HIERARCHICAL STOCHASTIC BLOCK MODEL FOR COMMUNITY DETECTION IN MULTIPLEX NETWORKS

MARINA S. PAEZ, ARASH A. AMINI, AND LIZHEN LIN

Abstract. Multiplex networks have become increasingly more prevalent in many fields, and have emerged as a powerful tool for modeling the complexity of real networks. There is a critical need for developing inference models for multiplex networks that can take into account potential dependencies across different layers, particularly when the aim is community detection. We add to a limited literature by proposing a novel and efficient Bayesian model for community detection in multiplex networks. A key feature of our approach is the ability to model varying communities at different network layers. In contrast, many existing models assume the same communities for all layers. Moreover, our model automatically picks up the necessary number of communities at each layer (as validated by real data examples). This is appealing, since deciding the number of communities is a challenging aspect of community detection, and especially so in the multiplex setting, if one allows the communities to change across layers. Borrowing ideas from hierarchical Bayesian modeling, we use a hierarchical Dirichlet prior to model community labels across layers, allowing dependency in their structure. Given the community labels, a stochastic block model (SBM) is assumed for each layer. We develop an efficient slice sampler for sampling the posterior distribution of the community labels as well as the link probabilities between communities. In doing so, we address some unique challenges posed by coupling the complex likelihood of SBM with the hierarchical nature of the prior on the labels. An extensive empirical validation is performed on simulated and real data, demonstrating the superior performance of the model over single-layer alternatives, as well as the ability to uncover interesting structures in real networks.

Keywords: Community detection; Hierarchical Stochastic block model (HSBM); Multiplex networks; Hierarchical Dirichlet Process; Random partition.

1. Introduction

Networks, which are used to model interactions among a set of entities, have emerged as one of the most powerful tools for modern data analysis. The last few decades have witnessed explosions in the development of models, theory, and algorithms for network analysis. Modern network data are often complex and heterogeneous. To model such heterogeneity, *multiplex networks*, which also go by various other qualifiers (multiple-layer, multiple-slice, multi-array, multi-relational), have arisen as a useful representation of complex networks.

A multiplex network typically consists of a fixed set of nodes but multiple types of edges, often representing heterogeneous relationships as well as the dynamic nature of the edges. These networks have become increasingly prevalent in many applications. Typical examples include various types of dynamic networks [4, 23] including temporal networks, dynamic social networks [9] and dynamic biological networks such as protein

networks [37]. An example of a multiplex social network is the Twitter network where different edge information is available such as "mention", "follow" and "retweet" [17]. Another example is the co-authorship network where the authors are recorded on relationships such as "co-publishes", "convenes" and "cites" [18]. Many other examples can be found in economy [31], neuroscience [13], biology [14] and so on.

Due to the ubiquity of multiplex network data, there is a critical need for developing realistic statistical models that can handle their complex structures. There has already been growing efforts in extending static (or single-layered) statistical network models to the multiplex setting [26]. Many statistical metrics have already been extended from basic notions such as degree and node centrality [7, 2] to the clustering coefficient and modularity [10, 6]. Popular latent space models [32] have also been extended to dynamic [34] and multiplex [16, 33] networks. Based on such models, one can perform learning tasks such as link probability estimation or link prediction.

Another task that has been extensively studied in single-layer networks is community detection, that is, clustering nodes into groups or communities. There have already been a few works proposing models or algorithms for community detection in multiplex and multilayer networks [26, 3, 22, 38, 12, 5]. The general strategy adopted is to either transform a multiplex network into a static one and then apply the existing community detection algorithms, or extend a model from static networks to the multiplex setting. One of the key shortcomings of many existing methods is that communities across different layers are assumed to be the same, which is clearly restrictive and often unrealistic. Instead, there is interest in monitoring or exploring how the communities vary across different layers or evolve across different time points.

We propose a novel Bayesian community detection model, namely, the hierarchical stochastic block model (HSBM) for multiplex networks that is fundamentally different from the existing approaches. Specifically, we impose a random partition prior based on the hierarchical Dirichlet process [36] on communities (or partitions) across different layers. Given these communities, a stochastic block model is assumed for each layer. One of the appealing features of our model is allowing the communities to vary across different layers of the network while being able to incorporate potential dependency across them. The hierarchical aspect of HSBM allows for effective borrowing of information or strength across different layers for improved estimation. Our approach has the added advantage of being able to handle a broad class of networks by allowing the number of nodes to vary, and not necessarily imposing the nodes to be fixed across layers. In addition, HSBM inherits the desirable property of hierarchical Dirichlet priors that allow for automatic and adaptive selection of the number of communities (for each layer) based on the data.

Our simulation study in Section 4 confirms that HSBM significantly outperforms a model that assumes independence of layers, especially when it comes to matching community labels across layers. In particular, the superior performance of our model over its independent single-layer counterpart is manifested by the significant improvement in the slicewise or aggregate NMI measures.

Although a hierarchical Dirichlet process is a natural choice for modeling dependent partition structures, extending the ideas from simple mixture models to community structured network models is not that straightforward. In particular, leveraging ideas from one of our recent works [1], we develop a new and efficient slice sampler for exact

sampling of the posterior of HSBM for inference. We will discuss some of the technical difficulties in Remark 1.

The rest of the paper is organized as follow. Section 2 introduces the HSBM in details. Section 3 is devoted to describing a novel MCMC algorithm for inference of HSBM. Simulation study and real data analysis are presented in Sections 4 and 5. The code for all the experiments is available at [27]. We conclude with a discussion in Section 6.

2. Hierarchical Stochastic block model (HSBM)

Consider a multiplex network with T layers (or T types of edge relations) and n_t nodes with labels in $[n_t] = \{1, \ldots, n_t\}$ for each layer $t = 1, \ldots, T$. Denote by \mathbf{A}_t the adjacency matrix of the network at layer t, so that an observed multiplex network consists of the collection $\mathbf{A}_t \in \{0, 1\}^{n_t \times n_t}$ for $t = 1, \ldots, T$. We let \mathbf{A} denote this collection and view it as a (partial or irregular) adjacency tensor. That is, $\mathbf{A} = (A_{tij}, t \in [T], i, j \in [n_t])$ and $\mathbf{A}_t = (A_{tij}, i, j \in [n_t])$ where $A_{tij} = 1$ if nodes i and j in layer t are connected. Our goal is to estimate the clustering or community structure of the nodes in each layer, given \mathbf{A} .

Specifically, to each node i, at each layer t, we assign a community label, encoded in a variable $\mathbf{Z} = (z_{ti}) \in \mathbb{N}^{T \times n_t}$, where $\mathbb{N} = \{1, 2, \dots, \}$ is the set of natural numbers. Let

(2.1)
$$\boldsymbol{\eta}_t = (\eta_{txy})_{x,y} \in [0,1]^{\mathbb{N} \times \mathbb{N}}, \quad \eta_{txy} \equiv \eta_t(x,y),$$

be a matrix of link probabilities between communities, indexed by \mathbb{N}^2 , for t = 1, ..., T. At times, we will use the equivalent notation $\eta_t(x, y) \equiv \eta_{txy}$ to increase readability.

For example, we interpret $\eta_{t12} = \eta_t(1,2)$ as the probability of an edge being formed between a node from community 1 and a node from community 2 at layer t. Note that we have assumed that the total number of community labels is infinite. However, for a given adjacency matrix A_t observed at layer t, the number of community labels is finite, unknown, and will be denoted by K_t .

For each layer t, we model the distribution of the adjacency matrix $\mathbf{A}_t \in \{0, 1\}^{n_t \times n_t}$ as a stochastic block model (SBM) with membership vector $\mathbf{z}_t = (z_{ti}) \in \mathbb{N}^{n_t}$, and edge probability matrix $\boldsymbol{\eta}_t$, that is,

$$A_t \mid z_t, \eta_t, \sim \text{SBM}(z_t; \eta_t) \iff A_{tij} \stackrel{\text{iid}}{\sim} \text{Ber}(\eta_t(z_{ti}, z_{tj})), \quad 1 \leq i < j \leq n_t.$$

In an SBM, the link probability between nodes i and j is uniquely determined by which communities these nodes belong to. In our notation, at a layer t, the link probability between nodes i and j is given by $\eta_t(z_{ti}, z_{tj})$. Note that our SBM notation is slightly different from the traditional stochastic block model where the set of community labels is random. In writing SBM $(z; \eta)$, we assume that z is given and nonrandom.

If there is belief or prior information that the community structure in one layer of the network is independent of the others, independent stochastic block models (SBM) could be assumed for the A_t 's. This assumption is, however, too restrictive when we are dealing with networks of different kinds of relations, but among the same set of nodes. The other extreme is to assume that all layers in the network share the same partition—meaning that $z_t = z$ for all t and the A_t 's are conditionally independent given z.

We believe, however, that a model that can incorporate various dependencies between these two extremes is more appropriate in many applications. In other words, it may be desirable to allow some change in the partition structure between layers, but also impose some kind of dependence among them. Here, we propose a model that achieves this goal by allowing the community structures at various layers to be different but dependent, using a hierarchical specification for the distribution of the partitions.

2.1. Hierarchical community prior. Before stating the prior on the labels, let us introduce a simplification, namely, that η_t does not vary by layer. Therefore, we drop index t, and denote the matrix of link probabilities by η and its elements by $\eta_{xy} \equiv \eta(x,y)$. This assumption is not necessarily restrictive since, as will become clear shortly, our model allows for an infinite number of communities. By imposing this restriction, we are simply stating that cluster i at a certain layer corresponds to the same cluster i at every other layer.

