

Infinite-degree-corrected stochastic block modelTue Herlau,^{*} Mikkel N. Schmidt,[†] and Morten Mørup[‡]*Section for Cognitive Systems, DTU Compute, Technical University of Denmark, DK-2800 Kongens Lyngby, Denmark*

(Received 9 December 2013; published 29 September 2014)

In stochastic block models, which are among the most prominent statistical models for cluster analysis of complex networks, clusters are defined as groups of nodes with statistically similar link probabilities within and between groups. A recent extension by Karrer and Newman [Karrer and Newman, *Phys. Rev. E* **83**, 016107 (2011)] incorporates a node degree correction to model degree heterogeneity within each group. Although this demonstrably leads to better performance on several networks, it is not obvious whether modeling node degree is always appropriate or necessary. We formulate the degree corrected stochastic block model as a nonparametric Bayesian model, incorporating a parameter to control the amount of degree correction that can then be inferred from data. Additionally, our formulation yields principled ways of inferring the number of groups as well as predicting missing links in the network that can be used to quantify the model's predictive performance. On synthetic data we demonstrate that including the degree correction yields better performance on both recovering the true group structure and predicting missing links when degree heterogeneity is present, whereas performance is on par for data with no degree heterogeneity within clusters. On seven real networks (with no ground truth group structure available) we show that predictive performance is about equal whether or not degree correction is included; however, for some networks significantly fewer clusters are discovered when correcting for degree, indicating that the data can be more compactly explained by clusters of heterogeneous degree nodes.

DOI: [10.1103/PhysRevE.90.032819](https://doi.org/10.1103/PhysRevE.90.032819)

PACS number(s): 89.75.Hc, 89.75.Da

I. INTRODUCTION

The stochastic block model (SBM) [1–3] has become a prominent tool for modeling group structure in complex networks [4]. However, as pointed out by Karrer and Newman [5], the stochastic block model has a tendency to group nodes according to their degree such that high-degree nodes group together even though their patterns of interactions with the remaining network may differ. This grouping thus reflects aspects of node degree rather than overall statistical patterns in the network. To alleviate this issue, Karrer and Newman introduced the degree-corrected stochastic block model (DCSBM) [5]. In their model, additional parameters modeling node degree heterogeneity are introduced, allowing nodes of varying degree to be clustered together, and they demonstrate that including this degree correction reduces the tendency to group nodes according to their degree distribution [5]. The parameters in the DCSBM are inferred using maximum likelihood (ML) estimation and since closed-form expressions for the ML estimates of the additional degree-correction parameters are available, the computational complexity of the inference procedure is similar to inference in the SBM.

Although Karrer and Newman demonstrate on several network data sets that degree correction leads to better performance [5], it is not obvious whether including a degree correction is always appropriate on real network data. Furthermore, the number of groups used in the analysis is likely to influence the results since groups of heterogeneous node degree can be reasonably modeled by a number of homogeneous subgroups. Not handling this issue in a principled manner

could potentially confound the results. Finally, an important subject of network modeling is validation. Although many real networks are hypothesized to possess group structure, no ground truth clustering is available, which makes it difficult to assess the goodness of the obtained clustering. A popular alternative is to measure the predictive performance on held-out links in the network. In order to do this in a principled manner the methods must be able to handle missing entries in the network data as well as define a predictive distribution over the missing entries.

In this paper we address these three important challenges when modeling network data by the DCSBM.

(i) Can we infer the extent to which degree correction is necessary?

(ii) How can we determine the number of components?

(iii) How can we predict links in the DCSBM?

In particular, we formulate a nonparametric Bayesian generative model for the DCSBM. The number of components is inferred using the Chinese Restaurant Process, which has previously been used to determine the number of components in stochastic block models [6,7]. Our generative model is characterized by admitting a simple inference procedure in which both the degree parameter and group interactions can be analytically marginalized such that inference reduces to estimating the assignments of nodes to clusters as for the DCSBM. We address the link-prediction problem using Markov chain Monte Carlo (MCMC) imputation. By inferring the hyperparameter in the prior distribution of the parameters that account for heterogeneous node degree our model is able to determine the extent to which a degree correction is necessary, possibly reducing to an uncorrected stochastic block model. On synthetic as well as seven real networks, we demonstrate the utility of our proposed model for determining the number of components, link prediction, and inferring the magnitude of the parameter controlling degree correction.

^{*}tuhe@dtu.dk[†]mnsc@dtu.dk[‡]mmor@dtu.dk

Past work on the SBM and DCSBM has not treated the problem of inferring components, the presence of degree heterogeneity, and link prediction under one unified framework. Although Bayesian approaches to inferring components and link prediction have a long history for the SBM [4,6,7], most work on the DCSBM has been focused on other inference methods. As noted, Karrer and Newman [5] treated the problem of inference in the DCSBM from a ML perspective. A related approach was taken by Peixoto [8], who considered degree correction as constraints on a block-model ensemble and derived an entropy-based cost function. For the SBM, a method relying on a minimum-description-length-based approach to learning has been proposed, giving rise to an efficient maximization procedure [9]. The minimum-description-length approach by Rosvall and Bergstrom [10] allows degree correction, but is otherwise analytically different from the DCSBM. For the DCSBM minimum-description-length-based procedures were considered by Peixoto [11] to give an efficient MCMC-based inference procedure (see also [12] for additional discussion of this approach and an application to the problem of estimating the number of components). The belief propagation method of Decelle *et al.* [13,14] may also be applied to the DCSBM. More related to our approach is that of Yan *et al.* [15], who consider the problem of inferring the number of groups in the DCSBM from a model-selection perspective. While these approaches represent important contributions to the problem of jointly modeling degree heterogeneity and block structure, none of the current proposals are based on a Bayesian generative model and allow joint inference of degree correction, number of components, and missing links using a MCMC-based approach.

