# Analysis and Control of Beliefs in Social Networks

Tian Wang, Hamid Krim, *Fellow, IEEE,* Yannis Viniotis

*Abstract*—In this paper, we investigate the problem of how beliefs diffuse among members of social networks. We propose an *information flow model* (IFM) of belief that captures how interactions among members affect the diffusion and eventual convergence of a belief. The IFM model includes a generalized Markov Graph (GMG) model as a social network model, which reveals that the diffusion of beliefs depends heavily on two characteristics of the social network characteristics, namely degree centralities and clustering coefficients. We apply the IFM to both converged belief estimation and belief control strategy optimization. The model is compared with an IFM including the Barabási-Albert model, and is evaluated via experiments with published real social network data.

*Index Terms*—Complex Networks, Information Flow, Machine Learning

## I. INTRODUCTION

A Social network, as an abstract model of a social environment, consists of a set of nodes, which could be a set of individuals or a set of groups of individuals, and a set of relationships with specific characteristics among these nodes. There are numerous types of existing social networks in our daily lives; examples include on-line social networks such as Facebook or LinkedIn, email networks and alumni networks. Possibilities for forming new ones abound; for example, social networks among patients with a rare disease. Compared with non-social types of networks (e.g., a sensor network), social networks exhibit certain specific characteristics in quantities such as degree distributions, clustering coefficient distributions, etc. [2], [16], [18], which allow one to further analyse and control them.

Within a social network, each node has a certain belief representing its current status. The belief of a node may be influenced by other nodes connected to it, and could be changed in time. The belief may be, for example, an opinion regarding the quality of a restaurant, or the preference to attend a school [4]. In a different example, the belief may be the opinion of a patient regarding the curability of their illness. A node can propagate its belief to other connected nodes via the network links by performing certain activities, such as dining in a certain restaurant or attending a particular school, or simply by informing the other nodes. Such exposure and exchange of beliefs is actually a flow of information in a social network.

The beliefs in a social network may have value for members of the network or even outsiders. For example, knowing members' opinions about a certain product can help a product manufacturer better predict the market and form merchandising strategies. In a clinical study, a doctor may be interested in researching ways to influence patients' behaviour by "facilitating" interactions among patients. Consequently, prediction or even control of the beliefs in a social network can be

an important and interesting problem. As such, it requires a mathematical model to simulate and analyse the flow of beliefs in a social network.

We call this model an *Information Flow Model* (IFM). Briefly speaking, IFM includes two main parts: a description of how beliefs are updated in time, and a description of the network structure. The first part describes how information of a member's belief flows in the network and influences the other members. The second one describes the paths inside the network over which the beliefs can be transmitted. Details about both parts are given in Section II.

To date, research in IFMs has been mainly conducted on social learning [1], [11], [4], [7], which focuses on the first part. The effects of the second part have not received much attention, primarily due to the complexity of social network structure. Early studies on IFM had assumed acyclic networks [4], [9], [12]. More recent works adapt simple descriptions of social networks, including some limited connectivity properties [1], [11], [3], [20], [7]. Real social networks, however, have much more complicated structure, that is better captured by special network models. There are a lot of candidates of social entwork models, such as the Chung-Lu model [5], the Sznajd model [14], the Barabási-Albert (BA) model [2] ,and the Generalized Markov Graph (GMG) model [16]. The GMG models can reveal intrinsic statistical properties of a social network (e.g., connectivity patterns, measured via distributions of centralities) and thus help analyze the information flow model.

Such intrinsic properties can (at least in theory) be exploited when one studies prediction or control problems in social networks. The properties of the network could be effectively used to determine, for example, which people should be chosen to spread the information to others in order to minimize the number of such people, maximize the number of people with the desired behaviour in the network, which is novel compared to other IFM on social networks [8], [14]. In another example, in a medical study with a limited budget, simulations of the social network could be used to select the cohort of patients in the hope of an expedited and less costly study.

In this paper, we propose an IFM with enhanced (that is, BA or GMG) models for describing the social network properties. Our motivation for proposing the specific IFM is the desire to design strategies that can control beliefs effectively, without excessive overheads in computation time or memory requirements.

The contributions of this paper are the following. First, we propose an IFM adapted to a realistic social network without excessive overheads. Second, we develop three methods to analyse beliefs in a social network; the methods can be tuned to trade off accuracy for overhead. Third, we develop strategies to control beliefs. Partial results of the paper also appeared in

a brief conference paper [17]. In addition to the contents in the conference version, experiments to verify basic assumptions of models, techniques to estimate the converged beliefs as well as proofs for all the theorems and detailed description of models are added in this paper.

The paper is organized as follows. In Section II, we define the concepts of belief, control strategy and social network, and introduce the IFM formally. We introduce the notion of "control power" as a metric that can be used to compare control strategies. In Section III, we elaborate on the BA and GMG models and show how one can analytically calculate the respective control power metric. In order to verify the proposed model, we use real network data in Section IV, to test the fundamental properties of the model, and specific control strategies. We finally conclude with some remarks and future planned work in Section V.

## II. INFORMATION FLOW MODEL

### A. Basic concepts in IFM

The basic elements that comprise an IFM are: the social network, the belief and the control strategy.

*1) Social Network:* For a social network $G$ with $N$ nodes, we use indices $i \in \{1, 2, ..., N\}$ to represent the nodes; the set of nodes is defined as: $node_G = \{1, 2, ..., N\}$. The set of edges, $edge_G$, includes all pairs of connected nodes in the network, $\{i, j\}$, where $i \in node_G, j \in node_G$; the social network is thus defined as: $G = \{node_G, edge_G\}$. The social network can also be represented by its adjacency matrix $\mathbf{A}$, whose elements are $A_{ij}$, where $A_{ij} = 1$ if $\{i, j\} \in edge_G$; $A_{ij} = 0$ otherwise.

*2) Belief:* We employ two kinds of beliefs in the IFM model: private and updated belief. The former is unchanged and taken as an input of the model. The latter, however, is updated at each time step.

*a) Private belief:* Private beliefs abstract the intrinsic characteristics of nodes in a network. They will not be changed during the process of information flow. In this model, we assume all nodes in the network have same probabilistic distribution of private beliefs, which means node $i$ in the network takes the private belief as a random number $w_i \in [-1, 1]$ with distribution $p(w_i)$. And we denote the private belief vector as $\mathbf{w}$, whose elements are $w_i$. The distribution $p(\mathbf{w})$ is common knowledge to everyone in the network.

*b) Current belief:* A current belief $B_{i,T}$ describes the current opinion of node $i$ in a network at time step $T \in \mathbb{N}$. It lies in the range $[-1, 1]$. $\mathbf{B}(T)$ is a vector whose elements are $B_{i,T}$. $\mathbf{B}(T)$ will be updated at each time step, and will converge to a limit $\mathbf{B}(T)$ in certain networks, as will be explained in Section III.

*3) Control Strategy:* To control the overall behaviour of the network, we propose a control strategy which chooses certain nodes in the network, the so called control nodes, and asks them to broadcast certain beliefs to their neighbours.

*a) Control Set:* The set of $c$ control nodes is defined as $\mathbf{C} = \{\theta_1, \theta_2, ..., \theta_c\}$, where $\theta_i$, $1 \leq i \leq c$, are indices of control nodes. The uncontrolled nodes thus belong to set $^\dagger C = \{\theta_{c+1}, ..., \theta_N\}$. And the belief chosen to be broadcast

by the $i^{th}$ control node, is set to a controlled belief $B_i^*$, where $B_i^* \in [-1, 1]$, such that:

$$w_{\theta_i} = B_i^*, B_{\theta_i, T} = B_i^*, \tag{1}$$

for $\theta_i \in \mathbf{C}$ and any value of $T$. A control strategy is specified by the control set $\mathbf{C}$ with the corresponding controlled beliefs $\mathbf{B}^*$.

*b) Control Power:* Control power is used to measure how much the beliefs in a network have changed from their initial status. Control power for an arbitrary node $i$ is defined as the expectation of the difference between the "final" belief $B_{i,\infty}$ and the initial belief $w_i$:

$$cp_i = E[B_{i,\infty} - w_i]. \tag{2}$$

The averaged $cp_i$ values over all nodes in the network is called the *network control power*:

$$cp = \Sigma_{i=1}^N E[B_{i,\infty} - w_i]/N. \tag{3}$$

### B. Information Flow Model

In an information flow model, $B_{i,0}$ is initialized to $w_i$. The value of $B_{i,T}$ is then updated at each time step and determined as the average of the current beliefs of the neighbours of node $i$, and the private belief of node $i$, $w_i$, under the influence of a control strategy. The average process is specified by an adjusted private belief vector and an adjusted adjacency matrix. The adjusted private belief vector is denoted by $\mathbf{w}^*$, with elements $w_i/(d_i + 1)$, where $d_i = \Sigma_{j=1}^N A_{ij}$. And the adjusted adjacency matrix $\mathbf{A}^*$ contains elements $A_{i,j}^* = A_{i,j}/(1 + d_j)$. The control strategy is specified by a control matrix and a control vector, which are determined by the control set $\mathbf{C}$ and the corresponding controlled belief $\mathbf{B}^*$. The control matrix is defined as $\mathbf{M}$, where $M_{i,i} = 1$ if $i \notin \mathbf{C}$ and $M_{i,j} = 0$ otherwise. The control vector is $\mathbf{V}$, where $V_{\theta_i} = B_i^*$ if $i \leq c$, and $V_{\theta_i} = 0$ otherwise. If $\mathbf{C} = \emptyset$, the control strategy is trivial. The updating process of the current belief vector $\mathbf{B}(T)$ is shown in Figure (1).
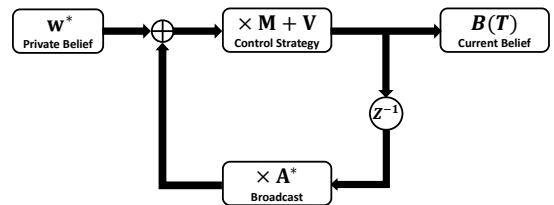


Fig. 1: Process for Information Update.