Our key idea is to impose a dependent random partition prior on the membership labels at different layers, i.e., $z_t, t \in [T]$. We adopt the prior based on a hierarchical Dirichlet process (HDP). More precisely, we have the following specification for the overall model:

(2.2)
$$\pi \mid \gamma_0 \sim \text{GEM}(\gamma_0)$$

(2.3)
$$\boldsymbol{w}_t \mid \alpha_0, \boldsymbol{\pi} \sim \mathrm{DP}(\alpha_0, \boldsymbol{\pi})$$

$$(2.4) z_{ti} \mid \boldsymbol{w}_t \sim \boldsymbol{w}_t$$

(2.5)
$$\eta = (\eta_{xy}) \stackrel{\text{iid}}{\sim} \text{Beta}(\alpha_{\eta}, \beta_{\eta})$$

$$(2.6) A_t \mid \boldsymbol{z}_t, \boldsymbol{\eta} \sim \mathrm{SBM}(\boldsymbol{z}_t, \boldsymbol{\eta})$$

for $i=1,\ldots,n_t$ and $t=1,\ldots,T$, where the draws across i and t are independent. GEM stands for Griffiths, Engen and McCloskey [30]; it is a distribution for a random measure on \mathbb{N} which has the well-known stick-breaking construction [35]. DP stands for the Dirichlet process [15] and Beta for the usual Beta distribution. Both π and \mathbf{w}_t are random measures supported on \mathbb{N} and all \mathbf{w}_t s are controlled by same π which acts as the base measure of the Dirichlet process generating them. Note that we have set

$$\mathbf{z}_t = (z_{ti}, i \in [n_t]), \quad \mathbf{A}_t = (A_{tij}) \in \{0, 1\}^{n_t \times n_t}.$$

The last line, (2.6), can be written explicitly as

$$A_{tij} \mid \boldsymbol{z}_t, \boldsymbol{\eta} \sim \operatorname{Ber}(\eta_{z_{ti}, z_{tj}}), \quad 1 \leq i < j \leq n_t,$$

Figure 1 provides a graphical illustration of the hierarchical structure of the model.

The first three equations, (2.2)–(2.4), in the model above are basically similar to the ones in equation (19) of the original HDP paper [36]. These three equations constitute the label part of the model which can be equivalently described as

(2.7)
$$\boldsymbol{\pi} \mid \gamma_0 \sim \operatorname{GEM}(\gamma_0), \quad \boldsymbol{\pi} = (\pi_k)$$

$$\boldsymbol{\gamma}_t \mid \alpha_0 \sim \operatorname{GEM}(\alpha_0), \quad \boldsymbol{\gamma}_t = (\gamma_{tg})$$

$$k_{tg} \mid \boldsymbol{\pi} \sim \boldsymbol{\pi}, \quad g \in \mathbb{N},$$

$$g_{ti} \mid \boldsymbol{\gamma}_t \sim \boldsymbol{\gamma}_t, \quad i = 1, \dots, n_t$$

$$z_{ti} \mid \boldsymbol{g}_t, \boldsymbol{k}_t = k_{t,g_{ti}}.$$

We assume that the reader is familiar with the HDP and its various interpretations, and in particular, the Chinese Restaurant Franchise process (CRF); see for example [36]. In

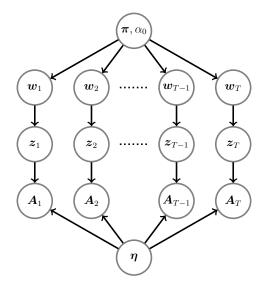


FIGURE 1. Hierarchical Stochastic Block Model

the CRF interpretation, each layer t of the network corresponds to a restaurant, and nodes are gathered around tables, or groups, in each restaurant. The variable g_{ti} denotes the group of node i in restaurant (i.e., layer) t. All the nodes in the same group g, at the same layer t, share the same dish (i.e., community) which is denoted by k_{tg} .

Thus, given all the groups $\mathbf{g}_t = (g_{ti})_i$ and group-communities $\mathbf{k}_t = (k_{tg})_g$, the label of node i (at layer t) is uniquely determined, as in the last line of (2.7). For more details on the equivalence of (2.7) and the model described by (2.2)–(2.4), as well as the interpretation of latent variable γ_t , we refer to [1].

Using the equivalent representation (2.7), the joint density for the label part of the model can be expressed as:

(2.8)
$$p(\boldsymbol{g}, \boldsymbol{k}, \boldsymbol{\gamma}', \boldsymbol{\pi}') = \prod_{t=1}^{T} \left[p(\boldsymbol{g}_{t} | \boldsymbol{\gamma}_{t}) p(\boldsymbol{\gamma}'_{t}) p(\boldsymbol{k}_{t} | \boldsymbol{\pi}) \right] p(\boldsymbol{\pi}'),$$

where $p(\boldsymbol{g}_t|\boldsymbol{\gamma}_t) = \prod_{i=1}^{n_t} \gamma_{t,g_{ti}}$ and $p(\boldsymbol{k}_t|\boldsymbol{\pi}) = \prod_{g=1}^{\infty} \pi_{k_{tg}}$. The new variables $\boldsymbol{\gamma}_t'$ and $\boldsymbol{\pi}'$ are related to $\boldsymbol{\gamma}_t$ and $\boldsymbol{\pi}$ via the stick-breaking construction for the GEM distribution [35, 19].

The idea behind this construction is to imagine a stick with length 1, which will be successively broken into smaller pieces. Let $F:[0,1]^{\mathbb{N}}\to [0,1]^{\mathbb{N}}$ be given by

(2.9)
$$[F(\boldsymbol{x})]_1 := x_1, \quad [F(\boldsymbol{x})]_j := x_j \prod_{\ell=1}^{j-1} (1 - x_\ell),$$

where $\mathbf{x} = (x_j, j \in \mathbb{N})$. The idea is that $[F(\mathbf{x})]_j$ is the length of the piece broken at iteration j after successive fractions x_1, x_2, \ldots, x_j are broken off. Both $\boldsymbol{\pi}$ and $\boldsymbol{\gamma}_j$ have stick-breaking representations of this form:

$$\gamma'_{tg} \sim \text{Beta}(1, \alpha_0), \quad \pi'_k \sim \text{Beta}(1, \gamma_0),$$

$$\gamma_t = F(\gamma'_t), \qquad \boldsymbol{\pi} = F(\boldsymbol{\pi}'),$$

where $\gamma_t' = (\gamma_{tg}')$ and $\beta' = (\beta_k')$. Let $b_{\alpha}(\cdot)$ be the density of Beta $(1, \alpha)$, that is, $b_{\alpha}(x) \propto (1-x)^{\alpha-1}$. Then, we can write:

(2.10)
$$p(\boldsymbol{g}, \boldsymbol{k}, \boldsymbol{\gamma}', \boldsymbol{\pi}') = \prod_{t=1}^{T} \left(\prod_{i=1}^{n_t} \gamma_{t, g_{ti}} \prod_{g=1}^{\infty} b_{\alpha_0}(\gamma'_{tg}) \prod_{g=1}^{\infty} \pi_{k_{tg}} \right) \prod_{k=1}^{\infty} b_{\gamma_0}(\pi'_k).$$

Adding the network part, we have the full joint density:

$$(2.11) p(\boldsymbol{A}, \boldsymbol{\eta}, \boldsymbol{g}, \boldsymbol{k}, \boldsymbol{\gamma}', \boldsymbol{\pi}') = p(\boldsymbol{A} \mid \boldsymbol{g}, \boldsymbol{k}, \boldsymbol{\eta}) \cdot p(\boldsymbol{g}, \boldsymbol{k}, \boldsymbol{\gamma}', \boldsymbol{\pi}') \cdot p(\boldsymbol{\eta})$$

$$= \prod_{t=1}^{T} \left(\prod_{1 \leq i < j \leq n_{t}} L\left(\eta(k_{t,g_{ti}}, k_{t,g_{tj}}); A_{tij}\right) \prod_{i=1}^{n_{t}} \gamma_{t,g_{ti}} \prod_{g=1}^{\infty} b_{\alpha_{0}}(\gamma'_{tg}) \prod_{g=1}^{\infty} \pi_{k_{tg}} \right) \times \prod_{k=1}^{\infty} b_{\gamma_{0}}(\pi'_{k}) \prod_{1 \leq k \leq \ell < \infty} b_{\alpha_{\eta},\beta_{\eta}}(\eta_{k\ell}),$$

where $L(p; a) = p^a (1 - p)^{1-a}$ is the Bernoulli likelihood and $b_{\alpha,\beta}(\cdot)$ is the density of Beta (α, β) . Note that we have used the alternative notation $\eta(x, y) = \eta_{xy}$ for readability. In the next section, we derive a novel MCMC algorithms for sampling the posterior distribution of our model.

3. SLICE SAMPLING FOR HSBM

We propose a slice sampler for HSBM, based on a slice sampling algorithm we recently developed for HDP [1]. Recall that in slice sampling from a density f(x), we introduce the nonnegative variable u, and look at the joint density $g(x,u) = 1\{0 \le u \le f(x)\}$ whose marginal over x is f(x). Then, we perform Gibbs sampling on the joint g. In the end, we only keep samples of x and discard those of g. This idea has been employed in [20] to sample from the classical DP mixture and extended in [1] to sample HDPs.

