
Model-Based Clustering of Large Random Graphs With High-Dimensional Predictors

Duy Q. Vu

Department of Mathematics and Statistics
University of Melbourne
duy.vu@unimelb.edu.au

Michael Schweinberger

Department of Statistics
Rice University
m.s@rice.edu

Abstract

Model-based clustering of random graphs has attracted much attention in machine learning and statistics. We improve the predictive power of model-based clustering by exploiting high-dimensional predictors and encouraging sparsity by using a penalized Bayesian approach. We estimate the high-dimensional posterior by using variational methods. To implement variational methods, we exploit the minorization-maximization (MM) principle to construct minorizing functions which separate the parameters and facilitate the exploration of the high-dimensional posterior. The resulting Bayesian Variational GEM-MM algorithm is less prone to be trapped at local maxima than conventional Variational EM algorithms. We demonstrate the success of the model and method by using two real-world networks with more than 10,000 nodes and 15,000 predictors.

1 Introduction

Data which can be represented by random graphs arise in computer science (e.g., social networks, World Wide Web), engineering (e.g., power networks), the life sciences (e.g., biological networks), the health sciences (e.g., sexual networks), and the social sciences (e.g., terrorist networks). An attractive approach to modeling random graphs model is based on the assumption that edge variables are exchangeable. An extension of de Finetti’s theorem to random graphs suggests that exchangeable edge variables can be modeled by mixture models [1]. A simple mixture model is a finite mixture model, which assumes that edges are governed by a finite mixture of distributions. Such models are known as stochastic block models [2, 3].

Scalable methods to estimate stochastic block models are variational methods, which give rise to approximate maximum likelihood estimates [4, 5, 6, 7]. An alternative is based on spectral clustering [8, 9, 10, 11], but spectral clustering neither allows users to specify the features of the data on the basis of which nodes are to be clustered nor allows users to incorporate available predictors. In contrast, variational methods can be used to estimate a wide range of models capturing interesting features of real-world networks and incorporating available predictors. The consistency and asymptotic normality of variational estimators was established by [12, 13]. An important, remaining challenge stems from the fact that the high-dimensional posterior and its variational approximations tend to be non-concave and multimodal and algorithms tend to be trapped at local maxima.

We improve the predictive power of model-based clustering by introducing a novel model which incorporates available predictors and address the non-concave and multimodal nature of the high-dimensional posterior and its variational approximations. In most applications, additional information in the form of covariates is available, but—despite the fact that covariates may improve the predictive power of models—it is common practice to discard them. A notable exception is [14] and related work on latent variable models by [15, 16, 17], but the work of [14] and others relies on Bayesian Markov chain Monte Carlo methods which are not scalable and therefore have been

limited to small networks with up to 200 nodes. We introduce a novel model which incorporates available covariates and allows covariate parameters to depend on blocks, thus allowing a covariate to be active in some blocks while being inactive in others and varying in strength across the blocks in which it is active. We address the non-concave and multimodal nature of the high-dimensional posterior and its variational approximations by using the MM principle [18, 19] to derive a Bayesian Variational GEM-MM algorithm. We demonstrate the success of the model and method by using two data sets with more than 10,000 nodes and 15,000 block-dependent covariate terms.

The paper is structured as follows. We introduce a novel model with predictors in Section 2, discuss a penalized Bayesian approach in Section 3, and introduce a novel Bayesian Variational GEM-MM in Section 4. Experiments are presented in Section 5.

2 Data-generating process

We extend stochastic block models to include predictors and model both the sequence of degrees and block-dependent predictors.

We assume that the set of nodes is partitioned into K subsets, called blocks, and that the block memberships are unobserved. We assume that node-dependent categorical covariates are available and denote by $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ the set of covariate vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ of nodes $1, \dots, n$. We denote by $Y_{i,j}$ undirected edge variables, where $Y_{i,j}$ means that there is an undirected edge between nodes i and j and $Y_{i,j} = 0$ otherwise. Since edge variables are undirected and self-edges are meaningless, we assume that $Y_{i,j} = Y_{j,i}$ and $Y_{i,i} = 0$ with probability 1.

The data-generating process can be described as follows.

1. Generate mixing probabilities:

$$\alpha_1, \dots, \alpha_K \mid \omega_1, \dots, \omega_K \sim \text{Dirichlet}(\omega_1, \dots, \omega_K). \quad (1)$$

2. Generate block membership indicators:

$$Z_{i1}, \dots, Z_{iK} \mid \alpha_1, \dots, \alpha_K \stackrel{\text{iid}}{\sim} \text{Multinomial}(1; \alpha_1, \dots, \alpha_K). \quad (2)$$

3. Generate edge variables:

$$Y_{ij} \mid Z_{ik} = 1, Z_{jl} = 1, \mathbf{x} \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\pi_{y_{ij};kl}(\boldsymbol{\delta}, \boldsymbol{\beta}, \mathbf{x})), \quad (3)$$

where

$$\log \left[\frac{\pi_{y_{ij};kl}(\boldsymbol{\delta}, \boldsymbol{\beta}, \mathbf{x})}{1 - \pi_{y_{ij};kl}(\boldsymbol{\delta}, \boldsymbol{\beta}, \mathbf{x})} \right] = \delta_{kl} + \boldsymbol{\beta}_{kl}^\top \mathbf{f}(x_i, x_j),$$

where δ_{kl} are block-dependent degree parameters, $\boldsymbol{\beta}_{kl}$ are block-dependent covariate parameters, and \mathbf{f} are functions of covariates.