In Figure (1), the adjusted private belief $\mathbf{w}^*$ and adjusted adjacency matrix $\mathbf{A}^*$ are for calculating an averaged belief for each node. The control matrix $\mathbf{M}$ and control vector $\mathbf{V}$ are used to set the belief of control nodes as their corresponding controlled belief $\mathbf{B}^*$. And $Z^{-1}$ means the beliefs are updated based on the information of the previous time step. Equation (5) shows the formula to calculate $\mathbf{B}(T)$:

$$\mathbf{B}(T) = [\mathbf{w}^* + \mathbf{B}(T-1) \times A^*] \times M + V. \qquad (4)$$

Based on Equation (4), we can derive a compact sigma notation of current belief as shown in Equation (5):

$$\mathbf{B}(T) = [\mathbf{w}^* \times M + V] \times [\Sigma_{t_1=0}^{T-1}(A^* \times M)^{t_1}]. \qquad (5)$$

assuming that the summation $\Sigma_{t_1=0}^{T}(A^* \times M)^{t_1}$ in Equation (5) converges to a finite matrix when $T$ approaches $\infty$, the converged belief $\mathbf{B}(\infty)$ can be represented as:

$$\mathbf{B}(\infty) = [\mathbf{w}^* \times M + V] \times [I - A^* \times M]^{-1}. \qquad (6)$$

From Equation (6) and Equation (3), we obtain the control power metric when using a control strategy with parameters $\mathbf{M}$ and $\mathbf{V}$ as[1]:

$$cp = \frac{1}{N}||E\left[[\mathbf{w}^* \times M + V] \times [I - A^* \times M]^{-1} - \mathbf{w}\right]||_1. \quad (7)$$

### C. An example

In certain applications, a specific direction of change of public beliefs may be of interest, and hence a desire to control the sign of $cp$, and maximize its magnitude. For example, in a social network $G$ the set $node_G$ may represent patients with certain disease, and the set $edge_G$ may be determined by the email and phone communications between these patients. Each patient has his/her own private belief $w_i$ about their disease. A doctor may want to influence the beliefs $B(T)_i$ of his/her patients and steer their final beliefs $B(\infty)_i$ towards a desired belief (e.g., "following your suggested treatment will help you"). To do so, the doctor may choose several patients in the social network as control nodes $\mathbf{C}$, and ask them to broadcast certain controlled beliefs $\mathbf{B}^*$ about taking the suggested treatment. However, the number of controlled patients $c$ may be limited by budget constraints. To reach maximum performance regarding influencing the patients beliefs, the doctor needs an optimized control strategy, as we describe later.

### III. ANALYSIS OF CONTROL POWER

In this section, we show how we can efficiently calculate the measure of control power in Equation (7) and select an appropriate control strategy for social networks via two social network models, the BA model and the GMG model, each of which has its own advantage in performance or resources in computation and information. We start each model by introducing the fundamental assumptions and the network synthesis processes, which will be used to verify the assumptions. The main results regarding calculation of control power are Equation (12) and Equation (18). Theorem III.3 and III.5 are the main results regarding optimization of the control power, as they provide optimized control strategies for each model.

According to Equation (7), to calculate the control power $cp$, the complete information about $\mathbf{A}$ is needed. In addition,

[1]$||\mathbf{x}||_1$ is the $l_1 \, norm$ of vector $\mathbf{x}$ with N elements: $||\mathbf{x}||_1 = \Sigma_{i=1}^N x_i$

the computational cost is the inversion of matrix $I - \mathbf{A}^* \times \mathbf{M}$. Such an exact solution does not shed any particular light on the choice of control set $\mathbf{C}$, or on the convergence speed of $\mathbf{B}(T)$ towards $\mathbf{B}(\infty)$. In order to reduce the information needed to predict the control power $cp$, as well as provide a detailed analysis about the control strategy and convergence speed, network models are required. In particular, the social network models are meant to reveal the intrinsic properties of $\mathbf{A}$, as the other elements in the proposed information flow model are well understood.

In practice, information about the network may not be complete, which means $\mathbf{A}$ is not always available. Furthermore, the network may contain a large number of nodes or edges, which requires significant computational power to process. To solve such problems, network models are necessary. A good network model can help calculate the converged beliefs using less information than $\mathbf{A}$, and more efficiently. Two important network models, the BA model [2] and the GMG model [16], will be introduced and applied in the analysis of the information flow model. The reason for choosing these two models is that they both provide probabilistic properties about the element $A_{i,j}$ of the adjacency matrix $\mathbf{A}$. The BA model assumes that $A_{i,j}$ only depends on the degree of nodes $i$ and $j$, and thus requires less information. The GMG model, on the other hand, extends the dependence of $A_{i,j}$ to both degree and clustering coefficient of nodes $i$ and $j$, and thus provides better accuracy.

### A. Barabási-Albert Model

*1) Basic assumption:* A BA model describes the growth of a network. We can, however, model a static network of interest with size $N$ as one evolves from an initial network of small size until the the number of nodes reaches $N$. Then the network growth is freeze and we perform information flow on it. The purpose of introducing the BA model is to reveal the statistical characteristic of the adjacency matrix of a social network. The BA model, however, receives different judgements from academics [13] [6], which leaves room for improvement such as the proposed GMG model.

One of the basic assumptions of the BA model is that the probability of a node $i$ attached by a new edge is proportional to its degree $d_i$ [2]:

$$\frac{\partial d_i}{\partial t} \sim d_i. \qquad (8)$$

Based on this assumption, we can derive the probability $P_{i,j}$ of an edge established between nodes $i$ and $j$ when the size of the network grows to a certain size, as shown in Theorem III.1. In this information flow model, $A_{i,j}$ is binary, so that $^1P_{i,j}$ is also the expected value of $A_{i,j}$ as: $\overline{A}_{i,j} = 1 \times {}^1 P_{i,j} + 0 \times (1 - {}^1 P_{i,j})$. Moreover, the summation of $\overline{A}_{i,j} = P_{i,j}$ over all choices for $i$ and $j$ is $\Sigma_{k=1}^N d_k$, which is the summation of all degrees and is $\Sigma_{i,j}A_{i,j}$. $^1P_{i,j}$ plays an important role in the analysis of control power estimation and of control strategy, as it represents the information about the matrix $\mathbf{A}$.

**Theorem III.1.** *In a BA model, the probability $^1P_{i,j}$ of two nodes $i$ and $j$ being connected in network $G$ with $N$ nodes*

*and a fixed degree sequence is:*

$$^1P_{i,j} = \frac{d_i d_j}{\Sigma_{k=1}^N d_k}, \tag{9}$$

*where $d_k$ is the degree corresponding to node $k$.*

The proof of Theorem III.1 is provided in Appendix C.

*2) Network Synthesis:* In order to verify the correctness of Equation (9), we need to generate sample social networks according to the basic assumption of the BA model. The network synthesis of a BA model is discussed in Barabási(2002) [2]. The input of the synthesis process is the total number of nodes $N$ and the average number of edges attached to each new incoming node $m$. The detailed process is shown in Appendix A.

*3) Calculation of Control Power:* Theorem III.1 has revealed the statistical properties of adjacency matrix **A**. If we take **A** as a random matrix, combined with the formula to calculate a converged belief, as shown in Equation (5) and Equation (6), we are able to obtain the expected value of converged beliefs for all the nodes in the network, which is shown in Theorem III.2.

**Theorem III.2.** *In a BA model, the expected value of converged belief $^1B_{i,\infty}$ of a non-controlled node $i$, $i \notin \mathbf{C}$, is:*

$$\overline{^1B}_{i,\infty} = \frac{1}{\Sigma_{k=1}^N d_k} \frac{d_i}{1+d_i} \frac{\Sigma_{j=1}^c B^*_{\theta_j} d_{\theta_j} + \Sigma_{j=c+1}^N \frac{\overline{w}_{\theta_j}}{1+d_{\theta_j}} d_{\theta_j}}{1-\beta_1}, \tag{10}$$

*where $w_{\theta_i}$ is the private belief of node $\theta_i$, $m$ is the average number of edges per node in a network $G$ with $N$ nodes, $d_i$ is the degree corresponding to node $i$, $B^*_{\theta_j}$ is the controlled belief of control node $j$, $c$ is the number of control nodes, $\theta_i \in \mathbf{C}$ for $i \leq c$, and $\beta_1$ is a constant which is smaller than 1:*

$$\beta_1 = \Sigma_{k=c+1}^N \frac{d_{\theta_k}{}^2}{1+d_{\theta_k}} / \Sigma_{k=1}^N d_k. \tag{11}$$

The proof of Theorem III.2 is in Appendix D.

Plugging Equation (10) into Equation (3), we obtain the control power:

$$^1cp = \Sigma_{i=1}^N \left( \frac{d_i(\Sigma_{j=1}^c B^*_{\theta_j} d_{\theta_j} + \Sigma_{j=c+1}^N \frac{\overline{w}_{\theta_j}}{1+d_{\theta_j}} d_{\theta_j})}{N\Sigma_{k=1}^N d_k(1+d_i)(1-\beta_1)} - \overline{w}_i \right) \tag{12}$$

The information needed for the calculation in Equatoin(12) is the degree list of network $G$, which is far less than the information of adjacency matrix **A**. In addition, degree lists follow the power-law distribution in most social networks [2], which means the degree list could be sampled from the network. The computational cost of such calculation is $O(N)$, which is much more efficient than the matrix inverse calculation required by Equation (7).

*4) Optimization of Control Strategy:* According to Equation (12), we can see that the control strategy, as well as the degrees of the control nodes, have a direct impact on the control power. Without loss of generality, we set the preferred sign of beliefs positive. If controlled beliefs $B^*_i$, $i \leq c$, are maximized to be 1, and the private belief $w_i$ has zero mean, then, as shown in Theorem III.3, the maximization of $^1cp$ requires the selection of a control group **C** to include nodes with highest degrees in the network $G$.