In order to perform the slice sampling for HSBM, we introduce (independent) variables $\mathbf{u} = (u_{ti})$ and $\mathbf{v} = (v_{tq})$ so that the augmented joint density for (2.8) is

$$p(\boldsymbol{g}, \boldsymbol{k}, \boldsymbol{\gamma}', \boldsymbol{\pi}', \boldsymbol{u}, \boldsymbol{v}) = p(\boldsymbol{g}, \boldsymbol{u} \mid \boldsymbol{\gamma}) \cdot p(\boldsymbol{\gamma}') \cdot p(\boldsymbol{k}, \boldsymbol{v} \mid \boldsymbol{\pi}) \cdot p(\boldsymbol{\pi}'),$$

where for example

$$p(\boldsymbol{g}, \boldsymbol{u} \mid \boldsymbol{\gamma}) = \prod_{t=1}^{T} \prod_{i=1}^{n_t} 1\{0 \le u_{ti} \le \gamma_{t, g_{ti}}\},$$

and similarly for $p(\mathbf{k}, \mathbf{v} \mid \boldsymbol{\pi})$. Note that marginalizing out \mathbf{u} from $p(\mathbf{g}, \mathbf{u} \mid \boldsymbol{\gamma})$ gives back $p(\mathbf{g} \mid \boldsymbol{\gamma}) = \prod_t \prod_i \gamma_{t,g_{ti}}$ as before. The full augmented joint density is now (3.1)

$$\begin{split} p(\boldsymbol{A}, \boldsymbol{\eta}, \boldsymbol{g}, \boldsymbol{k}, \boldsymbol{\gamma}', \boldsymbol{\pi}', \boldsymbol{u}, \boldsymbol{v}) &= p(\boldsymbol{A} \mid \boldsymbol{g}, \boldsymbol{k}, \boldsymbol{\eta}) \cdot p(\boldsymbol{g}, \boldsymbol{k}, \boldsymbol{\gamma}', \boldsymbol{\pi}', \boldsymbol{u}, \boldsymbol{v}) \cdot p(\boldsymbol{\eta}) \\ &= \prod_{t=1}^{T} \left(\prod_{1 \leq i < j \leq n_{t}} L\Big(\eta(k_{t, g_{ti}}, k_{t, g_{tj}}); A_{tij} \Big) \prod_{i=1}^{n_{t}} 1\{u_{ti} \leq \gamma_{t, g_{ti}}\} \prod_{g=1}^{\infty} b_{\alpha_{0}}(\gamma'_{tg}) \prod_{g=1}^{\infty} 1\{v_{tg} \leq \pi_{k_{tg}}\} \right) \times \\ &\prod_{k=1}^{\infty} b_{\gamma_{0}}(\pi'_{k}) \prod_{1 \leq k \leq \ell < \infty} b_{\alpha_{\eta}, \beta_{\eta}}(\eta_{k\ell}), \end{split}$$

with the support understood to be restricted to $u \geq 0$ and $v \geq 0$. We then perform block Gibbs sampling on the augmented density. Note that marginalizing variables (u_t) and (v_t) out, we get back original joint density (2.11). The idea is to sample (γ', u) jointly given the rest of variables, and similarly for (π', v) . The updates for variables u, γ', v and π' are similar to those in [1]. However, the updates for the underlying latent groups g and group-communities k require some care due to the coupling introduced by the SBM likelihood. As can be seen from the derivation below, these updates will be quite nontrivial in the case of SBM relative to case where the bottom layer is a simple mixture model.

3.1. Sampling $(\boldsymbol{u}, \boldsymbol{\gamma}') \mid \cdots$. First, we sample $(\boldsymbol{u} \mid \boldsymbol{\gamma}', \Theta_{-\boldsymbol{u}\boldsymbol{\gamma}'})$, where $\Theta_{-\boldsymbol{u}\boldsymbol{\gamma}'}$ denotes all variables except \boldsymbol{u} and $\boldsymbol{\gamma}'$. This density factorizes and coordinate posteriors are $p(u_{ti} \mid \boldsymbol{\gamma}', \Theta_{-\boldsymbol{u}\boldsymbol{\gamma}'}) \propto 1\{u_{ti} \leq \gamma_{t,q_{ti}}\}$, that is

$$u_{ti} \mid \boldsymbol{\gamma}', \Theta_{-\boldsymbol{u}\boldsymbol{\gamma}'} \sim \text{Unif}(0, \gamma_{t,q_{ti}}).$$

Next, we sample from $(\gamma' \mid \Theta_{-u\gamma'})$. To do this, we first marginalize out u in (3.1) which gives back (2.11). The corresponding posterior is, thus, proportional to (2.11) viewed only as a function of γ' . The posterior factorizes over t and g and we have [1, Lemma 1]

(3.2)
$$\gamma_{tg}' \mid \Theta_{-u\gamma'} \sim \text{Beta}\left(n_g(\boldsymbol{g}_t) + 1, n_{>g}(\boldsymbol{g}_t) + \alpha_0\right),$$

where $n_g(\mathbf{g}_t) = |\{i: g_{ti} = g\}|$ and $n_{>g}(\mathbf{g}_t) = |\{i: g_{ti} > g\}|$.

3.2. Sampling $(\boldsymbol{v}, \boldsymbol{\pi}') \mid \cdots$. First, we sample $(\boldsymbol{v} \mid \boldsymbol{\pi}', \Theta_{-\boldsymbol{v}\boldsymbol{\pi}'})$ which factorizes and coordinate posteriors are $p(v_{tg} \mid \boldsymbol{\pi}', \Theta_{-\boldsymbol{v}\boldsymbol{\pi}'}) \propto 1\{v_{tg} \leq \beta_{k_{tg}}\}$, that is

$$v_{tg} \mid \boldsymbol{\pi}', \Theta_{-\boldsymbol{v}\boldsymbol{\pi}'} \sim \text{Unif}(0, \pi_{k_{tg}}).$$

Next, we sample from $(\pi' \mid \Theta_{-v\pi'})$. As in the case of γ' , we first marginalize v which leads to the usual block Gibbs sampler updates: The posterior factorizes over k, and

(3.3)
$$\pi'_{k} \mid \Theta_{-\boldsymbol{v}\boldsymbol{\pi}'} \sim \operatorname{Beta}\left(n_{k}(\boldsymbol{k}) + 1, n_{>k}(\boldsymbol{k}) + \gamma_{0}\right),$$

where $n_k(\mathbf{k}) = |\{(t,g): k_{tg} = k\}|$ and similarly for $n_{>k}(\mathbf{k})$.

3.3. Sampling $g \mid \cdots$. This posterior factorizes over t (but not over i). From (3.1), we have

$$(3.4) \qquad \mathbb{P}(g_{ti} = g \mid \boldsymbol{g}_{-ti}, \Theta_{-\boldsymbol{g}}) \propto \prod_{j \in [n_t] \setminus \{i\}} L\left(\eta(k_{tg}, k_{t,g_{tj}}); A_{tij}\right) 1\{u_{ti} \leq \gamma_{tg}\}.$$

Let $G_{ti} := G(\gamma_t; u_{ti}) := \sup\{g : u_{ti} \leq \gamma_{tg}\}$. According to the above g_{ti} given everything else will be distributed as

$$g_{ti} \mid \cdots \sim (\rho_{ti}(g))_{g \in [G_{ti}]},$$

where we have defined

$$\rho_{ti}(g) := \prod_{j \in [n_t] \setminus \{i\}} L\left(\eta(k_{tg}, k_{t,g_{tj}}); A_{tij}\right) \\
= \prod_{j \in [n_t] \setminus \{i\}} \prod_{\ell} \left[L\left(\eta(k_{tg}, \ell); A_{tij}\right)\right]^{1\{k_{t,g_{tj}} = \ell\}} \\
= \prod_{j \in [n_t] \setminus \{i\}} \prod_{\ell} \left[\eta(k_{tg}, \ell)^{A_{tij}} [1 - \eta(k_{tg}, \ell)]^{1 - A_{tij}}\right]^{1\{k_{t,g_{tj}} = \ell\}} \\
= \prod_{\ell} \eta(k_{tg}, \ell)^{\tau_{ti\ell}} [1 - \eta(k_{tg}, \ell)]^{m_{ti\ell} - \tau_{ti\ell}},$$

where

(3.5)
$$\tau_{ti\ell} := \sum_{j \in [n_t] \setminus \{i\}} A_{tij} 1\{k_{t,g_{tj}} = \ell\}, \quad m_{ti\ell} := \sum_{j \in [n_t] \setminus \{i\}} 1\{k_{t,g_{tj}} = \ell\}.$$

These expressions can be simplified by noting that $k_{t,g_{tj}} = z_{tj}$.

3.4. Sampling $k \mid \cdots$. This posterior also factorizes over t (but not over g). First, note that since we are conditioning on g, we can simplify as

$$\prod_{1 \leq i < j \leq n_{t}} L\left(\eta(k_{t,g_{ti}}, k_{t,g_{tj}}); A_{tij}\right) = \prod_{1 \leq i < j \leq n_{t}} \prod_{g,g'=1}^{\infty} \left[L\left(\eta(k_{tg}, k_{tg'}); A_{tij}\right)\right]^{1\{g_{ti}=g, g_{tj}=g'\}}$$

$$= \prod_{1 \leq i < j \leq n_{t}} \prod_{g,g'=1}^{\infty} \left[\eta(k_{tg}, k_{tg'})^{A_{tij}} \left[1 - \eta(k_{tg}, k_{tg'})\right]^{1 - A_{tij}}\right]^{1\{g_{ti}=g, g_{tj}=g'\}}$$

$$= \prod_{g,g'=1}^{\infty} \eta(k_{tg}, k_{tg'})^{\xi_{tgg'}} \left[1 - \eta(k_{tg}, k_{tg'})\right]^{O_{tgg'} - \xi_{tgg'}},$$

where

$$(3.6) \quad \xi_{tgg'} := \sum_{1 \le i < j \le n_t} A_{tij} 1\{g_{ti} = g, g_{tj} = g'\}, \quad O_{tgg'} := \sum_{1 \le i < j \le n_t} 1\{g_{ti} = g, g_{tj} = g'\}.$$

Note that $\xi_{tqq'}$ is not symmetric in g and g'. Let us define

$$h_{tgg'}(k,\ell) := \eta(k,\ell)^{\xi_{tgg'}} [1 - \eta(k,\ell)]^{O_{tgg'} - \xi_{tgg'}},$$

so that we have

$$\prod_{1 \le i < j \le n_t} L\Big(\eta(k_{t,g_{ti}}, k_{t,g_{tj}}); A_{tij}\Big) = \prod_{g,g'=1}^{\infty} h_{tgg'}(k_{tg}, k_{tg'}).$$