The model, in its full generality, is not an attractive model of large random graphs for at least two reasons.

First, estimating and interpreting $O(K^2)$ block-dependent degree parameters δ_{kl} is challenging when K is large. It is desirable to impose constraints to facilitate estimation and interpretation [7]. We consider here constraints of the form

$$\delta_{kl} = \delta_k + \delta_l, \quad (4)$$

where δ_k and δ_l are the degree parameters of blocks k and l , respectively, implying that there are $O(K)$ rather than $O(K^2)$ degree parameters.

Second, and worse, the number of block-dependent covariate parameters $\boldsymbol{\beta}_{kl}$ is $O(p K^2)$, where p is the dimension of the covariate vectors. To deal with the large number of block-dependent covariate parameters, we encourage sparsity by using a penalized Bayesian approach.

3 Penalized Bayesian approach

A Bayesian approach has advantages over non-Bayesian approaches, being able to deal with parameter uncertainty and average with respect to nuisance parameters, such as the mixing proportions $\alpha_1, \dots, \alpha_K$.

We discuss here a penalized Bayesian approach which encourages sparsity. In Bayesian fashion, the penalty is imposed through the prior. We encourage most covariate-related parameters to be 0 by using a spike-and-slab prior of the following form:

- Block-dependent degree parameters: $\delta_k \mid \mu_d, \sigma_d^2 \stackrel{\text{iid}}{\sim} N(\mu_d, \sigma_d^2)$.
- Block-dependent covariate parameters: $\beta_{kl}^s \mid \tau^2, I_{kl}^s \stackrel{\text{iid}}{\sim} [N(0, \tau^2)]^{1-I_{kl}^s} + [N(0, g\tau^2)]^{I_{kl}^s}$, where $g \gg 1$ is given.
- Block-dependent activeness indicators: $I_{kl}^s \mid \gamma_{kl} \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\gamma_{kl})$.
- Block-dependent activeness expectation: $\gamma_{kl} \mid a_\gamma, b_\gamma \stackrel{\text{iid}}{\sim} \text{Beta}(a_\gamma, b_\gamma)$.
- Variance parameter: $\tau^2 \sim \text{Inverse-Gamma}(a_{\tau^2}, b_{\tau^2}^{-1})$.

That is, the covariate parameters are governed by a two-component mixture distribution, where both component distributions are Gaussians centered at 0, one having a small variance τ^2 (spike at 0) and the other having a large variance $g\tau^2$ with $g \gg 1$ (slab). The indicator I_{kl}^s is 1 if covariate s is active in pair of blocks $\{k, l\}$ and 0 otherwise. The idea is that the posterior mass of inactive covariate parameters is concentrated in a small neighborhood of 0. The smaller γ_{kl} , the more covariates parameters are expected to be inactive. As a result, a covariate can be active in some blocks while being inactive in others and varying in strength across the blocks in which it is active.

Last, since there is uncertainty about γ_{kl} , we place a prior on γ_{kl} . Likewise, we place a prior on the variance parameter τ^2 .

4 Bayesian Variational GEM-MM algorithm

We explore the posterior by variational methods. While variational methods are scalable methods, implementing variational methods in high-dimensional settings can be challenging, as discussed in Sections 4.1 and 4.2. We address these problems by exploiting the MM principle [18, 19] to derive a Bayesian Variational GEM-MM algorithm which implements both the E-step and M-step of a Variational EM algorithm by MM algorithms.

We denote the parameters by $\theta = (\alpha, \mathbf{Z}, \delta, \tau^2, \gamma, \mathbf{I}, \beta)$. To keep the notation manageable, we denote the probability mass function of \mathbf{y} by $p(\mathbf{y} \mid \theta)$ and suppress the notational dependence on constants, including covariates and hyperparameters. Bayesian inference is based on the posterior density of θ of the form

$$p(\theta \mid \mathbf{y}) = \frac{p(\mathbf{y} \mid \theta) p(\theta)}{p(\mathbf{y})}, \quad (5)$$

where $p(\theta)$ is the prior density of θ and $p(\mathbf{y})$ is the marginal likelihood of \mathbf{y} . The posterior $p(\theta \mid \mathbf{y})$ is intractable, because the marginal likelihood $p(\mathbf{y})$ is intractable. A tractable lower bound on the intractable marginal likelihood $p(\mathbf{y})$ can be derived by using Jensen's inequality, giving

