

A nonparametric view of network models and Newman–Girvan and other modularities

Peter J. Bickel^{a,1} and Aiyou Chen^b

^aUniversity of California, Berkeley, CA 94720; and ^bAlcatel-Lucent Bell Labs, Murray Hill, NJ 07974

Edited by Stephen E. Fienberg, Carnegie Mellon University, Pittsburgh, PA, and approved October 13, 2009 (received for review July 2, 2009)

Prompted by the increasing interest in networks in many fields, we present an attempt at unifying points of view and analyses of these objects coming from the social sciences, statistics, probability and physics communities. We apply our approach to the Newman–Girvan modularity, widely used for “community” detection, among others. Our analysis is asymptotic but we show by simulation and application to real examples that the theory is a reasonable guide to practice.

modularity | profile likelihood | ergodic model | spectral clustering

The social sciences have investigated the structure of small networks since the 1970s, and have come up with elaborate modeling strategies, both deterministic, see Doreian et al. (1) for a view, and stochastic, see Airoldi et al. (2) for a view and recent work. During the same period, starting with the work of Erdős and Rényi (3), a rich literature has developed on the probabilistic properties of stochastic models for graphs. A major contribution to this work is Bollobás et al. (4). On the whole, the goals of the analyses of ref. 4, such as emergence of the giant component, are not aimed at the statistical goals of the social science literature we have cited.

Recently, there has been a surge of interest, particularly in the physics and computer science communities in the properties of networks of many kinds, including the Internet, mobile networks, the World Wide Web, citation networks, email networks, food webs, and social and biochemical networks. Identification of “community structure” has received particular attention: the vertices in networks are often found to cluster into small communities, where vertices within a community share the same densities of connecting with vertices in their own community as well as different ones with other communities. The ability to detect such groups can be of significant practical importance. For instance, groups within the worldwide Web may correspond to sets of web pages on related topics; groups within mobile networks may correspond to sets of friends or colleagues; groups in computer networks may correspond to users that are sharing files with peer-to-peer traffic, or collections of compromised computers controlled by remote hackers, e.g. botnets (5). A recent algorithm proposed by Newman and Girvan (6), that maximizes a so-called “Newman–Girvan” modularity function, has received particular attention because of its success in many applications in social and biological networks (7).

Our first goal is, by starting with a model somewhat less general than that of ref. 4, to construct a nonparametric statistical framework, which we will then use in the analysis, both of modularities and parametric statistical models. Our analysis is asymptotic, letting the number of vertices go to ∞ . We view, as usual, asymptotics as being appropriate insofar as they are a guide to what happens for finite n . Our models can, on the one hand, be viewed as special cases of those proposed by ref. 4, and on the other, as encompassing most of the parametric and semiparametric models discussed in Airoldi et al. (2) from a statistical point of view and in Chung and Lu (8) for a probabilistic one. An advantage of our framework is the possibility of analyzing the properties of the Newman–Girvan modularity, and the reasons for its success and occasional failures. Our approach suggests an alternative modularity which is, in

principle, “fail-safe” for rich enough models. Moreover, our point of view has the virtue of enabling us to think in terms of “strength of relations” between individuals not necessarily clustering them into communities beforehand.

We begin, using results of Aldous and Hoover (9), by introducing what we view as the analogues of arbitrary infinite population models on infinite unlabeled graphs which are “ergodic” and from which a subgraph with n vertices can be viewed as a piece. This development of Aldous and Hoover can be viewed as a generalization of deFinetti’s famous characterization of exchangeable sequences as mixtures of i.i.d. ones. Thus, our approach can also be viewed as a first step in the generalization of the classical construction of complex statistical models out of i.i.d. ones using covariates, information about labels and relationships.

It turns out that natural classes of parametric models which approximate the nonparametric models we introduce are the “blockmodels” introduced by Holland, Laskey and Leinhardt ref. 10; see also refs. 2 and 11, which are generalizations of the Erdős–Rényi model. These can be described as follows.

In a possibly (at least conceptually) infinite population (of vertices) there are K unknown subcommunities. Unlabeled individuals (vertices) relate to each other through edges which for this paper we assume are undirected. This situation leads to the following set of probability models for undirected graphs or equivalently the corresponding adjacency matrices $\{A_{ij} : i, j \geq 1\}$, where $A_{ij} = 1$ or 0 according as there is or is not an edge between i and j .

1. Individuals independently belong to community j with probability π_j , $1 \leq j \leq K$, $\sum_{j=1}^K \pi_j = 1$.
2. A symmetric $K \times K$ matrix $\{P_{kl} : 1 \leq k, l \leq K\}$ of probabilities is given such that P_{ab} is the probability that a specific individual i relates to individual j given that $i \in a, j \in b$. The membership relations between individuals are established independently. Thus $1 - \sum_{1 \leq a, b \leq K} \pi_a \pi_b P_{ab}$ is the probability that there is no edge between i and j .

The Erdős–Rényi model corresponds to $K = 1$.

We proceed to define Newman–Girvan modularity and an alternative statistically motivated modularity. We give necessary and sufficient conditions for consistency based on the parameters of the block model, properties of the modularities, and average degree of the graph. By consistency we mean that the modularities can identify the members of the block model communities perfectly. We also give examples of inconsistency when the conditions fail. We then study the validity of the asymptotics in a limited simulation and apply our approach to a classical small example, the Karate Club and a large set of Private Branch Exchange (PBX) data. We conclude with a discussion and some open problems.

Author contributions: P.J.B. and A.C. performed research and analyzed data.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

¹To whom correspondence should be addressed. E-mail: bickel@stat.berkeley.edu.

This article contains supporting information online at www.pnas.org/cgi/content/full/0907096106/DCSupplemental.

Random Graph Models

Consider any probability distribution \mathbb{P} on an infinite undirected graph, or equivalently a probability distribution on the set of all matrices $\|A_{ij} : i, j \geq 1\|$ where $A_{ij} = 1$ or 0 , $A_{ij} = A_{ji}$ for all i, j pairs, and $A_{ii} = 0$ for all i , thus excluding self relation. If the graph is unlabeled, it is natural to restrict attention to \mathbb{P} such that $\|A_{\sigma_i \sigma_j}\| \sim \mathbb{P}$ for any permutation σ of $\{1, 2, 3, \dots\}$. Hoover (see ref. 9) has shown that all such probability distributions can be represented as,

$$A_{ij} = g(\alpha, \xi_i, \xi_j, \lambda_{ij})$$

where $\alpha, \{\xi_i\}$ and $\{\lambda_{ij}\}$ are i.i.d. $U(0, 1)$ variables and

$$g(u, v, w, z) = g(u, w, v, z)$$

for all u, v, w, z . The variables ξ correspond to latent variables, λ being completely individual specific, ξ generating relations between individuals and α a mixture variable which is unidentifiable even for an infinite graph. Note that g is unidentifiable and the ξ and λ could be put on another scale, e.g. Gaussian. We note that, this point of departure was also recently proposed by Hoff (12) but was followed to a different end.

It is clear that the distributions representable as,

$$A_{ij} = g(\xi_i, \xi_j, \lambda_{ij}) \quad [1]$$

where $\lambda_{ij} = \lambda_{ji}$, are the extreme points of this set and play the same role as sequences of i.i.d. variables play in de Finetti's theorem. Since given ξ_i and ξ_j , the λ_{ij} are i.i.d., these distributions are naturally parametrized by the function

$$h(u, v) \equiv \mathbb{P}[A_{ij} = 1 | \xi_i = u, \xi_j = v].$$

As Diaconis and Janson (13) point out $h(\cdot, \cdot)$ does not uniquely determine \mathbb{P} but if h_1 and h_2 define the same \mathbb{P} , then there exists $\varphi : [0, 1] \rightarrow [0, 1]$ which is: measure preserving, i.e. such that $\varphi(\xi_1)$ has a $U(0, 1)$ distribution; and $h_1(u, v) = h_2(\varphi(u), \varphi(v))$.

Given any h corresponding to \mathbb{P} , let

$$\mathbb{P}[X_{ij} = 1 | \xi_i = u] = g(u) = \int_0^1 h(u, v) dv.$$

It is well known (see section 10 of ref. 14) that there exists a measure preserving φ_g such that, $g(\varphi_g(v))$ is monotone non decreasing.

Define

$$h_{CAN}(u, v) = h(\varphi_g(u), \varphi_g(v)).$$

We claim that

$$g_{CAN}(u) \equiv \int_0^1 h_{CAN}(u, v) dv = F^{-1}(u)$$

where F is the cdf of $g_{CAN}(\xi_i)$, and h_{CAN} is unique up to sets of measure 0. To see this note that if h corresponds to \mathbb{P} and $g(u) \equiv \int_0^1 h(u, v) dv$ is non decreasing, then since F is determined by \mathbb{P} only, $g(u) = F^{-1}(u)$. But $g(\varphi_g(u)) = g_{CAN}(u)$ and $\varphi_g(u) = g^{-1}g_{CAN}(u) = u$.