It is not hard to see that the posterior of $k \mid \cdots$ factorizes over t, and for any fixed t,

(3.7)
$$p(\mathbf{k}_t \mid \mathbf{k}_{-t}, \Theta_{-\mathbf{k}}) \propto \prod_{g,g'=1}^{\infty} h_{tgg'}(k_{tg}, k_{tg'}) \prod_{g=1}^{\infty} 1\{v_{tg} \leq \pi_{k_{tg}}\}.$$

It then follows that for any fixed t and g:

(3.8)

$$p(k_{tg} \mid \mathbf{k}_{-tg}, \Theta_{-\mathbf{k}}) \propto 1\{v_{tg} \leq \pi_{k_{tg}}\}h_{tgg}(k_{tg}, k_{tg}) \prod_{g': g' \neq g} [h_{tgg'}(k_{tg}, k_{tg'}) h_{tg'g}(k_{tg'}, k_{tg})]$$

By symmetry of $\eta(k,\ell)$ in its two argument, $h_{tgg'}(k,\ell)$ is also symmetric in (k,ℓ) , hence

(3.9)
$$h_{tgg'}(k,\ell) h_{tg'g}(\ell,k) = \eta(k,\ell)^{\widetilde{\xi}_{tgg'}} [1 - \eta(k,\ell)]^{\widetilde{O}_{tgg'} - \widetilde{\xi}_{tgg'}}$$
$$=: \widetilde{h}_{tgg'}(k,\ell),$$

where

$$\widetilde{\xi}_{tqq'} = \xi_{tqq'} + \xi_{tq'q}, \quad \widetilde{O}_{tqq'} = O_{tqq'} + O_{tq'q}.$$

We note that

$$\widetilde{\xi}_{tgg'} := \sum_{1 \le i \ne j \le n_t} A_{tij} 1\{g_{ti} = g, \ g_{tj} = g'\}, \quad \widetilde{O}_{tgg'} := \sum_{1 \le i \ne j \le n_t} 1\{g_{ti} = g, \ g_{tj} = g'\},$$

It follows that

(3.10)
$$\mathbb{P}(k_{tg} = k \mid \mathbf{k}_{-tg}, \Theta_{-\mathbf{k}}) \propto 1\{v_{tg} \leq \pi_k\} h_{tgg}(k, k) \prod_{g': g' \neq g} \widetilde{h}_{tgg'}(k, k_{tg'})$$

Let us simplify the last factor further. Writing $\widetilde{h}_{tgg'}(k, k_{tg'}) = \prod_{\ell} [\widetilde{h}_{tgg'}(k, \ell)]^{1\{k_{tg'}=\ell\}}$, we have, using (3.9),

$$\prod_{g': g' \neq g} \widetilde{h}_{tgg'}(k, k_{tg'}) = \prod_{\ell} \prod_{g': g' \neq g} [\widetilde{h}_{tgg'}(k, \ell)]^{1\{k_{tg'} = \ell\}}
= \prod_{\ell} \eta(k, \ell)^{\zeta_{tg\ell}} [1 - \eta(k, \ell)]^{R_{tg\ell} - \zeta_{tg\ell}},$$

where

$$(3.11) \quad \zeta_{tg\ell} := \sum_{g': g' \neq g} \widetilde{\zeta}_{tgg'} 1\{k_{tg'} = \ell\} = \sum_{g': g' \neq g} \sum_{1 \leq i \neq j \leq n_t} A_{tij} 1\{g_{ti} = g, \ g_{tj} = g', \ k_{tg'} = \ell\},$$

and similarly

$$(3.12) \quad R_{tg\ell} := \sum_{g': g' \neq g} \widetilde{O}_{tgg'} 1\{k_{tg'} = \ell\} = \sum_{g': g' \neq g} \sum_{1 \leq i \neq j \leq n_t} 1\{g_{ti} = g, \ g_{tj} = g', \ k_{tg'} = \ell\},$$

Let $K_{tg} := K(\boldsymbol{\pi}; v_{tg}) := \sup\{k : v_{tg} \leq \pi_k\}$. According to the above, k_{tg} given everything else will be distributed as

$$k_{tg} \mid \cdots \sim (\delta_{tg}(k))_{k \in [K_{tg}]},$$

where

$$\delta_{tg}(k) = \prod_{\ell} \eta_{k\ell}^{\zeta_{tg\ell}} [1 - \eta_{k\ell}]^{R_{tg\ell} - \zeta_{tg\ell}}.$$

3.5. Sampling $\eta \mid \cdots$. Recalling $z_{ti} = k_{t,g_{ti}}$, the relevant part of (3.1) is

$$\left[\prod_{t=1}^{T}\prod_{1\leq i< j\leq n_{t}}L\left(\eta(z_{ti},z_{tj});A_{tij}\right)\right]\prod_{1\leq k\leq \ell<\infty}b_{\alpha_{\eta},\beta_{\eta}}(\eta_{k\ell}).$$

Using the fact that $\eta(k,\ell) = \eta_{k\ell}$ is symmetric in its two arguments (that is, we treat $\eta_{k\ell}$ and $\eta_{\ell k}$ as the same variable), we have

$$p(\eta_{k\ell} \mid \cdots) \propto \eta_{k\ell}^{\lambda_{k\ell}} (1 - \eta_{k\ell})^{N_{k\ell} - \lambda_{k\ell}} b_{\alpha_{\eta}, \beta_{\eta}}(\eta_{k\ell}),$$

where for $k \neq \ell$.

(3.13)

$$\lambda_{k\ell} = \sum_{t=1}^{T} \sum_{1 \le i \ne j \le n_t} A_{tij} 1\{z_{ti} = k, z_{tj} = \ell\}, \quad N_{k\ell} = \sum_{t=1}^{T} \sum_{1 \le i \ne j \le n_t} 1\{z_{ti} = k, z_{tj} = \ell\},$$

and

$$(3.14) \quad \lambda_{kk} = \sum_{t=1}^{T} \sum_{1 \le i < j \le n_t} A_{tij} 1\{z_{ti} = z_{tj} = k\}, \quad N_{kk} = \sum_{t=1}^{T} \sum_{1 \le i < j \le n_t} 1\{z_{ti} = z_{tj} = k\}.$$

Recalling that $b_{\alpha_{\eta},\beta_{\eta}}(\eta_{k\ell}) \propto \eta_{k\ell}^{\alpha_{\eta}-1}(1-\eta_{k\ell})^{\beta_{\eta}-1}$ we conclude that

$$\eta_{k\ell} \mid \cdots \sim \text{Beta} (\lambda_{k\ell} + \alpha_{\eta}, N_{k\ell} - \lambda_{k\ell} + \beta_{\eta}).$$

Remark 1. Comparing with the updates in [1], we observe that, under HSBM, the updates for k and g (which is equivalent to t in [1]) and η (which is equivalent to parameters of f in [1]) are much more complex. For the usual HDP, the bottom layer is a simple mixture model where observations are independent given the labels, and each, only depends on its own label. This causes the posterior for k, g and the parameters of the mixture components to factorize over their coordinates. In contrast, the SBM likelihood makes each observation A_{tij} dependent on two labels z_{ti} and z_{tj} . This causes the posterior for k, g and η to remain coupled under HSBM. Nevertheless, the Gibbs sampling scheme above allows us to effectively sample from these coupled multivariate posteriors.

3.6. Summary. A summary of the overall slice sampler for HSBM is presented in Appendix A. For comparison, we have also included the corresponding slice sampler for a single-layer DP-SBM, which can be thought of as HSBM with a single layer and with w_1 set equal to π . This model is essentially the same as the one considered in [25], but with a different sampling algorithm. Comparing the two algorithms reveals the complexity of the HSBM model relative the DP-SBM, hence the possibility of incorporating more dependencies when applied to real multiplex data.

4. Simulation study

In this section, we use simulated data to demonstrate the applicability of our model and the proposed algorithm. First, we simulate a hypothetical multilayer network, with each layer representing the friendship network among contestants in a competition between U.S. states (each state represented by a single person). For this simulation, cluster assignments were done independently. For a second illustration, a multilayer network was

simulated, but this time with dependence among cluster assignments. In this example, we also study the effect of increasing the number of layers on the performance of the algorithms.

We fit both the HSBM and the single-layer version, DP-SBM. The latter is fit separately to each layer. The two algorithms are summarized in Appendix A. In each case, we run the corresponding MCMC for ITRmax iterations, throw away the first burnin $(=0.6 \cdot \text{ITRmax} \text{ or } = \text{ITRmax}/2, \text{ usually})$ and use the rest of the samples for inference. As we will see, in these simulations, our model clearly outperforms the models that assume independence across layers.

Measuring performance. After obtaining the MCMC samples, the aim is usually to calculate posterior summaries and, in a simulation study, compare them to the true parameters (cluster assignments and link probabilities in our case). With Bayesian mixture models, however, we face the well-known label switching problem. A K-component mixture is invariant to permutations of the labels of the components, and, as a consequence, the posterior of the mixture parameters will have K! modes—making it impossible to perform component specific inference directly from the MCMC draws.

To get around this issue, we measure the accuracy of the estimated cluster labels using the normalized mutual information (NMI), a well-known measure of similarity between two cluster assignments. NMI takes values in [0,1] where 1 corresponds to a perfect match. A random assignment against the truth is guaranteed to map to NMI \approx 0. The NMI penalizes mismatch quite aggressively. In our setting, an NMI \approx 0.5 corresponds to a roughly 90% match.

In the multiplex setting, we can compute at least two NMIs. For the first one, which we refer to as the *slicewise NMI*, we compute the NMI between the cluster assignments separately for each layer. We then either report these individually or report their average, referred to as the *average slicewise NMI*. For the second one, which we refer to as the *aggregate NMI*, we consider the labels for all the layers together and compute a single NMI between the competing label assignments. Achieving a high aggregate NMI is more challenging since it requires consistency both within and across layers.

