

# Community extraction for social networks

Yunpeng Zhao, Elizaveta Levina<sup>1</sup>, and Ji Zhu

Department of Statistics, University of Michigan, Ann Arbor, MI 48109

Edited by Stephen E. Fienberg, Carnegie Mellon University, Pittsburgh, PA, and approved March 15, 2011 (received for review May 17, 2010)

**Analysis of networks and in particular discovering communities within networks has been a focus of recent work in several fields and has diverse applications. Most community detection methods focus on partitioning the entire network into communities, with the expectation of many ties within communities and few ties between. However, many networks contain nodes that do not fit in with any of the communities, and forcing every node into a community can distort results. Here we propose a new framework that extracts one community at a time, allowing for arbitrary structure in the remainder of the network, which can include weakly connected nodes. The main idea is that the strength of a community should depend on ties between its members and ties to the outside world, but not on ties between nonmembers. The proposed extraction criterion has a natural probabilistic interpretation in a wide class of models and performs well on simulated and real networks. For the case of the block model, we establish asymptotic consistency of estimated node labels and propose a hypothesis test for determining the number of communities.**

Understanding and modeling network structures have been a focus of attention in a number of diverse fields, including physics, biology, computer science, statistics, and social sciences. Applications of network analysis include friendship and social networks, marketing and recommender systems, the World Wide Web, disease models, and food webs, among others. A fundamental problem in the study of networks is community detection (see refs. 1–3 for comprehensive recent reviews). We focus here on undirected networks  $N = (V, E)$ , where  $V$  is the set of nodes and  $E$  is the set of edges, possibly weighted. The community detection problem is typically formulated as finding the partition  $V = V_1 \cup \dots \cup V_K$ , which gives the “best” communities in some suitable sense. The node sets  $V_1, \dots, V_K$  are usually taken to be disjoint, although there is some recent work on detecting overlapping communities (4, 5).

The extensive physics and computer science literature on networks typically thinks of communities as tightly knit groups with many connections between the group members and relatively few connections between groups. Thus detection methods focus on maximizing links within communities while minimizing links between communities. This can be achieved either implicitly through an algorithmic approach (6) or explicitly by optimizing a criterion that measures the quality of a proposed partition over all possible partitions. These criteria include ratio cuts (7), normalized cuts (8), spectral clustering (9), and modularity (10); see ref. 3 for a review. All of these are symmetric criteria, in the sense that all potential communities play the same role. There are many examples of networks where such a requirement makes sense, for example, the college football games network (11), and yet some commonly studied networks clearly do not fit this framework. One such example is when there are nodes without strong connections to any communities, such as in the high school friendship network of ref. 12 discussed later in the paper. In such cases, the partitioning methods typically split up weakly connected nodes and group them together with tighter communities. There is not much work in the networks literature focusing on such a structure, even though in traditional multivariate clustering there are methods that allow for a diffuse “background” cluster [e.g., DBSCAN (13) and DenClue (14)].

Another class of community detection methods relies on a statistical model for the network to estimate the partition, typically by maximizing some form of the likelihood directly or employing Gibbs sampling. The models used for partitioning include the stochastic block model (15–17), a mixture model (18), univariate (19) and multivariate (20) latent variable models, and latent feature models (21); for a comprehensive review of statistical models of networks, see ref. 2. In particular, the stochastic block model, described in detail later in the paper, allows for any density of connections within communities and can in principle handle any combination of sparse and tight communities. The block model assumes a uniform distribution of edges within a block, and an extension was recently proposed to accommodate arbitrary degree distributions within the blocks [the degree-corrected block model (22)].

In this paper, we propose a framework for community extraction that looks for one community at a time. Intuitively, our extraction criterion, like other partition methods, looks for a “tight” group with more links within itself than to the rest of the network; but, unlike partition methods, it allows for the rest of the network to include an arbitrary mixture of tight and weak communities or sparsely connected background nodes. Mathematically, our criterion matches the definition of community in a large class of probability models on networks, outlined below, which can be thought of as a generalization of the block model with some parameter constraints. Our goal is to extract the tightest community, which we do by focusing on the edges within the candidate community and edges connecting it to the rest of the network, and ignoring edges within the rest of the network. The key feature of the extraction criterion is that it is not symmetric in the two sets into which we are splitting the network. In practice, such extraction can be used on its own or in conjunction with graph partitioning, for example, to identify community cores.