$$\log p(\mathbf{y}) = \log \int \frac{p(\mathbf{y} \mid \theta) p(\theta)}{q(\theta)} q(\theta) d\theta \geq \int \left[\log \frac{p(\mathbf{y} \mid \theta) p(\theta)}{q(\theta)} \right] q(\theta) d\theta, \quad (6)$$

where Q is an auxiliary distribution with density q with the same support as the posterior. We denote the lower bound LB_1 in (6) by

$$LB_1(q) = \mathbb{E}_Q [\log p(\mathbf{y} \mid \theta) p(\theta)] - \mathbb{E}_Q [\log q(\theta)] \quad (7)$$

and assume that the auxiliary distribution belongs to a family of fully factorized distributions of the form

$$q(\boldsymbol{\theta} \mid \tilde{\boldsymbol{\theta}}) = q(\boldsymbol{\alpha} \mid \tilde{\boldsymbol{\alpha}}) \prod_{i=1}^n q(\mathbf{Z}_i \mid \tilde{\mathbf{Z}}_i) \prod_{k=1}^K q(\delta_k \mid \tilde{\mu}_{\delta_k}, \tilde{\sigma}_{\delta_k}^2) q((\tau^2)^{-1} \mid \tilde{a}_{\tau^2}, \tilde{b}_{\tau^2}) \\ \prod_{k=1}^K \prod_{l=1}^k q(\gamma_{kl} \mid \tilde{a}_{\gamma_{kl}}, \tilde{b}_{\gamma_{kl}}) \prod_{k=1}^K \prod_{l=1}^k \prod_{s=1}^p q(I_{kl}^s \mid \tilde{I}_{kl}^s) \prod_{k=1}^K \prod_{l=1}^k \prod_{s=1}^p q(\beta_{kl}^s \mid \tilde{\mu}_{\beta_{kl}^s}, \tilde{\sigma}_{\beta_{kl}^s}^2),$$

where $\tilde{\boldsymbol{\theta}} = (\tilde{\boldsymbol{\alpha}}, \tilde{\mathbf{Z}}, \tilde{\boldsymbol{\mu}}_{\delta}, \tilde{\boldsymbol{\sigma}}_{\delta}^2, \tilde{\boldsymbol{a}}_{\gamma}, \tilde{\boldsymbol{b}}_{\gamma}, \tilde{a}_{\tau^2}, \tilde{b}_{\tau^2}, \tilde{\mathbf{I}}, \tilde{\boldsymbol{\mu}}_{\beta}, \tilde{\boldsymbol{\sigma}}_{\beta}^2)$ denotes the parameters of the auxiliary distribution.

To make the lower bound as tight as possible—which is equivalent to minimizing the Kullback-Leibler divergence from $q(\boldsymbol{\theta} \mid \tilde{\boldsymbol{\theta}})$ to $p(\boldsymbol{\theta} \mid \mathbf{y})$ —we maximize the lower bound with respect to $\tilde{\boldsymbol{\theta}}$.

To maximize the lower bound with respect to $\tilde{\boldsymbol{\theta}}$, we need closed-form expressions of the expectations in (7) under Q . However, the expectation

$$\mathbb{E}_Q [\log \pi_{y_{ij};kl}(\boldsymbol{\delta}, \boldsymbol{\beta}, \mathbf{x})] = \left[\tilde{\mu}_{\delta_k} + \tilde{\mu}_{\delta_l} + \sum_{s=1}^p \tilde{\mu}_{\beta_{kl}^s} f_s(x_i, x_j) \right] y_{ij} \\ + \mathbb{E}_Q [-\log (1 + \exp[\delta_k + \delta_l + \boldsymbol{\beta}_{kl}^\top \mathbf{f}(x_i, x_j)])]$$

is intractable, because the second expectation on the right-hand side is intractable. We replace the intractable expectation by a tractable lower bound obtained by Jensen's inequality, which implies that $\mathbb{E}_Q [-\log(\cdot)] \geq -\log[\mathbb{E}_Q(\cdot)]$ and therefore that $\mathbb{E}_Q [-\log(\cdot)]$ is bounded below by $-\log[\Phi_{ij;kl}(\cdot)]$, where

$$\Phi_{ij;kl}(\mathbf{x}) = 1 + \exp \left[\tilde{\mu}_{\delta_k} + \frac{1}{2} \tilde{\sigma}_{\delta_k}^2 + \tilde{\mu}_{\delta_l} + \frac{1}{2} \tilde{\sigma}_{\delta_l}^2 + \sum_{s=1}^p \tilde{\mu}_{\beta_{kl}^s} f_s(x_i, x_j) + \frac{1}{2} \sum_{s=1}^p \tilde{\sigma}_{\beta_{kl}^s}^2 f_s^2(x_i, x_j) \right].$$

Maximizing the resulting lower bound, denoted by LB_2 , is feasible, but the maximization problem is high-dimensional. We facilitate the high-dimensional maximization problem by exploiting the MM principle. The MM principle reduces the high-dimensional maximization problem to low-dimensional maximization problems which are straightforward to solve. The resulting Bayesian Variational GEM-MM algorithm increases the lower bound LB_2 at each iteration. We sketch the Bayesian Variational GEM-MM algorithm in Table 1 and provide details in Sections 4.1 and 4.2.

