

# A Probabilistic Framework for Structural Analysis in Directed Networks

Cheng-Shang Chang, Duan-Shin Lee, Li-Heng Liou, Sheng-Min Lu, and Mu-Huan Wu  
Institute of Communications Engineering, National Tsing Hua University  
Hsinchu 30013, Taiwan, R.O.C.

Email: cschang@ee.nthu.edu.tw; lds@cs.nthu.edu.tw; dacapo1142@gmail.com;  
s103064515@m103.nthu.edu.tw; u9661106@oz.nthu.edu.tw;

**Abstract**—In our recent works, we developed a probabilistic framework for structural analysis in *undirected networks*. The key idea of that framework is to sample a network by a *symmetric* bivariate distribution and then use that bivariate distribution to formerly define various notions, including centrality, relative centrality, community, and modularity. The main objective of this paper is to extend the probabilistic framework to *directed networks*, where the sampling bivariate distributions could be *asymmetric*. Our main finding is that we can relax the assumption from *symmetric* bivariate distributions to bivariate distributions that have the same marginal distributions. By using such a weaker assumption, we show that various notions for structural analysis in directed networks can also be defined in the same manner as before. However, since the bivariate distribution could be *asymmetric*, the community detection algorithms proposed in our previous work cannot be directly applied. For this, we show that one can construct another sampled graph with a *symmetric* bivariate distribution so that for any partition of the network, the modularity index remains the same as that of the original sampled graph. Based on this, we propose a hierarchical agglomerative algorithm that returns a partition of communities when the algorithm converges.

**keywords:** centrality, community, modularity, PageRank

## I. INTRODUCTION

As the advent of on-line social networks, structural analysis of networks has been a very hot research topic. There are various notions that are widely used for structural analysis of networks, including centrality, relative centrality, similarity, community, modularity, and homophily (see e.g., Chapters 7 and 11 of the book by Newman [1]). In order to make these notions more mathematically precise, we developed in [2], [3] a probabilistic framework for structural analysis of *undirected networks*. The key idea of the framework is to “sample” a network to generate a bivariate distribution  $p(v, w)$  that specifies the probability that a pair of two nodes  $v$  and  $w$  are selected from a sample. The bivariate distribution  $p(v, w)$  can be viewed as a normalized *similarity* measure [5] between the two nodes  $v$  and  $w$ . A graph  $G$  associated with a bivariate distribution  $p(\cdot, \cdot)$  is then called a *sampled graph*.

In [2], [3], the bivariate distribution is assumed to be *symmetric*. Under this assumption, the two marginal distributions of the bivariate distribution, denoted by  $p_V(\cdot)$  and  $p_W(\cdot)$ , are the same and they represent the probability that a particular node is selected in the sampled graph. As such, the marginal

distribution  $p_V(v)$  can be used for defining the *centrality* of a node  $v$  as it represents the probability that node  $v$  is selected. The relative centrality of a set of nodes  $S_1$  with respect to another set of nodes  $S_2$  is then defined as the *conditional* probability that one node of the selected pair of two nodes is in the set  $S_1$  given that the other node is in the set  $S_2$ . Based on the probabilistic definitions of centrality and relative centrality in the framework, the *community strength* for a set of nodes  $S$  is defined as the difference between its relative centrality with respect to itself and its centrality. Moreover, a set of nodes with a *nonnegative* community strength is called a *community*. In the probabilistic framework, the *modularity* for a partition of a sampled graph is defined as the average community strength of the community. As such, a high modularity for a partition of a graph implies that there are communities with strong community strengths. It was further shown in [3] that the Newman modularity in [6] and the stability in [7], [8] are special cases of the modularity for certain sampled graphs.

The main objective of this paper is to extend the probabilistic framework in [2], [3] to *directed networks*, where the sampling bivariate distributions could be *asymmetric*. Our main finding is that we can relax the assumption from *symmetric* bivariate distributions to bivariate distributions that have the same marginal distributions. By using such a weaker assumption, we show that the notions of centrality, relative centrality, community and modularity can be defined in the same manner as before. Moreover, the equivalent characterizations of a community still hold. Since the bivariate distribution could be *asymmetric*, the agglomerative community detection algorithms in [2], [3] cannot be directly applied. For this, we show that one can construct another sampled graph with a *symmetric* bivariate distribution so that for any partition of the network, the modularity index remains the same as that of the original sampled graph. Based on this, we propose a hierarchical agglomerative algorithm that returns a partition of communities when the algorithm converges.