As an illustration, consider this toy example: Out of  $n = 60$  nodes, 15 belong to a community where links between members form independently with probability 0.5. The links between members and the other 45 nodes and links between the other 45 nodes all form independently with probability 0.1. A partition into two communities using modularity and our community extraction method are shown in Fig. 1. Modularity has to balance tightness of the two communities, and as a result includes a number of background nodes in the community. Extraction, on the other hand, separates out the true community perfectly.

Finally, we briefly mention other related work. The core-periphery partition methods (23, 24) use a different criterion to separate a tight “core” from a sparse “periphery,” whereas our criterion is designed to extract the tightest community regardless of whether the rest of the network is sparse or contains other communities. Local community detection (25, 26) looks for the tightest community around a given node, but not the globally tightest community. Finally, the hierarchical network model

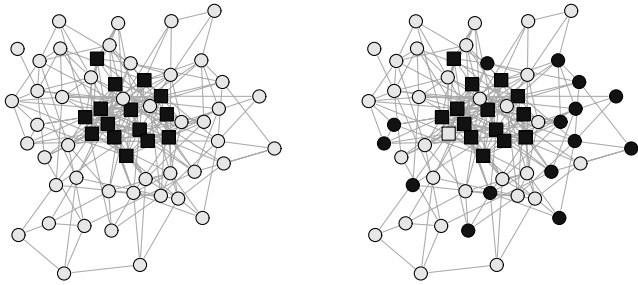
Author contributions: Y.Z., E.L., and J.Z. designed research; Y.Z., E.L., and J.Z. performed research; Y.Z. analyzed data; and Y.Z., E.L., and J.Z. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

<sup>1</sup>To whom correspondence should be addressed. E-mail: elevina@umich.edu.

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1006642108/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1006642108/-DCSupplemental).



**Fig. 1.** Toy example: shapes represent the truth and colors represent results using our extraction method (Left) and modularity (Right).

based on ensembles of trees proposed in ref. 27 also has the feature that the strength of a community does not depend on unrelated nodes.

### The Community Extraction Framework

First we introduce basic notation. A network  $N = (V, E)$  with  $|V| = n$  nodes can be represented by an  $n \times n$  adjacency matrix  $\mathbf{A} = [A_{ij}]$ , where  $A_{ij} > 0$  if there is an edge between nodes  $i$  and  $j$  and  $A_{ij} = 0$  otherwise. If the edges have weights, the positive  $A_{ij}$ 's are the weights; if not, they are set to 1. For undirected networks,  $\mathbf{A}$  is symmetric.

For simplicity, we start with partitioning into two sets  $V_1, V_2$ , where  $V_1 \cap V_2 = \emptyset$  and  $V = V_1 \cup V_2$ . A naive way to partition a network is to minimize the total weight  $R$  of edges connecting  $V_1$  and  $V_2$  (the min-cut method),  $R = \sum_{i \in V_1, j \in V_2} A_{ij}$ . However, minimizing  $R$  yields a trivial solution of  $V_1 = V$ , so various adjustments have been proposed. For example, the ratio cut (7) avoids the trivial solution by minimizing  $R/(|V_1| \cdot |V_2|)$ , where  $|V_1|$  and  $|V_2|$  are the sizes of the two groups. The important point for us here is that in all these criteria the sets  $V_1$  and  $V_2$  can be interchanged.