4.1 E-step: MM algorithm

The E-step maximizes LB_2 with respect to $\tilde{\mathbf{Z}}_1, \dots, \tilde{\mathbf{Z}}_n$, where

$$LB_2(q) = \sum_{i=1}^n \sum_{k=1}^K \tilde{Z}_{ik} \left[\Psi(\tilde{\alpha}_k) - \Psi \left(\sum_{l=1}^K \tilde{\alpha}_l \right) - \log \tilde{Z}_{ik} \right] + \sum_{i < j}^N \sum_{k=1}^K \sum_{l=1}^K \tilde{Z}_{ik} \tilde{Z}_{jl} \Upsilon_{y_{ij};kl}(\mathbf{x}) + \text{const},$$

where $\Upsilon_{y_{ij};kl}(\cdot)$ is the lower bound obtained of $\mathbb{E}_Q [\log \pi_{y_{ij};kl}(\cdot)]$ discussed above.

A simple approach to maximizing LB_2 with respect to $\tilde{\mathbf{Z}}_i$ is based on fixed-point (FP) updates of the form

$$\log \tilde{Z}_{ik}^{(t+1)} \propto \Psi \left(\tilde{\alpha}_k^{(t)} \right) - \Psi \left(\sum_{l=1}^K \tilde{\alpha}_l^{(t)} \right) + \sum_{j \neq i}^N \sum_{l=1}^K \tilde{Z}_{jl}^{(t)} \Upsilon_{y_{ij};kl}^{(t)}(\mathbf{x}),$$

where t denotes the iteration number. However, FP updates can be slow and are prone to be trapped at local maxima. Therefore, we leverage the MM principle [18, 19], which is less prone to be trapped at local maxima [7].

The MM principle exploits properties of the objective function to derive a simple-to-maximize minorizing function and then maximizes the minorizing function. Here, the objective function is the lower bound LB_2 , which we want to maximize as a function of $\tilde{\mathbf{Z}}_1, \dots, \tilde{\mathbf{Z}}_n$. By exploiting the log

Table 1: Sketch of Bayesian Variational GEM-MM algorithm

Minimize the Kullback-Leibler divergence from $q(\boldsymbol{\theta} \mid \tilde{\boldsymbol{\theta}})$ to $p(\boldsymbol{\theta} \mid \mathbf{y})$ by cycling through the components of $\tilde{\boldsymbol{\theta}}$ as follows:

E-step:

E.1 Update $\tilde{\mathbf{Z}}_1, \dots, \tilde{\mathbf{Z}}_n$ using the E-step MM algorithm described in Section 4.1.

E.2 Update $\tilde{\alpha}_1, \dots, \tilde{\alpha}_K$:

$$\tilde{\alpha}_k = \omega_k + \sum_{i=1}^n \tilde{Z}_{ik}.$$

M-step:

M.1 Update $\tilde{\boldsymbol{\mu}}_\delta, \tilde{\sigma}_\delta^2, \tilde{\boldsymbol{\mu}}_\beta, \tilde{\sigma}_\beta^2$ using the M-step MM algorithm described in Section 4.2.

M.2 Update $\tilde{a}_{\tau^2}, \tilde{b}_{\tau^2}$:

$$\begin{aligned} \tilde{a}_{\tau^2} &= a_{\tau^2} + \frac{p K (K+1)}{4} \\ \tilde{b}_{\tau^2} &= b_{\tau^2} + \frac{1}{2} \sum_{k=1}^K \sum_{l=1}^k \sum_{s=1}^p (1 - \tilde{I}_{kl}^s) (\tilde{\mu}_{\beta_{kl}^s}^2 + \tilde{\sigma}_{\beta_{kl}^s}^2) + \frac{1}{2g} \sum_{k=1}^K \sum_{l=1}^k \sum_{s=1}^p \tilde{I}_{kl}^s (\tilde{\mu}_{\beta_{kl}^s}^2 + \tilde{\sigma}_{\beta_{kl}^s}^2). \end{aligned}$$

M.3 Update $\tilde{\mathbf{a}}_\gamma, \tilde{\mathbf{b}}_\gamma$:

$$\begin{aligned} \tilde{a}_{\gamma_{kl}} &= a_\gamma + \sum_{s=1}^p \tilde{I}_{kl}^s \\ \tilde{b}_{\gamma_{kl}} &= b_\gamma + \sum_{s=1}^p (1 - \tilde{I}_{kl}^s). \end{aligned}$$