**Theorem III.3.** *Consider a social network $G$ of $N$ nodes, with degree list $\{d_i\}$, $i = 1, \ldots, N$, a private belief **w** with zero mean and maximized control beliefs $B^*_i = 1$, $i \leq c$. Suppose further that the number of control nodes, $c$, is fixed. The control set $\mathbf{C_o} = \{\theta_{o1}, \theta_{o2}, \ldots, \theta_{oc}\}$, where $d_{\theta_{o_i}} \geq d_{\theta_{o_j}}$ if $i \leq j$, $1 \leq i, j \leq N$, maximizes the control power $^1cp$:*

$$\mathbf{C_o} = \arg\max_{\mathbf{C}} {}^1cp(\mathbf{C}) \tag{13}$$

The proof of Theorem III.3 is provided in Appendix F.

*B. Generalized Markov Graph Model*

*1) Basic assumption:* In [16], the BA model is shown to be a special case of a Markov Graph model [10]. A Markov Graph model is based on the dependence between pairs of nodes. The basic assumption of the BA model depends however on the degree of nodes, which is a specific description of pairwise relationship between nodes. It is thus natural to extend the BA model to the GMG model to analyse the property of adjacency matrix **A**.

In a GMG Model, the probabilistic dependence of an edge is extended from the other attached edges to attached triangles. As degree is used to describe the dependence on attached edges, a clustering coefficient, which is related to both edges and triad relational structures, is added to the description of dependence in a GMG model. The assumption about the probability of a node $i$ attached by a new edge in a GMG model then becomes:

$$\frac{\partial d_i}{\partial t} \sim d_i(1+\gamma_i)^\alpha, \tag{14}$$

where $d_i$ is the degree of node $i$, $\gamma_i$ is the clustering coefficient of node $i$, and $\alpha$, which is called the clustering weight, is determined by the property of the network $G$.

The range of the clustering coefficient lies in $[0, 1]$. If there is no triangle attached to a node $i$, the clustering coefficient $\gamma_i$ will be zero. To ensure a non-zero probability of a node getting an edge, on account of a zero clustering coefficient, we adopt a $(1+\gamma_i)$ in our model.

In an arbitrary network $G$, the influence of degree $d_i$ and clustering coefficient $\gamma_i$ on the probability of node $i$ obtaining a new edge is not equivalent. Parameter $\alpha$ is thus used to adjust the importance of clustering coefficient versus degree. The value of $\alpha$ changes for different types of networks and different applications.

Based on Theorem III.1, we add the influence of clustering coefficients, and to make the summation of $^2P_{i,j}$ still be the sum of degrees $\Sigma_{k=1}^N d_k$, the probability $^2P_{i,j}$ becomes:

$$^2P_{i,j} = \frac{d_i(1+\gamma_i)^\alpha d_j(1+\gamma_j)^\alpha}{\eta} \Sigma_{k=1}^N d_k, \tag{15}$$

where $\eta$ is:

$$\eta = \Sigma_{i=1}^N \Sigma_{j=1, j\neq i}^N d_i(1+\gamma_i)^\alpha d_j(1+\gamma_j)^\alpha.$$

*2) Network Synthesis:* In order to verify the correctness of Equation (15), we need to generate sample social networks according to the basic assumption of the GMG model. In addition, the synthesized social networks can be used to estimate the clustering weight $\alpha$. In a GMG model [18], the input of the synthesis process is: the total number of nodes $N$, the average number of edges attached to each node $m$, and the clustering weight $\alpha$. The detailed process is shown in Appendix B.

*3) Calculation of Control Power:* Since the GMG model makes use of more information than the BA model, the probability $^2P_{i,j}$ should be more accurately describing the statistical property of adjacency matrix $\mathbf{A}$ than $^1P_{i,j}$. Based on this idea, we develop the converged belief for a GMG model, as shown in Theorem III.4.

**Theorem III.4.** *Define a constant $\beta_2$ as,*

$$\beta_2 = \Sigma_{k=c+1}^N \frac{(d_{\theta_k}(1+\gamma_{\theta_k})^\alpha)^2}{(1+d_{\theta_k})\eta} \Sigma_{k=1}^N d_k. \qquad (16)$$

*If $|\beta_2| < 1$, the expected value of converged belief $^2B_{i,\infty}$ of a non-controlled nodes $i$, $i \notin \mathbf{C}$, in a GMG model, is,*

$$\overline{^2B_{i,\infty}} = \frac{\Sigma_{k=1}^N d_k}{\eta} \frac{d_i(1+\gamma_i)^\alpha}{1+d_i}$$
$$\frac{\Sigma_{j=1}^c B_{\theta_j}^* d_{\theta_j}(1+\gamma_{\theta_j})^\alpha + \Sigma_{j=c+1}^N \frac{\overline{w_{\theta_j}}}{1+d_{\theta_j}} d_{\theta_j}(1+\gamma_{\theta_j})^\alpha}{1-\beta_2}, \qquad (17)$$

*where $w_j$ is the private belief of node $i$ in a network $G$ with $N$ nodes, $d_i$ is the degree corresponding to node $i$, $\gamma_i$ is the clustering coefficient of node $i$, $B_{\theta_j}^*$ is the controlled belief of control node $j$, $c$ is the number of control nodes, $\theta_j \in \mathbf{C}$ if $j \le c$, $\alpha$ is the clustering weight for network $G$, $\eta$ is defined in Equation (15).*

The proof of Theorem III.4 is in Appendix E.

The constant $\beta_2$, however, will not be guaranteed to be strictly smaller than 1, as $\beta_1$ is in a BA model, which is shown in Appendix D. The value of $\beta_2$ will be determined by the degree list, the clustering coefficient list and the clustering weight $\alpha$ of the network $G$.

Plugging Equation (17) into Equation (3), we obtain the control power:

$$^2cp = \frac{1}{N}\Sigma_{i=1}^N \left( \frac{\Sigma_{k=1}^N d_k}{\eta} \frac{d_i(1+\gamma_i)^\alpha}{1+d_i} \times \right.$$
$$\left. \frac{\Sigma_{j=1}^c B_{\theta_j}^* d_{\theta_j}(1+\gamma_{\theta_j})^\alpha + \Sigma_{j=c+1}^N \frac{\overline{w_{\theta_j}}}{1+d_{\theta_j}} d_{\theta_j}(1+\gamma_{\theta_j})^\alpha}{1-\beta_2} - \overline{w_i} \right). \qquad (18)$$

The calculation of a control power in Equation(18) requires the information of the degree list, the clustering coefficient list and the clustering weight $\alpha$ of network $G$. The clustering weight $\alpha$ is obtained by a learning process, and is introduced in Section IV. When calculating the control power, the information needed by the GMG model is still far less than the information of adjacency matrix $\mathbf{A}$. And due to the fact that $\eta$ could be rewritten as:

$$\eta = (\Sigma_{i=1}^N d_i(1+\gamma_i)^\alpha)^2 - \Sigma_{i=1}^N d_i(1+\gamma_i)^\alpha,$$

the computational cost of such a calculation is $O(N)$, which is the same as that in the BA model.

*4) Optimization of Control Strategy:* According to Equation (18), we can see that, in addition to the control strategy and the degrees of the control nodes, the clustering coefficients of the control nodes also have a direct impact on the control power. Without loss of generality, we set the preferred sign of beliefs positive. If controlled beliefs $B_i^*$, $i \le c$, are maximized to be 1, and the private belief $w_i$ has a zero mean, then, as we show in Theorem III.5, the maximization of $^2cp$ requires the selection of a control group $\mathbf{C}$ to include nodes with highest $d(1+\gamma)^\alpha$ values in the network $G$.

**Theorem III.5.** *Consider a social network $G$ of $N$ nodes, with degree list $\{d_i\}$, clustering weight $\alpha$, clustering coefficient list $\{\gamma_i\}$, $i = 1, \ldots, N$, and maximized control beliefs $B_i^* = 1$, $i \le c$. Suppose further that the private belief $\mathbf{w}$ has zero mean and the number of control nodes $c$ is fixed. Then, the control set $^2\mathbf{C_o} = \{\theta_{o1}^\dagger, \theta_{o2}^\dagger, \ldots, \theta_{oc}^\dagger\}$, where $d_{\theta_{oi}^\dagger}(1+\gamma_{\theta_{oi}^\dagger})^\alpha \ge d_{\theta_{oj}^\dagger}(1+\gamma_{\theta_{oj}^\dagger})^\alpha$ if $i \le j$, $1 \le i, j \le N$, maximizes the control power $^2cp$*

$$^2\mathbf{C_o} = \arg\max_{\mathbf{C}} {}^2cp(\mathbf{C}), \qquad (19)$$

*under the condition that:*

$$\frac{1}{\beta_2} > 1 + max(1, 2^\alpha) \frac{\Sigma_{j=1}^c d_{\theta_{oj}^\dagger}(1+\gamma_{\theta_{oj}^\dagger})^\alpha}{\Sigma_{k=c+1}^N \frac{(d_{\theta_{ok}^\dagger}(1+\gamma_{\theta_{ok}^\dagger})^\alpha)^2}{1+d_{\theta_{ok}^\dagger}}}. \qquad (20)$$

The proof of Theorem III.5 is provided in Appendix G.

The GMG model covers more types of social networks due to the presence of the clustering weight $\alpha$, an additional free parameter. Note that the range of $\alpha$ is constrained only by the learning data set, not by the model. The specific condition in Theorem III.5 means that the optimized control strategy is appropriate only for certain types of social networks where the beliefs converge fast, as denoted by $\beta_2$.

