
Bayesian latent structure discovery from multi-neuron recordings

Scott W. Linderman
Columbia University
swl2133@columbia.edu

Ryan P. Adams
Harvard and Twitter Cortex
rpa@seas.harvard.edu

Jonathan W. Pillow
Princeton University
pillow@princeton.edu

Abstract

Neural circuits contain heterogeneous groups of neurons that differ in type, location, connectivity, and basic response properties. However, traditional methods for dimensionality reduction and clustering are ill-suited to recovering the structure underlying the organization of neural circuits. In particular, they do not take advantage of the rich temporal dependencies in multi-neuron recordings and fail to account for the noise in neural spike trains. Here we describe new tools for inferring latent structure from simultaneously recorded spike train data using a hierarchical extension of a multi-neuron point process model commonly known as the generalized linear model (GLM). Our approach combines the GLM with flexible graph-theoretic priors governing the relationship between latent features and neural connectivity patterns. Fully Bayesian inference via Pólya-gamma augmentation of the resulting model allows us to classify neurons and infer latent dimensions of circuit organization from correlated spike trains. We demonstrate the effectiveness of our method with applications to synthetic data and multi-neuron recordings in primate retina, revealing latent patterns of neural types and locations from spike trains alone.

1 Introduction

Large-scale recording technologies are revolutionizing the field of neuroscience [e.g., 1, 5, 16]. These advances present an unprecedented opportunity to probe the underpinnings of neural computation, but they also pose an extraordinary statistical and computational challenge: how do we make sense of these complex recordings? To address this challenge, we need methods that not only capture variability in neural activity and make accurate predictions, but also expose meaningful structure that may lead to novel hypotheses and interpretations of the circuits under study. In short, we need exploratory methods that yield interpretable representations of large-scale neural data.

For example, consider a population of distinct retinal ganglion cells (RGCs). These cells only respond to light within their small receptive field. Moreover, decades of painstaking work have revealed a plethora of RGC types [17]. Thus, it is natural to characterize these cells in terms of their type and the location of their receptive field center. Rather than manually searching for such a representation by probing with different visual stimuli, here we develop a method to automatically discover this structure from correlated patterns of neural activity.

Our approach combines latent variable network models [6, 11] with generalized linear models of neural spike trains [12, 20, 14, 21] in a hierarchical Bayesian framework. The network serves as a bridge, connecting interpretable latent features of interest to the temporal dynamics of neural spike trains. Unlike many previous studies [e.g., 2, 3, 18], our goal here is not necessarily to recover true synaptic connectivity, nor is our primary emphasis on prediction. Instead, our aim is to explore and compare latent patterns of functional organization, integrating over possible networks. To do so, we develop an efficient Markov chain Monte Carlo (MCMC) inference algorithm by leveraging Pólya-gamma augmentation to derive collapsed Gibbs updates for the network. We illustrate the

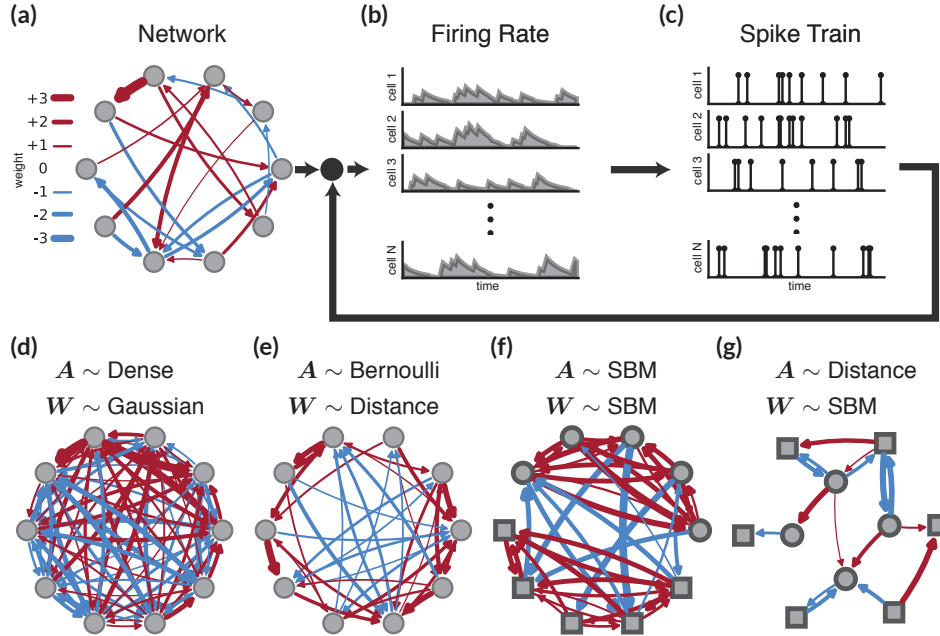


Figure 1: Components of the generative model. **(a)** Neurons influence one another via a sparse weighted network of interactions. **(b)** The network parameterizes an autoregressive model with a time-varying activation. **(c)** Spike counts are randomly drawn from a discrete distribution with a logistic link function. Each spike induces an impulse response on the activation of downstream neurons. **(d)** Standard GLM analyses correspond to a fully-connected network with Gaussian or Laplace distributed weights, depending on the regularization. **(e-g)** In this work, we consider structured models like the stochastic block model (SBM), in which neurons have discrete latent types (e.g. *square* or *circle*), and the latent distance model, in which neurons have latent locations that determine their probability of connection, capturing intuitive and interpretable patterns of connectivity.

robustness and scalability of our algorithm with synthetic data examples, and we demonstrate the scientific potential of our approach with an application to retinal ganglion cell recordings, where we recover the true underlying cell types and locations from spike trains alone, without reference to the stimulus.