In this paper, we also address two methods for sampling a directed network with a bivariate distribution that has the same marginal distributions : (i) PageRank and (ii) random walks with self loops and backward jumps. Experiments show that sampling by a random walk with self loops and backward jumps performs better than that by PageRank for community

detection. This might be due to the fact that PageRank adds weak links in a network and that changes the topology of the network and thus affects the results of community detection.

## II. SAMPLING NETWORKS BY BIVARIATE DISTRIBUTIONS WITH THE SAME MARGINAL DISTRIBUTIONS

In [3], a probabilistic framework for network analysis for undirected networks was proposed. The main idea in that framework is to characterize a network by a *sampld graph*. Specifically, suppose a network is modelled by a graph  $G(V_g, E_g)$ , where  $V_g$  denotes the set of vertices (nodes) in the graph and  $E_g$  denotes the set of edges (links) in the graph. Let  $n = |V_g|$  be the number of vertices in the graph and index the  $n$  vertices from  $1, 2, \dots, n$ . Also, let  $A = (a_{ij})$  be the  $n \times n$  adjacency matrix of the graph, i.e.,

$$a_{ij} = \begin{cases} 1, & \text{if there is an edge from vertex } i \text{ to vertex } j, \\ 0, & \text{otherwise.} \end{cases}$$

A sampling bivariate distribution  $p(\cdot, \cdot)$  for a graph  $G$  is the bivariate distribution that is used for *sampling* a network by randomly selecting an ordered pair of two nodes  $(V, W)$ , i.e.,

$$P(V = v, W = w) = p(v, w). \quad (1)$$

Let  $p_V(v)$  (resp.  $p_W(w)$ ) be the marginal distribution of the random variable  $V$  (resp.  $W$ ), i.e.,

$$p_V(v) = P(V = v) = \sum_{w=1}^n p(v, w), \quad (2)$$

and

$$p_W(w) = P(W = w) = \sum_{v=1}^n p(v, w). \quad (3)$$

**Definition 1: (Sampled graph)** A graph  $G(V_g, E_g)$  that is sampled by randomly selecting an ordered pair of two nodes  $(V, W)$  according to a specific bivariate distribution  $p(\cdot, \cdot)$  in (1) is called a *sampld graph* and it is denoted by the two tuple  $(G(V_g, E_g), p(\cdot, \cdot))$ .

For a given graph  $G(V_g, E_g)$ , there are many methods to generate sampled graphs by specifying the needed bivariate distributions. In [3], the bivariate distributions are all assumed to be *symmetric* and that limits its applicability to *undirected* networks. One of the main objectives of this paper is to relax the *symmetric* assumption for the bivariate distribution so that the framework can be applied to *directed* networks. The key idea of doing this is to assume that the bivariate distribution has the same marginal distributions, i.e.,

$$p_V(v) = p_W(v), \quad \text{for all } v. \quad (4)$$

Note that a symmetric bivariate distribution has the same marginal distributions and thus the assumption in (4) is much more general.

### A. PageRank

One approach for sampling a network with a bivariate distribution that has the same marginal distributions is to sample a network by an *ergodic Markov chain*. From the Markov chain theory (see e.g., [9]), it is well-known that an ergodic Markov chain converges to its steady state in the long run. Hence, the joint distribution of two successive steps of a *stationary and ergodic* Markov chain can be used as the needed bivariate distribution. Specifically, suppose that a network  $G(V_g, E_g)$  is sampled by a stationary and ergodic Markov chain  $\{X(t), t \geq 0\}$  with the state space  $\{1, 2, \dots, n\}$  being the  $n$  nodes in  $V_g$ . Let  $P = (p_{ij})$  be the  $n \times n$  transition probability matrix and  $\pi = (\pi_1, \pi_2, \dots, \pi_n)$  be the steady state probability vector of the stationary and ergodic Markov chain. Then we can choose the bivariate distribution

$$\begin{aligned} P(V = v, W = w) &= p(v, w) \\ &= P(X(t) = v, X(t+1) = w). \end{aligned} \quad (5)$$

As the Markov chain is stationary, we have

$$P(X(t) = v) = P(X(t+1) = w) = p_V(v) = p_W(w). \quad (6)$$

It is well-known that a random walk on the graph induces a Markov chain with the state transition probability matrix  $P = (p_{ij})$  with

$$p_{ij} = \frac{a_{ij}}{k_i^{\text{out}}}, \quad (7)$$

where

$$k_i^{\text{out}} = \sum_{j=1}^n a_{ij}, \quad (8)$$

is the number of outgoing edges from vertex  $i$ . In particular, if the graph is an undirected graph, i.e.,  $a_{ij} = a_{ji}$ , then the induced Markov chain is reversible and the steady state probability of state  $i$ , i.e.,  $\pi_i$ , is  $k_i/2m$ , where  $m = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n a_{ij}$  is the total number of edges of the undirected graph.