There is a reparametrization of h_{CAN} (we drop the CAN subscript in the future) which enables us to think of our model in terms more familiar to statisticians.

Let

$$\rho = \mathbb{P}(\text{Edge}) = \int_0^1 \int_0^1 h(u, v) du dv.$$

Then the conditional density of (ξ_i, ξ_j) given that there is an edge between i and j is $w(u, v) = \rho^{-1}h(u, v)$. This parametrization also permits us to decouple $\rho \propto \mathbb{E}(\text{Degree})$ of the graph from the inhomogeneity structure. It is natural finally to let ρ depend on n but

$w(\cdot, \cdot)$ to be fixed. If $\lambda_n \equiv \mathbb{E}(\text{Degree}) \rightarrow \infty$, we have what we may call the “dense graph” limit. If $\lambda_n = \Omega(1)$, we are in the case most studied in probability theory where, for instance, $\lambda_n = 1$ is the threshold at which the so called “giant component” appears. This is the situation Bollobas et al. focus on.

As we have noted, block models are of this type. Here we can think of the reparametrization as being $\rho, \pi, \|S_{ab}\| = \|\rho^{-1}P_{ab}\|$, or $\|W_{ab}\| \equiv \|\rho^{-1}P_{ab}\pi_a\pi_b\|$. The models studied by Chung and others (8) given by,

$$h(u, v) \propto a(u)a(v) \quad [2]$$

also fall under our description. The mixture model of Newman and Leicht (15) where, given communities $1, \dots, K$,

$$\mathbb{P}[X_{ij} = 1 | i \in s, j \in r] = \theta_{sr}\theta_{sj}$$

is not of this type, since it is not invariant under permutations. It can be made invariant by summing over all permutations of $\{1, \dots, n\}$, but is then generally not ergodic. Such models can be developed from our framework by permitting covariates Z_i depending on vertex identity or Z_{ij} depending on edge identity. Newman and Leicht's example where the communities are WEB pages falls under this observation. From a statistical point of view, these models bear the same relation to our models as regression models do to single population models.

Block models or models where

$$h(u, v, \theta) \propto \sum_{k=1}^{K_n} \theta_k a_k(u) a_k(v)$$

for known functions $\{a_k\}$ can be used to approximate general h . The latent eigenvalue model of Hoff (12) is of this type, but with a_k which are extremely rough and unidentifiable since the $a_j(\xi)$ are independent, and for which no unique choice exists. We can think of the canonical version of the block model as corresponding to a labeling $1, \dots, K$ of the communities in the order $W_1 \leq \dots \leq W_K$ where $W_j = \sum_k W_{jk}$, which is proportional to the expected degree of a member of community j . The function $h(\cdot, \cdot)$ then takes value P_{ab} on the (a, b) block of the product partition in which each axis is divided into consecutive intervals, of lengths π_1, \dots, π_K . Each corresponding vertical slice exhibits the relation pattern for that community with the diagonal block identifying the members of the community. The nonparametric $h(\cdot, \cdot)$ gives the same intuitive picture on an arbitrarily fine scale. We note that, as in nonparametric statistics, to estimate h or w , regularization is needed. That is, we need to consider $K_n \rightarrow \infty$ at rates which depend on n and λ_n to obtain good estimates of h or w by using estimates of θ above or of block model parameters. We will discuss this further later.

Newman–Girvan and Likelihood Modularities

The task of determining K communities corresponds to finding a good assignment for the vertices $\mathbf{e} \equiv \{e_1, \dots, e_n\}$ where $e_j \in \{1, \dots, K\}$. There are K^n such assignments. Suppose that the distribution of A follows a K block model with parameters $\pi = (\pi_1, \dots, \pi_K)$ and $P = \|P_{ab}\|_{K \times K}$. The observed A is a consequence of a realization $\mathbf{c} = (c_1, \dots, c_n)$ of n independent $\text{Multinomial}(1, \pi)$ variables. Evidently we can measure the adequacy of an assignment through the matrix

$$R(\mathbf{c}, \mathbf{e}) = \|R_{ab}\|_{K \times K},$$

where $R_{ab} = n^{-1} \sum_{i=1}^n I(c_i = b, e_i = a)$, the fraction of b members classified as a members if we use \mathbf{e} . It is natural to ask for consistency of an assignment \mathbf{e} , that is, $\mathbf{e} = \mathbf{c}$, i.e.

$$R(\mathbf{c}, \mathbf{e}) = \text{diag}(\mathbf{f}(\mathbf{c}))$$

where $\mathbf{f}_a(\mathbf{e}) = n^{-1}n_a(\mathbf{e})$ and $n_a(\mathbf{e}) = \sum_{i=1}^n I(e_i = a)$. The Newman–Girvan modularity $Q_{NG}(\mathbf{e}, A)$ is defined as follows. Let $\{i : e_i = k\}$ denote e-community k , i.e. as estimated by \mathbf{e} . Define

$$O_{kl}(\mathbf{e}, A) = \sum_{1 \leq i, j \leq n} A_{ij} I(e_i = k, e_j = l)$$

to be the block sum of A . Obviously, O_{kk} is twice the number of edges among nodes in the k -th e-community and for $k \neq l$, O_{kl} is the number of edges between nodes in the k -th e-community and nodes in the l -th e-community. Let $D_k(\mathbf{e}, A) = \sum_{l=1}^K O_{kl}(\mathbf{e}, A)$. It is easy to verify that D_k is the sum of degrees for nodes in the k -th e-community. Let $L = \sum_{k=1}^K D_k$ be twice the number of edges among all nodes. Then the Newman–Girvan modularity is defined by

$$Q_{NG}(\mathbf{e}, A) = \sum_{k=1}^K \frac{O_{kk}}{L} - \left(\frac{D_k}{L} \right)^2.$$

The Newman–Girvan algorithm then searches for the membership assignment vector \mathbf{e} that maximizes Q_{NG} . Notice that if edges are randomly generated uniformly among all pairs of nodes with given node degrees, then the number of edges between the k th e-community and l th e-community is expected to be $L^{-1}D_k D_l$. Therefore, the Newman–Girvan modularity measures the fraction of the edges on the graph that connect vertices of the same type (i.e. within-community edges) minus the expected value of the same quantity on a graph with the same community divisions but random connections between the vertices (6). Newman (16) contrasts and compares his modularity with spectral clustering, another common “community identification” method which we will also compare to the likelihood modularity below. We seek conditions under which the official N-G assignment

$$\hat{\mathbf{c}} = \arg \max Q_{NG}(\mathbf{e}, A)$$

is consistent with probability tending to 1. Before doing so we consider alternative modularities.

For fixed \mathbf{e} , the conditional, given $\{n_k(\mathbf{e})\}_1^K$, log-likelihood of A is $\frac{1}{2} \sum_{1 \leq a, b \leq K} (O_{ab} \log(P_{ab}) + (n_{ab} - O_{ab}) \log(1 - P_{ab}))$, where $n_{ab} = n_a n_b$ if $a \neq b$, $n_{aa} = n_a(n_a - 1)$. If we maximize over P , we obtain by letting $\tau(x) = x \log x + (1 - x) \log(1 - x)$,

$$Q_{LM}(\mathbf{e}, A) = \frac{1}{2} \sum_{a, b} n_{ab} \tau \left(\frac{O_{ab}}{n_{ab}} \right)$$

which we call the likelihood modularity. This is not a true likelihood but a profile likelihood where we treat \mathbf{e} as an unknown parameter. We will argue below that the profile likelihood is optimal in the usual parametric sense if $\frac{\lambda_n}{\log n} \rightarrow \infty$. But so are all other consistent modularities as defined below. However, we expect that if $\frac{\lambda_n}{\log n}$ is bounded and certainly if $\lambda_n = \Omega(1)$, the most important case, this is false. We are deriving optimal and computationally implementable procedures for this case.

We write general modularities in the form,

$$Q(\mathbf{e}, A) = F_n \left(\frac{O}{\mu_n}, \frac{L}{\mu_n}, \mathbf{f}(\mathbf{e}) \right)$$

where $\mu_n = E(L) = n(n - 1)\rho_n$ and the matrix $O(\mathbf{e}, A)$ is defined by its elements $O_{ab}(\mathbf{e}, A)$, $\mathbf{f}(\mathbf{e}) = (\frac{n_1(\mathbf{e})}{n}, \dots, \frac{n_K(\mathbf{e})}{n})^T$, and $F_n : \mathcal{M} \times \mathcal{R}^+ \times \mathcal{G} \rightarrow \mathcal{R}$ where \mathcal{M} is the set of all nonnegative $K \times K$ symmetric matrices and \mathcal{G} is the K simplex. Note that both Q_{NG}

and Q_{LM} can be written as such up to a proportionality constant. It is easy to see that if the K block model holds,

$$\frac{\mathbb{E}(O(\mathbf{e}, A) | \mathbf{c})}{E(L)} = R(\mathbf{c}, \mathbf{e}) S R^T(\mathbf{c}, \mathbf{e}) \quad [3]$$

where, by definition, $R^T \mathbf{1} = \mathbf{f}(\mathbf{c})$, $R \mathbf{1} = \mathbf{f}(\mathbf{e})$, $\mathbf{1} = (1, \dots, 1)^T$, and $S_{ab} = \rho_n^{-1} \mathbb{P}(A_{12} = 1 | c_1 = a, c_2 = b)$. Note that $W \equiv \mathcal{D}(\pi) S \mathcal{D}(\pi)$, where $\mathcal{D}(v) \equiv \text{diag}(v)$ for $v \in K \times 1$.