The criterion we propose extracts one community at a time by looking for a set with a large number of links within itself and a small number of links to the rest of the network. The links within the complement of this set do not affect the value of the criterion. To emphasize the lack of symmetry, we denote the community to be extracted by  $S$  and its complement by  $S^c$  (rather than  $V_1$  and  $V_2$ ). Then we maximize the following extraction criterion over all possible  $S$ :

$$W(S) = \frac{O(S)}{|S|^2} - \frac{B(S)}{|S||S^c|}, \quad [1]$$

where

$$O(S) = \sum_{i,j \in S} A_{ij}, \quad B(S) = \sum_{i \in S, j \in S^c} A_{ij}.$$

The term  $O(S)$  is twice the weight of the edges within  $S$ , and  $B(S)$  represents connections between  $S$  and the rest of the network. Each term is normalized by the total number of possible edges in each case, which gives these quantities a natural interpretation as probability estimates, discussed further below. Note that we normalize the first term by  $|S|^2$  rather than  $|S|(|S| - 1)$ , thus not explicitly excluding self-loops in order to be consistent with the probability models discussed below, but in practice this makes little difference. Subject to this small difference, our criterion can be described as the intracluster density minus the intercluster density; it is also related to conductance (3).

Criterion 1, like the graph min-cut, does not explicitly guard against splitting off small communities. The trivial solution does not maximize  $W$ , but in a large sparse network a very small community can give a high value of  $W$ , because the second term will be made negligible by the large  $|S^c|$ . To avoid this situation, we

make an adjustment in the spirit of the ratio cut and maximize the following criterion instead:

$$\tilde{W}(S) = |S||S^c| \left[ \frac{O(S)}{|S|^2} - \frac{B(S)}{|S||S^c|} \right]. \quad [2]$$

Because  $|S||S^c|$  is maximized at  $|S| = n/2$ , this factor penalizes very small and very large communities and produces more balanced solutions. Empirically, we found the adjustment helps in sparse networks, but plays no role in dense networks. Later we show that asymptotically both criteria are consistent under an appropriate probability model.

**Probabilistic Interpretation.** The criterion (Eq. 1) is motivated by the intuitive definition of community as a tightly knit group rather than by fitting a probability model to the network. However, it has a clear interpretation in the context of probability models on graphs. Consider a probability distribution  $P$  on symmetric adjacency matrices  $A$  that satisfies the following assumptions:

1. Each node  $i$  is associated with latent variables  $c_i$  and  $\theta_i$ , where  $c_i$  is the community label and  $\theta_i$  can contain any other node information. The labels  $c_i$  are independent and have a multinomial distribution with parameter  $\pi = (\pi_1, \dots, \pi_K)^T$ .
2. For any pair  $(i, j)$  and  $(i', j')$  that have no nodes in common,  $A_{ij}$  and  $A_{i'j'}$  are independent given the labels  $\mathbf{c}$ .
3.  $E(A_{ij} | \mathbf{c}) = P_{c_i c_j}$  for all  $i, j$ .

Assumption 1 is similar to the block model, except it allows for additional latent variables  $\theta_i$ . An assumption related to assumption 2 was proposed in ref. 28, stipulating independence conditional on other edges. In assumption 3, the expectation integrates out  $\theta$ , so the left-hand side is always a function of  $\mathbf{c}$ ; the assumption stipulates that it can depend only on the two labels  $c_i$  and  $c_j$ , as in the block model.

Let  $\mathbf{s}$  be an arbitrary label assignment, let  $O_{ab}(\mathbf{s}) = \sum_{ij} A_{ij} I(s_i = a, s_j = b)$ , and  $R_{ab}(s, \mathbf{c}) = n^{-1} \sum_{i=1}^n I(s_i = a, c_i = b)$ . Suppose  $R(s, \mathbf{c}) \xrightarrow{P} R$  as  $n \rightarrow \infty$ . Note that  $R$  satisfies  $\mathbf{1}^T R = \pi$ . Then the following holds:

**Theorem 1.** (a) Under assumptions 1–3, as  $n \rightarrow \infty$ ,

$$\frac{O_{ab}(\mathbf{s})}{n^2} \xrightarrow{P} (RPR^T)_{ab}.$$

(b) Assume that  $K = 2$ , and, without loss of generality,  $P_{11} \geq P_{22}$ . Then  $W(S) \xrightarrow{P} f(R, P)$ ,  $n^{-2} \tilde{W}(S) \xrightarrow{P} \tilde{f}(R, P)$ . Further, if  $P_{11} + P_{22} > 2P_{12}$ , then both  $f$  and  $\tilde{f}$  are maximized by  $R = \text{diag}(\pi)$ , under the constraint  $\mathbf{1}^T R = \pi$ . The proof and the expressions for functions  $f$  and  $\tilde{f}$  can be found in the *SI Text*. The theorem establishes that our extraction criterion is a natural data-based approximation to a population criterion that is maximized by the correct label assignment. Note that even though in part (b) we have  $K = 2$  because it applies to a single extraction step, this does not mean that extraction fails if  $K > 2$  (see more on this in Example 4). Next, we give several examples of models that satisfy our conditions.

**Example 1:** The stochastic block model corresponds to  $\theta_i = \text{const}$  and independent edges. Under the block model, each node is first assigned to one of the  $K$  blocks independently of other nodes. Then, conditional on  $\mathbf{c}$ , edges are generated independently with probabilities  $P[A_{ij} = 1 | c_i = a, c_j = b] = P_{ab}$ , which evidently satisfies assumptions 2 and 3. The block model is fully specified by its

parameters  $\pi$  and the  $K \times K$  symmetric matrix  $P$ . Thus part (a) of the theorem applies. Part (b) applies if we additionally assume  $K = 2$  and  $P_{11} + P_{22} > 2P_{12}$ .

**Example 2:** The degree-corrected block model. This generalization of the block model has been recently introduced to allow for different expected node degrees within blocks (22). There the labels  $c$  are treated as fixed and fixed node-specific parameters  $\theta$  are introduced to reflect the nodes' individual tendencies to form ties. We use exactly the same model but treat  $c$  and  $\theta$  as random, with  $c$  satisfying our assumption 1. Conditional on  $c$  and  $\theta$ ,  $A_{ij}$  are independent Poisson random variables with mean  $\theta_i \theta_j P_{c_i c_j}$ . The relaxation of the Bernoulli distribution to the Poisson is primarily for ease of technical derivations in ref. 22 and has few practical consequences. In either distributional form, the model satisfies assumptions 1–3 and part (a) holds. The model can be further constrained to satisfy conditions for (b).

**Example 3:** General exchangeable distributions. It is well known that any probability distribution invariant to node permutations on the matrix  $A$  can be written in the form  $A_{ij} = h(\mu, \xi_i, \xi_j, \lambda_{ij})$ , where  $\mu$ ,  $\xi_i$ 's, and  $\lambda_{ij}$  are independent and identically distributed (i.i.d.), and  $h$  is symmetric in its second and third arguments (for details, see, e.g., ref. 29). The equivalent of i.i.d. sequences in this class has the form  $A_{ij} = h(\xi_i, \xi_j, \lambda_{ij})$ . Let  $\xi_i = (c_i, \theta_i)$ , where  $c_i$  is the class label, and  $\theta_i$  is other node information that can be correlated with  $c_i$ . If the distribution of i.i.d. pairs  $(c_i, \theta_i)$  is such that the marginal distribution of  $c$  is multinomial, this model satisfies conditions 1–3. One concrete example of such a model is the latent position cluster model of ref. 20. Again, the general model can be further constrained to satisfy (b).

**Example 4:** Consider a block model with true  $K = 3$  where two blocks have been merged under one label. Let  $c_i$  (merged labels) take values 1 and 2 with probabilities  $\pi_c$  and  $1 - \pi_c$ , and let  $\theta_i$  (true labels) have the distribution  $P(\theta_i = 1 | c_i = 1) = 1$ ,  $P(\theta_i = 2 | c_i = 2) = \pi_\theta$ , and  $P(\theta_i = 3 | c_i = 2) = 1 - \pi_\theta$ . If  $Q$  is the  $3 \times 3$  matrix corresponding to the true block model, and  $\Pi$  is the  $2 \times 3$  matrix giving the joint probability distribution of the pair  $(c, \theta)$ , let  $P = \Pi Q \Pi^T$ . Then the population version of our criterion is maximized by extracting the first community as long as  $P_{11} > P_{12}$ ,  $P_{11} > P_{22}$ ,  $P_{11} + P_{22} > 2P_{12}$ , which ensures the condition of part (b) holds. This is true, for example, if  $Q_{11} > Q_{ij}$  for all  $(i, j) \neq (1, 1)$  and  $Q_{ij} \geq Q_{1k}$  for all  $1 < i, j, k \leq 3$ . In other words, if the ties within the first community are stronger than its ties to the mixture of second and third, the criterion will extract it correctly. Note that this is exactly the situation in the counterexample to consistency of modularity given by ref. 29. In that example, both profile likelihood and extraction (in two stages) are consistent, but modularity is not.