One problem for sampling a *directed* network by a simple random walk is that the induced Markov chain may not be ergodic even when the network itself is weakly connected. One genuine solution for this is to allow random jumps from states to states in a random walk. PageRank [10], proposed by Google, is one such example that has been successfully used for ranking web pages. The key idea behind PageRank is to model the behavior of a web surfer by a random walk (the random surfer model) and then use that to compute the steady state probability for a web surfer to visit a specific web page. Specifically, suppose that there are  $n$  web pages and a web surfer uniformly selects a web page with probability  $1/n$ . Once he is on a web page, he continues web surfing with probability  $\lambda$ . This is done by selecting *uniformly* one of the hyperlinks in that web page. On the other hand, with probability  $1 - \lambda$  he starts a new web page *uniformly* among all the  $n$  web pages. The transition probability from state  $i$  to state  $j$  for the induced Markov chain is then

$$p_{ij} = (1 - \lambda) \frac{1}{n} + \lambda \frac{a_{ij}}{k_i^{\text{out}}}, \quad (9)$$

where  $a_{ij} = 1$  if there is a hyperlink pointing from the  $i^{th}$  web page to the  $j^{th}$  web page and  $k_i^{out} = \sum_{j=1}^n a_{ij}$  is the total number of hyperlinks on the  $i^{th}$  web page. Let  $\pi_i$  be steady probability of visiting the  $i^{th}$  web page by the web surfer. It then follows that

$$\pi_i = (1 - \lambda) \frac{1}{n} + \lambda \sum_{j=1}^n \frac{a_{ji}}{k_j^{out}} \pi_j. \quad (10)$$

PageRank then uses  $\pi_i$  as the centrality of the  $i^{th}$  web page and rank web pages by their centralities. Unlike the random walk on an undirected graph, the steady state probabilities in (10) cannot be explicitly solved and it requires a lot of computation to solve the system of linear equations.

The sampled graph  $(G(V_g, E_g), p(\cdot, \cdot))$  by using PageRank then has the following bivariate distribution

$$p(v, w) = \pi_v p_{vw}, \quad (11)$$

where  $p_{vw}$  is defined in (9) and  $\pi_v$  is the solution of (10).

### B. Random walks with self loops and backward jumps

Another way to look at the Markov chain induced by PageRank in (9) is that it is in fact a random walk on a different graph with the adjacency matrix  $\tilde{A}$  that is constructed from the original graph with additional edge weights, i.e.,

$$\tilde{A} = (1 - \lambda) \frac{1}{n} \mathbf{1} + \lambda D^{-1} A, \quad (12)$$

where  $\mathbf{1}$  is an  $n \times n$  matrix with all its elements being 1 and  $D = (d_{ij})$  is the diagonal matrix with  $d_{ii} = k_i^{out}$  for all  $i = 1, 2, \dots, n$ .

In view of (12), another solution for the ergodic problem is to consider a random walk on the graph with the adjacency matrix

$$\hat{A} = \lambda_0 \mathbf{I} + \lambda_1 A + \lambda_2 A^T, \quad (13)$$

where  $\mathbf{I}$  is the  $n \times n$  identity matrix and  $A^T$  is the transpose matrix of  $A$ . The three parameters  $\lambda_0, \lambda_1, \lambda_2$  are positive and

$$\lambda_0 + \lambda_1 + \lambda_2 = 1.$$

A random walk on the graph with the adjacency matrix  $\hat{A}$  induces an ergodic Markov chain if the original graph is weakly connected. Also, with the additional edges from the identity matrix and the transpose matrix, such a random walk can be viewed as a random walk on the original graph with self loops and backward jumps.

## III. THE FRAMEWORK FOR DIRECTED NETWORKS

### A. Centrality and relative centrality

Centrality [11], [12], [1] is usually used as a measure for ranking the importance of a set of nodes in a (social) network. Under the assumption in (4), such a concept can be directly mapped to the probability that a node is selected as in [3].