We define asymptotic consistency of a sequence of assignments $\hat{\mathbf{c}}$ by

$$\mathbb{P}[\hat{\mathbf{c}} = \mathbf{c}] \rightarrow 1 \quad [4]$$

as $n \rightarrow \infty$.

We will assume that there exists a function $F : \mathcal{M} \times \mathcal{R}^+ \times \mathcal{G} \rightarrow \mathcal{R}$ such that F_n is approximated by F evaluated at the conditional expectation given \mathbf{c} of the argument of F_n . Suppose first $F_n \equiv F$. It is intuitively clear from [3] that if $\hat{\mathbf{c}} \rightarrow \mathbf{c}$, then $R(\mathbf{c}, \hat{\mathbf{c}}) \rightarrow \mathcal{D}(\pi)$. Then, since $\mathbf{f}(\mathbf{c}) \rightarrow \pi$, the following condition is natural.

- I. $F(RS R^T, \mathbf{1}, R \mathbf{1})$ is uniquely maximized over $\mathcal{R} = \{R : R \geq 0, R^T \mathbf{1} = \pi\}$ by $R = \mathcal{D}(\pi)$, for all (π, S) in an open set Θ .

This means \mathbf{c} with $\mathbf{f}(\mathbf{c}) = \pi$ is the right assignment for the limiting problem. Note that since F is not concave in R , this is a strong condition.

For (π, S) to be identifiable uniquely, we clearly also need that:

- II. S does not have two identical columns (Two communities cannot have identical probabilities of being related to other communities and within themselves) and π has all entries positive (Each community has some members).

We also need a few more technical conditions of a standard type.

- III. a) F is Lipschitz in its arguments. b) The directional derivatives $\frac{\partial^2 F}{\partial \epsilon^2} (M_0 + \epsilon(M_1 - M_0), r_0 + \epsilon(r_1 - r_0), \mathbf{t}_0 + \epsilon(\mathbf{t} - \mathbf{t}_0))|_{\epsilon=0+}$ are continuous in (M_1, r_1, \mathbf{t}) for all (M_0, r_0, \mathbf{t}_0) in a neighborhood of $(W, \mathbf{1}, \pi)$. c) Let $G(R, S) = F(RS R^T, \mathbf{1}, R \mathbf{1})$. Assume that on \mathcal{R} , $\frac{\partial G((1-\epsilon)\mathcal{D}(\pi) + \epsilon R, S)}{\partial \epsilon}|_{\epsilon=0+} < -C < 0$ for all $(\pi, S) \in \Theta$.

Theorem 1. Suppose F, S and π satisfy I–III and $\hat{\mathbf{c}}$ is the maximizer of $Q(\mathbf{e}, A)$. Suppose $\frac{\lambda_n}{\log n} \rightarrow \infty$. Then, for all $(\pi, S) \in \Theta$,

$$\overline{\lim}_{n \rightarrow \infty} \frac{\log \mathbb{P}(\hat{\mathbf{c}} \neq \mathbf{c})}{\lambda_n} \leq -s_Q(\pi, S) < 0.$$

The proof is given in *SI Appendix*.

We note that Snijders and Nowicki (11) established a related result, exponential convergence to 0 of the mis-classification probability for $\lambda_n = \Omega(n)$ using node degree K-means clustering for $K = 2$.

Let $F(\mathbf{c}, \mathbf{e}) \equiv F(R(\mathbf{c}, \mathbf{e}) S R^T(\mathbf{c}, \mathbf{e}), \mathbf{f}^T(\mathbf{c}) S \mathbf{f}(\mathbf{e}), \mathbf{f}(\mathbf{e}))$. In the general case it suffices to show,

$$\overline{\lim}_{n \rightarrow \infty} \frac{\mathbb{P}(\sup\{|Q(\mathbf{e}, A) - F(\mathbf{c}, \mathbf{e})| : \mathbf{e}\} \leq \delta \Delta_n)}{\lambda_n} \leq -\gamma \quad [5]$$

for all $\delta > 0$, some $\gamma > 0$, where

$$\Delta_n = \inf\{|F(\mathbf{c}, \mathbf{e}) - F(\mathbf{c}, \mathbf{c})| : |\mathbf{e} - \mathbf{c}| \geq 1\}$$

and $|\mathbf{e} - \mathbf{c}| = \sum_{i=1}^n 1(e_i \neq c_i)$. We show in *SI Appendix* that Q_{NG} and Q_{LM} satisfy I and III and Eq. 5 for selected (π, S) for Q_{NG} and all (π, S) for Q_{LM} .

An immediate consequence of the theorem is:

Corollary 1: If the conditions of Theorem 1 hold and,

$$\hat{W} \equiv L^{-1}O(\hat{\mathbf{c}}, A),$$

$$\hat{\pi} = \left\{ \frac{1}{n} \sum_{i=1}^n I(\hat{\mathbf{c}}_i = a) : a = 1, \dots, K \right\},$$

then,

$$\sqrt{n}(\hat{\pi} - \pi) \Rightarrow \mathcal{N}(0, \mathcal{D}(\pi) - \pi\pi^T),$$

$$\sqrt{n}(\hat{W} - W) \Rightarrow S \cdot (\pi\eta^T + \eta\pi^T) - 2(\pi^T S \eta)W,$$

$$\eta = \mathcal{N}(0, \mathcal{D}(\pi) - \pi\pi^T),$$

with $A \cdot B$ denoting point-wise product. The limiting variances are what we would get for maximum likelihood estimates if $\hat{\mathbf{c}} = \mathbf{c}$, i.e. we knew the assignment to begin with. So consistent modularities lead to efficient estimates of the parameters.

This follows since with probability tending to 1, $\hat{\mathbf{c}} = \mathbf{c}$.

To estimate $w(\cdot, \cdot)$ in the nonparametric case we need $K \rightarrow \infty$, and $w(\cdot, \cdot)$ and $\pi(\cdot)$ smooth. We approximate by $W_K \sim K^{-2} \|w(aK^{-1}, bK^{-1})\|$, $\pi_K(a) \sim K^{-1} \pi(aK^{-1})$, where $w(\cdot, \cdot)$, W_K are canonical and the modularity defining F_K , $F_K(\cdot, \cdot, \cdot)$ is of order K^{-2} . We have preliminary results in that direction but their formulation is complicated and we do not treat them further here.

Consistency of N-G, L-M

We show in *SI Appendix* using the appropriate F_{NG} , F_{LM} that the likelihood modularity is always consistent while the Newman–Girvan is not. This is perhaps not surprising since N-G focuses on the diagonal of O . In fact, we would hope that N-G is consistent under the submodel $\{(\rho, \pi, W) : W_{aa} > \sum_{b \neq a} W_{ab} \text{ for all } a\}$, which corresponds to Newman and Girvan’s motivation. We have shown this for $K = 2$ but it surprisingly fails for $K > 2$. Here is a counterexample. Let $K = 3$, $\pi = (1/3, 1/3, 1/3)^T$ and

$$P = \begin{bmatrix} .06 & .04 & 0 \\ .04 & .12 & .04 \\ 0 & .04 & .66 \end{bmatrix}.$$

As $n \rightarrow \infty$, with true labeling, Q_{NG} approaches 0.033. However, the maximum Q_{NG} , about 0.038, is achieved by merging the first two communities. That is, two sparser communities are merged. This is consistent with an observation of Fortunato and Barthelemy (17).

If for the profile likelihood we maximize only over \mathbf{e} such that $\hat{W}_{aa}(\mathbf{e}) > \sum_{b \neq a} \hat{W}_{ab}(\mathbf{e})$ for all a , we obtain $\hat{\mathbf{c}}$ which is consistent under the submodel above, and in the Karate Club example performs like N-G.

Computational Issues

Computation of optimal assignments using modularities is, in principle, NP hard. However, although the modularity is multimodal, in the examples we have considered and generally when the signal is strong, optimization from several starting points using a label switching algorithm (19) works well.