MAP estimate. We also compute the maximum a posteriori (MAP) label assignment, that is, for each node we find the label which is most likely according to the posterior: $\operatorname{argmax}_k \mathbb{P}(z_{ti} = k \mid \cdots)$. Associated with the MAP estimate, there is a confidence which is the value of the posterior probability. To compute the MAP estimate we use the posterior estimate given by $\frac{1}{\operatorname{ITRmax-burnin}} \sum_{j=\operatorname{burnin}+1}^{\operatorname{ITRmax}} 1(\hat{z}_{ti}^{(j)} = k)$ where $\hat{z}^{(j)} = (\hat{z}_{it}^{(j)})$ is the label assignment at MCMC iteration j. Here, we ignore the potential mismatch between $\hat{z}^{(j)}$ and $\hat{z}^{(j')}$ due to the potential label-switching. In practice, all labels after MCMC convergence are coming from a single mode of the posterior as can be verified by computing the NMI between consecutive samples $z^{(j)}$ and $z^{(j+1)}$.

4.1. A simulated multilayer friendship network. Consider, as an illustration, the group of young women who run, every year, for the Miss America title. Each state holds a preliminary competition to choose their delegate for the main race, and the winners hold the title of Miss from their home state.

Suppose that we are interested in identifying groups of competitors by their personalty type: extrovert, introvert or ambivert (a balance between the former two). Let us also suppose that around 40% of the female population in America can be classed as extrovert, 35% introvert, and 25% is in the ambivert group. Assume that the distribution of personality types among the contests of each edition of the Miss America competition is similar to that of the general population.

After meeting each other in the competition, some of the girls tend to bond and become friends at a social networking website. It is reasonable to assume that there is a higher chance that the extrovert girls befriend themselves and others. On the other hand, there is a lower chance that the introvert girls bond to other contestants. For this simulation exercise, we consider the probabilities of friendship between and within groups listed in Table 1.

Table 1. Hypothetical probability of friendship between young females in the U.S.A

	Extrovert	Ambivert	Introvert
Extrovert	90%	75%	50%
Ambivert	75%	60%	25%
Introvert	50%	25%	10%

In this example, each edition of the competition is a layer and the competitors of each state represent the nodes. Note, however, that given the distribution of the personality types, the cluster assignments across layers are independent. That makes sense, since the node for the state of Texas, for example, is represented by different girls in each edition.

The probabilities within and between clusters characterize the clusters. The cluster weights (i.e., proportion of nodes in each cluster) can vary, but are somewhat similar to those of the general population of young females in America.

The observed data are the binary matrices of friendship between fifty girls ($n_t = 50$), for each one of T = 5 editions of the competition (say from year 2014 to 2018). Our aim in this exercise is to be able to correctly classify the contenders into subgroups based on their friendship matrices. Figure 2 shows the simulated networks and their adjacency matrices, respectively.

Results. The algorithms were run for ITRmax = 5000 and burnin = 3000. Figure 3 shows both the slicewise NMI and the aggregate NMI for the two algorithms HSBM and DP-SBM, the latter considered the single-layer version of HSBM. The NMIs are computed between the MAP estimates in each case and the true labels; the boxplots are obtained over 20 replications. As the figure clearly shows, HSBM improves on DP-SBM even for the slicewise NMI which is computed separately for each layer. This consistent advantage, even though the cluster assignments are conditionally independent, is due to the ability of the proposed multilayer model to use all the layers to estimate the probability of friendship within and between clusters. The multilayer model also takes into account similarities between the cluster weights across different layers.

For the aggregate NMI, the advantage of HSBM is much more pronounced. This is expected since DP-SBM does not have any means of matching the cluster assignments

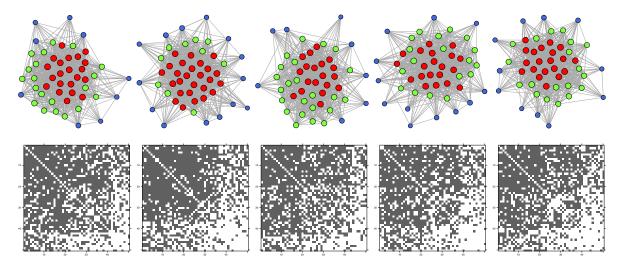


FIGURE 2. The five simulated networks and their corresponding adjacency matrices for the Miss America example. Red nodes correspond to extrovert girls, green corresponds to ambivert, and blue to introvert.

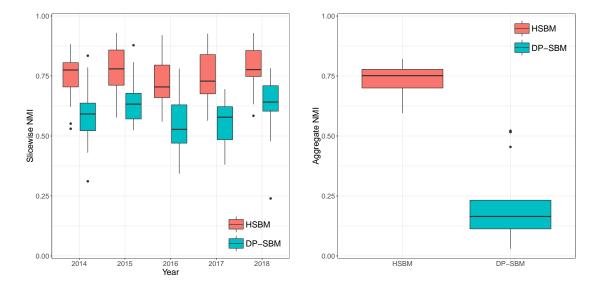


FIGURE 3. Normalized mutual information (NMI) plots for the Miss America example. Left: Slicewise NMI which is computed separately for each layer. Right: Aggregate NMI, computed by considering all the layers together.

across layers. The aggregate NMI plot shows that HSBM indeed succeeds at obtaining consistent cluster assignments across layers.

Figure 4 compares the posterior NMIs for the two methods. The posterior NMI is computed by considering the NMI of each MCMC iteration (after the burnin period) relative to the true label, for a single realization of the model. The plots, thus, show the concentration of the posteriors around the true labeling (as measured by the NMI). A similar behavior as in Figure 3 is observed, with visible improvements for slicewise NMIs and a significant improvement for the aggregate NMI.

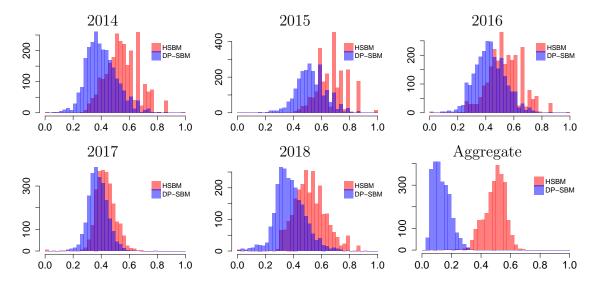


FIGURE 4. Plots of the histogram of the posterior NMI for each layer and in aggregate, for the Miss America example. The plots are obtained by evaluating the (slicewise or aggregate) NMI of each iteration of the algorithm after the burn-in period, relative to the true labels. The plots are for a single realization.

It is worth nothing that Figures 3 and 4 are not showing the same information in different formats; the former shows the frequentist variation of the mode of the posterior (over 20 trials) while the latter shows the variation of a single realization of the posterior. For example, Figure 3 shows that even individual samples from the posterior are performing well (i.e., have high NMI $\gtrsim 0.5$) without the aggregating effect of the MAP procedure.

4.2. Simulated multiplex network with dependent layers. As a second example, we simulate a sequence of layers, indexed by time, such that a given node is more likely to belong to the same community in two consecutive layers. The labels of the first layer are drawn from a k-component mixture with weights $w = (w_1, w_1, \ldots, w_k)'$. Then, given the layers up to time t-1, the community labels for time t layer are sampled in the following Markovian fashion: Each node i retains its label $z_{t-1,i}$ with probability p or is assigned a random label from $[k] \setminus \{z_{t-1,i}\}$ with probability 1-p.

Here, we simulate a multiplex network with $n_t = 50$ nodes, and k = 3 components, for all layers t. We consider different number of layers T = 2, 4, 8 and 12. The labels of the first layer are drawn from a mixture of three communities with equal weights. The label retention probability is set to p = 0.9 for subsequent layers. The connection probability matrix is

$$\eta^* = \begin{pmatrix} 0.8 & 0.1 & 0.3 \\ 0.1 & 0.9 & 0.2 \\ 0.3 & 0.2 & 0.7 \end{pmatrix}.$$

Results. As in the previous example, we compare the performance of the proposed model (HSBM) to the single layer models (DP-SBM) using the NMI for the MAP labels. The algorithms were run for ITRmax = 2000 and burnin = ITRmax/2. Figure 5 shows

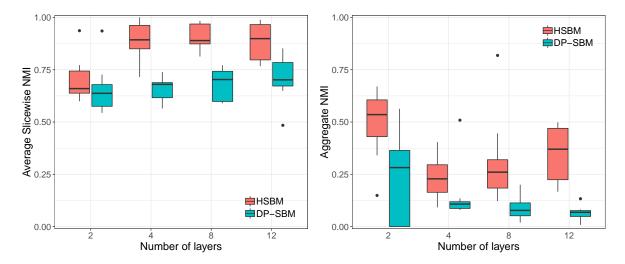


FIGURE 5. Normalized mutual information (NMI) versus the number of layers for the dependent (Markovian) example. Left: Average slicewise NMI for which the NMI is computed separately for each layer and then averaged. Right: Aggregate NMI, computed by considering all the layers together.

the plots of the average slicewise NMI (left) and the aggregate NMI (right), versus the number of layers (T), over 8 replications. Note that we are generating from four models with different number of layers.

The proposed HSBM clearly outperforms the single-layer model. Again, the gain is due to using more information to estimate the cluster structure in each layer. This, for example, is clear from the slicewise NMI plots, where for T=2 the performance of the proposed model is closer to the single-layer model, and the gap increases as T is increased. We also note the significant boost HSBM provides in terms of the aggregate NMI for large T.