II. METHODS

Let A be the adjacency matrix of an undirected observed network of n nodes such that A_{ij} is the number of links between nodes i and j . We allow a positive number of self-links A_{ii} in our model definition (note that in the original formulation of the DCSBM [5] A_{ii} is defined as twice the number of self-links). The DCSBM [5] for an undirected graph assumes that the links between nodes i and j follow a Poisson distribution

$$A_{ij} \sim \mathcal{P}(\theta_i \eta_{z_i z_j} \theta_j) \quad \text{for } i \neq j. \quad (1)$$

The parameter $\eta_{\ell m}$ controls the probability of links between nodes in groups ℓ and m , $z_i = \ell$ indicates that node i is assigned to group ℓ , and θ_i is a node-specific parameter that regulates this link probability and thus accounts for heterogeneous node degrees. The model is subject to the constraint that $\sum_i \delta_{z_i \ell} \theta_i = 1$ for all groups ℓ , i.e., the sum of the θ_i within each group is one.

We presently propose a nonparametric Bayesian generative model that extends the DCSBM dubbed the infinite-degree-corrected stochastic block model (IDCSBM). Like the DCSBM, we also maintain node weights θ_i to control the degree, however, to arrive at a Bayesian formulation we assume that the weights within each group are drawn from a Dirichlet distribution. More precisely, for each group ℓ containing n_ℓ nodes, we introduce an n_ℓ -dimensional vector of weights $(\phi_i)_{z_i=\ell}$ drawn from a Dirichlet distribution and define $\theta_i = n_\ell \phi_i$ in Eq. (1).

The scaling by n_ℓ makes the average degree of any given node independent of the size of the group the node belongs to. The full model now consists of (i) generating a random partition, (ii) generating the interaction between each group of the partition $\eta_{\ell m}$ from a gamma distribution, (iii) for each group, generating $(\phi_i)_{z_i=\ell}$ from a Dirichlet distribution and rescaling with n_ℓ , and finally (iv) using Eq. (1) to generate the number of links A_{ij} between node $i \neq j$.

The full model is given generatively below. The symbol \mathcal{D} denotes the Dirichlet distribution and \mathcal{G} the gamma distribution. For analytical convenience the model assumes a particular parametrization of the self-links A_{ii} , a point we will return to later:

$$\mathbf{z} \sim \mathcal{C}(\alpha) \quad (\text{clusters}), \quad (2)$$

$$(\phi_i)_{z_i=\ell} \sim \mathcal{D}(\gamma \mathbf{1}_{(n_\ell)}) \quad \text{for } \ell \geq 0, \quad (3)$$

$$\theta_i = n_{z_i} \phi_i \quad (\text{relative degree}), \quad (4)$$

$$\eta_{\ell m} \sim \mathcal{G}(\kappa, \lambda) \quad \text{for } \ell \leq m \text{ (link rate)}, \quad (4)$$

$$A_{ij} \sim \mathcal{P}(\theta_i \eta_{z_i z_j} \theta_j) \quad \text{for } i < j \text{ (link weight)}, \quad (5)$$

$$A_{ii} \sim \mathcal{P}\left(\frac{1}{2} \theta_i^2 \eta_{z_i z_i}\right) \quad \text{for } i = j.$$

In the above $\mathbf{1}_{(n_\ell)}$ is a vector of ones with length n_ℓ , $N = \sum_{\ell=1}^L n_\ell$ is the total number of nodes, and L is the number of groups. As a prior over the node partition \mathbf{z} we use $\mathcal{C}(\alpha)$, the Chinese restaurant process (CRP) parametrized by a single parameter α controlling the distribution of group size [16]. A potential advantage of the CRP over, for instance, a uniform prior over partitions is that the CRP is consistent under projections whereas the uniform prior is not. The simplest example is the case where \mathbf{z} is a partition of two nodes assigned to the same group (i.e., $z_1 = z_2 = 1$) and we consider a partition obtained by including a third node. In this case for the CRP it holds: $p(z_1 = z_2 = 1 | \alpha) = p(z_1 = z_2 = 1, z_3 = 1 | \alpha) + p(z_1 = z_2 = 1, z_3 = 2 | \alpha)$, however for the uniform prior the left-hand side is $\frac{1}{2}$ and the right-hand side $\frac{2}{5}$.