2 Probabilistic Model

Figure 1 illustrates the components of our framework. We begin with a prior distribution on networks that generates a set of weighted connections between neurons (Fig. 1a). A directed edge indicates a functional relationship between the spikes of one neuron and the activation of its downstream neighbor. Each spike induces a weighted impulse response on the activation of the downstream neuron (Fig. 1b). The activation is converted into a nonnegative firing rate from which spikes are stochastically sampled (Fig. 1c). These spikes then feed back into the subsequent activation, completing an autoregressive loop, the hallmark of the GLM [12, 20]. Models like these have provided valuable insight into complex population recordings [14]. We detail the three components of this model in the reverse order, working backward from the observed spike counts through the activation to the underlying network.

2.1 Logistic Spike Count Models

Generalized linear models assume a stochastic spike generation mechanism. Consider a matrix of spike counts, $\mathbf{S} \in \mathbb{N}^{T \times N}$, for T time bins and N neurons. The expected number of spikes fired by the n -th neuron in the t -th time bin, $\mathbb{E}[s_{t,n}]$, is modeled as a nonlinear function of the instantaneous activation, $\psi_{t,n}$, and a static, neuron-specific parameter, ν_n . Table 1 enumerates the three spike count models considered in this paper, all of which use the logistic function, $\sigma(\psi) = e^\psi / (1 + e^\psi)$, to rectify the activation. The Bernoulli distribution is appropriate for binary spike counts, whereas the binomial and negative binomial have support for $s \in [0, \nu]$ and $s \in [0, \infty)$, respectively. Notably

Distribution	$p(s \psi, \nu)$	Standard Form	$\mathbb{E}[s]$	$\text{Var}(s)$
Bern($\sigma(\psi)$)	$\sigma(\psi)^s \sigma(-\psi)^{1-s}$	$\frac{(e^\psi)^s}{1+e^\psi}$	$\sigma(\psi)$	$\sigma(\psi) \sigma(-\psi)$
Bin($\nu, \sigma(\psi)$)	$\binom{\nu}{s} \sigma(\psi)^s \sigma(-\psi)^{\nu-s}$	$\binom{\nu}{s} \frac{(e^\psi)^s}{(1+e^\psi)^\nu}$	$\nu \sigma(\psi)$	$\nu \sigma(\psi) \sigma(-\psi)$
NB($\nu, \sigma(\psi)$)	$\binom{\nu+s-1}{s} \sigma(\psi)^s \sigma(-\psi)^\nu$	$\binom{\nu+s-1}{s} \frac{(e^\psi)^s}{(1+e^\psi)^{\nu+s}}$	νe^ψ	$\nu e^\psi / \sigma(-\psi)$

Table 1: Table of conditional spike count distributions, their parameterizations, and their properties.

lacking from this list is the Poisson distribution, which is not directly amenable to the augmentation schemes we derive below; however, both the binomial and negative binomial distributions converge to the Poisson under certain limits. Moreover, these distributions afford the added flexibility of modeling under- and over-dispersed spike counts, a biologically significant feature of neural spiking data [4]. Specifically, while the Poisson has unit dispersion (its mean is equal to its variance), the binomial distribution is always under-dispersed, since its mean always exceeds its variance, and the negative binomial is always over-dispersed, with variance greater than its mean.

Importantly, all of these distributions can be written in a standard form, as shown in Table 1. We exploit this fact to develop an efficient Markov chain Monte Carlo (MCMC) inference algorithm described in Section 3.

2.2 Linear Activation Model

The instantaneous activation of neuron n at time t is modeled as a linear, autoregressive function of preceding spike counts of neighboring neurons,

$$\psi_{t,n} \triangleq b_n + \sum_{m=1}^N \sum_{\Delta t=1}^{\Delta t_{\max}} h_{m \rightarrow n}[\Delta t] \cdot s_{t-\Delta t,m}, \quad (1)$$

where b_n is the baseline activation of neuron n and $h_{m \rightarrow n} : \{1, \dots, \Delta t_{\max}\} \rightarrow \mathbb{R}$ is an impulse response function that models the influence spikes on neuron m have on the activation of neuron n at a delay of Δt . To model the impulse response, we use a spike-and-slab formulation [8],

$$h_{m \rightarrow n}[\Delta t] = a_{m \rightarrow n} \sum_{k=1}^K w_{m \rightarrow n}^{(k)} \phi_k[\Delta t]. \quad (2)$$