Simulation

We generate random matrices A and maximize Q_{NG} , Q_{LM} to obtain node labels respectively, where Q_{LM} is maximized using a label switching algorithm. To make a fair comparison, the initial labeling for Q_{NG} and Q_{LM} is to randomly choose 50% of the nodes with correct labels and the other 50% with random labels. For spectral clustering, we adopt the algorithm of (18) by using the first K eigenvectors of $\mathcal{D}(\mathbf{d})^{-1/2} A \mathcal{D}(\mathbf{d})^{-1/2}$, where $\mathbf{d} = (d_1, \dots, d_n)^T$ and d_i is the degree of the i -th node. We generate the P matrix randomly by forcing symmetry and then add a constant to diagonal entries

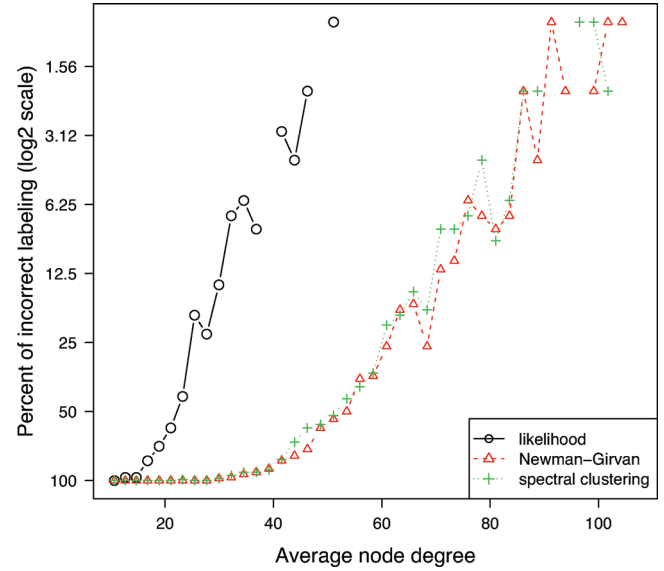


Fig. 1. Empirical comparison of Newman–Girvan, likelihood modularities and spectral clustering (18), where $K = 3$, the number of nodes n varies from 200 to 1500, and the percent of correct labeling is computed from 100 replicates of each simulation case. Here π, P are given in the text.

such that I holds. The π is generated randomly from the simplex. To be precise, the values for Fig. 1 are $\pi = (.203, .286, .511)^T$ and

$$P = bn^{-1} \log n \cdot \begin{bmatrix} .43 & .06 & .13 \\ .06 & .34 & .17 \\ .13 & .17 & .40 \end{bmatrix},$$

where n varies from 200 to 1,500 and b varies from 10 to 100. Obviously, Fig. 1 says that the likelihood method exhibits much less incorrect labeling than Newman–Girvan and spectral clustering. This is consistent with theoretical comparison.

Data Examples

We compare the L-M and N-G modularity algorithms below with applications to two real data sets. To deal with the issue of non-convex optimization, we simply use many restarting points.

Zachary’s “Karate Club” Network. We first compare L-M and N-G with the famous “Karate Club” network of ref. 20, from the social science literature, which has become something of a standard test for community detection algorithms. The network shows the patterns of friendship between the members of a karate club at a US university in the 1970s. The example is of particular interest because shortly after the observation and construction of the network, the club in question split into two components separated by the dashed line as shown in Figs. 2 and 3 as a result of an internal dispute. Fig. 2 *Left* shows two communities identified by maximizing the likelihood modularity where the shapes of the vertices denote the membership of the corresponding individuals, and similarly the right panel shows communities identified by N-G. Obviously, the N-G communities match the two sub-divisions identified by the split save for one mis-classified individual. The L-M communities are quite different, and obviously one community consists of five individuals with central importance that connect with many other nodes while the other community consists of the remaining individuals. Although not reflecting the split this corresponds to other plausible distinguishing characteristics of the individuals. However, if we force the constraint that within-community density is no less than the density of relationship to all other communities, the submodel we discussed, then we obtain two L-M communities that match the split perfectly. The same partitions as ours with and without constraint have also been reported

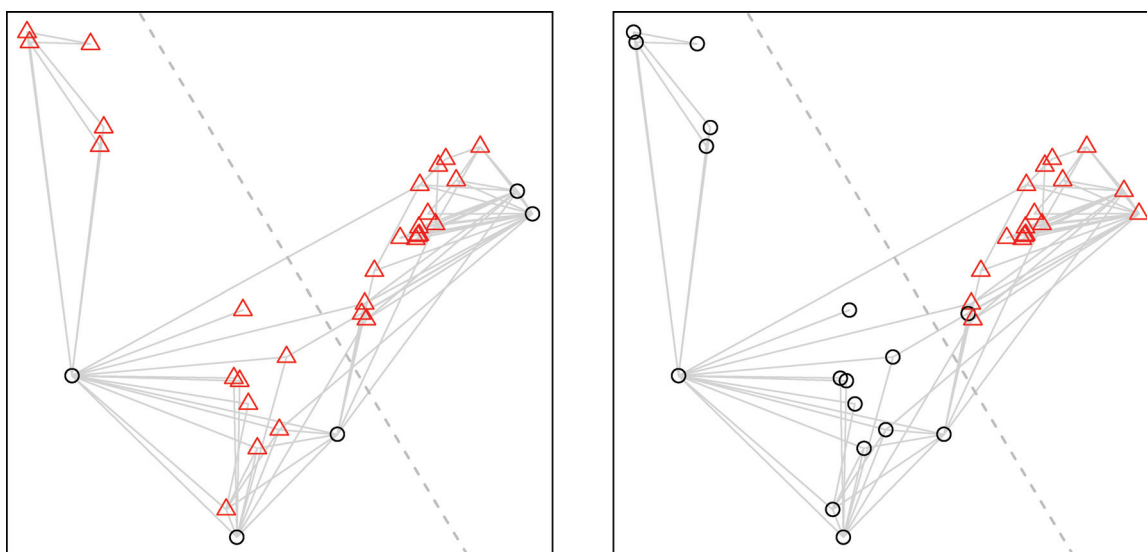


Fig. 2. Zachary's karate club network. Communities were identified by maximizing the likelihood modularity (*Left*) and by maximizing the Newman–Girvan modularity with $K = 2$ (*Right*), where the shapes of vertices indicate the membership of the corresponding individuals. The dashed line cuts the nodes into two groups which are the “known” communities that the club was split into.

by Rosvall and Bergstrom using a data compression criterion (21), which is closely related to L-M. We note that, as is usual in clustering, there is no ground truth, only features which can be validated expost fact. It is interesting to note that, if instead of $K = 2$, we put $K = 4$, as in Fig. 3, it is evident for both modularities that merging the communities on either side of the eigenvector split, gives the “correct” Karate Club split. This suggests that the standard policy mentioned by Newman (16) of increasing the number of communities by splitting is not necessarily ideal since in this case the “misclassified” individual of Fig. 2 would never be “correctly” classified.

Private Branch Exchange. Our second example is of a telephone communication network where connections are made among the internal telephones of a private business organization, a so called PBX. PBXs are differentiated from “key systems” in that users of key systems usually select their own outgoing lines, while PBXs select the outgoing line automatically. Our data contains 621

individuals. Fig. 4 *Left* shows the results of community detection by L-M, where the adjacency matrix is plotted but the nodes are sorted according to the membership of the corresponding individuals identified by maximizing the likelihood modularity. Similarly, the right panel of Fig. 4 shows the communities identified by N-G, where the maximum Newman–Girvan modularity is 0.4217. Note that the identified communities by L-M have sizes 323, 81, 78, 97, 41, and 1, respectively. The communities are ordered simply by their average node degrees, essentially the order for h_{CAN} . Interestingly, the last L-M community has only one node that communicates with almost everyone else, nodes in the second community only communicate with internal nodes, nodes in the fourth community and the sixth community, but not with others; Similarly, the third community only communicates with the fifth and sixth communities, and so on. In other words, communication between communities is sparse. However, the communities identified by N-G are quite different with only the fifth community heavily overlapping with a community identified by L-M. This

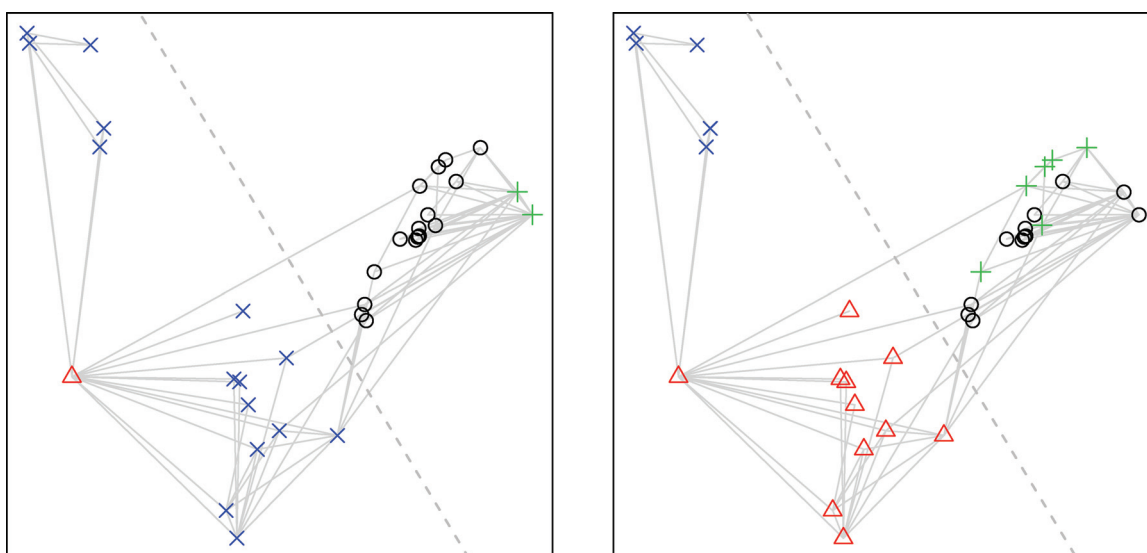


Fig. 3. Zachary's karate club network. Communities were identified by maximizing the likelihood modularity (*Left*) and by maximizing the Newman–Girvan modularity with $K = 4$ (*Right*), where the shapes of vertices indicate the membership of the corresponding individuals. The dashed line cuts the nodes into two groups which are the known communities that the club was split into.