## IV. EXPERIMENTS ON VERIFICATION AND APPLICATIONS OF INFORMATION FLOW MODELS

Three sets of experiments are described in this section. In the first set, presented in Section IV-A, we verify the basic assumptions about $P_{i,j}$, the probabilistic nature of the two network models, as showed in Theorem III.1 and Equation (15). In the second set, presented in Section IV-B, we test the performance of the two social network models when the private beliefs of nodes in a social network have zero-mean. In Section IV-C, the last set of experiments compares the effects of control strategies on maximizing the control power of two different models.

Real network data [15] is used in all three types of experiments. There are 3 different types of social networks: online social networks, p2p transmission networks and physicist
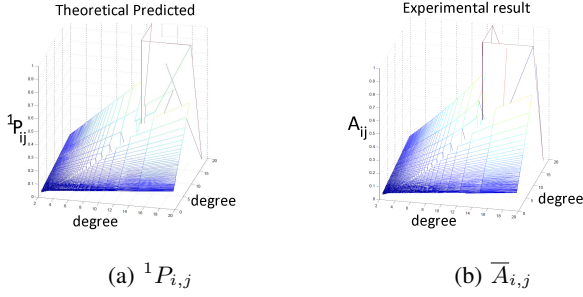
(a) $^1P_{i,j}$       (b) $\overline{A}_{i,j}$

Fig. 2: Comparison between $^1P_{i,j}$ and $\overline{A}_{i,j}$.

| m | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| Relative Error | 4.32e-2 | 4.59e-2 | 3.75e-2 | 4.12e-2 | 3.61e-2 |

TABLE I: Relative error between $^1P_{i,j}$ and $\overline{A}_{i,j}$ on synthesized data.

collaboration networks. Each of these three different social networks includes several subtypes of networks. On-line social networks include Slashdot network data of August 2008 and of February 2009, Wiki-vote network data and Epinions network data. P2p transmission networks include Gnutella network data at five different times. Physicist collaboration networks include collaboration networks of physicists studying astrophysics, condensed Matter Physics, theoretical high-energy Physics, experimental high-energy Physics and general relativity. All of these effectively constitute 14 subtypes of social networks and will be denoted by indices: $1, 2, \ldots, 14$. From each of these 14 networks, we sampled 50 sub-networks using the same sampling method.

### A. Basic Assumption Verification

The first type of experiments is designed to test the fundamental assumptions of the BA model, as shown in Section III-A1, and GMG model, as shown in Section III-B1. Both assumptions are tested on the synthesized data to show the consistence of the assumptions with the derived formula of $P_{i,j}$ in both models. The real network data are then used to test how $P_{i,j}$ computed in these two models fit into real applications. The performance of the models on synthesized data is better than that on real data, since the former were used to generate the data in the first place. The additional performance difference is due to data sampling noise as well as to the intrinsic unaccounted characteristics of the network itself.

*1) Barabási-Albert model:*

*a) Verification of Theorem III.1 on synthetic data:* In order to test the correctness of the basic assumption of the BA model, $^1P_{i,j}$ is computed to compare with $A_{i,j}$ averaged across 1,000 realizations of networks synthesized by the algorithm introduced in III-A2. An example of experimental results is shown in Figure (2) with networks containing 100 nodes and 300 edges. The average number of edges per node is thus $m = 3$. In this example, the average degree list is averaged across the degree lists of the 10,000 realizations of networks. The averaged degree list is then used to calculate $^1P_{i,j}$ as shown in Equation (9). The experiment is then repeated for different values of $m$, and the relative error between $^1P_{i,j}$ and $\overline{A}_{i,j}$ averaged across all pairs of nodes is shown in Table I.

From Table I, we can see that Equation (9) can be used to calculate the value of $^1P_{i,j}$, which is the expected value of $A_{i,j}$ for different choices of $m$ value. This experiment thus verifies the result of Theorem III.1.

*b) Performance of Theorem III.1 on real network data:* In this experiment, real network samples are used to validate Theorem III.1. For each subcategory of network samples, the number of nodes $N$ is chosen to be the minimal number of nodes of all samples, which ensures the same number of samples being used for the calculation of the average number of every element in the adjacency matrix $A$. The degree list, $N$, and $m$ are then averaged across all samples of the same subcategory. The value of $N$, the averaged degree list and the averaged $m$ are plugged into Equation (9) to calculate $^1P_{i,j}$. And the value of $\overline{A}_{i,j}$ is the average of $A_{i,j}$ from all 50 samples of networks belonging to the same subcategory. The relative errors between $^1P_{i,j}$ and $\overline{A}_{i,j}$ for 14 different types of network are shown in Table II.

| On-line | | P2p | | Collaboration | |
|---|---|---|---|---|---|
| Network | Error | Network | Error | Network | Error |
| 1 | 8.31e-2 | 5 | 22.13e-2 | 10 | 7.28e-2 |
| 2 | 8.57e-2 | 6 | 19.36e-2 | 11 | 9.15e-2 |
| 3 | 7.46e-2 | 7 | 24.18e-2 | 12 | 7.16e-2 |
| 4 | 6.59e-2 | 8 | 22.40e-2 | 13 | 8.52e-2 |
| | | 9 | 20.25e-2 | 14 | 9.73e-2 |

TABLE II: Relative error between $^1P_{i,j}$ and $\overline{A}_{i,j}$ on real network data.

From Table II, we can see that the value of $^1P_{i,j}$, according to Equation (9), is close to the expected value of $A_{i,j}$ in on-line social network samples and Collaboration networks. However, in p2p transmission networks, Theorem III.1 doesn't provide a good estimate of $\overline{A}_{i,j}$. The main reason for the inaccurate estimation of $A_{i,j}$ in this case is that the degree distribution of such a network does not strictly follow the power law distribution, which means the structure of these networks is different from what the BA model can describe. This indicates that a better model is needed for a more accurate description of the real network data.

*2) Generalized Markov Graph Model:*

*a) Verification of basic assumption on synthetic data:* The basic assumption of the GMG model is tested as follows. The number of nodes $N$ is 100. Five different values of $m$: $1, 2, 3, 4, 5$, and 4 different values of $\alpha$: $-2, -1, 1, 2$, are chosen as the inputs of the network synthesis process introduced in III-B2. For each combination of $m$ and $\alpha$, $^2P_{i,j}$ is calculated and compared with $A_{i,j}$ averaged across 1,000 realizations of synthesized networks. The relative error

between $^2P_{i,j}$ and $\overline{A}_{i,j}$ averaged across all pairs of nodes is shown in Table III.

| m<br>$\alpha$ | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| -2 | 3.25e-2 | 4.82e-2 | 4.26e-2 | 3.19e-2 | 3.55e-2 |
| -1 | 4.09e-2 | 4.91e-2 | 4.92e-2 | 3.31e-2 | 4.94e-2 |
| 1 | 4.91e-2 | 3.97e-2 | 4.60e-2 | 3.28e-2 | 3.84e-2 |
| 2 | 4.83e-2 | 4.58e-2 | 4.91e-2 | 4.31e-2 | 3.07e-2 |

TABLE III: Relative error between $^2P_{i,j}$ and $\overline{A}_{i,j}$ on synthetic data.

From Table III, we can see that Equation (15) can be used to calculate the value of $^2P_{i,j}$, which is the expected value of $A_{i,j}$ for different choices of $m$ and $\alpha$.

*b) Verification of basic assumptions on real network data:* In this experiment, real network samples are used to validate Equation (15). For each subcategory of the network samples, the number of nodes $N$ is chosen to be the minimal number of nodes of all 50 samples. The degree list of length $N$, the clustering coefficient list of length $N$, and the average number of edges $m$ are then averaged across all samples of the same subcategory. Then 25 of randomly chosen sampled networks for each sub-category of networks are used as a training set to learn the value $\alpha$. The rest of sampled networks are used as a test set. For each sub-category of networks, $\alpha$ is determined as the value that generates $^2P_{i,j}$ closest to $A_{i,j}$ averaged across the same 25 chosen sampled networks. Next, for each sub-category of networks, averaged degree lists, average clustering coefficient lists and averaged $m$ of the other 25 sampled networks, along with the learned $\alpha$, are used in Equation (15) to compute $^2P_{i,j}$. The value of $\overline{A}_{i,j}$ is the average of $A_{i,j}$ over the test set for each subcategory of networks. The relative error between $^2P_{i,j}$ and $\overline{A}_{i,j}$ for 14 different types of network is shown in Table IV.

| On-line | | P2p | | Collaboration | |
|---|---|---|---|---|---|
| Network | Error | Network | Error | Network | Error |
| 1 | 5.14e-2 | 5 | 8.68e-2 | 10 | 5.45e-2 |
| 2 | 6.23e-2 | 6 | 7.51e-2 | 11 | 6.69e-2 |
| 3 | 4.75e-2 | 7 | 8.64e-2 | 12 | 6.28e-2 |
| 4 | 5.86e-2 | 8 | 9.13e-2 | 13 | 4.19e-2 |
| | | 9 | 7.26e-2 | 14 | 6.94e-2 |

TABLE IV: Relative error between $^2P_{i,j}$ and $\overline{A}_{i,j}$ on real network data.

From Table IV, we can see that the value of $^2P_{i,j}$ from Equation (9) is close to the expected value of $A_{i,j}$ for all three different categories of networks. Although in p2p transmission networks, Equation (9) doesn't give as good an estimate of $\overline{A}_{i,j}$ as in the other two categories, the overall performance of the GMG model is better than that of the BA model. Note for p2p transmission networks, whose degree distribution does not strictly follow the power law distribution, the GMG model could greatly increase the accuracy of $A_{i,j}$ estimate.

## B. Converged Belief Estimation

In this experiment, we evaluate the accuracy of the converged belief $\mathbf{B}(\infty)$ estimate for both models of interest. The control strategy is set to push neutral public opinions towards positive opinions. The private beliefs $w_i$ of nodes in the network are set to be neutral, i.e, they obey a uniform distribution on $[-1, 1]$. For both the BA and the GMG models, 25 randomly chosen networks are selected from each sub-category of networks, and used as a test set. The other 25 networks for each sub-category of networks are used as a training set for the GMG model to learn the clustering weight $\alpha$.