**Definition 2: (Centrality)** For a sampled graph  $(G(V_g, E_g), p(\cdot, \cdot))$  with the bivariate distribution  $p(\cdot, \cdot)$  that has the same marginal distributions in (4), the *centrality*

of a set of nodes  $S$ , denoted by  $C(S)$ , is defined as the probability that a node in  $S$  is selected, i.e.,

$$C(S) = P(V \in S) = P(W \in S). \quad (14)$$

As a generalization of centrality, relative centrality in [3] is a (probability) measure that measures how important a set of nodes in a network is with respect to another set of nodes.

**Definition 3: (Relative centrality)** For a sampled graph  $(G(V_g, E_g), p(\cdot, \cdot))$  with the bivariate distribution  $p(\cdot, \cdot)$  that has the same marginal distributions in (4), the *relative centrality* of a set of nodes  $S_1$  with respect to another set of nodes  $S_2$ , denoted by  $C(S_1|S_2)$ , is defined as the conditional probability that the randomly selected node  $W$  is inside  $S_1$  given that the random selected node  $V$  is inside  $S_2$ , i.e.,

$$C(S_1|S_2) = P(W \in S_1|V \in S_2) \quad (15)$$

We note that if we choose  $S_2 = V_g$ , then the relative centrality of a set of nodes  $S_1$  with respect to  $V_g$  is simply the *centrality* of the set of nodes  $S_1$ .

**Example 4: (Relative PageRank)** PageRank described in Section II-A has been commonly used for ranking the importance of nodes in a directed network. Here we can use Definition 3 to define relative PageRank that can be used for ranking the relative importance of a set of nodes to another set of nodes in a directed network. Specifically, let  $\pi$  be the PageRank for node  $i$  in (10) and  $p_{i,j}$  be the transition probability from state  $i$  to state  $j$  for the induced Markov chain in (9). Then the relative PageRank of a set  $S_1$  with respect to another set  $S_2$  is

$$\begin{aligned} C(S_1|S_2) &= P(W \in S_1|V \in S_2) \\ &= \frac{P(W \in S_1, V \in S_2)}{P(V \in S_2)} = \frac{\sum_{i \in S_2} \sum_{j \in S_1} \pi_i p_{ij}}{\sum_{i \in S_2} \pi_i}. \end{aligned} \quad (16)$$

Analogous to the relative centrality in [3], there are also several properties of relative centrality in Definition 3. However, the reciprocity property in Proposition 5(iv) is much weaker than that in [3]. The proof of Proposition 5 is omitted due to space limitation.

**Proposition 5:** For a sampled graph  $(G(V_g, E_g), p(\cdot, \cdot))$  with the bivariate distribution  $p(\cdot, \cdot)$  that has the same marginal distributions in (4), the following properties for the relative centrality defined in Definition 3 hold.

- (i)  $0 \leq C(S_1|S_2) \leq 1$  and  $0 \leq C(S_1) \leq 1$ . Moreover,  $C(V_g|S_2) = 1$  and  $C(V_g) = 1$ .
- (ii) (Additivity) If  $S_1$  and  $S_2$  are two disjoint sets, i.e.,  $S_1 \cap S_2$  is an empty set, then for an arbitrary set  $S_3$ ,

$$C(S_1 \cup S_2|S_3) = C(S_1|S_3) + C(S_2|S_3). \quad (17)$$

In particular, when  $S_3 = \{1, 2, \dots, n\}$ , we have

$$C(S_1 \cup S_2) = C(S_1) + C(S_2). \quad (18)$$

- (iii) (Monotonicity) If  $S_1$  is a subset of  $S'_1$ , i.e.,  $S_1 \subset S'_1$ , then  $C(S_1|S_2) \leq C(S'_1|S_2)$  and  $C(S_1) \leq C(S'_1)$ .

(iv) (Reciprocity) Let  $S^c = V_g \setminus S$  be the set of nodes that are not in  $S$ .

$$C(S)C(S^c|S) = C(S^c)C(S|S^c).$$

### B. Community strength and communities

The notions of community strength and modularity in [3] generalizes the original Newman's definition [13] and unifies various other generalizations, including the stability in [7], [8]. In this section, we further extend these notions to directed networks.

**Definition 6: (Community strength and communities)** For a sample graph  $(G(V_g, E_g), p(\cdot, \cdot))$  with a bivariate distribution  $p(\cdot, \cdot)$  that has the same marginal distributions in (4), the *community strength* of a subset set of nodes  $S \subset V_g$ , denoted by  $Str(S)$ , is defined as the difference of the relative centrality of  $S$  with respect to itself and its centrality, i.e.,

$$Str(S) = C(S|S) - C(S). \quad (19)$$

In particular, if a subset of nodes  $S \subset V_g$  has a nonnegative community strength, i.e.,  $Str(S) \geq 0$ , then it is called a *community*.