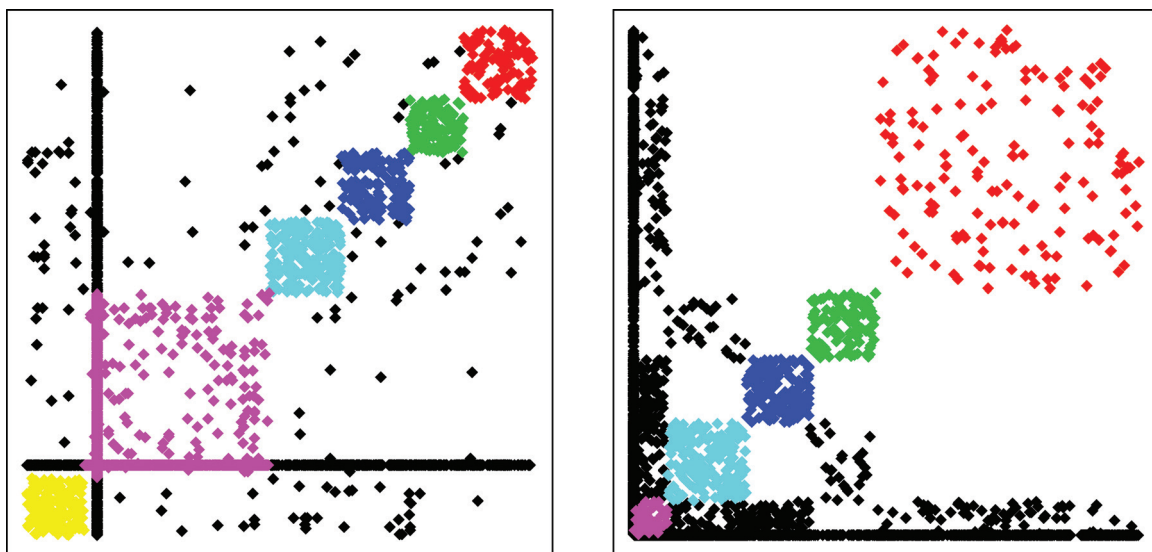


Fig. 4. Private branch exchange data. (Left) The adjacency matrix where the nodes are sorted according to the membership of the corresponding individuals identified by maximizing L-M. (Right) Same as Left, but with individuals identified by maximizing N-G with $K = 6$. The colors on the within-community edges are used to differentiate the communities for both L-M and N-G.

difference appears to be caused by the nodes in the 5th and 6th L-M communities which have many more between-community connections than within-community connections while N-G more or less maximizes within-community connections. We have verified that the group communicating with all others is a service group.

Discussion

1. As we noted, under our conditions the usual statistical goal of estimating the parameters π and P is trivial, since, once we have assigned individuals to the K communities consistently, the natural estimates, \hat{W} and $\hat{\pi}$, are not just consistent but efficient. However, in the more realistic case where $\lambda_n = \Omega(1)$, or even just $\lambda_n = \Omega(\log n)$, this is no longer true. Elsewhere, we shall show that, indeed, estimation of parameters by maximum likelihood and Bayes classification of individuals (no longer perfect) is optimal.
2. A difficulty faced by all these methods, modularities or likelihoods, is that if K is large, searching over the space of classifications becomes prohibitively expensive. In subsequent work we intend to show that this difficulty may

be partly overcome by using the method of moments to first estimate π and P , and then study the likelihood in a neighborhood of the estimated values.

Open Problems

1. A fundamental difficulty not considered in the literature is the choice of K . From our nonparametric point of view, this can equally well be seen as, how to balance bias and variance in the estimation of $w(\cdot, \cdot)$. We would like to argue that, as in nonparametric statistics, estimating $w(\cdot, \cdot)$ without prior prejudices on its structure is as important an exploratory step in this context as, using histograms in ordinary statistics.
2. The linking of this framework to covariates depending on vertex or edge identity is crucial, permitting relationship strength to be assessed as a function of vector variables.
3. The links of our approach to spectral graph clustering and more generally clustering on the basis of similarities seem intriguing.

ACKNOWLEDGMENTS. We thank Tin K Ho for help in obtaining the PBX data and for helpful discussions. We also thank the referees, whose references and comments improved this article immeasurably.

1. Doreian P, Batagelj V, Ferligoj A (2005) *Generalized Blockmodeling* (Cambridge Univ Press, Cambridge, UK).
2. Airoldi EM, Blei DM, Fienberg SE, Xing XP (2008) Mixed-membership stochastic blockmodels. *J Machine Learning Res* 9:1981–2014.
3. Erdős P, Rényi A (1960) On the evolution of random graphs. *Publ Math Inst Hungar Acad Sci* 5:17–61.
4. Bollobas B, Janson S, Riordan O (2007) The phase transition in inhomogeneous random graphs. *Random Struct Algorithms* 31:3–122.
5. Zhao Y, et al. (2009) Botgraph: Large scale spamming botnet detection. *Proceedings of the 6th USENIX Symposium on Networked Systems Design and Implementation* (USENIX, Berkeley, CA), pp 321–334.
6. Newman MEJ, Girvan M (2004) Finding and evaluating community structure in networks. *Phys Rev E* 69:026113.
7. Guimera R, Amaral LAN (2005) Functional cartography of complex metabolic networks. *Nature* 433:895–900.
8. Chung FRK, Lu L (2006) *Complex Graphs and Networks*. CBMS Regional Conference Series in Mathematics (Am Math Soc, Providence, RI).
9. Kallenberg O (2005) Probabilistic symmetries and invariance principles. *Probability and Its Application* (Springer, New York).
10. Holland PW, Laskey KB, Leinhardt S (1983) Stochastic blockmodels: First steps. *Soc Networks* 5:109–137.
11. Snijders T, Nowicki K. (1997) Estimation and prediction for stochastic block-structures for graphs with latent block structure. *J Classification* 14:75–100.
12. Hoff PD (2008) Modeling homophily and stochastic equivalence in symmetric relational data. *Advances in Neural Information Processing Systems*, eds Platt J, Koller D, Roweis S (MIT Press, Cambridge, MA) Vol 20, pp 657–664.
13. Diaconis P, Janson S (2008) Graph limits and exchangeable random graphs. *Rendiconti di Matematica* 28:33–61.
14. Hardy G, Littlewood J, Polya G (1988) *Inequalities*. (Cambridge Univ Press, Cambridge, UK), 2nd Ed.
15. Newman MEJ, Leicht EA (2007) Mixture models and exploratory analysis in networks. *Proc Natl Acad Sci USA* 104:9564–9569.
16. Newman MEJ (2006) Finding community structure in networks using the eigenvectors of matrices. *Phys Rev E* 74:036104.
17. Fortunato S, Barthélemy M (2007) Resolution limit in community detection. *Proc Natl Acad Sci USA* 104:36–41.
18. Ng AY, Jordan MI, Weiss Y (2001) On spectral clustering: Analysis and an algorithm. *Advances in Neural Information Processing Systems* (MIT Press, Cambridge, MA) Vol 14, pp 849–856.
19. Stephens M (2000) Dealing with label-switching in mixture models. *J R Stat Soc B* 62:795–809.
20. Zachary W (1977) An information flow model for conflict and fission in small groups. *J Anthropol Res* 33:452–473.
21. Rosvall M, Bergstrom C (2007) An information-theoretic framework for resolving community structure in complex network. *Proc Natl Acad Sci USA* 104:7327–7331.

SUPPLEMENTARY MATERIAL FOR ‘A NONPARAMETRIC VIEW OF NETWORK MODELS AND NEWMAN-GIRVAN AND OTHER MODULARITIES’

PETER BICKEL AND AIYOU CHEN

1. PROOF OF THEOREM 1

We establish Theorem 1. The generalization is immediate.