M.4 Update $\tilde{\mathbf{I}}$:

$$\tilde{I}_{kl}^s = \left\{ 1 + \exp \left[\psi(\tilde{b}_{\gamma_{kl}}) - \psi(\tilde{a}_{\gamma_{kl}}) - \frac{1}{2} \frac{\tilde{a}_{\tau^2}}{\tilde{b}_{\tau^2}} (\tilde{\mu}_{\beta_{kl}^s}^2 + \tilde{\sigma}_{\beta_{kl}^s}^2) \left(1 - \frac{1}{g} \right) + \frac{1}{2} \log g \right] \right\}^{-1}.$$

concavity and arithmetic-geometric mean inequality [19], we construct a minorizing function for LB_2 of the form:

$$\begin{aligned} M_1(\tilde{\mathbf{Z}}_1, \dots, \tilde{\mathbf{Z}}_n) &= \sum_{i=1}^n \left[\sum_{k=1}^K \tilde{Z}_{ik} \left(-\frac{1}{\tilde{Z}_{ik}^{(t)}} + \sum_{j \neq i} \sum_{l=1}^K \frac{\tilde{Z}_{jl}^{(t)}}{2 \tilde{Z}_{ik}^{(t)}} \Upsilon_{y_{ij};kl}^{(t)}(\cdot) \right) \right. \\ &\quad \left. + \sum_{k=1}^K \tilde{Z}_{ik} \left(\Psi(\tilde{\alpha}_k^{(t)}) - \Psi(\sum_{k'} \tilde{\alpha}_{k'}^{(t)}) - \log \tilde{Z}_{ik}^{(t)} + 1 \right) \right]. \end{aligned}$$

As a result, the high-dimensional maximization problem is replaced by K -dimensional quadratic programming problems which are straightforward to maximize [20].

4.2 M-step: MM algorithm

The M-step maximizes the lower bound LB_2 with respect to the parameters $\tilde{\boldsymbol{\mu}}_\delta, \tilde{\sigma}_\delta^2, \tilde{\boldsymbol{\mu}}_\beta$, and $\tilde{\sigma}_\beta^2$, where

$$\begin{aligned} LB_2(q) &= \sum_{k=1}^K \left(\frac{1}{2} \log \tilde{\sigma}_{\delta_k}^2 - \frac{\tilde{\mu}_{\delta_k}^2 + \tilde{\sigma}_{\delta_k}^2 - 2\tilde{\mu}_{\delta_k} \mu_d}{2\sigma_d^2} \right) + \sum_{k \leq l} \sum_{s=1}^p \frac{1}{2} \log \tilde{\sigma}_{\beta_{kl}^s}^2 \\ &\quad + \sum_{k \leq l} \sum_{s=1}^p \frac{1}{2} \frac{\tilde{a}_{\tau^2}}{\tilde{b}_{\tau^2}} (\tilde{\mu}_{\beta_{kl}^s}^2 + \tilde{\sigma}_{\beta_{kl}^s}^2) \left[\tilde{I}_{kl}^s \left(1 - \frac{1}{g} \right) - 1 \right] + \sum_{k=1}^K \tilde{\mu}_{\delta_k} \sum_{i=1}^n \tilde{Z}_{ik} d_i(\mathbf{y}) \\ &\quad + \sum_{1 \leq i < j \leq N} \sum_{k=1}^K \sum_{l=1}^K \tilde{Z}_{ik} \tilde{Z}_{jl} \left(\sum_{s=1}^p \tilde{\mu}_{\beta_{kl}^s} y_{ij} f_s(x_i, x_j) - \log \Phi_{ij;kl}(\mathbf{x}) \right) + \text{const}, \end{aligned}$$

where $d_1(\mathbf{y}), \dots, d_n(\mathbf{y})$ denote the degrees of nodes $1, \dots, n$, respectively.

To implement the M-step, we exploit the MM principle, because the large number of parameters makes gradient ascent and Newton's method problematic: The matrix inversion required by Newton's methods is expensive, while coordinate ascent methods may neither be stable nor straightforward to implement [19]. In contrast, MM updates do not require matrix inversion, because all parameters are separated, and are robust and straightforward to implement. A second reason is that increasing rather than maximizing objective functions can speed up convergence [21, 19],