In the following theorem, we show various equivalent statements for a set of nodes to be a community. The proof of Theorem 7 is omitted due to space limitation.

**Theorem 7:** Consider a sample graph  $(G(V_g, E_g), p(\cdot, \cdot))$  with a bivariate distribution  $p(\cdot, \cdot)$  that has the same marginal distributions in (4), and a set  $S$  with  $0 < C(S) < 1$ . Let  $S^c = V_g \setminus S$  be the set of nodes that are not in  $S$ . The following statements are equivalent.

- (i) The set  $S$  is a community, i.e.,  $Str(S) = C(S|S) - C(S) \geq 0$ .
- (ii) The relative centrality of  $S$  with respect to  $S$  is not less than the relative centrality of  $S$  with respect to  $S^c$ , i.e.,  $C(S|S) \geq C(S|S^c)$ .
- (iii) The relative centrality of  $S^c$  with respect to  $S$  is not greater than the centrality of  $S^c$ , i.e.,  $C(S^c|S) \leq C(S^c)$ .
- (iv) The relative centrality of  $S$  with respect to  $S^c$  is not greater than the centrality of  $S$ , i.e.,  $C(S|S^c) \leq C(S)$ .
- (v) The set  $S^c$  is a community, i.e.,  $Str(S^c) = C(S^c|S^c) - C(S^c) \geq 0$ .
- (vi) The relative centrality of  $S^c$  with respect to  $S^c$  is not less than the relative centrality of  $S^c$  with respect to  $S$ , i.e.,  $C(S^c|S^c) \geq C(S^c|S)$ .

### C. Modularity and community detection

As in [3], we define the modularity index for a partition of a network as the average community strength of a randomly selected node in Definition 8.

**Definition 8: (Modularity)** Consider a sampled graph  $(G(V_g, E_g), p(\cdot, \cdot))$  with a bivariate distribution  $p(\cdot, \cdot)$  that has the same marginal distributions in (4). Let  $\mathcal{P} = \{S_c, c = 1, 2, \dots, C\}$ , be a partition of  $\{1, 2, \dots, n\}$ , i.e.,  $S_c \cap S_{c'} = \emptyset$  for  $c \neq c'$  and  $\cup_{c=1}^C S_c = \{1, 2, \dots, n\}$ . The modularity index  $Q(\mathcal{P})$  with respect to the partition  $S_c$ ,

$c = 1, 2, \dots, C$ , is defined as the weighted average of the community strength of each subset with the weight being the centrality of each subset, i.e.,

$$Q(\mathcal{P}) = \sum_{c=1}^C C(S_c) \cdot Str(S_c). \quad (20)$$

We note the modularity index in (20) can also be written as follows:

$$\begin{aligned} Q(\mathcal{P}) &= \sum_{c=1}^C P(V \in S_c, W \in S_c) - P(V \in S_c)P(W \in S_c) \\ &= \sum_{c=1}^C \sum_{v \in S_c} \sum_{w \in S_c} (p(v, w) - p_V(v)p_W(w)). \end{aligned} \quad (21)$$

As the modularity index for a partition of a network is the average community strength of a randomly selected node, a good partition of a network should have a large modularity index. In view of this, one can then tackle the community detection problem by looking for algorithms that yield large values of the modularity index. For sampled graphs with *symmetric* bivariate distributions, there are already various community detection algorithms in [2], [3] that find local maxima of the modularity index. However, they cannot be directly applied as the bivariate distributions for sampling directed networks could be *asymmetric*. For this, we show in the following lemma that one can construct another sampled graph with a *symmetric* bivariate distribution so that for any partition of the network, the modularity index remains the same as that of the original sampled graph. The proof of Lemma 9 is omitted due to space limitation.

**Lemma 9:** Consider a sampled graph  $(G(V_g, E_g), p(\cdot, \cdot))$  with a bivariate distribution  $p(\cdot, \cdot)$  that has the same marginal distributions in (4). Construct the sampled graph  $(G(V_g, E_g), \tilde{p}(\cdot, \cdot))$  with the symmetric bivariate distribution

$$\tilde{p}(v, w) = \frac{p(v, w) + p(w, v)}{2}. \quad (22)$$

Let  $Q(\mathcal{P})$  (resp.  $\tilde{Q}(\mathcal{P})$ ) be the modularity index for the partition  $\mathcal{P} = \{S_c, c = 1, 2, \dots, C\}$  of the sampled graph  $(G(V_g, E_g), p(\cdot, \cdot))$  (resp. the sampled graph  $(G(V_g, E_g), \tilde{p}(\cdot, \cdot))$ ). Then