Define $D(\mathbf{c}) = \text{diag}(\mathbf{f}(\mathbf{c}))$, $W(\mathbf{c}) = D(\mathbf{c})SD(\mathbf{c})$ and $V(\mathbf{c}, \mathbf{e}) = R(\mathbf{c}, \mathbf{e})(D(\mathbf{c}))^{-1}$. Then

$$\begin{aligned} RPR^T &= \rho_n RSR^T \\ &= \rho_n VW(\mathbf{c})V^T \end{aligned}$$

and $V^T \mathbf{1} = \mathbf{1}$, $V\mathbf{f}(\mathbf{c}) = \mathbf{f}(\mathbf{e})$, where $V \equiv V(\mathbf{c}, \mathbf{e})$ and $O \equiv O(\mathbf{e}, A)$.

Let $\mu_n = n(n-1)\rho_n = (n-1)\lambda_n$, and define the matrix valued process,

$$X_n(\mathbf{e}) = (\mu_n^{-1}O - VW(\mathbf{c})V^T).$$

Let $|\mathbf{e} - \mathbf{c}| = \sum_{i=1}^n 1(e_i \neq c_i)$.

Lemma 1.1. *The following inequalities hold.*

$$(1) \quad \mathbb{P}(\sup\{||X_n(\mathbf{e})|| : |\mathbf{e}| \geq \epsilon\}) \leq 2K^{n+2} \exp\left(-\frac{1}{2}\epsilon\mu_n\right)$$

for $\epsilon \geq \mu_n^{-1/2}$,

$$(2) \quad \mathbb{P}(\sup\{||X_n(\mathbf{e}) - X_n(\mathbf{c})|| : |\mathbf{e} - \mathbf{c}| \leq m\} \geq \epsilon) \leq 2\binom{n}{m} K^{m+2} \exp\left(-\frac{1}{2}\epsilon\mu_n\right)$$

for $\epsilon \geq \frac{2m}{n}$.

Proof. We prove (2). Then

$$\mu_n^{-1}O_{ab}(\mathbf{e}, A) = \mu_n^{-1} \sum_{i,j} A_{ij} I(e_i = a, e_j = b)$$

where the A_{ij} are conditionally independent Bernoulli, $P(A_{ij} = 1|\mathbf{c}) = P_{c_i c_j}$.

$$\begin{aligned} \mu_n^{-1}\mathbb{E}(O_{ab}|\mathbf{c}) &= \mu_n^{-1} \sum_{i,j} P_{c_i c_j} I(e_i = a, e_j = b) \\ &= \rho_n^{-1}(RPR^T)_{ab} \\ &= VW(\mathbf{c})V^T \end{aligned}$$

where R is as defined in the text

$$R_{ab} = n^{-1} \sum_{i=1}^n I(e_i = a, c_i = b).$$

Thus, if $e_1 \neq c_1, \dots, e_m \neq c_m, e_{m+1} = c_{m+1}, \dots, e_n = c_n$,

$$\begin{aligned} [X_n(\mathbf{e}) - X_n(\mathbf{c})]_{ab} &= \mu_n^{-1} \sum_{i=1}^m \sum_{j=1}^n \sum_{u \neq a} (A_{ij} - P_{ub}) I(c_i = u, e_i = a) I(c_j = e_j = b) \\ &\quad + \mu_n^{-1} \sum_{i=1}^m \sum_{j=1}^n \sum_{u \neq a} (A_{ij} - P_{bu}) I(c_i = e_i = b) I(c_j = u, e_j = a). \end{aligned}$$

Bernstein's inequality (c.f. [1]) states: If Y_i are independent, $|Y_i| \leq 1$, $EY_i = 0$, $1 \leq i \leq I$, $S_I = \sum_{i=1}^I Y_i$, then

$$\mathbb{P}(|S_I| \geq w) \leq \exp\left(-\min\left(\frac{w^2}{4\text{var}(S_I)}, \frac{w}{2}\right)\right).$$

We identify Y_i with $(A_{ij} - P_{ub})I(c_i = u, c_j = b, e_i = a, e_j = b)$, $I = 2mn$. Then for $|\mathbf{e} - \mathbf{c}| \leq m$,

$$\begin{aligned} \text{var}(\mu_n(X_n(\mathbf{e}) - X_n(\mathbf{c}))|\mathbf{c}) &\leq 2mn\rho_n \\ &= 2m\lambda_n \end{aligned}$$

since $P_{uv} \leq \rho_n$. Put $w = \epsilon\mu_n$ and obtain the bound (2). \square

For future reference, we note the inequality,

$$(3) \quad \mathbb{P}(\|\mathbf{f}(\mathbf{c}) - \pi\| \geq \epsilon) \leq 2Ke^{-2n\epsilon^2}.$$

Lemma 1.2. *The following inequalities hold.*

$$(4) \quad \mathbb{P}(|\mu_n^{-1}L - \mathbf{1}^T W(\mathbf{c})\mathbf{1}| \geq \epsilon) \leq \exp\left(-\frac{1}{2}\mu_n\epsilon\right)$$

and for $\epsilon \geq 4$,

$$(5) \quad \mathbb{P}(\|W(\mathbf{c}) - W\|_o \geq \epsilon) \leq \exp\left(-\frac{n\lambda\delta(\epsilon)}{4\|W\|_o}\right),$$

where $\lambda = \min_j \pi_j$, $\|W\|_o$ is the operator norm of W , and $\delta(\epsilon) = \sqrt{1 + \epsilon} - 1$.

Proof. Since

$$\mu_n^{-1}E(L|\mathbf{c}) = \mathbf{1}^T W(\mathbf{c})\mathbf{1},$$

and L is conditionally twice the sum of $\binom{n}{2}$ independent $A_{ij} \sim \text{Bernoulli}(P_{c_i c_j})$, (4) is immediate from Bernstein. For (5) write $W(\mathbf{c}) = D(\mathbf{c})\mathcal{D}^{-1}(\pi)W\mathcal{D}^{-1}(\pi)D(\mathbf{c})$ and apply the bound

$$\|AJA - BJB\|_o \leq \|A - B\|_o(\|A\|_o + \|B\|_o)\|J\|_o$$

with $A = D(\mathbf{c})\mathcal{D}^{-1}(\pi)$, $B = I$, $J = W$, to get,

$$\|W(\mathbf{c}) - W\|_o \leq 2 \max_j \frac{|\mathbf{f}_j(\mathbf{c}) - \pi_j|}{\pi_j} + \max_j \frac{|\mathbf{f}_j(\mathbf{c}) - \pi_j|^2}{\pi_j^2}$$

and

$$\mathbb{P}(\|W(\mathbf{c}) - W\|_o \geq \epsilon) \leq \mathbb{P}\left(\max_j |\mathbf{f}_j(\mathbf{c}) - \pi_j| \geq \frac{\lambda\delta(\epsilon)}{2\|W\|_o}\right)$$

since $2x + x^2 \geq \epsilon$ is equivalent to $x \geq \delta(\epsilon)$. Apply (3). \square

Now we proceed to prove Theorem 1 for $F_n = F$.

Proof. By III a),

$$(6) \quad \left| F\left(\frac{O(\mathbf{c}, \mathbf{e})}{\mu_n}, \frac{L}{\mu_n}, \mathbf{f}(\mathbf{e})\right) - F(VW(\mathbf{c})V^T, \mathbf{1}^T W(\mathbf{c})\mathbf{1}, V\pi) \right| \leq M(|X_n(\mathbf{e})| + |\mu_n^{-1}L - \mathbf{1}^T W(\mathbf{c})\mathbf{1}|)$$

since $V(\mathbf{c}, \mathbf{e})\pi = \mathbf{f}(\mathbf{e})$. Similarly,

$$(7) \quad \begin{aligned} |F(VW(\mathbf{c})V^T, \mathbf{1}^T W(\mathbf{c})\mathbf{1}, V\pi) - F(VWV^T, \mathbf{1}, V\pi)| &\leq M(\|V\|_o^2 \|W(\mathbf{c}) - W\|_o + K\|W(\mathbf{c}) - W\|_o) \\ &\leq 2MK\|W(\mathbf{c}) - W\|_o \end{aligned}$$

since $\|V\|_o^2 \leq \text{Trace}(V^2) \leq K$, $\|\mathbf{1}\|^2 = K$.

Since $(V, \mathbf{f}(\mathbf{c})) \rightarrow F(VW(\mathbf{c})V^T, \mathbf{1}, V\pi)$ is continuous, we conclude from (6) and (7) that for some $\delta_n \downarrow 0$,

$$(8) \quad \mathbb{P}\left(\max\left\{F\left(\frac{O(\mathbf{c}, \mathbf{e})}{\mu_n}, \frac{L}{\mu_n}, \mathbf{f}(\mathbf{e})\right) : |V(\mathbf{c}, \mathbf{e}) - I| \geq \delta_n\right\} < F\left(\frac{O(\mathbf{c}, \mathbf{c})}{\mu_n}, \frac{L}{\mu_n}, \mathbf{f}(\mathbf{c})\right)\right) \geq 1 - e^{-\epsilon\mu_n}$$

for some $\epsilon > 0$.