By using log concavity and exponential convexity inequality [19], we construct a minorizing function for LB_2 of the form:

$$\begin{aligned} M_2(\tilde{\mu}_\delta, \tilde{\sigma}_\delta^2, \tilde{\mu}_\beta, \tilde{\sigma}_\beta^2) &= \sum_{k=1}^K M_{2, \tilde{\mu}_\delta}(\tilde{\mu}_{\delta_k}) + \sum_{k=1}^K M_{2, \tilde{\sigma}_\delta^2}(\tilde{\sigma}_{\delta_k}^2) \\ &+ \sum_{k \leq l}^K \sum_{s=1}^p M_{2, \tilde{\mu}_\beta}(\tilde{\mu}_{\beta_{kl}^s}) + \sum_{k \leq l}^K \sum_{s=1}^p M_{2, \tilde{\sigma}_\beta^2}(\tilde{\sigma}_{\beta_{kl}^s}^2) + \text{const}, \end{aligned}$$

where

$$\begin{aligned} Q_{\tilde{\mu}_\delta}(\tilde{\mu}_{\delta_k}) &= -\frac{\tilde{\mu}_{\delta_k}^2}{2\sigma_d^2} + \tilde{\mu}_{\delta_k} \left(\frac{\mu_d}{\sigma_d^2} + \sum_{i=1}^n \tilde{Z}_{ik} d_i(\mathbf{y}) \right) \\ &- \sum_{i=1}^n \sum_{j \neq i} \sum_{l=1}^K \tilde{Z}_{ik} \tilde{Z}_{jl} \frac{\Phi_{ij;kl}^{(t)}(\mathbf{x}) - 1}{(2p+4) \Phi_{ij;kl}^{(t)}(\mathbf{x})} \exp[(2p+4)(\tilde{\mu}_{\delta_k} - \tilde{\mu}_{\delta_k}^{(t)})] \\ Q_{\tilde{\sigma}_\delta^2}(\tilde{\sigma}_{\delta_k}^2) &= \frac{1}{2} \log \tilde{\sigma}_{\delta_k}^2 - \frac{\tilde{\sigma}_{\delta_k}^2}{2\sigma_d^2} \\ &- \sum_{i=1}^n \sum_{j \neq i} \sum_{l=1}^K \tilde{Z}_{ik} \tilde{Z}_{jl} \frac{\Phi_{ij;kl}^{(t)}(\mathbf{x}) - 1}{(2p+4) \Phi_{ij;kl}^{(t)}(\mathbf{x})} \exp[(2p+4) \frac{1}{2} (\tilde{\sigma}_{\delta_k}^2 - \tilde{\sigma}_{\delta_k}^{(t)2})] \\ Q_{\tilde{\mu}_\beta}(\tilde{\mu}_{\beta_{kl}^s}) &= \frac{1}{2} \frac{\tilde{a}_{\tau^2}}{\tilde{b}_{\tau^2}} \left[\tilde{I}_{kl}^s \left(1 - \frac{1}{g} \right) - 1 \right] \tilde{\mu}_{\beta_{kl}^s} + \tilde{\mu}_{\beta_{kl}^s} \sum_{i=1}^n \sum_{j \neq i} \tilde{Z}_{ik} \tilde{Z}_{jl} y_{ij} f_s(x_i, x_j) \\ &- \sum_{i=1}^n \sum_{j \neq i} \tilde{Z}_{ik} \tilde{Z}_{jl} \frac{\Phi_{ij;kl}^{(t)}(\mathbf{x}) - 1}{(2p+4) \Phi_{ij;kl}^{(t)}(\mathbf{x})} \exp[(2p+4) f_s(x_i, x_j) (\tilde{\mu}_{\beta_{kl}^s} - \tilde{\mu}_{\beta_{kl}^s}^{(t)})] \\ Q_{\tilde{\sigma}_\beta^2}(\tilde{\sigma}_{\beta_{kl}^s}^2) &= \frac{1}{2} \log \tilde{\sigma}_{\beta_{kl}^s}^2 + \frac{1}{2} \frac{\tilde{a}_{\tau^2}}{\tilde{b}_{\tau^2}} \left[\tilde{I}_{kl}^s \left(1 - \frac{1}{g} \right) - 1 \right] \tilde{\sigma}_{\beta_{kl}^s}^2 \\ &- \sum_{i=1}^n \sum_{j \neq i} \tilde{Z}_{ik} \tilde{Z}_{jl} \frac{\Phi_{ij;kl}^{(t)}(\mathbf{x}) - 1}{(2p+4) \Phi_{ij;kl}^{(t)}(\mathbf{x})} \exp \left[(2p+4) \frac{1}{2} f_s^2(x_i, x_j) (\tilde{\sigma}_{\beta_{kl}^s}^2 - \tilde{\sigma}_{\beta_{kl}^s}^{(t)2}) \right]. \end{aligned}$$

These one-dimensional minorizing functions can be maximized by one-step Newton updates [19].

5 Experiments

We use two blog networks, summarized in Table 2.

The first blog network is the Political Blogs network of [22], where a directed edge from blogger i to blogger j indicates that blogger i cited blogger j . Since the blog network is directed, we convert it into an undirected network by specifying an undirected edge between two blogs i and j if there is an edge from i to j or j to i or both. We use two covariates: The first one indicates whether both blogs are liberal, while the second one indicates whether both blogs are conservative.