*1) Barabási-Albert Model:* For each of the network samples in the test set of each sub-category of networks, a degree list is recorded. The control set $\mathbf{C}$ is set to be the top $5\%$ nodes with the highest degree in the network, thus fixing the number of control nodes to: $c = \lceil 5\% N \rceil$. The controlled beliefs $\mathbf{B}^*$ for nodes in the control set are set to be 1. For each network sample, the private belief list $w_i$ is generated 100 times according to a uniform distribution on $[-1, 1]$.

For each of the generated private belief lists, the value of control power $^1cp_i$ is calculated as in Section III-A3. The same network is used for calculating the exact value of the control power according to Equation (6). The relative error between the $^1cp_i$ and the exact value is then recorded as the relative error for this network sample under this belief list. The 100 relative errors for all generated private belief lists are then averaged and recorded as the relative error of this network sample. Next, the relative errors for all network samples in the test set of each sub-category of networks are averaged and recorded as the relative error of the BA model on this sub-category. The relative errors for 14 different sub-categories of networks are shown in Figure (3) and Figure (4) to compare with the results from the GMG model.

*2) Generalized Markov Graph Model:* In a GMG model, the clustering weight $\alpha$ plays a key role in describing the characteristics of the social network and has to be learned from the training data. We propose two ways to learn $\alpha$. The first one makes use of complete information of the adjacency matrix $A$ of the social network and provides better accuracy. The second one learns $\alpha$ using less information, namely only the degree list and clustering coefficient list of the social network. The learned $\alpha$ is then used to calculate the control power and the results are compared with the one from the BA model.

*a) learning $\alpha$ from $A$:* The first step in a GMG model is to learn the clustering weight $\alpha$ from the training sets for each sub-category of networks. For each training network sample, the degree list and clustering coefficient list are recorded. The control set $\mathbf{C}$ is set to be the top $5\%$ nodes with highest degree in the network, and thus the number of control nodes is set to be: $c = \lceil 5\% N \rceil$. The controlled beliefs $B_i^*$ for nodes in the control set are set to be 1, and 100 private belief lists sampled from a uniform distribution on $[-1, 1]$ are obtained.

The clustering weight $\alpha$ is chosen to minimize the difference between $cp$ and $^2cp$. For each network sample in the training set of each sub-category of networks, the control power $^2cp_i$ is calculated for each of the 100 private belief lists with an

arbitrary choice of $\alpha$. The exact solution of control power is also calculated for the 100 generated private belief lists. The relative error between the $^2cp_i$ and the exact value of the same private belief list is then calculated, averaged across all generated private belief lists, and recorded as the relative error for this network sample. For all 25 network samples in the training set, the relative errors are calculated and their average value is recorded as the relative error of this sub-category of networks. The clustering weight $\alpha$ for each sub-category of networks is chosen as the one that minimizes the relative error, as shown in Table V.

| On-line | | P2p | | Collaboration | |
|---------|---|-----|---|---------------|---|
| Network | $\alpha$ | Network | $\alpha$ | Network | $\alpha$ |
| 1 | -1.00 | 5 | -0.20 | 10 | 1.48 |
| 2 | -1.00 | 6 | 1.29 | 11 | -2.57 |
| 3 | -0.50 | 7 | 1.52 | 12 | -2.57 |
| 4 | -0.30 | 8 | 1.45 | 13 | -2.84 |
| | | 9 | 1.12 | 14 | 0.52 |

TABLE V: $\alpha$ learned from complete information of $\mathbf{A}$ from training set.

The learned $\alpha$ is subsequently used to test the performance of the GMG model on the test set. For each of the network samples in the test set, the degree list and the clustering coefficient list are recorded. The control set is similar to that in the training set. For each network sample, the private belief list $w_i$ is generated 100 times according to a uniform distribution on $[-1, 1]$. The relative errors are calculated on the test set by the same method as that used on the training set. The relative errors for 14 different sub-categories of networks are shown in Figure (3), together with the result from the BA model.
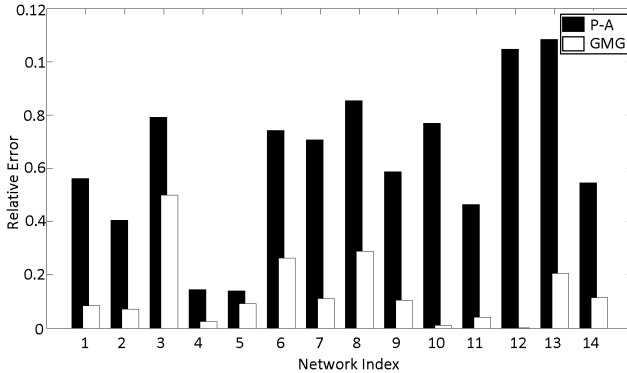


Fig. 3: Relative error of control power

*b) learning $\alpha$ from the degree list and the clustering coefficient list:* The second method of learning $\alpha$ only uses partial information of the networks in the training set, i.e, the degree list and the clustering coefficient list. For each of the network samples in the training set, the number of nodes $N$, the number of edges $N \times m$, the degree list and the clustering list are recorded. For an arbitrary value of $\alpha$, 100 networks with $N$ nodes and $N \times m$ edges are generated according to the network synthesis algorithm introduced in Section III-B2.

For each of the generated networks, the degree list and the clustering coefficient list are recorded and averaged across these 100 networks. The averaged degree list and clustering coefficient list are then compared with the one of the network sample in the training set. The difference between them is calculated for each of the network samples in the training set, and averaged across all 25 networks in the training set. The $\alpha$ that minimizes the average difference is selected as the clustering weight of this sub-category of networks. The results for all 14 sub-categories of networks are shown in Table VI.

| On-line | | P2p | | Collaboration | |
|---------|---|-----|---|---------------|---|
| Network | $\alpha$ | Network | $\alpha$ | Network | $\alpha$ |
| 1 | -0.40 | 5 | -0.15 | 10 | 0.58 |
| 2 | -0.30 | 6 | 0.29 | 11 | -2.30 |
| 3 | -0.40 | 7 | 0.50 | 12 | -1.50 |
| 4 | -0.20 | 8 | 0.40 | 13 | -1.70 |
| | | 9 | 0.10 | 14 | 0.50 |

TABLE VI: $\alpha$ learned from partial information of $\mathbf{A}$ from training set.

The learned $\alpha$ is then used to test the performance of the GMG model on the test set. The experiment is similar as in Section IV-B2a, and the relative errors of control power are computed for all 14 subcategories of networks, as shown in Figure (4).
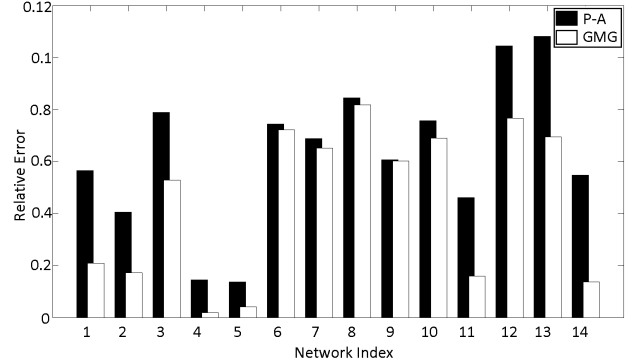


Fig. 4: Relative error of control power

*c) Discussion of the experiments:* Figures (3) and (4) indicate that the GMG model outperformed the BA model, as the former one yields smaller error for both $\alpha$ learned from adjacency matrix $A$ and from the degree list and the clustering list. When the complete information of $\mathbf{A}$ of the training set is used, the learned $\alpha$ will make the GMG model more accurate than the case when only partial information of $\mathbf{A}$ is learned. And the more accurate of the models, the more information is needed.

## C. Control strategy optimization

In the experiments described in this section, we are interested in comparing the performance of control strategies under the two network models.

*1) Barabási-Albert Model:* According to Theorem III.3, the control strategy of the BA model requires the control set $\mathbf{C}$ to include the $c$ nodes with the highest degrees in the network. The test set is the same as in Section IV-B. For each network in the test set, the adjacency matrices are recorded. The indices of the first $\lceil 5\%N \rceil$ nodes with the highest degrees are set as the first $\lceil 5\%N \rceil$ nodes, so that they constitute the control set. The private belief lists with uniform distribution on $[-1, 1]$ are sampled 100 times.

For the networks in the test set of each sub-category network, the converged beliefs and control power of each network are calculated according to Equation (6), based on each of the 100 private belief lists. The control power averaged on these 100 private belief lists and the 25 networks is then recorded as the average control power of this sub-category of networks. The results are shown in Figure (5) and Figure (6).

*2) Generalized Markov Graph Model:* In a GMG model, according to Theorem III.5, the control set should include nodes with highest $d(1 + \gamma)^\alpha$ in order to achieve maximum control power. In this experiment, $\alpha$ is learned from either complete or partial information of the training set. The test set is the same as that in Section IV-B. For each network in the test set, the adjacency matrix is recorded. The indices of the first $\lceil 5\%N \rceil$ nodes with highest $d(1 + \gamma)^\alpha$ are selected for the control set. In order to utilize all the sample networks in the test set, the ones with $\beta_2$ values that do not satisfy the condition in Theorem III.5 are also included. Note that this will yield suboptimal control power values for a GMG model. However, as the experiments show, the control strategy from the GMG model still outperforms the best one from the BA model.