By Lemma 1.1 and III a), we can write

$$(9) \quad F\left(\frac{O(\mathbf{c}, \hat{\mathbf{c}})}{\mu_n}, \frac{L}{\mu_n}, \mathbf{f}(\hat{\mathbf{c}})\right) - F\left(\frac{O(\mathbf{c}, \mathbf{c})}{\mu_n}, \frac{L}{\mu_n}, \mathbf{f}(\mathbf{c})\right) = F\left(\hat{V}W(\mathbf{c})\hat{V}^T, \mathbf{1}^T W(\mathbf{c})\mathbf{1}, \hat{V}\pi\right) \\ - F\left(W(\mathbf{c}), \mathbf{1}^T W(\mathbf{c})\mathbf{1}, \mathbf{f}(\mathbf{c})\right) + \Delta_n(\hat{\mathbf{c}}, \mathbf{c})$$

where $\hat{V} \equiv V(\mathbf{c}, \hat{\mathbf{c}})$ and $|\Delta_n(\mathbf{c}, \hat{\mathbf{c}})| \leq M(\|X_n(\hat{\mathbf{c}}) - X_n(\mathbf{c})\| + \|\mu_n^{-1}L - \mathbf{1}^T W(\mathbf{c})\mathbf{1}\|)$.

We now use III b), c) since $\mathbf{f}(\mathbf{c})$ is with probability greater or equal to $1 - e^{-n\epsilon}$ within ϵ of π . And by (7), $\hat{V} - I = o_P(1)$, we write,

$$F\left(\hat{V}W(\mathbf{c})\hat{V}^T, \mathbf{1}^T W(\mathbf{c})\mathbf{1}, \mathbf{f}(\hat{\mathbf{c}})\right) - F\left(W(\mathbf{c}), \mathbf{1}^T W(\mathbf{c})\mathbf{1}, \mathbf{f}(\mathbf{c})\right) = \sum_{a \neq b} c_{ab}(W(\mathbf{c}))\hat{V}_{ab} \\ + O_P(\|\hat{V} - I\|_1^2)$$

where $\|M\|_1 = \sum_{a,b} |M_{ab}|$ if $M = \|M_{ab}\|$. Next, note that

$$\frac{|\mathbf{c} - \mathbf{e}|}{n} = \frac{1}{n} \sum_{i=1}^n 1(c_i \neq e_i) \\ = 1 - \sum_{a=1}^K R_{aa}(\mathbf{c}, \mathbf{e}) \\ = 1 - \sum_{a=1}^K V_{aa}(\mathbf{c}, \mathbf{e})\mathbf{f}_a(\mathbf{c}) \\ = \sum_{a=1}^K \mathbf{f}_a(\mathbf{c})(1 - V_{aa}(\mathbf{c}, \mathbf{e})).$$

Since $\sum_{a=1}^K (1 - V_{aa}(\mathbf{c}, \mathbf{e})) = \frac{1}{2}\|V(\mathbf{c}, \mathbf{e}) - I\|_1$, we conclude that for some $\lambda > 0$,

$$(10) \quad \frac{\lambda}{2}\|V(\mathbf{c}, \mathbf{e}) - I\|_1 \leq \frac{1}{n}|\mathbf{c} - \mathbf{e}| \leq \frac{1}{2}\|V(\mathbf{c}, \mathbf{e}) - I\|_1.$$

Next note that by III c), $\sum_{a \neq b} c_{ab}(W)\hat{V}_{ab} \leq -\frac{C}{2}\|\hat{V} - I\|_1$ since V is a linear transform of R . Thus, the left hand side of (9) is less than 0 for $0 < \|\hat{V} - I\|_1 \leq \delta(F, W, \pi)$. Combining (8), (10), we conclude that to prove the theorem for $F_n = F$, we need to show,

$$(11) \quad \mathbb{P}\left(\frac{2}{n\lambda} \leq \|\hat{V} - I\|_1 \leq \delta_n\right) \geq 1 - Ce^{-\lambda_n \alpha}$$

for some $C, \alpha > 0$.

In view of (9), this will follow if

$$\mathbb{P}\left(\sup\left\{\frac{\|X_n(\mathbf{e}) - X_n(\mathbf{c})\|}{\|V(\mathbf{c}, \mathbf{e}) - I\|_1} : \delta_n \geq \|V(\mathbf{c}, \mathbf{e}) - I\|_1 \geq \frac{2}{n\lambda}\right\} \geq \frac{\nu}{4}\right) \leq Ce^{-\lambda_n \alpha}.$$

Substitute $\epsilon = \frac{2m}{n}$ into (2) and use $\binom{n}{m} \leq \left(\frac{n}{m}\right)^m (1 - \frac{m}{n})^{n-m}$. Put $m = \alpha n$, $\alpha \leq \frac{1}{2}$ and find the bound for the right hand side of (2), $e^{n\alpha \log \frac{1}{\alpha}} e^{n\alpha \log K} e^{-\alpha(n-1)\lambda_n} = e^{-\lambda_n(1+o(1))}$ if $\log \frac{1}{\alpha} \leq \log K = o(\lambda_n)$ which holds if $\alpha \geq \frac{1}{n}$. The theorem follows. \square

2. THE NEWMAN-GIRVAN MODULARITY

Q_{NG} evidently corresponds to

$$F_{NG}(M, r, \mathbf{t}) = F_{NG}(M, r) = \sum_a \left(\frac{M_{aa}}{r} - \left(\frac{\sum_b M_{ab}}{r} \right)^2 \right).$$

It is easy to check that I, III are satisfied if for all $a \neq b$, $\mathcal{E}_{aa} > 0$, $\mathcal{E}_{ab} < 0$ where $\mathcal{E} = W - (W\mathbf{1})(W\mathbf{1})^T$.

Our counter example shows inconsistency for other plausible (π, W) .

3. THE LM MODULARITY

If $\rho_n \equiv \rho > 0$, we can write

$$Q_{LM}(\mathbf{e}, A) = \frac{n^2}{2} F_{LM}^{(1)}\left(\frac{O}{\mu_n}, \frac{L}{\mu_n}, \mathbf{f}(\mathbf{e})\right)$$

where

$$F_{LM}^{(1)}(M, r, \mathbf{t}) = \sum_{a,b} t_a t_b \tau\left(\frac{\rho M_{ab}}{t_a t_b}\right)$$

with $\tau(x) = x \log x + (1-x) \log(1-x)$.

If $\rho_n \rightarrow 0$,

$$\tau(\rho_n x) = \rho_n(x \log x - x) + x \rho_n \log \rho_n + O(\rho_n^2 x^2).$$

Therefore

$$\begin{aligned} \frac{2}{n^2} Q_{LM}(\mathbf{e}, A) &= F_n\left(\frac{O}{\mu_n}, \frac{L}{\mu_n}, \mathbf{f}(\mathbf{e})\right) \\ (12) \quad &= \rho_n F_{LM}^{(2)}\left(\frac{O}{\mu_n}, \mathbf{f}(\mathbf{e})\right) + (\rho_n \log \rho_n) \frac{L}{\mu_n} + O(\rho_n^2) \end{aligned}$$

where

$$F_{LM}^{(2)}(M, \mathbf{t}) = \sum_{a,b} t_a t_b \sigma\left(\frac{M_{ab}}{t_a t_b}\right)$$

and $\sigma(x) = x \log x - x$.

In both cases, it is easy to check condition III. We now check consistency. Without loss of generality, we take $\mathbf{f}(\mathbf{c}) = \pi$ and $W(\mathbf{c}) = W$.

We need the following conjugate function lemma.

Lemma 3.1. (*Conjugate functions [2]*) *Let $\tau(s) = s \log s + (1-s) \log(1-s)$ and $\sigma(s) = s \log s - s$. Then*

$$\tau(s) = \inf_x (x^T s - \log(1 + e^x))$$

and

$$\sigma(s) = \inf_x (x^T s - e^x).$$

The following lemma is used for proving uniqueness of the maximizer.

Lemma 3.2. *Let S and R be two $K \times K$ matrices with nonnegative entries. Assume that*

- a) S is symmetric;
- b) S does not have two identical columns;
- c) there exists at least one nonzero entry in each column of R ;
- d) for $1 \leq a, b, k, l \leq K$, $S_{kl} = S_{ab}$ whenever $R_{ka} R_{lb} > 0$.

Then R is a diagonal matrix or a permutation of a diagonal matrix.

Proof. There are two possibilities for R and we prove the result in two steps as below.