The second blog network is the Catalog Blogs 3 network of [23], where an undirected edge between two bloggers i and j indicates friendship. We use observed group memberships as covariates. We construct 39 indicators, one for each of the 39 groups: The indicator of group s is 1 if both blogs are members of group s otherwise 0. Note that although the number of original predictors is small, the number of block-dependent parameters is $O(K^2)$. In all of the experiments, we use the constants $\omega_1 = \dots = \omega_K = 1$, $\mu_d = 0$, $\sigma_d^2 = 1$, $a_\gamma = 1$, $b_\gamma = 1$, $a_{\tau^2} = 1$, $b_{\tau^2} = 1$, and $g = 10,000$.

Data	Nodes	Edge variables	Edges	Covariates	Block-dependent predictors	
					$K = 5$	$K = 20$
P. Blogs	1,490	1,109,305	16,783	2	70	880
C. Blogs	10,312	53,163,516	333,983	39	1,180	16,420

Table 2: Summary statistics of two blog networks.

To demonstrate the advantage of the Bayesian Variational GEM-MM algorithm with MM updates over the Bayesian Variational EM algorithm with FP updates, we compare FP and MM updates using both blog networks with $K = 5$ and $K = 20$ blocks. In both cases, we use MM algorithms in the M-step and averaged the results over 50 runs with random starting values.

Figure 1 compares lower bounds of two update methods based FP and MM. Overall, the MM method is far superior to the FP method on both small (i.e., Political Blogs) and large (i.e., Catalog Blogs) data sets under $K = 5$ as well as $K = 20$. We speculate that the superior nature of MM updates stems from the fact that the MM principle separates $\tilde{Z}_1, \dots, \tilde{Z}_n$. The separation in turn weakens the dependence of updates of $\tilde{Z}_1, \dots, \tilde{Z}_n$.

(a)(d)
 PPCC.
 Political Blogs
 Catalog Blogs
 K=5
 K=20
 52020

Figure 1: Lower bounds based on FP and MM implementation of E-step with $K = 5$ and $K = 20$.

We select K based on the lower bound on the logarithm of the marginal likelihood $p(\mathbf{y})$, on which Bayesian model selection is based. For both data sets, we vary K from 1 to 20 and select the K corresponding to the best lower bound on $\log p(\mathbf{y})$. For the Political Blogs, the best lower bound is achieved at $K = 12$ as shown in Figure 2(a). Figure 2(b) shows a histogram of parameter means $\tilde{\mu}_\beta$ of all predictors and demonstrates that predictors are relevant across blocks. Figures 2(c) and 2(d) show that the covariate terms vary in strength across blocks.

(a)(d)
 LPPCC-
 Political Blogs
 Catalog Blogs
 K=12
 12

Figure 2: Political Blogs: (a) Lower bounds under $K = 1, \dots, 20$. The best lower bound is achieved at $K = 12$. (b) Histogram of parameter means $\tilde{\mu}_\beta$ of predictors. (c) and (d) Predictors 1 and 2 are active in most blocks.

In the case of the Catalog Blogs data set, the best lower bound is achieved at $K = 11$ (Figure 3(a)). In contrast to the Polical Blogs data set, there is evidence of sparsity (Figure 3(b)): The high peak at 0 indicates that a substantial proportion of covariate terms is inactive in some blocks. Figure 3(d) shows parameter estimates of predictor 12 across all blocks. Almost 50% of them are close to 0 which demonstrates that predictor 12 is inactive in most blocks. Observe that a predictor can act have both positive and negative effects on edge probabilities depending on the blocks involved, as shown in Figure 3(c).

(a) (b) (c) (d)
 Lower-
 bounds
 Histogram
 of
 parameter
 means
 of
 predictors
 212

Figure 3: Catalog Blogs: (a) Lower bounds under $K = 1, \dots, 20$. The best lower bound is achieved at $K = 11$. (b) Histogram of parameter means $\tilde{\mu}_\beta$ of predictors. The peak at 0 indicates that many block-dependent covariate parameters are inactive. (c) and (d) Predictor 2 is active in most blocks while predictor 12 is inactive in most.

To compare the predictive power of models, we construct a test data set by sampling 10% zero and 10% non-zero edges at random, then run models on the training data set where all test edges are set to zero plus other edges that are not in the test data set. We compare three models: the stochastic block model without predictors, the stochastic block model with block-invariant predictors, i.e., $\beta_{kl} = \beta$, and the stochastic block model with block-dependent predictors. The cross-validation procedure is repeated 10 times to obtain standard errors. We set $K = 10$ motivated by the results above. We use the predictive probabilities of all zero and non-zero test edges to compare three methods:

$$p(y_{ij} | \mathbf{y}) = \sum_{k=1}^K \sum_{l=1}^K p(y_{ij} | Z_{ik} = 1, Z_{jl} = 1, \tilde{\mu}_\delta, \tilde{\mu}_\beta) p(Z_{ik} = 1 | \tilde{\mathbf{Z}}_i) p(Z_{jl} = 1 | \tilde{\mathbf{Z}}_j). \quad (8)$$

Table 3 shows predictive probabilities of three models averaged over 10 training and test data sets. It is evident that incorporating predictors increases predictive power and that block-dependent predictors are best in terms of predictive power.