Here, $a_{m \rightarrow n} \in \{0, 1\}$ is a binary variable indicating the presence or absence of a connection from neuron m to neuron n , the weight $\mathbf{w}_{m \rightarrow n} = [w_{m \rightarrow n}^{(1)}, \dots, w_{m \rightarrow n}^{(K)}]$ denotes the strength of the connection, and $\{\phi_k\}_{k=1}^K$ is a collection of fixed basis functions. In this paper, we consider scalar weights ($K = 1$) and use an exponential basis function, $\phi_1[\Delta t] = e^{-\Delta t/\tau}$, with time constant of $\tau = 15\text{ms}$. Since the basis function and the spike train are fixed, we precompute the convolution of the spike train and the basis function to obtain $\hat{s}_{t,m}^{(k)} = \sum_{\Delta t=1}^{\Delta t_{\max}} \phi_k[\Delta t] \cdot s_{t-\Delta t,m}$. Finally, we combine the connections, weights, and filtered spike trains and write the activation as,

$$\psi_{t,n} = (\mathbf{a}_n \odot \mathbf{w}_n)^\top \hat{\mathbf{s}}_t, \quad (3)$$

where $\mathbf{a}_n = [1, a_{1 \rightarrow n} \mathbf{1}_K, \dots, a_{N \rightarrow n} \mathbf{1}_K]$, $\mathbf{w}_n = [b_n, \mathbf{w}_{1 \rightarrow n}, \dots, \mathbf{w}_{N \rightarrow n}]$, and $\hat{\mathbf{s}}_t = [1, \hat{s}_{t,1}^{(1)}, \dots, \hat{s}_{t,N}^{(K)}]$. Here, \odot denotes the Hadamard (elementwise) product and $\mathbf{1}_K$ is length- K vector of ones. Hence, all of these vectors are of size $1 + NK$. The difference between our formulation and the standard GLM is that we have explicitly modeled the sparsity of the weights in $a_{m \rightarrow n}$. In typical formulations [e.g., 14], all connections are present and the weights are regularized with ℓ_1 and ℓ_2 penalties to promote sparsity. Instead, we consider structured approaches to modeling the sparsity and weights.

2.3 Random Network Models

Patterns of functional interaction can provide great insight into the computations performed by neural circuits. Indeed, many circuits are informally described in terms of ‘‘types’’ of neurons that perform a particular role, or the ‘‘features’’ that neurons encode. Random network models formalize these intuitive descriptions. Types and features correspond to latent variables in a probabilistic model that governs how likely neurons are to connect and how strongly they influence each other.

Name	$\rho(\mathbf{u}_m, \mathbf{u}_n, \boldsymbol{\theta})$	$\boldsymbol{\mu}(\mathbf{v}_m, \mathbf{v}_n, \boldsymbol{\theta})$	$\boldsymbol{\Sigma}(\mathbf{v}_m, \mathbf{v}_n, \boldsymbol{\theta})$
Dense Model	1	$\bar{\boldsymbol{\mu}}$	$\bar{\boldsymbol{\Sigma}}$
Independent Model	$\bar{\rho}$	$\bar{\boldsymbol{\mu}}$	$\bar{\boldsymbol{\Sigma}}$
Stochastic Block Model	$\rho_{\mathbf{u}_m \rightarrow \mathbf{u}_n}$	$\boldsymbol{\mu}_{\mathbf{v}_m \rightarrow \mathbf{v}_n}$	$\boldsymbol{\Sigma}_{\mathbf{v}_m \rightarrow \mathbf{v}_n}$
Latent Distance Model	$\sigma(-\ \mathbf{u}_n - \mathbf{v}_m\ _2 + \gamma_0)$	$-\ \mathbf{v}_n - \mathbf{v}_m\ _2 + \mu_0$	η^2

Table 2: Random network models for the binary adjacency matrix or the Gaussian weight matrix.

Let $\mathbf{A} = \{\{a_{m \rightarrow n}\}\}$ and $\mathbf{W} = \{\{\mathbf{w}_{m \rightarrow n}\}\}$ denote the binary adjacency matrix and the real-valued array of weights, respectively. Now suppose $\{\mathbf{u}_n\}_{n=1}^N$ and $\{\mathbf{v}_n\}_{n=1}^N$ are sets of neuron-specific latent variables that govern the distributions over \mathbf{A} and \mathbf{W} . Given these latent variables and global parameters $\boldsymbol{\theta}$, the entries in \mathbf{A} are conditionally independent Bernoulli random variables, and the entries in \mathbf{W} are conditionally independent Gaussians. That is,

$$p(\mathbf{A}, \mathbf{W} \mid \{\mathbf{u}_n, \mathbf{v}_n\}_{n=1}^N, \boldsymbol{\theta}) = \prod_{m=1}^N \prod_{n=1}^N \text{Bern}(a_{m \rightarrow n} \mid \rho(\mathbf{u}_m, \mathbf{u}_n, \boldsymbol{\theta})) \times \mathcal{N}(\mathbf{w}_{m \rightarrow n} \mid \boldsymbol{\mu}(\mathbf{v}_m, \mathbf{v}_n, \boldsymbol{\theta}), \boldsymbol{\Sigma}(\mathbf{v}_m, \mathbf{v}_n, \boldsymbol{\theta})), \quad (4)$$

where $\rho(\cdot)$, $\boldsymbol{\mu}(\cdot)$, and $\boldsymbol{\Sigma}(\cdot)$ are functions that output a probability, a mean vector, and a covariance matrix, respectively. We recover the standard GLM when $\rho(\cdot) \equiv 1$, but here we can take advantage of structured priors like the stochastic block model (SBM) [10], in which each neuron has a discrete type, and the latent distance model [6], in which each neuron has a latent location. Table 2 outlines the various models considered in this paper.

We can mix and match these models as shown in Figure 1(d-g). For example, in Fig. 1g, the adjacency matrix is distance-dependent and the weights are block structured. Thus, we have a flexible language for expressing hypotheses about patterns of interaction. In fact, the simple models enumerated above are instances of a rich family of exchangeable networks known as Aldous-Hoover random graphs, which have been recently reviewed by Orbanz and Roy [11].