1). If there exists a permutation of the rows and columns of R such that its diagonals are all positive after permutation, i.e. $R_{bb} > 0$ for $b = 1, \dots, K$. Now if R is not diagonal, then there exists $k \neq a$ such that $R_{ka} > 0$, which gives $R_{ka} R_{bb} > 0$ and thus by d) $S_{kb} = S_{ab}$ for $b = 1, \dots, K$. But this contradicts with b). Thus R must be a diagonal matrix or a permutation of it.

2). The remaining possibility is that there does not exist a permutation such that diagonal entries of R are all positive. Then we can always permute rows and columns of R , such that for some $m \geq 1$, $R_{ij} = 0$ for $1 \leq i, j \leq m$, and $R_{bb} > 0$ for $b = m+1, \dots, K$. Note that there is at least one nonzero entry in each of the first m columns, say $R_{k_i i} > 0$ for $i = 1, \dots, m$ and some $k_i \in \{m+1, \dots, K\}$. So $R_{k_i i} R_{bb} > 0$ and $R_{k_i i} R_{k_j j} > 0$, by d), for $i, j = 1, \dots, m$ and $b = m+1, \dots, K$,

$$(13) \quad S_{k_i k_j} = S_{ij}, \quad \text{and} \quad S_{k_i b} = S_{ib}.$$

Combining with symmetry a), then for $1 \leq i \leq m$,

$$S_{1i} = S_{i1} = S_{k_1 k_1} = S_{i k_1} = S_{k_1 i}.$$

From (13), we also have $S_{1b} = S_{k_1 b}$ for $b \geq m+1$. Thus the 1-st row and k_1 -th row of S are identical. This contradicts b). Thus this possibility does not exist. The conclusion follows. \square

3.1. Case 1: $\rho_n \geq \epsilon > 0$. We consider $\rho_n \equiv \rho > 0$ for simplicity and take $F \equiv F_{LM}^{(1)}$.

Note that $VWV^T = RSR^T$. Letting $\mathbf{t} = V\pi$, we write

$$F(VWV^T, V\pi) = \sum_{k,l} t_k t_l \tau \left(\frac{(RSR^T)_{kl}}{t_k t_l} \right).$$

By Lemma 3.1, we have

$$\begin{aligned} (14) \quad F(VWV^T, V\pi) &\leq \sum_{k,l} t_k t_l \left(\frac{(RSR^T)_{kl}}{t_k t_l} \log \frac{S_{kl}}{1 - S_{kl}} - \log \left(1 + e^{\log \frac{S_{kl}}{1 - S_{kl}}} \right) \right) \\ &= \sum_{k,l} t_k t_l \left(\frac{(RSR^T)_{kl}}{t_k t_l} \log \frac{S_{kl}}{1 - S_{kl}} - \log \left(\frac{1}{1 - S_{kl}} \right) \right) \\ &= \sum_{k,l} ((RSR^T)_{kl} \log S_{kl} + (R(11^T - S)R^T)_{kl} \log(1 - S_{kl})) \\ &= \sum_{a,b} \left(\sum_{k,l} R_{ka} R_{lb} (S_{ab} \log S_{kl} + (1 - S_{ab}) \log(1 - S_{kl})) \right) \\ (15) \quad &\leq \sum_{a,b} \sum_{k,l} R_{ka} R_{lb} (S_{ab} \log S_{ab} + (1 - S_{ab}) \log(1 - S_{ab})) \\ &= \sum_{a,b} \pi_a \pi_b (S_{ab} \log S_{ab} + (1 - S_{ab}) \log(1 - S_{ab})) \equiv F(W, \pi). \end{aligned}$$

where the last inequality follows from the information inequality

$$p \log q + (1 - p) \log(1 - q) \leq p \log p + (1 - p) \log(1 - p).$$

Obviously, $R = \mathcal{D}(\pi)$ (or its permutation) is the maximizer of $F(RSR^T, R\mathbf{1})$ over \mathcal{R} . Below we show that this is the unique maximizer.

Suppose that $R \in \mathcal{R}$ is a maximizer. Note that for '=' to hold in (15), we have for $1 \leq a, b, k, l \leq K$,

$$R_{ka} R_{lb} \left(S_{ab} \log \left(\frac{S_{ab}}{S_{kl}} \right) + (1 - S_{ab}) \log \left(\frac{1 - S_{ab}}{1 - S_{kl}} \right) \right) = 0,$$

that is,

$$(16) \quad S_{kl} = S_{ab} \quad \text{whenever} \quad R_{ka} R_{lb} > 0.$$

Now by Lemma 3.2, R must be a diagonal matrix, i.e. $\mathcal{D}(\pi)$, or its permutation.

3.2. Case 2: $\rho_n \rightarrow 0$. We take $F \equiv F_{LM}^{(2)}$. By using $\mathbf{t} = V\pi$, we write

$$F(VWV^T, V\pi) = \sum_{k,l} t_k t_l \sigma \left(\frac{(RSR^T)_{kl}}{t_k t_l} \right).$$

The proof for this case is similar to the above using the dual inequality of $\sigma(\cdot)$ in Lemma 3.1 and Lemma 3.2. The details are omitted. Given $\frac{\lambda_n}{\log n} \rightarrow \infty$, (12) implies [5].

4. PROOF OF COROLLARY 2 AND THE LIMITING VARIANCE

Since $\hat{\mathbf{c}} = \mathbf{c}$ with probability approaching 1, it is enough to consider the case with $\hat{\mathbf{c}} = \mathbf{c}$. Then

$$O_{ab} = \sum_{i,j} A_{ij} I(c_i = a, c_j = b).$$

For $a \neq b$, then

$$O_{ab} | \mathbf{c} \sim \text{Binomial}(n^2 \mathbf{f}_a(\mathbf{c}) \mathbf{f}_b(\mathbf{c}), P_{ab}).$$

Thus

$$O_{ab} - E(O_{ab}) = (O_{ab} - E(O_{ab} | \mathbf{c})) + (E(O_{ab} | \mathbf{c}) - E(O_{ab})).$$

It is easy to check that

$$\begin{aligned} \text{var}(O_{ab} - E(O_{ab} | \mathbf{c})) &= E(\text{var}(O_{ab} | \mathbf{c})) \\ &= E(n^2 \mathbf{f}_a(\mathbf{c}) \mathbf{f}_b(\mathbf{c}) P_{ab} (1 - P_{ab})) \\ &= \mu_n S_{ab} (1 - \rho S_{ab}) \pi_a \pi_b (1 + \Omega(n^{-1/2})) \end{aligned}$$

and since

$$E(O_{ab} | \mathbf{c}) - E(O_{ab}) = n^2 P_{ab} (\mathbf{f}_a(\mathbf{c}) \mathbf{f}_b(\mathbf{c}) - E(\mathbf{f}_a(\mathbf{c}) \mathbf{f}_b(\mathbf{c}))),$$

then

$$\begin{aligned} \text{var}(E(O_{ab} | \mathbf{c}) - E(O_{ab})) &= \text{var}(n^2 \mathbf{f}_a(\mathbf{c}) \mathbf{f}_b(\mathbf{c}) P_{ab}) \\ &= \mu_n^2 S_{ab}^2 \text{var}(\mathbf{f}_a(\mathbf{c}) \mathbf{f}_b(\mathbf{c})) \\ &= n^{-1} \mu_n^2 S_{ab}^2 \text{var}(\eta_a \eta_b) (1 + \Omega(n^{-1})) \end{aligned}$$

where $(\eta_1, \dots, \eta_K) = \mathcal{N}(0, \mathcal{D}(\pi) - \pi \pi^T)$, since

$$\begin{aligned} (n \mathbf{f}_1(\mathbf{c}), \dots, n \mathbf{f}_K(\mathbf{c})) &= \text{Multinomial}(n, \pi) \\ &= n \pi + \sqrt{n} \eta (1 + \Omega_P(n^{-1/2})). \end{aligned}$$

Further

$$\sqrt{n} (\mathbf{f}_a(\mathbf{c}) \mathbf{f}_b(\mathbf{c}) - \pi_a \pi_b) = \pi_a \eta_b + \pi_b \eta_a + \Omega_P(n^{-1/2}).$$

Since $\frac{\mu_n}{n} \rightarrow \infty$, the second term $E(O_{ab} | \mathbf{c}) - E(O_{ab})$ dominates and thus

$$\sqrt{n} (\mu_n^{-1} O - W) \Rightarrow S \cdot (\pi \eta^T + \eta \pi^T).$$

Now $L = \mathbf{1}^T O \mathbf{1}$, thus

$$\sqrt{n} (\mu_n^{-1} L - 1) = 2 \pi^T S \eta + \Omega_P(n^{-1/2}).$$

Therefore,

$$\sqrt{n} (L^{-1} O - W) \Rightarrow S \cdot (\pi \eta^T + \eta \pi^T) - 2(\pi^T S \eta) W.$$

Here $(A \cdot B)_{ab} = A_{ab} B_{ab}$.

REFERENCES

- [1] G. Shorack and J. Wellner. *Empirical Processes with Applications to Statistics*. Wiley, New York, 1986.
- [2] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.