$$\tilde{Q}(\mathcal{P}) = Q(\mathcal{P}). \quad (23)$$

As  $\tilde{Q}(\mathcal{P}) = Q(\mathcal{P})$ , one can then use the community detection algorithms for the sampled graph  $(G(V_g, E_g), \tilde{p}(\cdot, \cdot))$  with the symmetric bivariate distribution to solve the community detection problem for the original sampled graph  $(G(V_g, E_g), p(\cdot, \cdot))$ . Analogous to the hierarchical agglomerative algorithms in [13], [14], in the following we propose a hierarchical agglomerative algorithm for community detection in directed networks. The idea behind this algorithm is

*modularity maximization.* For this, we define the correlation measure between two nodes  $v$  and  $w$  as follows:

$$q(v, w) = \tilde{p}(v, w) - \tilde{p}_V(v)\tilde{p}_W(w) \\ = \frac{p(v, w) - p_V(v)p_W(w) + p(w, v) - p_V(w)p_W(v)}{2}. \quad (24)$$

For any two sets  $S_1$  and  $S_2$ , define the correlation measure between these two sets as

$$q(S_1, S_2) = \sum_{v \in S_1} \sum_{w \in S_2} q(v, w). \quad (25)$$

Also, define the *average* correlation measure between two sets  $S_1$  and  $S_2$  as

$$\bar{q}(S_1, S_2) = \frac{1}{|S_1| \cdot |S_2|} q(S_1, S_2). \quad (26)$$

With this correlation measure, we have from Lemma 9, (21) and (25) that the modularity index for the partition  $\mathcal{P} = \{S_c, c = 1, 2, \dots, C\}$  is

$$Q(\mathcal{P}) = \tilde{Q}(\mathcal{P}) = \sum_{c=1}^C q(S_c, S_c), \quad (27)$$

Moreover, a set  $S$  is a community if and only if  $q(S, S) \geq 0$ .

**Algorithm 1: a hierarchical agglomerative algorithm for community detection in a directed network**

(P0) Input a sampled graph  $(G(V_g, E_g), p(\cdot, \cdot))$  with a bivariate distribution  $p(\cdot, \cdot)$  that has the same marginal distributions in (4).

(P1) Initially, there are  $n$  sets, indexed from 1 to  $n$ , with each set containing exactly one node. Specifically, let  $S_i$  be the set of nodes in set  $i$ . Then  $S_i = \{i\}$ ,  $i = 1, 2, \dots, n$ .

(P2) For all  $i, j = 1, 2, \dots, n$ , compute the correlation measures  $q(S_i, S_j) = q(\{i\}, \{j\})$  from (24).

(P3) If there is only one set left or there do not exist nonnegative correlation measures between two distinct sets, i.e.,  $q(S_i, S_j) < 0$  for all  $i \neq j$ , then the algorithm outputs the current sets.

(P4) Find two sets that have a nonnegative correlation measure. Merge these two sets into a new set. Suppose that set  $i$  and set  $j$  are grouped into a new set  $k$ . Then  $S_k = S_i \cup S_j$  and update

$$q(S_k, S_k) = q(S_i, S_i) + 2q(S_i, S_j) + q(S_j, S_j). \quad (28)$$

Moreover, for all  $\ell \neq k$ , update

$$q(S_k, S_\ell) = q(S_\ell, S_k) = q(S_i, S_\ell) + q(S_j, S_\ell). \quad (29)$$

(P5) Repeat from (P3).

The hierarchical agglomerative algorithm in Algorithm 1 has the following properties.

*Theorem 10:*

- (i) For the hierarchical agglomerative algorithm in Algorithm 1, the modularity index is non-decreasing in every iteration and thus converges to a local optimum.

- (ii) When the algorithm converges, every set returned by the hierarchical agglomerative algorithm is indeed a *community*.
- (iii) If, furthermore, we use the *greedy* selection that selects the two sets with the *largest* average correlation measure to merge in (P4) of Algorithm 1, then the average correlation measure of the two selected sets in each merge operation is non-increasing.

The proof of Theorem 10 is given in Appendix A. For (i) and (ii) of Theorem 10, it is not necessary to specify how we select a pair of two sets with a nonnegative correlation. One advantage of using the greedy selection in (iii) of Theorem 10 is the *monotonicity* property for the dendrogram produced by a greedy hierarchical agglomerative algorithm (see [15], Chapter 13.2.3). With such a monotonicity property, there is no *crossover* in the produced dendrogram.