3 Bayesian Inference

Generalized linear models are often fit via maximum *a posteriori* (MAP) estimation [12, 20, 14, 21]. However, as we scale to larger populations of neurons, there will inevitably be structure in the posterior that is not reflected with a point estimate. Technological advances are expanding the number of neurons that can be recorded simultaneously, but “high-throughput” recording of many individuals is still a distant hope. Therefore we expect the complexities of our models to expand faster than the available distinct data sets to fit them. In this situation, accurately capturing uncertainty is critical. Moreover, in the Bayesian framework, we also have a coherent way to perform model selection and evaluate hypotheses regarding complex underlying structure. Finally, after introducing a binary adjacency matrix and hierarchical network priors, the log posterior is no longer a concave function of model parameters, making direct optimization challenging (though see Soudry et al. [18] for recent advances in tackling similar problems). These considerations motivate a fully Bayesian approach.

Computation in rich Bayesian models is often challenging, but through thoughtful modeling decisions it is sometimes possible to find representations that lead to efficient inference. In this case, we have carefully chosen the logistic models of the preceding section in order to make it possible to apply the Pólya-gamma augmentation scheme [15]. The principal advantage of this approach is that, given the Pólya-gamma auxiliary variables, the conditional distribution of the weights is Gaussian, and hence is amenable to efficient Gibbs sampling. Recently, Pillow and Scott [13] used this technique to develop inference algorithms for negative binomial factor analysis models of neural spike trains. We build on this work and show how this conditionally Gaussian structure can be exploited to derive efficient, collapsed Gibbs updates.

3.1 Collapsed Gibbs updates for Gaussian observations

Suppose the observations were actually Gaussian distributed, i.e. $s_{t,n} \sim \mathcal{N}(\psi_{t,n}, \nu_n)$. The most challenging aspect of inference is then sampling the posterior distribution over discrete connections, \mathbf{A} . There may be many posterior modes corresponding to different patterns of connectivity.

Moreover, $a_{m \rightarrow n}$ and $w_{m \rightarrow n}$ are often highly correlated, which leads to poor mixing of naïve Gibbs sampling. Fortunately, when the observations are Gaussian, we may integrate over possible weights and sample the binary adjacency matrix from its collapsed conditional distribution.

We combine the conditionally independent Gaussian priors on $\{w_{m \rightarrow n}\}$ and b_n into a joint Gaussian distribution, $\mathbf{w}_n | \{\mathbf{v}_n\}, \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{w}_n | \boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$, where $\boldsymbol{\Sigma}_n$ is a block diagonal covariance matrix. Since $\psi_{t,n}$ is linear in \mathbf{w}_n (see Eq. 3), a Gaussian likelihood is conjugate with this Gaussian prior, given \mathbf{a}_n and $\widehat{\mathbf{S}} = \{\widehat{s}_t\}_{t=1}^T$. This yields the following closed-form conditional:

$$p(\mathbf{w}_n | \widehat{\mathbf{S}}, \mathbf{a}_n, \boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n) \propto \mathcal{N}(\mathbf{w}_n | \boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n) \prod_{t=1}^T \mathcal{N}(s_{t,n} | (\mathbf{a}_n \odot \mathbf{w}_n)^\top \widehat{\mathbf{s}}_t, \nu_n) \propto \mathcal{N}(\mathbf{w}_n | \widetilde{\boldsymbol{\mu}}_n, \widetilde{\boldsymbol{\Sigma}}_n),$$

$$\widetilde{\boldsymbol{\Sigma}}_n = \left[\boldsymbol{\Sigma}_n^{-1} + \left(\widehat{\mathbf{S}}^\top (\nu_n^{-1} \mathbf{I}) \widehat{\mathbf{S}} \right) \odot (\mathbf{a}_n \mathbf{a}_n^\top) \right]^{-1}, \quad \widetilde{\boldsymbol{\mu}}_n = \widetilde{\boldsymbol{\Sigma}}_n \left[\boldsymbol{\Sigma}_n^{-1} \boldsymbol{\mu}_n + \left(\widehat{\mathbf{S}}^\top (\nu_n^{-1} \mathbf{I}) \mathbf{s}_{:,n} \right) \odot \mathbf{a}_n \right].$$

Now, consider the conditional distribution of \mathbf{a}_n , integrating out the corresponding weights. The prior distribution over \mathbf{a}_n is a product of Bernoulli distributions with parameters $\boldsymbol{\rho}_n = \{\rho(\mathbf{u}_m, \mathbf{u}_n, \boldsymbol{\theta})\}_{m=1}^N$. The conditional distribution is proportional to the ratio of the prior and posterior partition functions,

$$p(\mathbf{a}_n | \widehat{\mathbf{S}}, \boldsymbol{\rho}_n, \boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n) = \int p(\mathbf{a}_n, \mathbf{w}_n | \widehat{\mathbf{S}}, \boldsymbol{\rho}_n, \boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n) d\mathbf{w}_n$$

$$= p(\mathbf{a}_n | \boldsymbol{\rho}_n) \frac{|\boldsymbol{\Sigma}_n|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \boldsymbol{\mu}_n^\top \boldsymbol{\Sigma}_n^{-1} \boldsymbol{\mu}_n \right\}}{|\widetilde{\boldsymbol{\Sigma}}_n|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \widetilde{\boldsymbol{\mu}}_n^\top \widetilde{\boldsymbol{\Sigma}}_n^{-1} \widetilde{\boldsymbol{\mu}}_n \right\}}.$$

Thus, we perform a joint update of \mathbf{a}_n and \mathbf{w}_n by collapsing out the weights to directly sample the binary entries of \mathbf{a}_n . We iterate over each entry, $a_{m \rightarrow n}$, and sample it from its conditional distribution given $\{a_{m' \rightarrow n}\}_{m' \neq m}$. Having sampled \mathbf{a}_n , we sample \mathbf{w}_n from its Gaussian conditional.