Notice the role played by γ in the Dirichlet distribution in Eq. (3). If $\gamma \rightarrow \infty$, we will have $\phi_i \rightarrow \frac{1}{n_\ell}$ for $z_i = \ell$ or simply $\theta_i \rightarrow 1$ for all i (the limits are understood in the distribution) and the model is thus independent of degree in Eq. (1). On the other hand, for $\gamma \rightarrow 0$, within each group ℓ a single node i^* will have mass $\theta_{i^*} = n_\ell$ and the network becomes very nearly entirely dominated by a few nodes. We return to the properties of the model in Sec. II B. The advantage of a Bayesian formulation is that we can infer not only θ_i , but also a distribution of the degree-correction variable γ representing the appropriateness of modeling degree heterogeneity for the network.

By collecting variables of the same type the joint density factorizes as

$$\begin{aligned} p(A, \phi, \eta, \mathbf{z} | \alpha, \gamma, \kappa, \lambda) \\ = p(A | \theta, \eta, \mathbf{z}) p(\eta | \kappa, \lambda) p(\phi | \mathbf{z}, \gamma) p(\mathbf{z} | \alpha). \end{aligned} \quad (6)$$

The model thus depend on parameters $(\alpha, \gamma, \kappa, \lambda)$. While one could fix these at a particular value, a more principled approach we have taken is to introduce vague uninformative priors and sample these as well [17]. Either choice has no effect on the following derivation below. In our notation the relevant

densities are

$$p(\mathbf{z}|\alpha) = \frac{\alpha^L \Gamma(\alpha)}{\Gamma(N + \alpha)} \prod_{\ell=1}^L \Gamma(n_\ell) \quad (\text{CRP}), \quad (7)$$

$$\mathcal{D}(\mathbf{x}|\gamma) = \frac{1}{B(\gamma)} \prod_i x_i^{\gamma_i-1}, \quad B(\gamma) = \frac{\prod_i \Gamma(\gamma_i)}{\Gamma(\sum_i \gamma_i)}, \quad (8)$$

$$\mathcal{G}(x|\kappa, \lambda) = \frac{1}{G(\kappa, \lambda)} x^{\kappa-1} e^{-\lambda x}, \quad G(\kappa, \lambda) = \lambda^{-\kappa} \Gamma(\kappa). \quad (9)$$

The advantage of the present formulation is the use of the Dirichlet distribution within each group and the particular parametrization of A_{ij} , which allows the node weights as well as group interactions to be integrated out analytically. To see this we introduce the shorthand notation for between- and within-group link counts

$$N_{\ell m}^+ = \begin{cases} \sum_{\substack{i: z_i = \ell, \\ j: z_j = m}} A_{ij}, & \ell \neq m \\ \sum_{\substack{i \leq j: \\ z_i = z_j = \ell}} A_{ij}, & \ell = m, \end{cases}$$

$$N_{\ell m} = \begin{cases} n_\ell n_m, & \ell \neq m \\ \frac{n_\ell n_\ell}{2}, & \ell = m. \end{cases}$$

as well as node degrees $k_i = \sum_j A_{ij}$ and $\hat{k}_i = k_i + A_{ii}$. It now follows by some algebra that

$$\begin{aligned} p(\mathbf{A}|\boldsymbol{\theta}, \boldsymbol{\eta}, \mathbf{z}) &= \prod_{i < j} \frac{(\theta_i \eta_{z_i z_j} \theta_j)^{A_{ij}}}{A_{ij}! e^{\theta_i \eta_{z_i z_j} \theta_j}} \prod_i \frac{\left(\frac{\theta_i^2 \eta_{z_i z_i}}{2}\right)^{A_{ii}}}{A_{ii}! e^{(1/2)\theta_i^2 \eta_{z_i z_i}}} \\ &= \frac{\prod_i 2^{-A_{ii}}}{\prod_{i < j} A_{ij}!} \prod_{\ell \leq m} \eta_{\ell m}^{N_{\ell m}^+} e^{-\eta_{\ell m} N_{\ell m}} \prod_i \theta_i^{k_i + A_{ii}} \\ &= \frac{\prod_i 2^{-A_{ii}}}{\prod_{i < j} A_{ij}!} \left[\prod_{\ell \leq m} \eta_{\ell m}^{N_{\ell m}^+} e^{-\eta_{\ell m} N_{\ell m}} \right] \prod_\ell n_\ell^{\hat{k}_\ell} \prod_{i: z_i = \ell} \phi_i^{\hat{k}_i}, \end{aligned} \quad (10)$$

$$p(\boldsymbol{\eta}|\kappa, \lambda) = \prod_{\ell \leq m} \frac{1}{G(\kappa, \lambda)} \eta_{\ell m}^{\kappa-1} e^{-\eta_{\ell m} \lambda}, \quad (11)$$

$$p(\boldsymbol{\phi}|\mathbf{z}, \gamma) = \prod_\ell \frac{1}{B(\gamma \mathbf{1}_{(n_\ell)})} \frac{\prod_{i: z_i = \ell} \Gamma(\gamma) \left(\frac{\theta_i}{n_\ell}\right)^{\gamma-1}}{n_\ell \Gamma(n_\ell \gamma)}. \quad (12)$$