The private belief lists with uniform distribution on $[-1, 1]$ are sampled 100 times. For each sub-category of networks, according to Equation (6), the converged beliefs and control power of each network in the test set are calculated based on each of the 100 private belief list for both $\alpha$ values as shown in Table V and Table VI. The control power averaged across these 100 private belief lists and the 25 networks is recorded as the average control power of this sub-category of networks for each value of $\alpha$. The results are shown in Figure (5) and Figure (6).
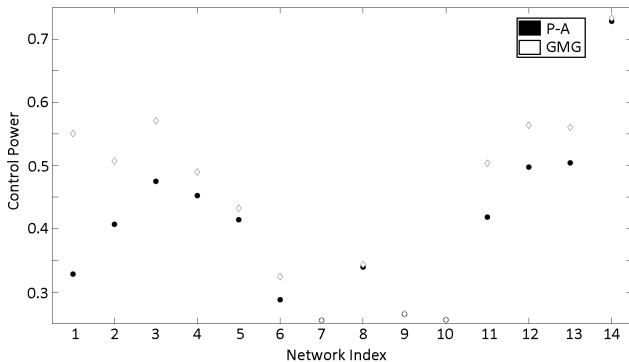


Fig. 5: Control Power of networks following control strategy of the BA model and the GMG model with complete information of training set.
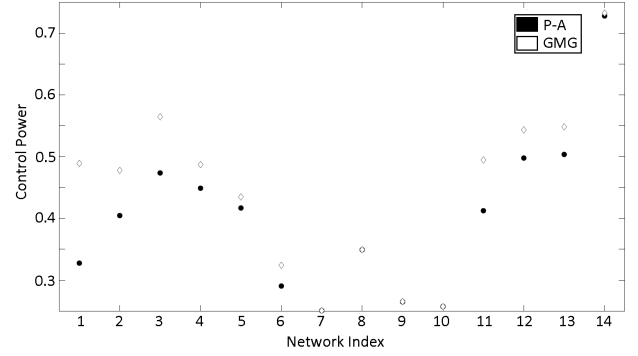


Fig. 6: Control Power of networks following control strategy of the BA model and the GMG model with partial information of training set.

From the experimental results, we can see that, with the same budget on the control set, the control strategy of the GMG model generates a higher control power than that of the BA model. Although the performance of the GMG model with partial information of the training set on control power estimation is significantly worse than the one with full information, their performance on control strategy are very close to each other.

## V. CONCLUSION AND FUTURE RESEARCH

In this paper, we introduced an information flow model to simulate the information flow in a social network. Two social network models are used to optimize the control strategy, and compute the converged beliefs of agents in a social network. Compared to a direct calculation of the converged beliefs, these two models use less information and require less computational power, but still provide results with good accuracy. In addition, the GMG Model outperforms the BA model both in belief estimation accuracy and belief control strategy, since it has a more realistic assumption and uses more information.

Future work includes better machine learning techniques for computing $\alpha$, potentially theoretical calculation for $\alpha$ from network models, and more complicated information flow models. Other future open questions include temporal variability of social network features and their influence on belief estimation.

## APPENDIX A
## NETWORK SYNTHESIS IN BARABÁSI-ALBERT MODEL

The input of the synthesis process in the BA model is the total number of nodes $N$ and the average number of edges attached to each new incoming node $m$. At the very beginning of the process, when the time step $t = 1$, there are $m+1$ nodes in the network and they form a fully connected network. The first node is then introduced to the network with $m$ edges attached. Each edge will choose a node to connect in the network based on its degree and clustering coefficient. At time step $t$, when the $\hat{m}^{th}$ edge of the $m + 1 + t$ node is seeking

another node to attach to, an arbitrary node $i$ in the network is chosen with probability $^1p_i(t, \hat{m})$:

$$^1p_i(t, \hat{m}) = \frac{d_i(a + \gamma_i)^\alpha}{\Sigma_{j=1}^{m+t} d_i(a + \gamma_i)^\alpha},$$

according to Equation (8). After all $m$ edges of the first node are connected to the nodes in the network, the time step $t$ changes to $t = 2$, and the second node is introduced. At an arbitrary time step $t$, where $1 < t \leq N - m$, a new node is introduced to the network with $m$ edges attached. Each edge will choose a node in the network with probability proportional to $^1p_i(t, \hat{m})$. The process ends when there are $N$ nodes in the network and, the edges brought by the last node are all connected with the nodes in the network. The degree distribution of networks generated by this process will follow a power law distribution [2]. The flow chart of this process is shown in Figure (7).
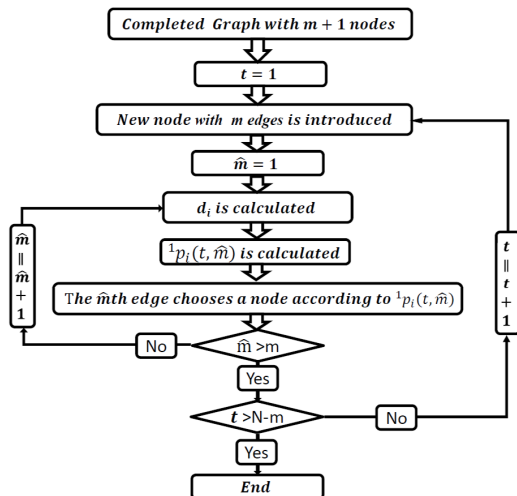


Fig. 7: Flow chart for network synthesis in a BA model.

## B  NETWORK SYNTHESIS IN A GENERALIZED MARKOV GRAPH MODEL

In a GMG model, the network synthesis includes the following steps [18]. The input of the synthesis process is: the total number of nodes $N$, the average number of edges attached to each node $m$, and the clustering weight $\alpha$. At the very beginning of the process, at time step $t = 1$, there are $m + 1$ nodes in the network and they form a fully connected network. The first node is then introduced to the network with $m$ edges attached. Each edge will choose a node in the network to connect to by its degree and clustering coefficient. At time step $t$, when the $\hat{m}^{th}$ edges of the $m + 1 + t$ node is seeking another node to attach to, an arbitrary node $i$ in the network is chosen with probability $^2p_i(t, \hat{m})$:

$$^2p_i(t, \hat{m}) = \frac{d_i(a + \gamma_i)^\alpha}{\Sigma_{j=1}^{m+t} d_i(a + \gamma_i)^\alpha},$$

according to Equation (14). After all $m$ edges of the first node are connected to the nodes in the network, the time step $t$ changes to $t = 2$ and the second node is introduced. At

an arbitrary time step $t$, where $1 < t \leq N - m$, a new node is introduced to the network with $m$ edges attached. The $\hat{m}^{th}$ edge will then choose a node in the network with probability $^2p_i(t, \hat{m})$. The process ends when there are $N$ nodes in the network and the edges brought by the last node are all connected to the nodes in the network. The degree distribution of networks generated by a GMG model can change according to the choice of value of $\alpha$ [19]. The flow chart of this process is shown in Figure (8).
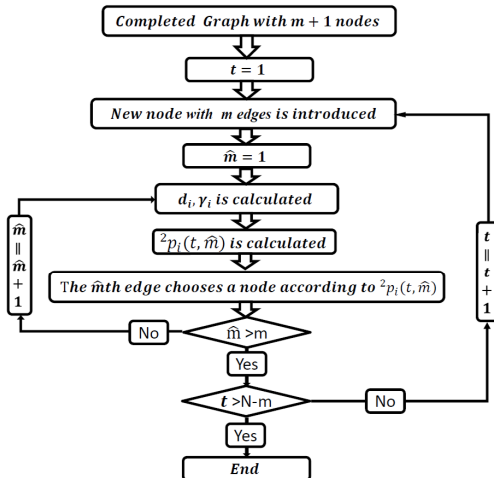


Fig. 8: Flow chart for network synthesis in a GMG model.

## C  PROOF FOR THEOREM III.1

The master equation and the solution of the master equation for the BA model is introduced in Barabási(2002) [2]. As a reminder, we restate it here.

In a BA model, at each time step $t$, where $t$ lies in the range $[1, N]$, and $N$ is the highest node index at the last time step in the network, a new node with $m$ edges is added to the network. We denote this node by its joining time $t$. The probability of an existing node $i$ in the network be connected to the new node $t$ is proportional to its degree $d_i(t)$ at time $t$. This means an increasing rate of degree of an existing node $i$ in the network is proportional to its degree, as shown in the master equation, Equation (21) [2]:

$$\frac{\partial d_i(t)}{\partial t} = m \frac{d_i(t)}{\Sigma_{j=1}^{t} d_j(t)}. \tag{21}$$

And the summation on the denominator of Equation (21) is known to be $2mt$, as in each time step $m$ edges are added to the network. So we can rewrite Equation (21) as,

$$\frac{\partial d_i(t)}{\partial t} = \frac{d_i(t)}{2t}, \tag{22}$$

and the solution to this equation is [2]:

$$d_i(t) = m\sqrt{t/i} \tag{23}$$

In the BA model, the probability of a node $i$ be connected with node $j$ is the probability of the following event happening: at time step $max(i,j)$, node $min(i,j)$ is chosen to be

attached to one of the $m$ edges brought by the new node. According to the basic assumption of the model and Equation (21), such a probability will be:

$$^1P_{i,j} = d_{min(i,j)}(max(i,j))/(2max(i,j)). \qquad (24)$$

According to Equation (23), the degree of node $min(i,j)$ at time step $max(i,j)$ is:

$$d_{min(i,j)} = m\sqrt{max(i,j)/min(i,j)}, \qquad (25)$$

Plugging Equation (25) into Equation (24), yields:

$$^1P_{i,j} = m/2\sqrt{1/(ij)}. \qquad (26)$$

At an arbitrary time step $t$, the degree of node $i$ and $j$ are: $m\sqrt{t/i}$ and $m\sqrt{t/j}$. The degree summation at time step $t$ is $\Sigma_{k=1}^t d_k(t) = 2mt$, which together with Equation (26) leads to:

$$^1P_{i,j} = d_i(t)d_j(t)/\Sigma_{k=1}^t d_k(t). \qquad (27)$$

If the number of nodes in the network reaches N, which means the time step is at: $t = N$, and we omit the time notation for degree $d_i(t)$, we get:

$$^1P_{i,j} = d_i d_j / \Sigma_{k=1}^N d_k. \qquad (28)$$

This concludes the proof. ∎

### D PROOF FOR THEOREM III.2

According to Equation (5) and the definition of control vector $\mathbf{V}$, in a network $G$ of $N$ nodes with a control strategy, the expected belief of an arbitrary non-control node $i$, where $i \notin \mathbf{C}$, at an arbitrary time step $T$, where $T \in \mathbb{Z}^+$, will be,

$$\overline{^1B}_{i,t} = \Sigma_{t=0}^{T-1}\Sigma_{j=1}^c B_{\theta_j}^* \overline{(A^* \times M)^t}_{\theta_j,i}$$
$$+\Sigma_{t=0}^{T-1}\Sigma_{j=c+1}^N \overline{w^*}_{\theta_j} \overline{(A^* \times M)^t}_{\theta_j,i}. \qquad (29)$$