3.2 Pólya-gamma augmentation for discrete observations

Now, let us turn to the non-conjugate case of discrete count observations. The Pólya-gamma augmentation [15] introduces auxiliary variables, $\omega_{t,n}$, conditioned upon which the discrete likelihood appears Gaussian and our collapsed Gibbs updates apply. The integral identity underlying this scheme is

$$c \frac{(e^\psi)^a}{(1 + e^\psi)^b} = c 2^{-b} e^{\kappa \psi} \int_0^\infty e^{-\omega \psi^2 / 2} p_{\text{PG}}(\omega | b, 0) d\omega, \quad (5)$$

where $\kappa = a - b/2$ and $p(\omega | b, 0)$ is the density of the Pólya-gamma distribution $\text{PG}(b, 0)$, which does not depend on ψ . Notice that the discrete likelihoods in Table 1 can all be rewritten like the left-hand side of (5), for some a, b , and c that are functions of s and ν . Using (5) along with priors $p(\psi)$ and $p(\nu)$, we write the joint density of (ψ, s, ν) as

$$p(s, \nu, \psi) = \int_0^\infty p(\nu) p(\psi) c(s, \nu) 2^{-b(s, \nu)} e^{\kappa(s, \nu) \psi} e^{-\omega \psi^2 / 2} p_{\text{PG}}(\omega | b(s, \nu), 0) d\omega. \quad (6)$$

The integrand of Eq. 6 defines a joint density on (s, ν, ψ, ω) which admits $p(s, \nu, \psi)$ as a marginal density. Conditioned on the auxiliary variable, ω , the likelihood as a function of ψ is,

$$p(s | \psi, \nu, \omega) \propto e^{\kappa(s, \nu) \psi} e^{-\omega \psi^2 / 2} \propto \mathcal{N}(\omega^{-1} \kappa(s, \nu) | \psi, \omega^{-1}).$$

Thus, after conditioning on s, ν , and ω , we effectively have a linear Gaussian likelihood for ψ .

We apply this augmentation scheme to the full model, introducing auxiliary variables, $\omega_{t,n}$ for each spike count, $s_{t,n}$. Given these variables, the conditional distribution of \mathbf{w}_n can be computed in closed form, as before. Let $\boldsymbol{\kappa}_n = [\kappa(s_{1,n}, \nu_n), \dots, \kappa(s_{T,n}, \nu_n)]$ and $\boldsymbol{\Omega}_n = \text{diag}([\omega_{1,n}, \dots, \omega_{T,n}])$. Then we have $p(\mathbf{w}_n | s_n, \widehat{\mathbf{S}}, \mathbf{a}_n, \boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n, \boldsymbol{\omega}_n, \nu_n) \propto \mathcal{N}(\mathbf{w}_n | \widetilde{\boldsymbol{\mu}}_n, \widetilde{\boldsymbol{\Sigma}}_n)$, where

$$\widetilde{\boldsymbol{\Sigma}}_n = \left[\boldsymbol{\Sigma}_n^{-1} + \left(\widehat{\mathbf{S}}^\top \boldsymbol{\Omega}_n \widehat{\mathbf{S}} \right) \odot (\mathbf{a}_n \mathbf{a}_n^\top) \right]^{-1}, \quad \widetilde{\boldsymbol{\mu}}_n = \widetilde{\boldsymbol{\Sigma}}_n \left[\boldsymbol{\Sigma}_n^{-1} \boldsymbol{\mu}_n + \left(\widehat{\mathbf{S}}^\top \boldsymbol{\kappa}_n \right) \odot \mathbf{a}_n \right].$$