**Maximizing the Extraction Criterion.** To maximize the extraction criteria, we use a local optimization technique based on label switching known as tabu search (30, 31). The key idea of tabu search is that once a node label has been switched, it cannot be switched again for the next  $T$  iterations (the node has “tabu” status) This guards against being trapped in a local maximum. The algorithm starts from an initial value and examines all current nontabu nodes in order. If the current value of the global maximum can be improved, the node label is switched, its status changed to tabu, and the algorithm returns to node 1. If no node can be switched to improve the global maximum, the node that gives the largest increase in the current criterion value is switched, and if no increase is possible, the node that gives the smallest decrease is switched. The algorithm is run for a prescribed number of iterations, and the best solution seen in the course of these iterations is taken to be the final solution. Note

that the value of  $\tilde{W}$  can be updated efficiently in  $O(n)$  operations for a single label switch. To help guard against local maxima, we run the algorithm for a number of random starting values and random orderings of nodes. Each run will converge to a local maximum, and although the algorithm is not guaranteed to find the global maximum, we have not encountered any problems with local maxima in either simulations or real data examples.

### The Stochastic Block Model Case

If we focus on the special case of the stochastic block model, we can obtain additional results on the properties of the extraction criterion. First, we show that the estimated node labels are asymptotically consistent, using the recent results of ref. 29. Second, we describe a hypothesis test that can be performed at every sequential split to determine whether the remainder after extraction contains any more communities.

**Asymptotic Consistency.** We consider asymptotic consistency of label assignments by the extraction method as the number of nodes  $n \rightarrow \infty$ . If  $P$  does not change with  $n$ , the network becomes very dense as  $n$  grows, so we allow  $P_n$  to depend on  $n$  and write  $P_n = \rho_n P$ , where  $\rho_n = P[A_{ij} = 1] \rightarrow 0$  is the probability of an edge between arbitrary nodes  $i$  and  $j$ . The expected node degree  $\lambda_n = n\rho_n$  becomes the natural parameter to control as  $n \rightarrow \infty$ .

Bickel and Chen (29) developed a general framework for checking whether a community-finding criterion can recover the true node labels as  $n \rightarrow \infty$ , under the assumptions of the block model and  $\lambda_n / \log n \rightarrow \infty$ . The latter may not be universally applicable, but in many examples the degree does grow with  $n$ , and faster than logarithmic growth is a very mild requirement. Further details are given in *SI Text*; briefly, the main condition is that the proposed criterion is maximized by the true label assignment when all the sample quantities in the criterion are replaced by their population equivalents. This can be viewed as a special case of the theory of minimum contrast estimation (32).

We focus on checking one-step consistency for the case  $K = 2$  (one extracted community plus the rest of the network). The matrix  $P$  is  $2 \times 2$  with three unique parameters  $P_{11}$ ,  $P_{22}$ ,  $P_{12}$ , and the vector of class probabilities  $\{\pi, 1 - \pi\}$  is determined by the single parameter  $\pi$ . Let  $\hat{c}^{(n)}$  be the maximizer of criterion 1, and  $\tilde{c}^{(n)}$  of criterion 2. It turns out that the adjustment factor of  $|S||S^c|$  has no effect in the limit, and both criteria are asymptotically consistent, as shown in the following theorem:

**Theorem 2.** Suppose  $\lambda_n / \log n \rightarrow \infty$ ,  $P_{11} > P_{12}$ ,  $P_{11} > P_{22}$  and  $P_{11} + P_{22} > 2P_{12}$ , and  $c$  are the true labels. Then

$$P[\hat{c}^{(n)} = c] \rightarrow 1 \quad \text{and} \quad P[\tilde{c}^{(n)} = c] \rightarrow 1.$$

Note that the simplest case of our toy example (one community with other weakly connected nodes,  $P_{12} = P_{22} = p$ ) is covered by the theorem as long as  $P_{11} > p$ . The proof is given in *SI Text*.