5. Data analysis

In this section, we illustrate the performance of our model and algorithm on two real data examples. First, we present an application on the Glasgow Teenage Friends and Lifestyle Study (TFLS) dataset [8, 24, 28, 29]. Then, the model is applied to the Krebs multiplex network of IT workers [21].

5.1. **TFLS multiplex network.** The Teenage Friends and Lifestyle Study was initially aimed to identify the processes by which attitudes towards smoking and smoking behavior change over adolescence (from early to mid teenage years). Students were followed over a two year period starting in February 1995, with 160 thirteen-year-old students, and ending in January 1997. From the original 160 students, only 135 were present at all the three measurement points. The friendship networks were formed by allowing the students to name up to six friends. Students were also asked about their lifestyle, including substance use and leisure activities. The school was representative of others in the region in terms of the social class composition [8].

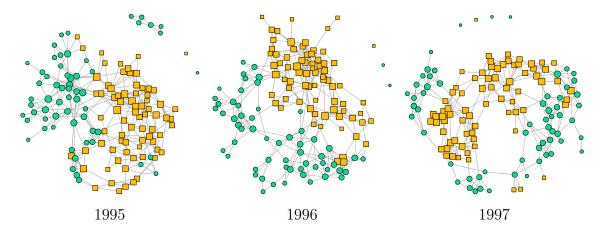


FIGURE 6. Three layers of the multiplex friendship network of the TFLS dataset. Each column corresponds to a different year, i.e., a different layer. The colors (and also shapes) denote the two genders, with green circles representing girls and yellow squares representing boys.

For this study, we considered the 135 students who were present at all the measurement points. Figure 6 shows the friendship networks using a Fruchterman-Reingold layout obtained by igraph [11]. The nodes are colored by gender making it easy to observe that kids of the same gender tend congregate together. Two different shapes (circles versus squares) are also used to indicate the gender of a node. Later, we keep these shapes to represent gender, while changing the color to represent the communities inferred by the algorithm (cf. Figure 7). In this data, the estimated probability of friendship between two boys is 5.54%, between two girls is 2.50% and between a boy and a girl is 2.57%. The average degrees of the three layers are 5.04, 4.93 and 4.96, respectively.

Results. We ran Algorithm 1 for 10000 iterations, and used the last 5000 for inference. The hyperparameters were set at $\alpha_0 = \gamma_0 = 5$. Figure 7 shows the results. The nodes are colored by the maximum a posteriori (MAP) cluster assignments for each layer. The shades indicate posterior confidence, i.e., the posterior probability of the MAP label: the darker the color, the higher the posterior probability. The shape of a node indicates the student's gender as in Figure 6.

For the top row in Figure 7, we fixed the layout to be that of year 1995 for all the layers. The bottom row shows the 1996 and 1997 layers with layouts that were free to adapt to the structure of the network.

It is interesting to note that the HSBM inference has converged to a 2-community structure for all the layers. The labels of the two communities are also the same across the layers, indicating that the model believes the same communities are present in all the layers (for example, color blue corresponds to the same community label in all the years). The structure of the communities however changes across layers. This is natural since the connections among the nodes change quite considerably from year to year, as can be seen from the fixed layout networks (top row in Figure 7).

On the other hand, as the free layout plots show (bottom row in Figure 7) the inferred community structure is very reasonable, bisecting each network into two pieces with few connections in between. Note that the inferred communities do not correspond to the

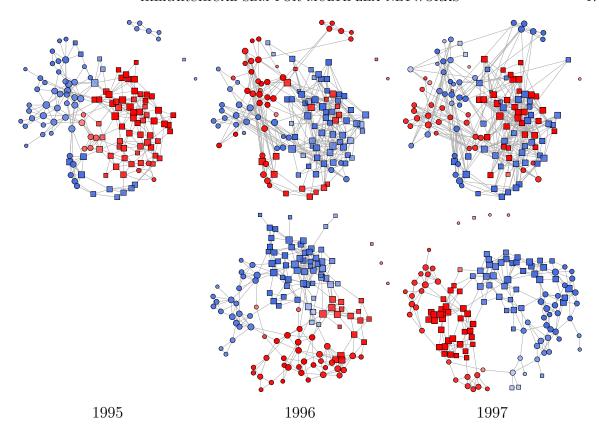


FIGURE 7. Multiplex friendship network of the TFLS dataset. Each column corresponds to a different year, i.e., a different layer. The nodes are colored according to their maximum a posteriori (MAP) community assignment by the HSBM algorithm. The shade (i.e., transparency) indicates posterior confidence, with darker colors representing a higher posterior value for the MAP label. The top row shows all the networks with the fixed layout of the 1995 network. The bottom row plots the 1996 and 1997 networks with layouts that were allowed to adapt to their structure.

bisection of the network by gender. Both communities contain nodes of different genders. It is also interesting to note that low confidence assignment occurs for some of the nodes at the boundary of the two communities in each layer.

5.2. Krebs multiplex network of IT workers. We next consider the data that was collected by Valdis Krebs in the IT Department of a Fortune 500 company. The data and the following description is available from [21] (the original source is unknown). Fifty-six individuals were surveyed. The first 7 were administrative staff, the next 12, 17 and 17 were together in three different departments and the last 3 were executives. In plotting the networks, we draw nodes in these groups with the following shapes: sphere, circle, square, triangle and star.

Each individual was surveyed on five kinds of interactions (links):

- (1) BUSINESS-1: With whom do you work with in Business Process 1?
- (2) BUSINESS-2: With whom do you work with in Business Process 2?

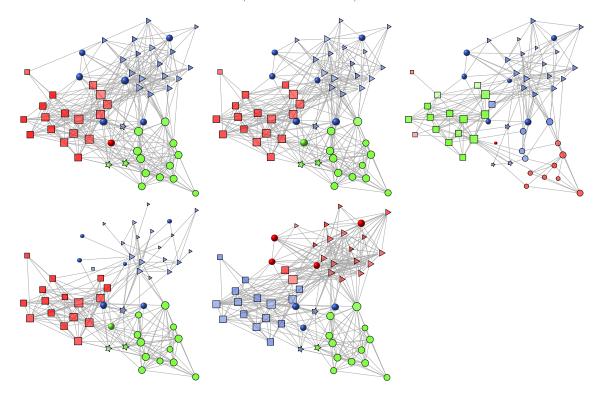


FIGURE 8. Krebs multiplex network, fixed layouts. The layouts are constrained to be the same as the first layer. From top-left to bottom-right, the layers are BUSINESS-1, BUSINESS-2, ADVICE, TECHNICAL and CUSTOMER. The shape of the nodes indicate their group: sphere = administrative, star = executive, and circle, square, and triangle each represents one of the three departments. The color of the nodes indicate the MAP communities estimated by HSBM.

(3) ADVICE: Who do you seek for advice before making a key decision?

- (4) TECHNICAL: Who do you seek for technical expertise in IT?
- (5) CUSTOMER: With whom do you discuss customer needs and issues?

They also reported the frequency of their interactions in each category on five levels (from yearly or less, to daily and more). The interactions were directed. We have ignored the frequency and direction and only considered the binary symmetric adjacency matrices representing whether there was any interaction.

Results. Figures 8 and 9 show the results of fitting HSBM to Krebs network. In Figure 8, the layouts are all fixed to be the same as the first layer, while in Figure 9, they are free to adapt to the structure of the network at each layer. The hyperparameters of HSBM were set to $\alpha_0 = \gamma_0 = 1$ and ITRmax = 2burnin = 10000. The color of the nodes denote the estimated MAP communities.

The model consistently recovers three communities in each of the layers. This is consistent with the existence of three big departments. (The administrative staff and executives are assigned to one of these communities in each layer.) A feature of the estimated labels is the matching between communities at different layers. For example,

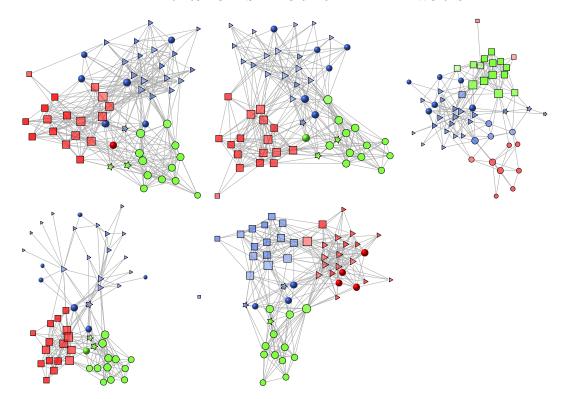


FIGURE 9. Krebs example, free layouts. The same as Figure 8 but with layouts allowed to adapt to the structure of the network in each layer.

according to HSBM the red community is the same in all layers meaning that it has the same signature connectivity vector to other communities. This is confirmed in Figure 10 where separate estimates of the connectivity matrix based on the labels of each layer are plotted.

Figure 8 reveals the following interesting point: Roughly the same nodes are assigned to the same communities in layers BUSINESS-1, BUSINESS-2 and TECHNICAL. This is very reasonable and indicates that the nature of technical collaboration is the same among the individuals. Note that the TECHNICAL layer is sparser, especially over the blue community, and without the signal from the other two layers, classifying the blue nodes into the same community would be difficult.

6. Discussion

In this work, we proposed a novel Bayesian model for community detection in multiplex networks by adopting the well-known HDP as a prior distribution for community assignments. Under the random partition prior, a block model is assumed. This model facilitates flexible modeling of community structure as well as link probabilities with its ability of incorporating potential dependency and borrowing strength among networks from different layers. For posterior inference of HSBM, we develop an efficient slicer sampler. The principles behind the slice sampler can be applied to developing sampling algorithms for many other models. Future work will be focused on developing models for

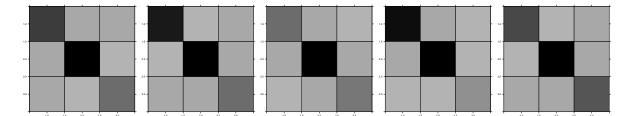


FIGURE 10. Estimated connectivity matrices for each layer of the Krebs network. We have used the labels of each layer to estimate the connectivity matrix η , separately for each layer, using the empirical averages. The plots confirm that communities 1 (red), 2 (green) and 3 (blue) play roughly the same roles in each layer.

community detection in networks with covariates, and for inference of network-valued objects.