Inserting into Eq. (6), collecting terms, and exploiting the conjugacy of the Dirichlet and gamma distributions to the Poisson distribution, we can analytically marginalize (i.e., collapse) $\boldsymbol{\phi}$ and $\boldsymbol{\eta}$ to obtain

$$\begin{aligned} p(\mathbf{A}, \mathbf{z}|\alpha, \gamma, \kappa, \lambda) &= \int d\boldsymbol{\eta} d\boldsymbol{\phi} p(\mathbf{A}|\boldsymbol{\theta}, \boldsymbol{\eta}, \mathbf{z}) p(\boldsymbol{\eta}|\kappa, \lambda) p(\boldsymbol{\phi}|\mathbf{z}, \gamma) p(\mathbf{z}|\alpha) \\ &= \frac{1}{\prod_{i < j} A_{ij}! \prod_i 2^{A_{ii}}} \prod_{\ell \leq m} \frac{G(N_{\ell m}^+ + \kappa, N_{\ell m} + \lambda)}{G(\kappa, \lambda)} \\ &\quad \times \left[\prod_\ell \frac{B(\gamma \mathbf{1}_{(n_\ell)} + (\hat{k}_i)_{i: z_i = \ell})}{B(\gamma \mathbf{1}_{(n_\ell)})} n_\ell^{\hat{k}_\ell} \right] \frac{\alpha^L \Gamma(\alpha)}{\Gamma(N + \alpha)} \prod_{\ell=1}^L \Gamma(n_\ell). \end{aligned} \quad (13)$$

In the above derivation we exploit that $\sum_{z_i = \ell} \theta_i = n_\ell$ and thus the derivation requires access to the entire network. As a result, the inference of our generative model is reduced to determining the posterior distribution of the assignment of nodes to groups \mathbf{z} .

The assignment matrix \mathbf{z} is inferred using standard Gibbs sampling [6] and using the Bayesian framework we can treat the hyperparameters γ , α , λ , and κ as random variables. In particular, we will invoke the uninformative prior $p(x) \propto x^{-1}$ for all four parameters and infer them using random-walk Metropolis updates of the form $x^* = \exp(\ln x + z)$, $z \sim N(0, \sigma = 0.1)$. For each Gibbs sweep over \mathbf{z} , we performed 20 Metropolis-Hastings updates of the hyperparameters. While Metropolis-Hastings updates with random proposals are not very computationally efficient, we noticed throughout the experiments that this step had a small computational cost compared to sampling \mathbf{z} .

A. Imputation and link prediction

Missing (unobserved) links commonly occur in the network and predicting missing links is an important goal of network modeling. Comparing the prediction of a model on unobserved data to the actual value is furthermore a popular way to validate a model. In addition, the self-links A_{ii} are often unknown or, if the network cannot contain self-links such as the case of a friendship network, they should be treated as auxiliary variables that are integrated out.

For the IDCSBM the (marginalized) expression for \mathbf{z} in Eq. (13) requires access to all entries in the adjacency matrix and so it is not possible to marginalize missing data simply by ignoring the corresponding terms in the likelihood function. To overcome this difficulty we marginalize missing entries by formulating a Markov chain Monte Carlo algorithm jointly over the parameters and the missing links. This is done by sampling \mathbf{z} and the hyperparameters using Gibbs sampling and random-walk Metropolis-Hastings updates and then conditionally on \mathbf{A} and \mathbf{z} drawing values of $\eta_{\ell m}$ and $(\phi_i)_{z_i = \ell}$ conditional on the full matrix \mathbf{A} and assignments \mathbf{z} and conditionally on these values draw the values of \mathbf{A} corresponding to the missing links from the Poisson distribution equation (5). This corresponds to imputing the missing values from their predictive distribution in each step of the MCMC algorithm and, assuming convergence of the Markov chain, is equivalent to marginalizing the missing links. We use this framework both to handle self-links and also for link prediction in general. Another popular method to predict missing data is simply replacing missing entries of \mathbf{A} with 0 [4, 5, 18]; however, as the diagonal of \mathbf{A} is often fully missing and the Poisson rate for A_{ii} is proportional to θ_i^2 , this approach would create an undesirable bias for θ_i .

B. Properties of the model

An important property of the model is that it can accurately determine the degree distribution of the data and the link density between the groups. Suppose \mathbf{A}_0 is an observed network and let \mathbf{z} be any fixed cluster. Conditional on \mathbf{A}_0 and \mathbf{z} , we may compute the posterior over $\boldsymbol{\eta}$ and $\boldsymbol{\theta}$ and check if these distributions accurately reflect relevant properties of

A_0 . Notice from Eq. (10) that the posterior distributions of η and θ are

$$p(\eta_{\ell m} | A_0, \mathbf{z}) = \mathcal{G}(\eta_{\ell m} | N_{\ell m}^+ + \kappa, N_{\ell m} + \lambda), \quad (14)$$

$$p\left[\left(\frac{\theta_i}{n_\ell}\right)_{z_i=\ell} \middle| A_0, \mathbf{z}\right] = \mathcal{D}\left[\left(\frac{\theta_i}{n_\ell}\right)_{z_i=\ell} \middle| \gamma \mathbf{1}_{n_\ell} + (\hat{k}_i)_{z_i=\ell}\right]. \quad (15)$$

Recall for two Poisson distributed random variables $X \sim \mathcal{P}(a)$ and $Y \sim \mathcal{P}(b)$ that their sum is Poissonian with rate $a + b$: $X + Y \sim \mathcal{P}(a + b)$. This, along with the derivation (10), allows us to compute various properties of the model.