**Determining the Number of Communities.** The full extraction procedure consists of sequentially applying criterion 2: We extract a community and apply the extraction again to its complement. Ideally, the user would have information on the true or desired number of communities  $K$  to be extracted. In the absence of such prior information, determining the number of communities in a network is an open problem, and a rigorous solution would require fully specifying a statistical model. Here we propose a hypothesis test that can be used under the assumption of the block model. In this case, we need to test the null hypothesis  $H_0: K = 1$  against the alternative  $H_a: K \geq 2$  for the subgraph induced by the nodes in the remainder.  $H_0$  means that the remainder is an Erdos–Renyi (ER) graph, where all edges form independently with the same probability. This approach has an







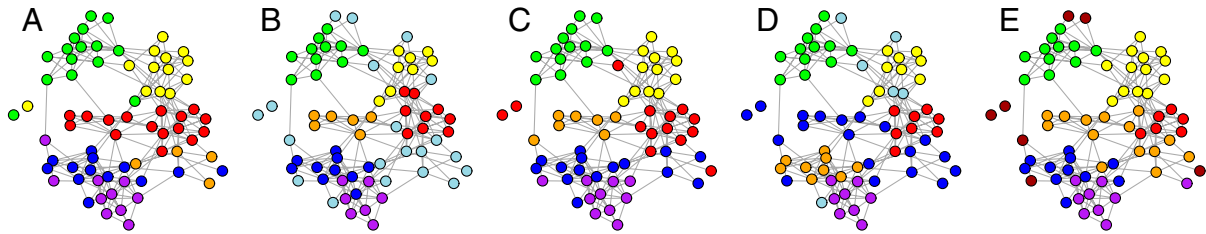


Fig. 6. The school friendship network. (A) grades, (B) extraction, (C) modularity, (D) the block model fitted via MCMC, and (E) the latent position cluster model.

## Acknowledgments

We thank the editor and two referees for many helpful comments and Mark Newman for constructive feedback and sharing his

code. This work is partially supported by National Science Foundation Grants DMS-0805798 (to E.L.) and DMS-0748389 (to J.Z.).