Data	without predictors	block-invariant predictors	block-dependent predictors
Political Blogs	-5,693 (5)	-4,925 (5)	-4,788 (6)
Catalog Blogs	-122,509 (42)	-121,321 (41)	-120,906 (41)

Table 3: Predictive probabilities of three models on test data sets. Mean test log likelihoods over 10 runs with random initializations are reported together with standard errors in parentheses.

6 Discussion

We believe that the MM principle is a promising approach to deal with the large number of local maxima encountered in applications of variational methods to statistical learning. We note that the convergence of variational GEM-MM algorithms can be sped up by (a) exploiting the separation of the parameters of the maximization problem and implementing parallel computing and by (b) using quasi-Newton methods.

References

- [1] P. J. Bickel and A. Chen. A nonparametric view of network models and Newman-Girvan and other modularities. In *PNAS*, volume 106, pages 21068–21073, 2009.
- [2] T. A. B. Snijders and K. Nowicki. Estimation and prediction for stochastic blockmodels for graphs with latent block structure. *Journal of Classification*, 14:75–100, 1997.
- [3] K. Nowicki and T. A. B. Snijders. Estimation and prediction for stochastic blockstructures. *Journal of the American Statistical Association*, 96(455):1077–1087, 2001.
- [4] J. J. Daudin, F. Picard, and S. Robin. A mixture model for random graphs. *Statistics and Computing*, 18(2):173–183, 2008.
- [5] H. Zanghi, F. Picard, V. Miele, and C. Ambroise. Strategies for online inference of model-based clustering in large and growing networks. *Annals of Applied Statistics*, 4:687–714, 2010.
- [6] P. Gopalan, D. Mimno, S. M. Gerrish, M. J. Freedman, and D. M. Blei. Scalable inference of overlapping communities. In *Neural Information Processing Systems*, 2012.
- [7] D. Q. Vu, D. R. Hunter, and M. Schweinberger. Model-based clustering of large networks. *Annals of Applied Statistics*, 7:1010–1039, 2013.
- [8] K. Rohe, S. Chatterjee, and B. Yu. Spectral clustering and the high-dimensional stochastic block model. *Annals of Statistics*, 39:1878–1915, 2011.
- [9] C. E. Priebe, D. L. Sussman, M. Tang, and J. T. Vogelstein. Statistical inference on errofully observed graphs. *Journal of the American Statistical Association*, 107:1119–1128, 2012.
- [10] P. Sarkar and P. Bickel. The role of normalization in spectral clustering of stochastic block-models. Technical report, 2013. arXiv preprint arxiv:1310.0532.
- [11] J. Lei and A. Rinaldo. Consistency of spectral clustering in sparse stochastic block models. Technical report, 2013. arXiv preprint arxiv:1312.2050.
- [12] A. Celisse, J. J. Daudin, and L. Pierre. Consistency of maximum-likelihood and variational estimators in the stochastic block model. *Electronic Journal of Statistics*, 6:1847–1899, 2012.
- [13] P. J. Bickel, D. Choi, X. Chang, and H. Zhang. Asymptotic normality of maximum likelihood and its variational approximation for stochastic blockmodels. *Annals of Statistics*, 41:1922–1943, 2013.
- [14] C. Tallberg. A Bayesian approach to modeling stochastic blockstructures with covariates. *Journal of Mathematical Sociology*, 29:1–23, 2005.
- [15] P. D. Hoff. Random effects models for network data. In *Dynamic Social Network Modeling and Analysis*, pages 303–312. National Academies Press, 2003.
- [16] M. A. J. van Duijn, T. A. B. Snijders, and B. J. H. Zijlstra. P2: a random effects model with covariates for directed graphs. *Statistica Neerlandica*, 58:234–254, 2004.
- [17] P. N. Krivitsky, M. S. Handcock, A. E. Raftery, and P. D. Hoff. Representing degree distributions, clustering, and homophily in social networks with latent cluster random effects models. *Social Networks*, 31:204–213, 2009.
- [18] D. R. Hunter and K. Lange. A tutorial on MM algorithms. *The American Statistician*, 58(1):30–38, 2004.
- [19] K. Lange. *Numerical Analysis for Statisticians*. Springer, New York, 2010.
- [20] S. M. Stefanov. Convex quadratic minimization subject to a linear constraint and box constraints. *Appl Math Res Express*, 2004(1):17–42, 2004.
- [21] R. M. Neal and G. E. Hinton. A new view of the EM algorithm that justifies incremental and other variants. In *Learning in Graphical Models*, pages 355–368, 1993.
- [22] Lada A. Adamic and Natalie Glance. The political blogosphere and the 2004 U.S. election: Divided they blog. In *Proceedings of LinkKDD 2005*, pages 36–43. ACM, 2005.
- [23] R. Zafarani and H. Liu. Social computing data repository at ASU, 2009.