence maximization results between the greedy and community method when the number of initial spreaders represents less than 1% of the network. However, the Nemenyi test indicates significant differences between the methods based on community centrality and those that consider random selection of spreaders or selection according to their centrality.

We verify in Figure 6 that for $|\mathbf{S}| > 1\%$, the methods based on greedy optimization and community centrality provide the highest number of informed nodes. Thus, we perform the statistical hypothesis test only on the methods based on optimization and community organization. For evaluation of these two algorithms in multiple data sets, we employ the Wilcoxon signed-rank test [45]. This statistical test is a non-parametric alternative to the paired t-test. We adopt the Wilcoxon test because it is less sensible to outliers and does not assume a particular population distribution [45]. For rejecting the null-hypothesis of similar performance, the W-value returned by test should be smaller than the corresponding critical W_c value of Wilcoxon test table.

We compute the Wilcoxon statistical test at 95 percentile for the greedy and community approach considering the number of initial spreaders as $|\mathbf{S}|=10\%|V|$ (Table III). As a result, we obtain a W-value = 8. The critical value for eight networks at p=0.05 is $W_c=3$. Therefore, the null-hypothesis of similar performance of the methods cannot be rejected. These results suggest that the fraction of informed nodes provided by the greedy optimization algorithm and the method based on community centrality is statistically similar. Therefore, since the community-based method is computationally faster than the greedy algorithm, it is more suitable to address the influence maximization problem in practice.

V. CONCLUSION

We have analyzed the role of degree-degree correlation in the influence maximization problem. To simulate the information spreading, we consider the rumor model proposed by Maki and Thompson [18], which is more suitable to represent the information dynamics in social networks [3]. We have proposed a method to maximize the influence transmission based on network community organization. This method has been analyzed by performing simulations on the top of eight real and six artificial complex networks. We have verified that our method is statistically similar, in terms of the information reach, to the approach based on greedy optimization, which is computationally expensive. Thus, our results suggest that our method is more suitable in practice, since it can provides similar results as the greedy approach, but it is less time consuming.

Our analysis can be extended with the consideration of patterns of connections inside networks (e.g. [46, 47]) to select the set of initial spreaders. The study of weighted [48], multilayer [49, 50] and dynamical networks [51] is also promising. In all these cases, general methods for community identification in networks are necessary.

VI. ACKNOWLEDGMENTS

D.A.VO acknowledges CNPq (grant 140688/2013-7). F.A.R. acknowledges CNPq (grant 305940/2010-4), FAPESP (Fundação de Amparo à Pesquisa do Estado de São Paulo, grants 2016/25682-5 and grants 2013/07375-0). L.F.C thanks CNPq (307333/2013-2), NAP-PRP-USP and FAPESP (11/50761-2) for the financial support. This researched is also supported by FAPESP (grant 2015/50122-0) and DFG-GRTK (grant 1740/2).

^[1] D. Kempe, J. Kleinberg, and É. Tardos, in *Proceedings* of the 9th ACM SIGKDD international conference on

- [2] D. Kempe, J. Kleinberg, and E. Tardos, Theory of Computing 11, 105 (2015).
- [3] R. Pastor-Satorras, C. Castellano, P. Van Mieghem, and A. Vespignani, Reviews of Modern Physics 87, 925 (2015).
- [4] G. F. de Arruda, A. L. Barbieri, P. M. Rodríguez, F. A. Rodrigues, Y. Moreno, and L. da F. Costa, Physical Review E 90, 032812 (2014).
- [5] T. Cao, X. Wu, S. Wang, and X. Hu, Expert Systems with Applications 38, 13128 (2011).
- [6] F. Morone and H. A. Makse, Nature (2015).
- [7] F. Radicchi and C. Castellano, Physical Review E 95, 012318 (2017).
- [8] I. Z. Kiss, D. M. Green, and R. R. Kao, Journal of The Royal Society Interface 5, 791 (2008).
- [9] M. Boguñá, R. Pastor-Satorras, and A. Vespignani, Physical Review Letters 90, 028701 (2003).
- [10] M. Bertotti, J. Brunner, and G. Modanese, Chaos, Solitons & Fractals 90, 55 (2016).
- [11] S. Fortunato and D. Hric, Physics Reports 659, 1 (2016).
- [12] A. Galstyan, V. Musoyan, and P. Cohen, Physical Review E 79, 056102 (2009).
- [13] Y. Wang, G. Cong, G. Song, and K. Xie, in Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '10 (ACM Press, New York, New York, USA, 2010) p. 1039.
- [14] L. Weng, F. Menczer, and Y.-Y. Ahn, Scientific Reports 3, 2522 (2013).
- [15] X. Zhang, J. Zhu, Q. Wang, and H. Zhao, Knowledge-Based Systems 42, 74 (2013).
- [16] M. Hosseini-Pozveh, K. Zamanifar, and A. R. Naghsh-Nilchi, Journal of Information Science, 0165551515621005 (2016).
- [17] D. J. Daley and D. G. Kendall, Nature **204**, 1118 (1964).
- [18] D. P. Maki and M. Thompson, Mathematical Models and Applications, with Emphasis on the Social, Life, and Management Sciences (Prentice-Hall, 1973).
- [19] D. H. Zanette, Phys. Rev. E 64, 050901 (2001).
- [20] Y. Moreno, M. Nekovee, and A. F. Pacheco, Physical Review E 69, 066130 (2004).
- [21] M. Nekovee, Y. Moreno, G. Bianconi, and M. Marsili, Physica A: Statistical Mechanics and its Applications 374, 457 (2007).
- [22] J. Borge-Holthoefer and Y. Moreno, Physical Review E 85, 026116 (2012).
- [23] M. Newman, Networks: an introduction (Oxford University Press, Inc., 2010).
- [24] M. Newman, Physical Review Letters 89, 208701 (2002).
- [25] T. K. D. Peron, P. Ji, F. A. Rodrigues, and J. Kurths, Physical Review E 91, 052805 (2015).
- [26] C. Castellano, S. Fortunato, and V. Loreto, Reviews of Modern Physics 81, 591 (2009).
- [27] D. A. Vega-Oliveros, L. da F Costa, and F. A. Rodrigues, Journal of Statistical Mechanics: Theory and Experiment 2017, 023401 (2017).
- [28] S. Stieglitz and L. Dang-Xuan, J. of Management Information Systems 29, 217 (2013).