First consider the total interaction strength between two groups ℓ and m . The interaction $\sum_{i \leq j} \delta_{z_i=\ell} \delta_{z_j=m} A_{ij}$, considered as a random variable, is then distributed as $\mathcal{P}(\eta_{\ell m} N_{\ell m})$. If $X \sim \mathcal{P}(\lambda)$ then $\mathbb{E}[X] = \lambda$ and so the average between-group interaction is [the expectation is with respect to $p(\cdot | A_0, \mathbf{z})$]

$$\mathbb{E}\left[\sum_{i \leq j} \delta_{z_i=\ell} \delta_{z_j=m} A_{ij}\right] = \mathbb{E}[N_{\ell m} \eta_{\ell m}] = \frac{N_{\ell m}(N_{\ell m}^+ + \kappa)}{N_{\ell m} + \lambda}. \quad (16)$$

For analytical simplicity, we will consider the degree plus the diagonal element. To this end we define the degree of node i as $d_i = \sum_j A_{ij} + A_{ii}$. Since each A_{ij} is Poisson distributed the degree too is a Poissonian random variable. If $z_i = \ell$ then d_i 's distribution is given by

$$d_i \sim \mathcal{P}\left(\sum_{j \neq i} \theta_i \eta_{\ell z_j} \theta_j + 2 \frac{\theta_i^2 \eta_{\ell \ell}}{2}\right) = \mathcal{P}\left(\theta_i \sum_m \eta_{\ell m} n_m\right). \quad (17)$$

We may now compute the average, again with respect to A_0 and fixed \mathbf{z} ,

$$\begin{aligned} \mathbb{E}[d_i] &= \mathbb{E}\left[\theta_i \sum_m \eta_{\ell m} n_m\right] \\ &= n_\ell \frac{\hat{k}_i + \gamma}{\sum_{j: z_j=\ell} \hat{k}_j + \gamma n_\ell} \sum_m \frac{N_{\ell m}^+ + \kappa}{N_{\ell m} + \lambda} n_m \\ &= (\hat{k}_i + \gamma) \sum_m \frac{N_{\ell m} 2^{\delta_{\ell m}}}{N_{\ell m} + \lambda} \frac{N_{\ell m}^+ + \kappa}{\sum_h N_{\ell h}^+ 2^{\delta_{\ell h}} + \gamma n_\ell}. \end{aligned} \quad (18)$$

Assuming that the groups are fairly large and in the low limit of the prior γ , the sum will be 1 to first order. The derivations (16) and (18) show in the limit of large systems that the relative influence of the prior terms will vanish and the model will accurately capture the between-group link density as well as the node degree.

III. RESULTS AND DISCUSSION

We analyze synthetic data sets generated from our model as well as seven real networks from the literature.

A. Synthetic data

In our synthetic simulation studies we generated networks of $N = 80$ nodes from our generative model with the parameters κ and α fixed at $\kappa = 0.5$ and $\alpha = 4$ and under different values of λ and γ . Each such network was analyzed using our IDCSBM as well as the corresponding ISBM without degree correction. In Fig. 1 the normalized mutual information (NMI), the ratio of the true number of components to the estimated number of components $L_{\text{frac}} = \langle \frac{L}{L_{\text{true}}} \rangle$, and the area under the curve (AUC) of the receiver operator characteristic are given (error bars indicate standard deviation of the mean where the deviation is computed over ten restarts of the sampler). In the analysis we ran the samplers for 1000 iterations and discarded the first half as burn-in. The AUC scores were computed by treating 5% of the links and a similar number of nonlinks as missing.

From the plot of the NMI we see that the degree-corrected model (IDCSBM) better recovers the true generated group structure than the degree-uncorrected model (ISBM) and as expected the performances of the two methods converge as γ increases, corresponding to networks that do not exhibit degree heterogeneity. Furthermore, the IDCSBM recovers the correct number of groups whereas the ISBM generates more than the true number of groups in order to account for the effect of a skewed degree distribution. The predictive performance as quantified by the AUC scores are more or less similar with a tendency of slightly better predictions for the IDCSBM. As expected, this is most notable for small values of γ . We further observe that structure is better recovered when the contrast in the interactions is high, as influenced by the values of λ . This too can be expected since very sparse networks presumably have little recoverable structure.

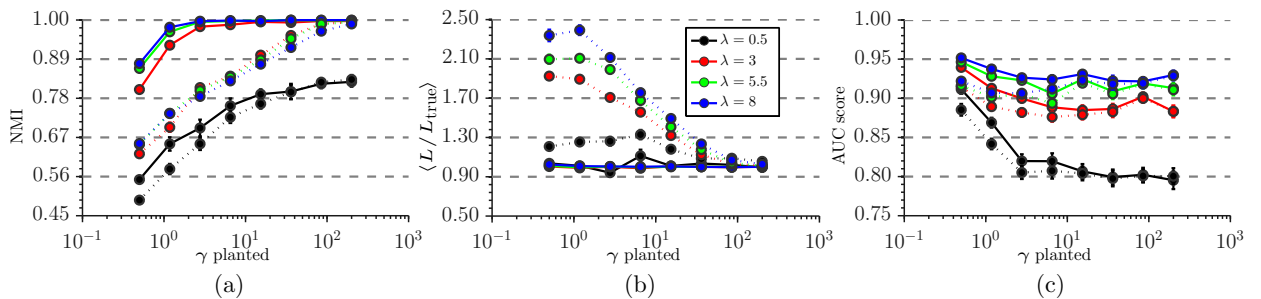


FIG. 1. (Color online) Plots of the IDCSBM and ISBM results on simulated networks, showing (a) the normalized mutual information, (b) the ratio of estimated to true number of components L_{frac} , and (c) the area under the curve of the receiver operator characteristics as computed by running the proposed methods on networks produced from the generative model of the IDCSBM with different values of λ and γ . The solid lines indicate results for the IDCSBM and the dotted lines indicate results for the ISBM.