The matrix elements $(A^* \times M^t)_{i,j}$, with $A_{i,j}^* = A_{i,j}/(1+d_j)$ and the definition of control matrix $\mathbf{M}$, can be written as:

$$(A^* \times M)_{i,j}^t$$
$$=\Sigma_{k_1=c+1}^N \cdots \Sigma_{k_{t-1}=c+1}^N \frac{A_{i,\theta_{k_1}}}{1+d_{\theta_{k_1}}}\frac{A_{\theta_{k_1},\theta_{k_2}}}{1+d_{\theta_{k_2}}}\cdots\frac{A_{\theta_{k_{t-1}},j}}{1+d_{\theta_{k_j}}}. \qquad (30)$$

To calculate the expected value of $(A^* \times M^t)_{i,j}$, we combine the expected value of $A_{i,j}$, which equates $P_{i,j}$ as shown in Equation (9), together with Equation (30) to yield

$$\overline{(A^* \times M)^t}_{i,j} = \frac{d_i d_j}{(1+d_j)(\Sigma_{n=1}^N d_n)}\left(\frac{\Sigma_{k=c+1}^N \frac{d_{\theta_k}^2}{1+d_{\theta_k}}}{\Sigma_{n=1}^N d_n}\right)^{t-1}. \qquad (31)$$

Let

$$\beta_1 = \Sigma_{k=c+1}^N \frac{d_{\theta_k}^2}{1+d_{\theta_k}}/\Sigma_{n=1}^N d_n.$$

Because degree $d_k$ is a positive integer: $d_k \geq 1$, we have,

$$d_k > \frac{d_k^2}{1+d_k},$$

and

$$\Sigma_{n=1}^N d_n > \Sigma_{k=c+1}^N d_{\theta_k}^2/(1+d_{\theta_k}).$$

So $\beta_1 < 1$, and the larger $c$ of the control group is, the smaller $\beta_1$ is.

Inserting Equation (31) into Equation (29), and seting $T = \infty$ yields

$$\overline{^1B}_{i,\infty} = \Sigma_{j=1}^c B_{\theta_j}^* d_i d_{\theta_j} \Sigma_{t=0}^\infty \beta_1^{t-1}/(1+d_i)/(\Sigma_{n=1}^N d_n)$$
$$+\Sigma_{j=c+1}^N w_{\theta_j} d_i d_{\theta_j} \Sigma_{t=0}^\infty \beta_1^{t-1}/((1+d_i)(1+d_{\theta_j})(\Sigma_{n=1}^N d_n)). \qquad (32)$$

As $\beta_1 < 1$, we can rewrite $\Sigma_{t=0}^\infty \beta_1^{t-1}$ as $\frac{1}{1-\beta_1}$ and Equation (32) becomes

$$\overline{^1B}_{i,\infty} = \frac{1}{\Sigma_{n=1}^N d_n}\frac{d_i}{1+d_i}\frac{\Sigma_{j=1}^c B_{\theta_j}^* d_{\theta_j} + \Sigma_{j=c+1}^N \frac{\overline{w}_{\theta_j}}{1+d_{\theta_j}}d_{\theta_j}}{1-\beta_1}. \qquad (33)$$

This concludes the proof. ∎

### E PROOF FOR THEOREM III.4

The proof of Theorem III.4 will be similar to that of Theorem III.2. Based on Equation (5) and the definition of control vector $\mathbf{V}$, in a network $G$ of $N$ nodes with control strategy, the expected belief of an arbitrary non-control node $i$, where $i \notin \mathbf{C}$, at an arbitrary time step $T$, where $T \in \mathbb{Z}^+$, will be:

$$\overline{^2B}_{i,t} = \Sigma_{t=0}^{T-1}\Sigma_{j=1}^c B_{\theta_j}^* \overline{(A^* \times M)^t}_{\theta_j,i}$$
$$+\Sigma_{t=0}^{T-1}\Sigma_{j=c+1}^N \overline{w^*}_{\theta_j} \overline{(A^* \times M)^t}_{\theta_j,i}. \qquad (34)$$

The matrix element $(A^* \times M^t)_{i,j}$, with $A_{i,j}^* = A_{i,j}/(1+d_j)$ and the definition of control matrix $\mathbf{M}$, may be written as:

$$(A^* \times M)_{i,j}^t$$
$$=\Sigma_{k_1=c+1}^N \Sigma_{k_2=c+1}^N \cdots \Sigma_{k_{t-1}=c+1}^N \frac{A_{i,k_1}}{1+d_{k_1}}\frac{A_{k_1,k_2}}{1+d_{k_2}}\cdots\frac{A_{k_{t-1},j}}{1+d_{k_j}}. \qquad (35)$$

$\overline{(A^* \times M)^t}_{i,j}$ is then calculated from the expected value of $A_{i,j}$, which equates $^2P_{i,j}$ as shown in Equation (15), and Equation (35):

$$\overline{(A^* \times M)^t}_{i,j}$$
$$= \frac{d_i(1+\gamma_i)^\alpha d_j(1+\gamma_j)^\alpha(\Sigma_{n=1}^N d_n)}{(1+d_j)\eta}$$
$$\left(\Sigma_{k=c+1}^N \frac{(d_{\theta_k}(1+\gamma_{\theta_k})^\alpha)^2}{(1+d_{\theta_k})\eta}\Sigma_{n=1}^N d_n\right)^{t-1}. \qquad (36)$$

Let

$$\beta_2 = \frac{\Sigma_{k=c+1}^N \frac{(d_{\theta_k}(1+\gamma_{\theta_k})^\alpha)^2}{1+d_{\theta_k}}}{\eta}\Sigma_{n=1}^N d_n.$$

Again inserting Equation (36) into Equation (34), and setting $T = \infty$ yields

$$\overline{^2B_{i,\infty}}$$
$$= \Sigma_{j=1}^c \frac{B_{\theta_j}^* d_i(1+\gamma_i)^\alpha d_{\theta_j}(1+\gamma_{\theta_j})^\alpha (\Sigma_{n=1}^N d_n)}{(1+d_i)\eta} \Sigma_{t=0}^\infty \beta_1^{t-1}$$
$$+ \Sigma_{j=1}^c \frac{\frac{\overline{w}_{\theta_j}}{1+d_{\theta_j}} d_i(1+\gamma_i)^\alpha d_{\theta_j}(1+\gamma_{\theta_j})^\alpha (\Sigma_{n=1}^N d_n)}{(1+d_i)\eta} \Sigma_{t=0}^\infty \beta_1^{t-1}.$$
$$(37)$$

If $\beta_2 < 1$, we can rewrite $\Sigma_{t=0}^\infty \beta_2^{t-1}$ as $\frac{1}{1-\beta_2}$, and Equation (37) becomes

$$\overline{^2B_{i,\infty}} = \frac{\Sigma_{n=1}^N d_n}{\eta} \frac{d_i(1+\gamma_i)^\alpha}{1+d_i}$$
$$\frac{\Sigma_{j=1}^c B_{\theta_j}^* d_{\theta_j}(1+\gamma_{\theta_j})^\alpha + \Sigma_{j=c+1}^N \frac{\overline{w}_{\theta_j}}{1+d_{\theta_j}} d_{\theta_j}(1+\gamma_{\theta_j})^\alpha}{1-\beta_2}$$
$$(38)$$

This concludes the proof. ∎

## F Proof for Theorem III.3

Consider an arbitrary control set $\mathbf{C}' = \{\theta_1', \theta_2', \ldots, \theta_c'\}$ with corresponding degrees $\{d_{\theta_1'}, d_{\theta_2'}, \ldots, d_{\theta_c'}\}$. Without loss of generality, we set $d_{\theta_1'} \geq d_{\theta_2'} \geq \cdots \geq d_{\theta_c'}$. For control set $\mathbf{C_o} = \{\theta_{o1}, \theta_{o2}, \ldots, \theta_{oc}\}$, the corresponding degrees satisfy $d_{\theta_{oi}} \geq d_{\theta_{oj}}$ if $i \leq j$, $1 \leq i, j \leq N$, so, for $1 \leq i \leq c$, we have

$$d_{\theta_{oi}} \geq d_{\theta_i'}. \tag{39}$$

Since $\overline{w}_i = 0$, for $i = 1, \ldots, N$, $B_j^* = 1$, for $1 \leq j \leq c$, according to Equation (12), $^1cp$ can be rewritten as:

$$^1cp = \Sigma_{i=1}^N \frac{d_i}{1+d_i}(\Sigma_{j=1}^c d_{\theta_j'})/(\Sigma_{k=1}^N d_k - \Sigma_{m=c+1}^N \frac{d_{\theta_m'}^2}{1+d_{\theta_m'}}). \tag{40}$$