#### IV. EXPERIMENTAL RESULTS

In this section, we compare the sampling methods by PageRank in Section II-A and random walks with self loops and backward jumps in Section II-B for community detection. We conduct various experiments based on the stochastic block model with two blocks. The stochastic block model, as a generalization of the Erdos-Renyi random graph, is a commonly used method for generating random graphs that can be used for benchmarking community detection algorithms. In a stochastic block model with two blocks (communities), the total number of nodes in the random graph are evenly distributed to these two blocks. The probability that there is an edge between two nodes within the same block is  $p_{in}$  and the probability that there is an edge between two nodes in two different blocks is  $p_{out}$ . These edges are generated independently. Let  $c_{in} = np_{in}$  and  $c_{out} = np_{out}$ .

In our experiments, the number of nodes  $n$  in the stochastic block model is 200 with 100 nodes in each of these two blocks. The average degree of a node is set to be 3. The values of  $c_{in} - c_{out}$  of these graphs are in the range from 2.5 to 5.9 with a common step of 0.1. We generate 100 graphs for each  $c_{in} - c_{out}$ . Isolated vertices are removed. Thus, the exact numbers of vertices used in this experiment might be slightly less than 200. For PageRank, the parameter  $\lambda$  is chosen to be 0.9. For the random walk with self loops and backward jumps, the three parameters are  $\lambda_0 = 0.05$ ,  $\lambda_1 = 0.85$  and  $\lambda_2 = 0.1$ . We run the greedy hierarchical agglomerative algorithm in Algorithm 1 until there are only two sets (even when there do not exist nonnegative correlation measures between two distinct sets). We then evaluate the overlap with the true labeling. In Figure 1, we show the experimental results, where each point is averaged over 100 random graphs from the stochastic block model. The error bars are the 95% confidence intervals. From Figure 1, one can see that the performance of random walks with self loops and backward jumps is better than that of PageRank. One reason for this is that PageRank uniformly adds an edge (with a small weight) between any two nodes and these added edges change the network topology. On the other hand, mapping by a random walk with backward

jumps in (13) does not change the network topology when it is viewed as an undirected network.

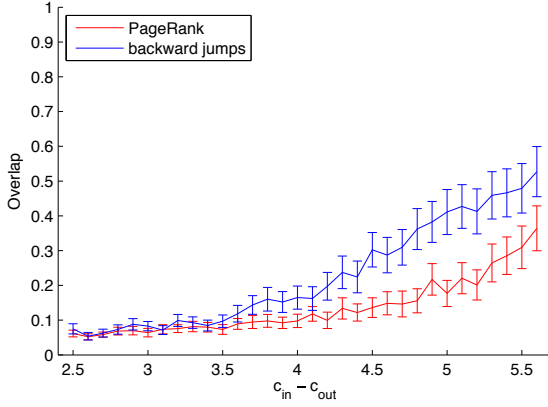


Fig. 1. Community detection of the stochastic block model by using PageRank in (12) and a random walk with self loops and backward jumps in (13).

## V. CONCLUSION

In this paper we extended our previous work in [2], [3] to *directed networks*. Our approach is to introduce bivariate distributions that have the same marginal distributions. By doing so, we were able to extend the notions of centrality, relative centrality, community strength, community and modularity to *directed networks*. For community detection, we propose a hierarchical agglomerative algorithm that guarantees every set returned from the algorithm is a community. We also tested the algorithm by using PageRank and random walks with self loops and backward jumps. The experimental results show that sampling by random walks with self loops and backward jumps perform better than sampling by PageRank for community detection. Further extensions and comparisons with existing clustering (community detection) algorithms in the literature are addressed in [4].