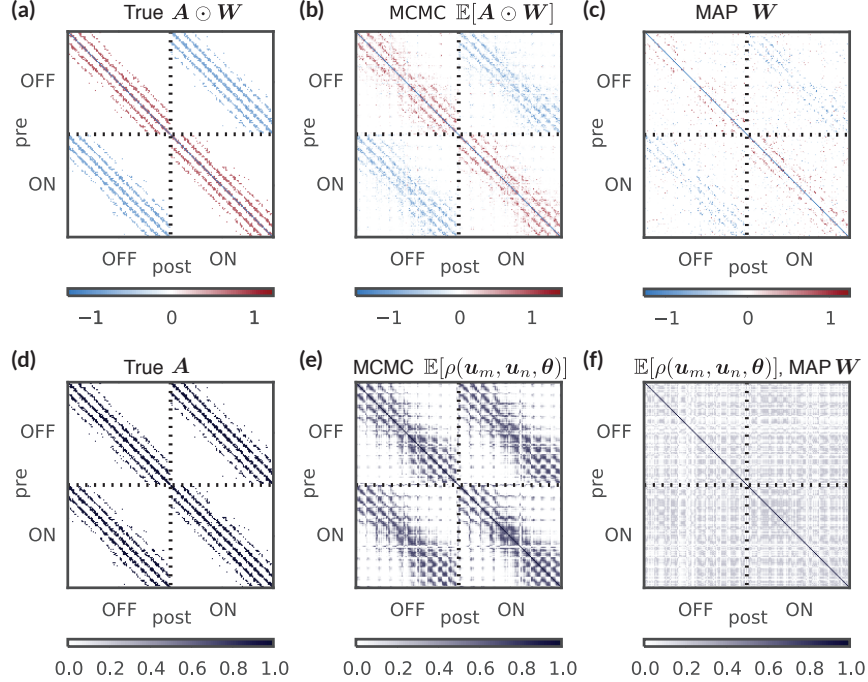


Figure 2: Weighted adjacency matrices showing inferred networks and connection probabilities for synthetic data. **(a,d)** True network. **(b,e)** Posterior mean using joint inference of network GLM. **(c,f)** MAP estimation.

Having introduced auxiliary variables, we must now derive Markov transitions to update them as well. Fortunately, the Pólya-gamma distribution is designed such that the conditional distribution of the auxiliary variables is simply a “tilted” Pólya-gamma distribution,

$$p(\omega_{t,n} | s_{t,n}, \nu_n, \psi_{t,n}) = p_{\text{PG}}(\omega_{t,n} | b(s_{t,n}, \nu_n), \psi_{t,n}).$$

These auxiliary variables are conditionally independent given the activation and hence can be sampled in parallel. Moreover, efficient algorithms are available to generate Pólya-gamma random variates [22]. Our Gibbs updates for the remaining parameters and latent variables (ν_n , \mathbf{u}_n , \mathbf{v}_n , and θ) are described in the supplementary material. A Python implementation of our inference algorithm is available at <https://github.com/slinderman/pyglm>.

4 Synthetic Data Experiments

The need for network models is most pressing in recordings of large populations where the network is difficult to estimate and even harder to interpret. To assess the robustness and scalability of our framework, we apply our methods to simulated data with known ground truth. We simulate a one minute recording (1ms time bins) from a population of 200 neurons with discrete latent types that govern the connection strength via a stochastic block model and continuous latent locations that govern connection probability via a latent distance model. The spikes are generated from a Bernoulli observation model.

First, we show that our approach of jointly inferring the network and its latent variables can provide dramatic improvements over alternative approaches. For comparison, consider the two-step procedure of Stevenson et al. [19] in which the network is fit with an ℓ_1 -regularized GLM and *then* a probabilistic network model is fit to the GLM connection weights. The advantage of this strategy is that the expensive GLM fitting is only performed once. However, when the data is limited, both the network and the latent variables are uncertain. Our Bayesian approach finds a very accurate network (Fig. 2b) by jointly sampling networks and latent variables. In contrast, the standard GLM does not account for latent structure and finds strong connections as well as spuriously correlated neurons (Fig. 2c). Moreover, our fully Bayesian approach finds a set of latent locations that mimics the true locations and therefore accurately estimates connection probability (Fig. 2e). In contrast, subsequently fitting a latent distance model to the adjacency matrix of a thresholded GLM network finds an embedding

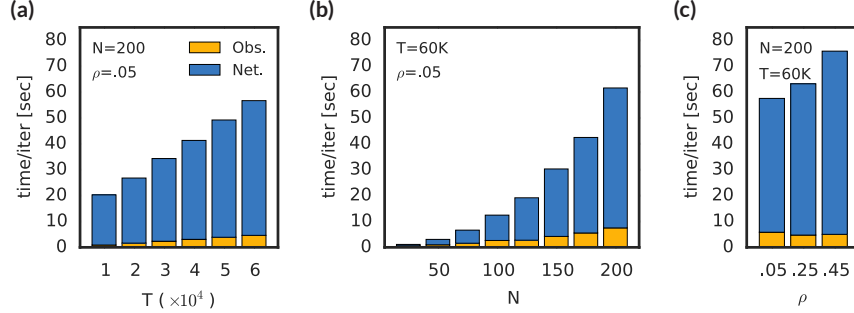


Figure 3: Scalability of our inference algorithm as a function of: **(a)** the number of time bins, T ; **(b)** the number of neurons, N ; and **(c)** the average sparsity of the network, ρ . Wall-clock time is divided into time spent sampling auxiliary variables (“Obs.”) and time spent sampling the network (“Net.”).

that has no resemblance to the true locations, which is reflected in its poor estimate of connection probability (Fig. 2f).

Next, we address the scalability of our MCMC algorithm. Three major parameters govern the complexity of inference: the number of time bins, T ; the number of neurons, N ; and the level of sparsity, ρ . The following experiments were run on a quad-core Intel i5 with 6GB of RAM. As shown in Fig. 3a, the wall clock time per iteration scales linearly with T since we must resample NT auxiliary variables. We scale at least quadratically with N due to the network, as shown in Fig. 3b. However, the total cost could actually be worse than quadratic since the cost of updating each connection could depend on N . Fortunately, the complexity of our collapsed Gibbs sampling algorithm only depends on the number of incident connections, d , or equivalently, the sparsity $\rho = d/N$. Specifically, we must solve a linear system of size d , which incurs a cubic cost, as seen in Fig. 3c.