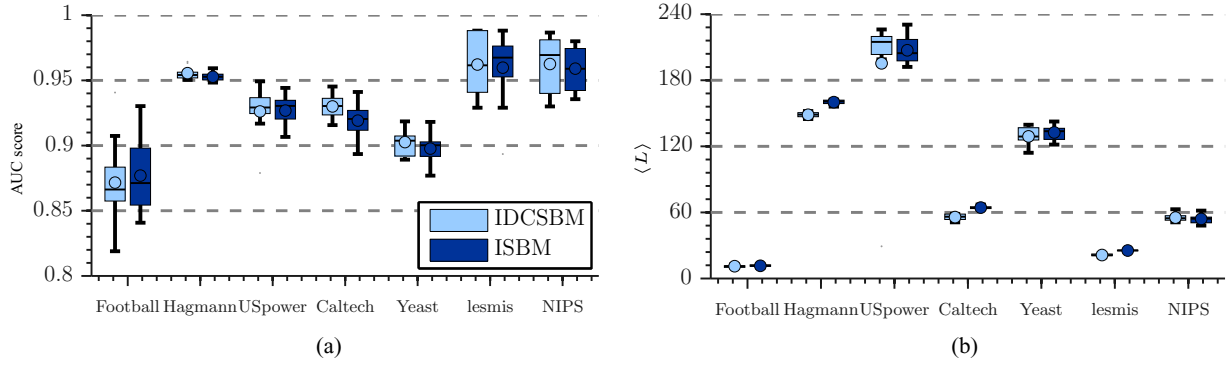


FIG. 2. (Color online) Plots of the IDCSBM and ISBM results on the seven real network for (a) AUC scores on held-out links and (b) the number of inferred groups L ; results are averaged over ten random restarts. The degree-corrected and -noncorrected methods perform roughly similarly with a tendency for the degree-corrected model to find fewer groups.

B. Real data

We analyzed the following seven networks.

(i) *Football*: undirected unweighted network of American football games between 115 Division IA colleges in the fall of 2000 [19].

(ii) *Hagmann*: undirected weighted network of the number of links between 998 brain regions as estimated by tractography from diffusion spectrum imaging across five subjects [20]. The graph of each subject has been symmetrized, thresholded at zero, and the five subject graphs added together.

(iii) *USPower*: undirected unweighted network of 4941 nodes representing the topology of the Western States Power Grid of the United States compiled in Ref. [21].

(iv) *Caltech*: the Caltech39 social network of 769 students from the Facebook100 data set [22].

(v) *Yeast*: the interaction network between 2361 proteins of yeast [23].

(vi) *Lesmis*: undirected and weighted graph of the coappearances of 77 characters in Les Misérables by Hugo [24].

(vii) *NIPS*: undirected weighted network of the number of coauthorships between 234 authors of papers presented at the Conference on Neural Information Processing Systems 1–12 [25].

Figure 2 shows the results for the IDCSBM and the ISBM on the seven networks in terms of AUC score treating 5% of the links (and a similar number of nonlinks) as missing. Furthermore, the numbers of components estimated by the two models are given. The samplers were run for 1000 iterations (half discarded as burn-in) and the results are averaged over ten restarts.

From Fig. 2 it can be seen that in general the performance in predicting a link as quantified by the AUC scores is on par for the IDCSBM and ISBM. However, as observed also in the synthetic study, the IDCSBM model extracts fewer components than the ISBM for the Hagmann, Caltech, and Lesmis networks. Thus, the model allocates fewer groups when compared to the ISBM, which allocates additional clusters in order to compensate for its lack of ability to explicitly account for degree.

Another way to examine this effect is to look at the degree distribution within each group. Since the groups have vastly different sizes it is hard to summarize this effect into a single

number; however, if we consider a fixed group structure \mathbf{z} and a single group ℓ of size n_ℓ we may compute the empirical mean $\mathbb{E}[k_\ell] = \frac{1}{n_\ell} \sum_{i:z_i=\ell} k_i$ and standard deviation $\mathcal{S}[k_\ell] = \sqrt{\frac{1}{n_\ell} \sum_{i:z_i=\ell} (k_i - \mathbb{E}[k_\ell])^2}$ of the degree within this group.

Plotted in Fig. 3 is the average of the empirical standard deviation of the degree distribution as a function of group size, that is, for each point (k, y) in Fig. 3, y is an estimate of $\mathbb{E}[\mathcal{S}[k_\ell]]$, where the expectation is conditional on $n_\ell = k$. This quantity is easily estimated based on the last 500 states of a MCMC chain. The error bars are the standard deviation of the mean of each point based on ten random restarts of the sampler.