- Newman MEJ (2006) Modularity and community structure in networks. *Proc Natl Acad Sci USA* 103:8577–8582.
- Goldenberg A, Zheng AX, Fienberg SE, Airoldi EM (2009) A survey of statistical network models. *Found Trends Mach Learn* 2:129–233.
- Fortunato S (2010) Community detection in graphs. *Phys Rep* 486:75–174.
- Palla G, Derenyi I, Farkas I, Vicsek T (2005) Uncovering the overlapping community structure of complex networks in nature and society. *Nature* 435:814–818.
- Airoldi EM, Blei DM, Fienberg SE, Xing EP (2008) Mixed membership stochastic blockmodels. *J Mach Learn Res* 9:1981–2014.
- Newman MEJ (2004) Detecting community structure in networks. *Eur Phys J B* 38:321–330.
- Wei Y-C, Cheng C-K (1989) Toward efficient hierarchical designs by ratio cut partitioning. *Proceedings of the IEEE International Conference on Computer Aided Design (IEEE, New York)*, pp 298–301.
- Shi J, Malik J (2000) Normalized cuts and image segmentation. *IEEE T Pattern Anal* 22:888–905.
- Ng A, Jordan M, Weiss Y (2001) On spectral clustering: Analysis and an algorithm. *Neural Information Processing Systems 14*, eds T Dietterich, S Becker, and Z Ghahramani (MIT Press, Cambridge, MA), pp 849–856.
- Newman MEJ, Girvan M (2004) Finding and evaluating community structure in networks. *Phys Rev E* 69:026113.
- Girvan M, Newman MEJ (2002) Community structure in social and biological networks. *Proc Natl Acad Sci USA* 99:7821–7826.
- Hunter DR, Goodreau SM, Handcock MS (2008) Goodness of fit of social network models. *J Am Stat Assoc* 103:248–258.
- Ester M, Kriegel H-P, Sander J, Xu X (1996) A density-based algorithm for discovering clusters in large spatial databases with noise. *Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining (KDD-96)* (AAAI Press, Menlo Park, CA), pp 226–231.
- Hinneburg A, Keim DA (1998) An efficient approach to clustering in large multimedia databases with noise. *Proceedings of the 4th International Conference on Knowledge Discovery and Data Mining (KDD-98)* (AAAI Press, Menlo Park, CA), pp 58–65.
- Holland PW, Laskey KB, Leinhardt S (1983) Stochastic blockmodels: first steps. *Soc Networks* 5:109–137.
- Snijders T, Nowicki K (1997) Estimation and prediction for stochastic block-structures for graphs with latent block structure. *J Classif* 14:75–100.
- Nowicki K, Snijders TAB (2001) Estimation and prediction for stochastic blockstructures. *J Am Stat Assoc* 96:1077–1087.
- Newman MEJ, Leicht EA (2007) Mixture models and exploratory analysis in networks. *Proc Natl Acad Sci USA* 104:9564–9569.
- Hoff PD, Raftery AE, Handcock MS (2002) Latent space approaches to social network analysis. *J Am Stat Assoc* 97:1090–1098.
- Handcock MD, Raftery AE, Tantrum JM (2007) Model-based clustering for social networks. *J R Stat Soc A-G* 170:301–354.
- Hoff PD (2007) Modeling homophily and stochastic equivalence in symmetric relational data. *Advances in Neural Information Processing Systems*, (MIT Press, Cambridge, MA), Vol 19.
- Karrer B, Newman MEJ (2011) Stochastic blockmodels and community structure in networks. *Phys Rev E* 83:016107.
- Borgatti SP, Everett MG (2000) Models of core/periphery structures. *Soc Networks* 21:375–395.
- Everett MG, Borgatti SP (2000) Peripheries of cohesive subsets. *Soc Networks* 21:397–407.
- Flake GM, Lawrence S, Giles CL, Coetzee FM (2002) Self-organization and identification of web communities. *IEEE Computer* 35:66–71.
- Clauset A (2005) Finding local community structure in networks. *Phys Rev E* 72:026132.
- Clauset A, Moore C, Newman MEJ (2008) Hierarchical structure and the prediction of missing links in networks. *Nature* 453:98–101.
- Frank O, Strauss D (1986) Markov graphs. *J Am Stat Assoc* 81:832–842.
- Bickel PJ, Chen A (2009) A nonparametric view of network models and Newman-Girvan and other modularities. *Proc Natl Acad Sci USA* 106:21068–21073.
- Beasley JE (1998) Heuristic algorithms for the unconstrained binary quadratic programming problem. Technical report, Management School, Imperial College, London, UK.
- Glover FW, Lagunas M (1997) *Tabu Search* (Kluwer Academic, Boston).
- Bickel PJ, Doksum KA (2007) *Mathematical Statistics: Basic Ideas and Selected Topics*, (Pearson Prentice Hall, Upper Saddle River, NJ), 2nd Ed, Vol 1.
- Tibshirani R, Walther G, Hastie T (2001) Estimating the number of clusters in a dataset via the gap statistic. *J R Stat Soc A* 63:411–423.
- Newman MEJ (2006) Finding community structure in networks using the eigenvectors of matrices. *Phys Rev E* 74:036104.
- Hubert L, Arabie P (1985) Comparing partitions. *J Classif* 2(1):193–218.
- Zachary WW (1977) An information flow model for conflict and fission in small groups. *J Anthropol Res* 33:452–473.
- Karrer B, Levina E, Newman MEJ (2008) Robustness of community structure in networks. *Phys Rev E* 77:046119.