- [29] M. Kitsak, L. K. Gallos, S. Havlin, F. Liljeros, L. Muchnik, H. E. Stanley, and H. A. Makse, Nature Physics 6, 888 (2010).
- [30] K. Kandhway and J. Kuri, IEEE Transactions on Systems, Man, and Cybernetics: Systems PP, 1 (2016).
- [31] S. Huang, T. Lv, X. Zhang, Y. Yang, W. Zheng, and C. Wen, PloS one 9, e103733 (2014).
- [32] D. A. Vega-Oliveros and L. Berton, in SIMBig, Vol. 1478 (2015) pp. 73–82.
- [33] M. E. J. Newman, Physical Review E 69, 66133 (2004).
- [34] R. Xulvi-Brunet and I. M. Sokolov, Acta Physica Polonica B 36, 1431 (2005).
- [35] J. McAuley and J. Leskovec, in Advances in Neural Information Processing Systems (2012) pp. 548–556.
- [36] M. E. J. Newman, "Personal Website Network Data: http://www-personal.umich.edu/~mejn/netdata/," (2013).
- [37] J. Leskovec, J. Kleinberg, and C. Faloutsos, in Proceedings of the Eleventh ACM SIGKDD International Conference on Knowledge Discovery in Data Mining, KDD '05 (ACM, New York, NY, USA, 2005) pp. 177–187.
- [38] P. Massa, M. Salvetti, and D. Tomasoni, in Proc, Int. Conf. Dependable, Automatic and Secure Computing (2009) pp. 658–663.
- [39] R. Guimera, L. Danon, A. Diaz-Guilera, F. Giralt, and A. Arenas, Physical Review E 68, 2003 (2003).
- [40] J. Kunegis, "Hamsterster full network dataset KONECT," (2014).
- [41] M. Boguñá, R. Pastor-Satorras, A. Díaz-Guilera, and A. Arenas, Phys. Rev. E 70, 056122 (2004).
- [42] M. E. J. Newman, in *Natl. Acad. Sci. USA*, 98 (2001) pp. 404 409.
- [43] N. X. Vinh, J. Epps, and J. Bailey, J. Mach. Learn. Res. 11, 2837 (2010).
- [44] T. M. Cover and J. A. Thomas, Elements of Information Theory (Wiley-Interscience, 2006).
- [45] J. Demšar, J. Mach. Learn. Res. 7, 1 (2006).
- [46] L. d. F. Costa and F. A. Rodrigues, EPL (Europhysics Letters) 85, 48001 (2009).
- [47] L. d. F. Costa, F. A. Rodrigues, C. C. Hilgetag, and M. Kaiser, EPL (Europhysics Letters) 87, 18008 (2009).
- [48] A. Barrat, M. Barthelemy, R. Pastor-Satorras, and A. Vespignani, Proceedings of the National Academy of Sciences of the United States of America 101, 3747 (2004).
- [49] M. Kivelä, A. Arenas, M. Barthelemy, J. P. Gleeson, Y. Moreno, and M. A. Porter, Journal of Complex Networks 2, 203 (2014).
- [50] G. F. de Arruda, E. Cozzo, T. P. Peixoto, F. A. Rodrigues, and Y. Moreno, Physical Review X 7, 011014 (2017).
- [51] P. Holme and J. Saramäki, Physics Reports 519, 97 (2012).