As can be seen, the IDCSBM reveals larger groups of nodes, confirming our previous findings in Fig. 2 and, more importantly, that the variance of the degree distribution within groups is larger than for the ISBM for all groups sizes. This shows that the compensation for degree heterogeneity affects not only a few large groups the IDCSBM lumps together and the ISBM splits apart, but groups of all sizes.

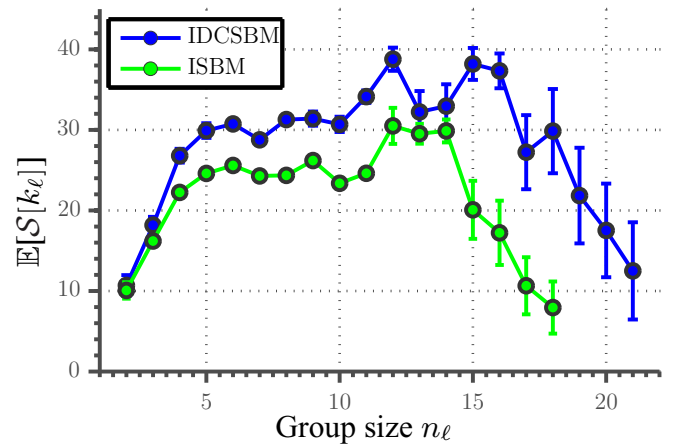


FIG. 3. (Color online) Variance of degree heterogeneity for the ISBM and IDCSBM for the Hagmann data set. Each point (k, y) is an estimate of the standard deviation of the degree distribution for nodes in a group ℓ of size $n_\ell = k$ (see main text for details).

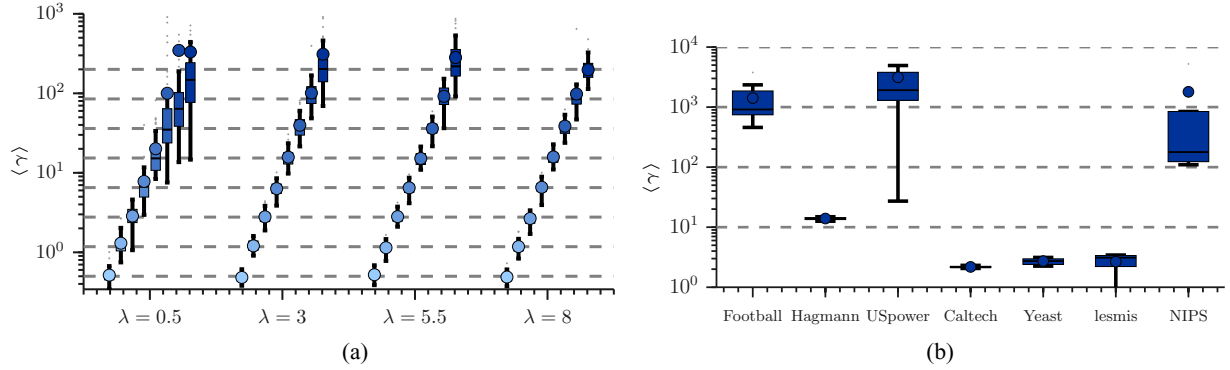


FIG. 4. (Color online) Inferred values of $\langle \gamma \rangle$ for (a) the artificial and (b) the real networks. The box plots show the inferred mean of γ for each of the 10 (or 50) MCMC chains (on artificial or real networks). For the artificial network (a), the networks are grouped according to the planted value of λ (controlling link density) and each of the eight boxes in a group corresponds to a planted value of γ , the planted values indicated by the horizontal lines. In the limit of good sampling the boxes should lie on the dotted lines. As shown, the sampler infers the correct value of degree correction for the artificial networks except for very sparse networks ($\lambda = 0.5$). For the real networks the model infers very different degrees of node heterogeneity.

To better understand the role of γ , we examined the behavior of the mean value of γ , $\langle \gamma \rangle$, across the random restarts of the chains for both the artificial and real data sets (see Fig. 4). For the artificial data sets [Fig. 4(a)] we grouped the networks according to the values of λ and γ used to generate the networks and plot the value of $\langle \gamma \rangle$ across the 50 restarts. Consistent with the other findings, the model has more difficulties recovering the true value of γ for very low link density ($\lambda = 0.5$) or when the planted value of γ is very high, here 200 as the highest value. The later finding may be related to this value not being favored by the prior. However, the sampler generally recovers the planted value of γ well across chains.

For the real networks [Fig. 4(b)], the recovered values of $\langle \gamma \rangle$ across chains show quite high variability for some of the larger networks, indicating that they may exhibit mixing times significantly longer than the 1000 iterations used here. Notice that since high values of γ are associated with a nearly vanishing effect of the degree, we see that the model correctly identifies the skewed degree distribution of the social networks Caltech and Yeast, while indicating that the effect of degree for the (very strongly) community-structured network Football and the spatially embedded USPower network is vanishing.