5 Retinal Ganglion Cells

Finally, we demonstrate the efficacy of this approach with an application to spike trains simultaneously recorded from a population of 27 retinal ganglion cells (RGCs), which have previously been studied by Pillow et al. [14]. Retinal ganglion cells respond to light shown upon their receptive field. Thus, it is natural to characterize these cells by the location of their receptive field center. Moreover, retinal ganglion cells come in a variety of types [17]. This population is comprised of two types of cells, *on* and *off* cells, which are characterized by their response to visual stimuli. *On* cells increase their firing when light is shone upon their receptive field; *off* cells decrease their firing rate in response to light in their receptive field. In this case, the population is driven by a binary white noise stimulus. Given the stimulus, the cell locations and types are readily inferred. Here, we show how these intuitive representations can be discovered in a purely unsupervised manner given one minute of spiking data alone and no knowledge of the stimulus.

Figure 4 illustrates the results of our analysis. Since the data are binned at 1ms resolution, we have at most one spike per bin and we use a Bernoulli observation model. We fit the 12 network models of Table 2 (4 adjacency models and 3 weight models), and we find that, in terms of predictive log likelihood of held-out neurons, a latent distance model of the adjacency matrix and SBM of the weight matrix performs best (Fig. 4a). See the supplementary material for a detailed description of this comparison. Looking into the latent locations underlying the adjacency matrix our network GLM (NGLM), we find that the inferred distances between cells are highly correlated with the distances between the true locations. For comparison, we also fit a 2D Bernoulli linear dynamical system (LDS) — the Bernoulli equivalent of the Poisson LDS [7] — and we take rows of the $N \times 2$ emission matrix as locations. In contrast to our network GLM, the distances between LDS locations are nearly uncorrelated with the true distances (Fig. 4b) since the LDS does not capture the fact that distance only affects the probability of connection, not the weight. Not only are our distances accurate, the inferred locations are nearly identical to the true locations, up to affine transformation. In Fig. 4c, semitransparent markers show the inferred *on* cell locations, which have been rotated and scaled to best align with the true locations shown by the outlined marks. Based solely on patterns of correlated spiking, we have recovered the receptive field arrangements.

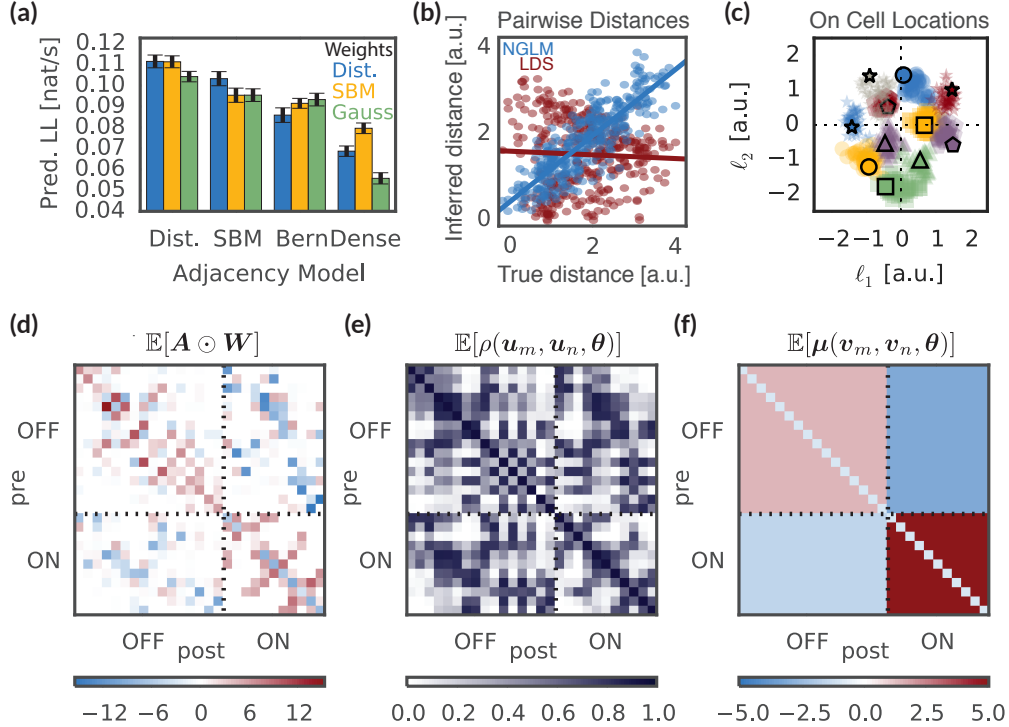


Figure 4: Using our framework, retinal ganglion cell types and locations can be inferred from spike trains alone. **(a)** Model comparison. **(b)** True and inferred distances between cells. **(c)** True and inferred cell locations. **(d-f)** Inferred network, connection probability, and mean weight, respectively. See main text for further details.

Fig. 4d shows the inferred network, $\mathbf{A} \odot \mathbf{W}$, under a latent distance model of connection probability and a stochastic block model for connection weight. The underlying connection probabilities from the distance model are shown in Fig. 4e. Finally, Fig. 4f shows that we have discovered not only the cell locations, but also their latent types. With an SBM, the mean weight is a function of latent type, and under the posterior, the neurons are clearly clustered into the two true types that exhibit the expected within-class excitation and between-class inhibition.