## REFERENCES

- [1] M. Newman, *Networks: an introduction*. OUP Oxford, 2009.
- [2] C.-S. Chang, C.-Y. Hsu, J. Cheng, and D.-S. Lee, "A general probabilistic framework for detecting community structure in networks," in *IEEE INFOCOM '11*, April 2011.
- [3] C.-S. Chang, C.-J. Chang, W.-T. Hsieh, D.-S. Lee, L.-H. Liou, and W. Liao, "Relative centrality and local community detection," *Network Science*, vol. FirstView, pp. 1–35, 9 2015.
- [4] C.-S. Chang, W. Liao, Y.-S. Chen and L.-H. Liou, "A mathematical theory for clustering in metric spaces," to appear in *IEEE Transactions on Network Science and Engineering*.
- [5] D. Liben-Nowell and J. Kleinberg, "The link prediction problem for social networks," in *Proceedings of the twelfth international conference on Information and knowledge management*. ACM, 2003, pp. 556–559.
- [6] M. E. Newman and M. Girvan, "Finding and evaluating community structure in networks," *Physical review E*, vol. 69, no. 2, p. 026113, 2004.
- [7] R. Lambiotte, "Multi-scale modularity in complex networks," in *Modeling and Optimization in Mobile, Ad Hoc and Wireless Networks (WiOpt), 2010 Proceedings of the 8th International Symposium on*. IEEE, 2010, pp. 546–553.
- [8] J.-C. Delvenne, S. N. Yaliraki, and M. Barahona, "Stability of graph communities across time scales," *Proceedings of the National Academy of Sciences*, vol. 107, no. 29, pp. 12 755–12 760, 2010.

- [9] R. Nelson, *Probability, stochastic processes, and queueing theory: the mathematics of computer performance modeling*. Springer Verlag, 1995.
- [10] S. Brin and L. Page, "The anatomy of a large-scale hypertextual web search engine," *Computer networks and ISDN systems*, vol. 30, no. 1, pp. 107–117, 1998.
- [11] L. C. Freeman, "A set of measures of centrality based on betweenness," *Sociometry*, pp. 35–41, 1977.
- [12] —, "Centrality in social networks conceptual clarification," *Social networks*, vol. 1, no. 3, pp. 215–239, 1979.
- [13] M. E. Newman, "Fast algorithm for detecting community structure in networks," *Physical review E*, vol. 69, no. 6, p. 066133, 2004.
- [14] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2008, no. 10, p. P10008, 2008.
- [15] S. Theodoridis and K. Koutroumbas, *Pattern Recognition*. Elsevier Academic press, USA, 2006.

## APPENDIX A

In this section, we prove Theorem 10.

(i) Since we choose two sets that have a nonnegative correlation measure, i.e.,  $q(S_i, S_j) \geq 0$ , to merge, it is easy to see from (28) and (27) that the modularity index is non-decreasing in every iteration.

(ii) Suppose that there is only one set left. Then this set is  $V_g$  and it is the trivial community. On the other hand, suppose that there are  $C \geq 2$  sets  $\{S_1, S_2, \dots, S_C\}$  left when the algorithm converges. Then we know that  $q(S_i, S_j) < 0$  for  $i \neq j$ .

Note from (24) and (25) that for any node  $v$ ,

$$q(\{v\}, V_g) = \sum_{w \in V_g} q(v, w) = 0. \quad (30)$$

Thus,

$$q(S_i, V_g) = \sum_{v \in S_i} q(\{v\}, V_g) = 0. \quad (31)$$

Since  $\{S_1, S_2, \dots, S_C\}$  is a partition of  $V_g$ , it then follows that

$$0 = q(S_i, V_g) = q(S_i, S_i) + \sum_{j \neq i} q(S_i, S_j). \quad (32)$$

Since  $q(S_i, S_j) < 0$  for  $i \neq j$ , we conclude that  $q(S_i, S_i) > 0$  and thus  $S_i$  is a community.

(iii) Suppose that  $S_i$  and  $S_j$  are merged into the new set  $S_k$ . According to the update rules in the algorithm and the symmetric property of  $q(\cdot, \cdot)$ , we know that

$$\begin{aligned} q(S_k, S_\ell) &= q(S_\ell, S_k) = q(S_i, S_\ell) + q(S_j, S_\ell) \\ &= q(S_i, S_\ell) + q(S_\ell, S_j), \end{aligned}$$

for all  $\ell \neq k$ . Thus,

$$\bar{q}(S_k, S_\ell) = \frac{|S_i|}{|S_i| + |S_j|} \bar{q}(S_i, S_\ell) + \frac{|S_j|}{|S_i| + |S_j|} \bar{q}(S_\ell, S_j).$$

Since we select the two sets with the *largest* average correlation measure in each merge operation, we have  $\bar{q}(S_i, S_\ell) \leq \bar{q}(S_i, S_j)$  and  $\bar{q}(S_\ell, S_j) \leq \bar{q}(S_i, S_j)$ . These then lead to

$$\bar{q}(S_k, S_\ell) \leq \bar{q}(S_i, S_j).$$

Thus,  $\bar{q}(S_i, S_j)$  is not less than the average correlation measure between any two sets after the merge operation. As such, the average correlation measure at each merge is non-increasing.