IV. CONCLUSION

In this paper we extended the degree-corrected stochastic block model [5] to a nonparametric Bayesian generative model (the IDCSBM). The advantage of the proposed model is that the number of blocks, i.e., the distribution of the number of groups, can be inferred, extending the model to an infinite representation similar to what has previously been done for

the regular stochastic block model [6,7]. By exploiting that the model is formulated generatively, we have derived a Markov chain Monte Carlo algorithm that handles missing links explicitly by marginalizing missing entries. We have further shown that we can learn the parameter γ in the process and thereby determine the extent to which networks can use the degree-correction parameter θ introduced in the degree-corrected stochastic block model. We have shown analytically that under a wide range of conditions the model will be able to accurately model between-group link density as well as node degree.

On synthetic and real networks, we demonstrated that the IDCSBM can result in a more compact representation of network structure. The IDCSBM also tends to use fewer components than the ISBM while accounting equally well for the networks as quantified by the AUC link prediction scores. On synthetic data with degree heterogeneity, we have shown that the proposed model, which corrects for degree skewness, is able to infer the parameters controlling degree heterogeneity correctly and obtain a more compact and accurate representation. As expected, this also translates in improved link prediction. On real network data, we have shown that a model that captures degree skewness does not dominate a model that does not in terms of prediction; however, the IDCSBM is able to consistently determine vastly different values of γ and thereby the presence or absence of degree heterogeneity.

ACKNOWLEDGMENT

This project was supported by the Lundbeck Foundation, Grant No. R105-9813.

- [1] H. C. White, S. A. Boorman, and R. L. Breiger, *Am. J. Sociol.* **81**, 730 (1976).
- [2] P. W. Holland, K. B. Laskey, and S. Leinhardt, *Soc. Networks* **5**, 109 (1983).

- [3] K. Nowicki and T. A. B. Snijders, *J. Am. Stat. Assoc.* **96**, 1077 (2001).
- [4] R. Guimer and M. Sales-Pardo, *Proc. Natl. Acad. Sci. USA* **106**, 22073 (2009).

- [5] B. Karrer and M. E. J. Newman, *Phys. Rev. E* **83**, 016107 (2011).
- [6] C. Kemp, J. B. Tenenbaum, T. L. Griffiths, T. Yamada, and N. Ueda, in *Proceedings of the 21st National Conference on Artificial Intelligence* (AAAI, Palo Alto, 2006), Vol. 3, p. 5.
- [7] Z. Xu, V. Tresp, K. Yu, and H.-P. Kriegel, in *Proceedings of the 22nd Annual Conference on Uncertainty in Artificial Intelligence*, edited by R. Dechter and T. Richardson (AUAI Press, Arlington, Virginia, 2006), pp. 544–551.
- [8] T. P. Peixoto, *Phys. Rev. E* **85**, 056122 (2012).
- [9] M. Rosvall and C. T. Bergstrom, *Proc. Natl. Acad. Sci. USA* **104**, 7327 (2007).
- [10] M. Rosvall and C. T. Bergstrom, *Proc. Natl. Acad. Sci. USA* **105**, 1118 (2008).
- [11] T. P. Peixoto, *Phys. Rev. Lett.* **110**, 148701 (2013).
- [12] T. P. Peixoto, *Phys. Rev. E* **89**, 012804 (2014).
- [13] A. Decelle, F. Krzakala, C. Moore, and L. Zdeborová, *Phys. Rev. Lett.* **107**, 065701 (2011).
- [14] A. Decelle, F. Krzakala, C. Moore, and L. Zdeborová, *Phys. Rev. E* **84**, 066106 (2011).
- [15] X. Yan, C. R. Shalizi, J. E. Jensen, F. Krzakala, C. Moore, L. Zdeborová, P. Zhang, and Y. Zhu, *J. Stat. Mech.* (2014) P05007.
- [16] D. J. Aldous, in *École d'Été de Saint-Flour XIII—1983*, Exchangeability and related topics, edited by P.-L. Hennequin, Lecture Notes in Math., Vol. 1117 (Springer-Verlag, Berlin, 1985), pp. 1–198.
- [17] E. Jaynes, *Probability Theory: The Logic of Science* (Cambridge University Press, Cambridge, 2003).
- [18] A. Clauset, C. Moore, and M. E. Newman, *Nature (London)* **453**, 98 (2008).
- [19] M. Girvan and M. E. Newman, *Proc. Natl. Acad. Sci. USA* **99**, 7821 (2002).
- [20] P. Hagmann, L. Cammoun, X. Gigandet, R. Meuli, C. J. Honey, V. J. Wedeen, and O. Sporns, *PLoS Biol.* **6**, e159 (2008).
- [21] D. J. Watts and S. H. Strogatz, *Nature (London)* **393**, 440 (1998).
- [22] Available at <https://archive.org/details/oxford-2005-facebook-matrix>.
- [23] D. Bu, Y. Zhao, L. Cai, H. Xue, X. Zhu, H. Lu, J. Zhang, S. Sun, L. Ling, N. Zhang *et al.*, *Nucleic Acids Res.* **31**, 2443 (2003).
- [24] D. E. Knuth, *The Stanford GraphBase: A Platform for Combinatorial Computing* (Addison-Wesley, Reading, 1993), Vol. 4.
- [25] Available at <http://www.cs.nyu.edu/~roweis/data.html>.