ACKNOWLEDGEMENT

The contribution of LL is funded by NSF grants IIS 1663870 and DMS 1654579, and a DARPA grant N66001-17-1-4041.

APPENDIX A. SUMMARY OF SAMPLING ALGORITHMS

Algorithm 1 summarizes the slicer sampler for HSBM presented in Section 3. For comparison, we have also included the slice sampler for a single-layer DP-SBM in Algorithm 2, obtained by the same principle as discussed in Section 3. DP-SBM can be considered a single layer (T=1) version of HSBM with $\mathbf{w}_1 = \boldsymbol{\pi}$.

A.1. Slice sampler for single-layer DP-SBM. For the sake of completeness, let us also derive the slice sampler for the single-layer version of our model which can be thought of as the usual SBM with a DP prior on the latent communities. In the single-layer case, with a single adjacency matrix A, the model simplifies to

(A.1)
$$\begin{aligned} \boldsymbol{\gamma} \mid \alpha_0 &\sim \operatorname{GEM}(\alpha_0), \\ z_i \mid \boldsymbol{\gamma} &\sim \boldsymbol{\gamma}, \quad i = 1, \dots, n. \\ \boldsymbol{\eta} &= (\eta_{xy}) \overset{\text{iid}}{\sim} \operatorname{Beta}(\alpha_{\eta}, \beta_{\eta}) \\ A \mid \boldsymbol{z}, \boldsymbol{\eta} &\sim \operatorname{SBM}(\boldsymbol{z}, \boldsymbol{\eta}). \end{aligned}$$

The full joint distribution of the model including the augmentation by variable $\mathbf{u} = (u_i)$ introduced for slice sampling is given by

(A.2)
$$p(\boldsymbol{A}, \boldsymbol{z}, \boldsymbol{\eta}, \boldsymbol{\gamma}', \boldsymbol{u}) = \prod_{1 \le i < j \le n} L(\eta(z_i, z_j); A_{ij}) \prod_{i=1}^{n} 1\{u_i \le \gamma_{z_i}\} \prod_{z=1}^{\infty} b_{\alpha_0}(\gamma_z') \prod_{1 \le k \le \ell < \infty} b_{\alpha_{\eta}, \beta_{\eta}}(\eta_{k\ell}).$$

Sampling $(u, \gamma') \mid \cdots$. As usual, we first we sample $(u \mid \gamma', \Theta_{-u\gamma'})$, where $\Theta_{-u\gamma'}$ collects all variables except u and γ' . This density factorizes and coordinate posteriors

Algorithm 1 Slice sampler for HSBM

```
1: Initialize G_t^{\text{cap}} and K^{\text{cap}} to pre-specified values (say 10).
  2: Initialize g_t and k_t to all-ones vectors.
  3: Initialize (u_{ti}) and (v_{tq}) to independent uniform variables.
       while not CONVERGED, nor maximum iterations reached do
              for t = 1, \dots, T do
  5:
                     Sample \gamma'_{tq} \sim \text{Beta}\left(n_g(\boldsymbol{g}_t) + 1, n_{>g}(\boldsymbol{g}_t) + \alpha_0\right) for all g \in [G_t^{\text{cap}}].
  6:
                     Let [\gamma_t]_{1:G_{\star}^{\operatorname{cap}}} \leftarrow [F(\gamma_t')]_{1:G_{\star}^{\operatorname{cap}}}.
  7:
                     Let G_{ti} \leftarrow \max\{g: u_{ti} \leq \gamma_{tg}\}, \forall i \in [n_t] \text{ and } G_t \leftarrow \max_{i=1,\dots,n_t} G_{ti}.
  8:
  9:
              end for
              Sample \pi'_k \sim \text{Beta}\left(n_k(\mathbf{k}) + 1, n_{>k}(\mathbf{k}) + \gamma_0\right) for all k \in [K^{\text{cap}}].
10:
              Let [\boldsymbol{\pi}]_{1:K^{\operatorname{cap}}} \leftarrow [F(\boldsymbol{\pi}')]_{1:K^{\operatorname{cap}}}.
11:
              Let K_{tg} \leftarrow \max\{k: v_{tg} \leq \pi_k\}, \forall g \in [G_t^{\text{cap}}], t \in [T].
12:
              Let K_t \leftarrow \max_q K_{tq}, and K \leftarrow \max_t K_t.
13:
              Let \lambda_{k\ell} and N_{k\ell} be as in (3.13) and (3.14).
14:
              Sample \eta_{k\ell} \sim \text{Beta}\left(\lambda_{k\ell} + \alpha_{\eta}, \ N_{k\ell} - \lambda_{k\ell} + \beta_{\eta}\right) \text{ for all } k, \ell \in [K^{\text{cap}}].
Let a_{k\ell} := a(k,\ell) \leftarrow \log \frac{\eta_{k\ell}}{1 - \eta_{k\ell}} \text{ and } b_{k\ell} := b(k,\ell) \leftarrow \log(1 - \eta_{k\ell}).
15:
16:
17:
              for t = 1, \ldots, T do
18:
                     Let \zeta_{tg\ell} and R_{tg\ell} be as in (3.11) and (3.12).
                     Compute \xi_{tgg} and O_{tgg} by setting g' = g in (3.6). Sample, for all g \in [G_t^{\text{cap}}],
19:
20:
                              k_{tg} \sim \left(\exp\left[a_{kk}\,\xi_{tgg} + b_{kk}O_{tgg} + \sum_{e} a_{k\ell}\,\zeta_{tg\ell} + b_{k\ell}\,R_{tg\ell}\right]\right)_{k \in [K_{tg}]}.
                     Sample v_{tg} \sim \text{Unif}(0, \pi_{k_{tg}}) \text{ for all } g \in [G_t^{\text{cap}}].
21:
22:
                     Let \tau_{ti\ell} and m_{ti\ell} be as in (3.5).
                     Sample g_{ti} \sim \left( \exp \left[ \sum_{\ell} a(k_{tg}, \ell) \tau_{ti\ell} + b(k_{tg}, \ell) m_{ti\ell} \right] \right)_{g \in [G_{ti}]} for all i \in [n_t].
23:
                     Sample u_{ti} \sim \text{Unif}(0, \gamma_{t,q_{ti}}) for all i \in [n_t].
24:
                     Set z_{ti} \leftarrow k_{t,q_{ti}}.
25:
              end for
26:
27: end while
```

are
$$p(u_i \mid \boldsymbol{\gamma}', \Theta_{-\boldsymbol{u}\boldsymbol{\gamma}'}) \propto 1\{u_i \leq \gamma_{z_i}\}$$
, that is $u_i \mid \boldsymbol{\gamma}', \Theta_{-\boldsymbol{u}\boldsymbol{\gamma}'} \sim \text{Unif}(0, \gamma_{z_i}).$

Next we sample from $(\gamma' \mid \Theta_{-u\gamma'})$. To do this, we marginalize out u from the joint distribution leading to the usual block sampler updates: The posterior factorizes over z and we have

(A.3)
$$\gamma'_z \mid \boldsymbol{\gamma}'_{-z}, \Theta_{-\boldsymbol{u}\boldsymbol{\gamma}'} \sim \operatorname{Beta}\left(n_z(\boldsymbol{z}) + 1, n_{>z}(\boldsymbol{z}) + \alpha_0\right).$$

Note that here $z \in \mathbb{N}$ is a dummy variable indexing $\gamma' = (\gamma'_z)$, whereas $z = (z_i)$ is the current vector of node labels.

Sampling $z \mid \cdots$. This posterior does not factorize but from (3.1), we have

(A.4)
$$\mathbb{P}(z_i = z \mid \boldsymbol{z}_{-i}, \Theta_{-\boldsymbol{z}}) \propto \prod_{j \in [n] \setminus \{i\}} L(\eta(z, z_j); A_{ij}) 1\{u_i \leq \gamma_z\}.$$

Algorithm 2 Slice sampler for single-layer DP-SBM

- 1: Initialize Z^{cap} to a pre-specified value (say 10).
- 2: Initialize z to the all-ones vector.
- 3: Initialize (u_i) independent uniform variables.
- while not CONVERGED, nor maximum iterations reached do
- Sample $\gamma'_z \sim \text{Beta}\left(n_z(z) + 1, n_{>z}(z) + \alpha_0\right)$ for all $z \in [Z^{\text{cap}}]$. 5:
- Let $[\gamma]_{1:Z^{\operatorname{cap}}} \leftarrow [F(\gamma')]_{1:Z^{\operatorname{cap}}}$. 6:
- Let $Z_i \leftarrow \max\{z : u_i \leq \gamma_z\}, \forall i \in [n] \text{ and } Z \leftarrow \max_{i=1,\dots,n} Z_i.$ 7:
- $\operatorname{doubling}(Z, Z^{\operatorname{cap}})$ 8:
- Sample $\eta_{k\ell} \sim \text{Beta}\left(\lambda_{k\ell} + \alpha_{\eta}, N_{k\ell} \lambda_{k\ell} + \beta_{\eta}\right)$ for all $k, \ell \in [Z^{\text{cap}}]$. Let $a_{k\ell} \leftarrow \log \frac{\eta_{k\ell}}{1 \eta_{k\ell}}$ and $b_{k\ell} \leftarrow \log(1 \eta_{k\ell})$. 9:
- 10:
- Let $\tau_{i\ell}$ and $m_{i\ell}$ be as in (A.5). 11:
- Sample $z_i \sim \left(\exp \left[\sum_{\ell} a_{z\ell} \tau_{i\ell} + b_{z\ell} m_{i\ell} \right] \right)_{z \in [Z_i]}$ for all $i \in [n]$. 12:
- Sample $u_i \sim \text{Unif}(0, \gamma_{z_i})$ for all $i \in [n]$. 13:
- 14: end while
- 15: macro doubling (K, K^{cap})
- if $K < K^{\text{cap}}$, then continue else $K^{\text{cap}} \leftarrow 1.5K^{\text{cap}}$ and go to the previous iteration. 16:

Let $Z_i := Z(z; u_i) := \sup\{z : u_i \leq \gamma_z\}$. According to the above z_i given everything else will be distributed as

$$z_i \mid \cdots \sim (\rho_i(z))_{z \in [Z_i]},$$

where we have defined

$$\rho_{i}(g) := \prod_{j \in [n] \setminus \{i\}} L(\eta(z, z_{j}); A_{ij})
= \prod_{j \in [n] \setminus \{i\}} \prod_{\ell} \left[L(\eta(z, \ell); A_{ij}) \right]^{1\{z_{j} = \ell\}}
= \prod_{j \in [n] \setminus \{i\}} \prod_{\ell} \left[\eta_{z\ell}^{A_{ij}} [1 - \eta_{z\ell}]^{1 - A_{ij}} \right]^{1\{z_{j} = \ell\}} = \prod_{\ell} \eta_{z\ell}^{\tau_{i\ell}} [1 - \eta_{z\ell}]^{m_{i\ell} - \tau_{i\ell}}$$

where

(A.5)
$$\tau_{i\ell} := \sum_{j \in [n] \setminus \{i\}} A_{ij} 1\{z_j = \ell\}, \quad m_{i\ell} := \sum_{j \in [n] \setminus \{i\}} 1\{z_j = \ell\}.$$

Sampling $\eta \mid \cdots$. These updates are exactly those of the multilayer case with T=1. The slice-sampler for the single-layer DP-SBM is summarized in Algorithm 2.

References

- [1] Arash A. Amini, Marina Paez, Lizhen Lin, and Zahra S. Razaee. Exact slice sampler for Hierarchical Dirichlet Processes. arXiv e-prints, page arXiv:1903.08829, Mar 2019.
- [2] Federico Battiston, Vincenzo Nicosia, and Vito Latora. Structural measures for multiplex networks. Phys. Rev. E, 89:032804, Mar 2014.

- [3] M. Berlingerio, M. Coscia, and F. Giannotti. Finding and characterizing communities in multidimensional networks. In 2011 International Conference on Advances in Social Networks Analysis and Mining, pages 490–494, July 2011.
- [4] Michele Berlingerio, Michele Coscia, Fosca Giannotti, Anna Monreale, and Dino Pedreschi. Evolving networks: Eras and turning points. *Intell. Data Anal.*, 17(1):27–48, January 2013.
- [5] Sharmodeep Bhattacharyya and Shirshendu Chatterjee. Spectral clustering for multiple sparse networks: I. arXiv preprint arXiv:1805.10594, 2018.
- [6] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008(10):P10008, 2008.
- [7] P. Brdka, K. Skibicki, P. Kazienko, and K. Musia. A degree centrality in multi-layered social network. In 2011 International Conference on Computational Aspects of Social Networks (CASoN), pages 237–242, Oct 2011.
- [8] H. Bush, P. West, and L. Michell. The role of friendship groups in the uptake and maintenance of smoking amongst pre-adolescent and adolescent children: Distribution of frequencies. MRC Medical Sociology Unit Glasgow, (62), 1997.
- [9] Kathleen M. Carley, Michael K. Martin, and Brian R. Hirshman. The etiology of social change. *Topics in Cognitive Science*, 1(4):621–650, 2009.
- [10] Emanuele Cozzo, Mikko Kivelä, Manlio De Domenico, Albert Solé, Alex Arenas, Sergio Gómez, Mason A Porter, and Yamir Moreno. Clustering coefficients in multiplex networks. arXiv preprint arXiv:1307.6780, 2013.
- [11] Gabor Csardi and Tamas Nepusz. The igraph software package for complex network research. *InterJournal, Complex Systems*, 1695(5):1–9, 2006.
- [12] Caterina De Bacco, Eleanor A Power, Daniel B Larremore, and Cristopher Moore. Community detection, link prediction, and layer interdependence in multilayer networks. *Physical Review E*, 95(4):042317, 2017.
- [13] Manlio De Domenico, Shuntaro Sasai, and Alex Arenas. Mapping multiplex hubs in human functional brain networks. *Frontiers in Neuroscience*, 10:326, 2016.
- [14] Gilles Didier, Christine Brun, and Anaïs Baudot. Identifying communities from multiplex biological networks. 3:e1525, 12 2015.
- [15] Thomas S. Ferguson. A bayesian analysis of some nonparametric problems. *Ann. Statist.*, 1(2):209–230, 03 1973.
- [16] Isabella Gollini and Thomas Brendan Murphy. Joint modeling of multiple network views. *Journal of Computational and Graphical Statistics*, 25(1):246–265, 2016.
- [17] Derek Greene and Pádraig Cunningham. Producing a unified graph representation from multiple social network views. In *Proceedings of the 5th Annual ACM Web Science Conference*, WebSci '13, pages 118–121, New York, NY, USA, 2013. ACM.
- [18] Manel Hmimida and Rushed Kanawati. Community detection in multiplex networks: A seed-centric approach. *Networks and Heterogeneous Media*, 10(1):71–85, 2015.
- [19] Hemant Ishwaran and Lancelot F. James. Gibbs sampling methods for stick-breaking priors. *Journal of the American Statistical Association*, 96(453):161–173, 3 2001.
- [20] Maria Kalli, Jim E. Griffin, and Stephen G. Walker. Slice sampling mixture models. *Statistics and Computing*, 21(1):93–105, 2011.
- [21] V. Krebs. Fortune 500 teams. http://moreno.ss.uci.edu/data.html#krebs. Accessed: 03-09-2019.
- [22] Zhana Kuncheva and Giovanni Montana. Community detection in multiplex networks using locally adaptive random walks. In *Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015*, ASONAM '15, pages 1308–1315, New York, NY, USA, 2015. ACM.
- [23] A. Majdandzic, B. Podobnik, S. V. Buldyrev, D. Y. Kenett, S. Havlin, and H. Eugene Stanley. Spontaneous recovery in dynamical networks. *Nature Physics*, 10:34–38, January 2014.

- [24] L. Michell and P. West. Peer pressure to smoke: the meaning depends on the method. *Health Education Research*, 11(1):39–49, 1996.
- [25] Morten Mørup and Mikkel N. Schmidt. Bayesian community detection. *Neural Comput.*, 24(9):2434–2456, September 2012.
- [26] Peter J Mucha, Thomas Richardson, Kevin Macon, Mason A Porter, and Jukka-Pekka Onnela. Community structure in time-dependent, multiscale, and multiplex networks. *Science*, 328(5980):876–878, 2010.
- [27] Marina Paez, Arash A. Amini, and Lizhen Lin. Hierarchical stochastic block model. https://github.com/aaamini/hsbm.
- [28] M. Pearson and L. Michell. Smoke rings: Social network analysis of friendship groups, smoking and drug-taking. *Drugs: education, prevention and policy*, 7:21–37, 2000.
- [29] M. Pearson and P. West. Drifting smoke rings: Social network analysis and markov processes in a longitudinal study of friendship groups and risk-taking. *Connections*, 25(2):59–76, 2003.
- [30] J. Picard and J. Pitman. Combinatorial Stochastic Processes: Ecole d'Eté de Probabilités de Saint-Flour XXXII - 2002. Lecture Notes in Mathematics. Springer Berlin Heidelberg, 2006.
- [31] Sebastian Poledna, José Luis Molina-Borboa, Serafín Martínez-Jaramillo, Marco Van Der Leij, and Stefan Thurner. The multi-layer network nature of systemic risk and its implications for the costs of financial crises. *Journal of Financial Stability*, 20:70–81, 2015.
- [32] A. E. Raftery, M. S. Handcock, and P. D. Hoff. Latent space approaches to social network analysis. Journal of the American Statistical Association, 15:460, 2002.
- [33] Michael Salter-Townshend and Tyler H. McCormick. Latent space models for multiview network data. *Ann. Appl. Stat.*, 11(3):1217–1244, 09 2017.
- [34] P. Sarkar and A. Moore. Dynamic social network analysis using latent space models. In *Advances in Neural Information Processing Systems*. 2005.
- [35] J. Sethuraman. A constructive definition of Dirichlet priors. Statistica Sinica, 4:639–650, 1994.
- [36] Yee Whye Teh, Michael I Jordan, Matthew J Beal, and David M Blei. Hierarchical dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- [37] Jianxin Wang, Xiaoqing Peng, Wei Peng, and Fang-Xiang Wu. Dynamic protein interaction network construction and applications. *PROTEOMICS*, 14(4-5):338–352, 2014.
- [38] James D. Wilson, John Palowitch, Shankar Bhamidi, and Andrew B. Nobel. Community extraction in multilayer networks with heterogeneous community structure. *J. Mach. Learn. Res.*, 18(1):5458–5506, January 2017.

DEPARTMENT OF STATISTICS, FEDERAL UNIVERSITY OF RIO DE JANEIRO, BRAZIL

 $E ext{-}mail\ address: marina@im.ufrj.br}$

DEPARTMENT OF STATISTICS, THE UNIVERSITY OF CALIFORNIA AT LOS ANGELES, LOS ANGELES, USA

E-mail address: aaamini@ucla.edu

DEPARTMENT OF APPLIED AND COMPUTATIONAL MATHEMATICS AND STATISTICS, THE UNIVERSITY OF NOTRE DAME, NOTRE DAME, USA

E-mail address: lizhen.lin@nd.edu