Since $\Sigma_{i=1}^N \frac{d_i}{1+d_i}$ and $\Sigma_{k=1}^N d_k$ are positive constants, the derivative of control power $^1cp$ with respect to the degree of a controlled node $d_{\theta_i'}$ is:

$$\frac{\partial^1 cp}{\partial d_{\theta_i'}} = D_1 \frac{\Sigma_{k=c+1}^N \frac{d_{\theta_k'}}{1+d_{\theta_k'}} + \frac{\Sigma_{j=1}^c d_{\theta_j'}}{(1+d_{\theta_i'})^2}}{D_2^2}, \tag{41}$$

where $D_1 = \Sigma_{j=1}^N \frac{d_j}{1+d_j}$, $D_2 = \Sigma_{k=1}^N \frac{d_k}{1+d_k} + \Sigma_{i=c+1}^N \frac{d_{\theta_i'}}{1+d_{\theta_i'}}$. As $D_1 > 0$, one can get:

$$\frac{\partial^1 cp}{\partial d_{\theta_i'}} > 0. \tag{42}$$

From Equations (39) and (42), one can get:

$$^1cp(\mathbf{C_o}) \geq {}^1cp(\mathbf{C}'). \tag{43}$$

This concludes the proof. ∎

## G Proof for Theorem III.5

Define the weight $\xi_i$ of node $i$ as: $\xi_i = d_i(1+\gamma_i)^\alpha$. Consider an arbitrary control set $\mathbf{C}' = \{\theta_1', \theta_2', \ldots, \theta_c'\}$ with corresponding weight $\{d_{\xi_1'}, d_{\xi_2'}, \ldots, d_{\xi_c'}\}$. Without loss of generality, we set $d_{\xi_1'} \geq d_{\xi_2'} \cdots \geq d_{\xi_c'}$. For control set $\mathbf{^2C_o} = \{\theta_{o1}^\dagger, \theta_{o2}^\dagger, \ldots, \theta_{oc}^\dagger\}$, the corresponding weights satisfy $\xi_{\theta_{oi}^\dagger} \geq \xi_{\theta_{oj}^\dagger}$ if $i \leq j$, $1 \leq i, j \leq N$, so, for $1 \leq i \leq c$, we have

$$\xi_{\theta_{oi}^\dagger} \geq \xi_{\theta_i'}. \tag{44}$$

Since $\overline{w}_i = 0$, for $i = 1, \ldots, N$, and $B_j^* = 1$, for $1 \leq j \leq c$, according to Equation (18) $^2cp$ can be rewritten as:

$$^2cp = K_1 \frac{\Sigma_{j=1}^c \xi_{\theta_j'}}{K_2 + K_3 \Sigma_{k=1}^c \frac{\xi_{\theta_k'}^2}{1+d_{\xi_{\theta_k'}}}}, \tag{45}$$

where $K_1 = \frac{1}{N}\Sigma_{i=1}^N \frac{\xi_i}{1+d_i}\Sigma_{j=1}^N d_j$,
$K_2 = \Sigma_{i=1}^N \xi_i^2 - \Sigma_{j=1}^N \xi_j - \Sigma_{k=1}^N d_k \Sigma_{m=1}^N \frac{\xi_m^2}{1+d_m}$,
$K_3 = \Sigma_{j=1}^N d_j$.
The derivative of control power $^2cp$ with respect to the degree of a controlled node $d_{\theta_i'}$ is:

$$\frac{\partial^2 cp}{\partial d_{\theta_i'}} = \frac{K_4}{K_5^2}[(\frac{1}{\beta_2}-1)\Sigma_{j=c+1}^N \frac{\xi_{\theta_j'}^2}{1+d_{\theta_j'}} - \Sigma_{k=1}^c \xi_{\theta_k'}(1+\gamma_{\theta_i'})^\alpha$$
$$+\Sigma_{m=1}^c \xi_{\theta_m'}\frac{(1+\gamma_{\theta_i'})^\alpha}{(1+d_{\theta_i'})^2}] \tag{46}$$

where $K_4 = \frac{1}{N}\Sigma_{i=1}^N \frac{\xi_i}{1+d_i}$,
$K_5 = \frac{\Sigma_{i=1}^N \xi_i^2 - \Sigma_{j=1}^N \xi_j}{K_3} - \Sigma_{k=1}^N \frac{\xi_k^2}{1+d_k}$.
Since $K_4 > 0$, we will have

$$\frac{\partial^2 cp}{\partial d_{\theta_i'}} > 0, \tag{47}$$

assuming that $\beta_2$ satisfies:

$$\frac{1}{\beta_2} > 1 + max(1, 2^\alpha)\frac{\Sigma_{j=1}^c d_{\theta_{oj}^\dagger}(1+\gamma_{\theta_{oj}^\dagger})^\alpha}{\Sigma_{k=c+1}^N \frac{(d_{\theta_{ok}^\dagger}(1+\gamma_{\theta_{ok}^\dagger})^\alpha)^2}{1+d_{\theta_{ok}^\dagger}}}. \tag{48}$$

Considering Equation (44), one can get:

$$^2cp(\mathbf{C_o^\dagger}) \geq {}^2cp(\mathbf{C}'), \tag{49}$$

if Equation (48) holds.
This concludes the proof. ∎

## References

[1] Daron Acemoglu, Munther A. Dahleh, Ilan Lobel, and Asuman Ozdaglar. Bayesian learning in social networks. *Review of Economic Studies, Oxford University Press*, 78(4):1201–1236, 2008.

[2] Réka Albert and Albert Barabási. Statistical mechanics of complex networks. *Reviews of Modern Physics*, 74:47–96, January 2002.

[3] Venkatesh Bala and Sanjeev Goyal. Learning from neighbours. *The Review of Economic Studies, Vol. 65, No. 3 (Jul., 1998), pp. 595-621*.

[4] Abhijit V. Banerjee. A simple model of herd behavior. *The Quarterly Journay of Economics*, CVII, Issue 3, August 1992.

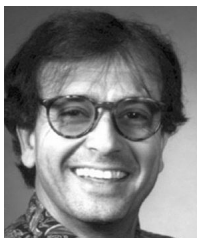[5] Fan Chung and Linyuan Lu. *Complex Graphs and Networks*. the AMS and CBMS, 2006.

[6] Aaron Clauset, Cosma Rohilla Shalizi, and M. E. J. Newman. Power-law distributions in empirical data. *SIAM REVIEW*, Vol. 51, No. 4, pp. 661-703, 2009.

[7] Peter Clifford and Aidan W Sudbury. A model for spatial conflict. *Biometrika*, 60: 581-588, 1973.

[8] M. H. DeGroot. Reaching a consensus. *Journal of the American Statistical Association*, 69:118-121, 1974.

[9] Glenn Ellison and Drew Fudenberg. Word-of-mouth communication and social learning. *The Quarterly Journal of Economics, Vol. 110, No. 1 (Feb., 1995), pp. 93-125.*

[10] Ove Frank and David Strauss. Markov graphs. *Journal of the American Statistical Association*, 81:832–842, September 1986.

[11] Douglas Gale and Shachar Kariv. Bayesian learning in social networks. *Games and Economic Behavior*, 45:329–346, November 2003.

[12] Usman A. Khan, Soummya Kar, and Jos M. F. Moura. Higher dimensional consensus: Learning in large-scale networks. *IEEE Transactions on Singal Processing, VOL. 58, NO. 5, MAY 2010.*

[13] Michael P. H. Stumpf and Mason A. Porter. Critical truths about power laws. *Science*, Vol. 335 no. 6069 pp. 665-666, February 2012.

[14] Katarzyna Sznajd-Weron. Sznajd model and its applications. *Acta Physica Polonica B, vol.36, no. 8 (2005).*

[15] StanFord University. Stanford large network dataset collection, retrieved from: http://snap.stanford.edu/data/, September 2011.

[16] Tian Wang, H. Krim, and Y. Viniotis. A generalized markov graph model: Application to social network analysis. *IEEE Journal of Selected Topics in Signal Processing*, 7:318–332, April 2013.

[17] Tian Wang and Hamid Krim. Control and prediction of beliefs on social network. *IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing. 2013.*

[18] Tian Wang and Hamid Krim. Statistical classification of social networks. *ICASSP*, 2012. ICASSP.

[19] Stanley Wasserman and Philippa Patrlson. Logit models and logistic regressions for social networks. *Psychometrika*, 61:401–425, September 1996.

[20] Duncan J. Watts. A simple model of global cascades on random networks. *Physical Sciences - Applied Mathematics*, 99 (9) 5766-5771, 2002.

**Yannis Viniotis** received his Ph.D. from the University of Maryland, College Park, in 1988 and is currently Professor, Department of Electrical and Computer Engineering at North Carolina State University. Dr. Viniotis is the author of over one hundred technical publications, including two engineering textbooks. He has served as the cochair of two international conferences in computer networking. His research interests include Service-Oriented Architectures, service engineering, and design and analysis of stochastic algorithms. Dr. Viniotis was the cofounder of Orologic, a successful startup networking company in Research Triangle Park, NC, that specialized in ASIC implementation of integrated traffic management solutions for high-speed networks.

**Tian Wang** received his B.S. (2008) degree in Physics from the University of Science and Technology of China and the M.Sc. (2011) degree in electrical engineering and Physics from North Carolina State University, where he been a Ph.D. student in Physics since 2008. His main research interest is in modelling and analysis of social network and information flow on social network.

**Hamid Krim** received his degrees in Electrical Engineering. As a member of technical staff at AT&T Bell Labs, he has worked in the area of telephony and digital communication systems/subsystems. In 1991 he became a NSF Post-doctoral scholar at Foreign Centers of Excellence (LSS Sup-elec/Univ. of Orsay, Paris, France). He subsequently joined the Laboratory for Information and Decision Systems, MIT, Cambridge, MA, as a Research Scientist performing/supervising research in his area of interest, and later as a faculty in the ECE dept. at North Car. State Univ. in Raleigh, N.C. in 1998. He is an original contributor and now an affiliate of the Center for Imaging Science sponsored by the Army. He also is a recipient of the NSF Career Young Investigator award. He is on the editorial board of the IEEE Trans. on SP and regularly contributes to the society in a variety of ways. His research interests are in statistical estimation and detection and mathematical modeling with a keen emphasis on applications.