6 Conclusion

Our results with both synthetic and real neural data provide compelling evidence that our methods can find meaningful structure underlying neural spike trains. Given the extensive work on characterizing retinal ganglion cell responses, we have considerable evidence that the representation we learn from spike trains alone is indeed the optimal way to summarize this population of cells. This lends us confidence that we may trust the representations learned from spike trains recorded from more enigmatic brain areas as well. While we have omitted stimulus from our models and only used it for confirming types and locations, in practice we could incorporate it into our model and even capture type- and location-dependent patterns of stimulus dependence with our hierarchical approach. Likewise, the network GLM could be combined with the PLDS as in Vidne et al. [21] to capture sources of low dimensional, shared variability.

Latent functional networks underlying spike trains can provide unique insight into the structure of neural populations. Looking forward, methods that extract interpretable representations from complex neural data, like those developed here, will be key to capitalizing on the dramatic advances in neural recording technology. We have shown that networks provide a natural bridge to connect neural types and features to spike trains, and demonstrated promising results on both real and synthetic data.

Acknowledgments. We thank E. J. Chichilnisky, A. M. Litke, A. Sher and J. Shlens for retinal data. SWL is supported by the Simons Foundation SCGB-418011. RPA is supported by NSF IIS-1421780 and the Alfred P. Sloan Foundation. JWP was supported by grants from the McKnight Foundation, Simons Collaboration on the Global Brain (SCGB AWD1004351), NSF CAREER Award (IIS-1150186), and NIMH grant MH099611.

References

- [1] M. B. Ahrens, M. B. Orger, D. N. Robson, J. M. Li, and P. J. Keller. Whole-brain functional imaging at cellular resolution using light-sheet microscopy. *Nature methods*, 10(5):413–420, 2013.
- [2] D. R. Brillinger, H. L. Bryant Jr, and J. P. Segundo. Identification of synaptic interactions. *Biological Cybernetics*, 22(4):213–228, 1976.
- [3] F. Gerhard, T. Kispersky, G. J. Gutierrez, E. Marder, M. Kramer, and U. Eden. Successful reconstruction of a physiological circuit with known connectivity from spiking activity alone. *PLoS Computational Biology*, 9(7):e1003138, 2013.
- [4] R. L. Goris, J. A. Movshon, and E. P. Simoncelli. Partitioning neuronal variability. *Nature Neuroscience*, 17(6):858–865, 2014.
- [5] B. F. Grewe, D. Langer, H. Kasper, B. M. Kampa, and F. Helmchen. High-speed in vivo calcium imaging reveals neuronal network activity with near-millisecond precision. *Nature methods*, 7(5):399–405, 2010.
- [6] P. D. Hoff. Modeling homophily and stochastic equivalence in symmetric relational data. *Advances in Neural Information Processing Systems 20*, 20:1–8, 2008.
- [7] J. H. Macke, L. Buesing, J. P. Cunningham, M. Y. Byron, K. V. Shenoy, and M. Sahani. Empirical models of spiking in neural populations. In *Advances in neural information processing systems*, pages 1350–1358, 2011.
- [8] T. J. Mitchell and J. J. Beauchamp. Bayesian Variable Selection in Linear Regression. *Journal of the American Statistical Association*, 83(404):1023—1032, 1988.
- [9] R. M. Neal. MCMC using Hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo*, pages 113–162, 2010.
- [10] K. Nowicki and T. A. B. Snijders. Estimation and prediction for stochastic blockstructures. *Journal of the American Statistical Association*, 96(455):1077–1087, 2001.
- [11] P. Orbanz and D. M. Roy. Bayesian models of graphs, arrays and other exchangeable random structures. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 37(2):437–461, 2015.
- [12] L. Paninski. Maximum likelihood estimation of cascade point-process neural encoding models. *Network: Computation in Neural Systems*, 15(4):243–262, Jan. 2004.
- [13] J. W. Pillow and J. Scott. Fully bayesian inference for neural models with negative-binomial spiking. In F. Pereira, C. Burges, L. Bottou, and K. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1898–1906. 2012.
- [14] J. W. Pillow, J. Shlens, L. Paninski, A. Sher, A. M. Litke, E. Chichilnisky, and E. P. Simoncelli. Spatio-temporal correlations and visual signalling in a complete neuronal population. *Nature*, 454(7207):995–999, 2008.
- [15] N. G. Polson, J. G. Scott, and J. Windle. Bayesian inference for logistic models using Pólya–gamma latent variables. *Journal of the American Statistical Association*, 108(504):1339–1349, 2013.
- [16] R. Prevedel, Y.-G. Yoon, M. Hoffmann, N. Pak, G. Wetzstein, S. Kato, T. Schrödel, R. Raskar, M. Zimmer, E. S. Boyden, et al. Simultaneous whole-animal 3d imaging of neuronal activity using light-field microscopy. *Nature methods*, 11(7):727–730, 2014.
- [17] J. R. Sanes and R. H. Masland. The types of retinal ganglion cells: current status and implications for neuronal classification. *Annual review of neuroscience*, 38:221–246, 2015.
- [18] D. Soudry, S. Keshri, P. Stinson, M.-h. Oh, G. Iyengar, and L. Paninski. A shotgun sampling solution for the common input problem in neural connectivity inference. *arXiv preprint arXiv:1309.3724*, 2013.
- [19] I. H. Stevenson, J. M. Rebesco, N. G. Hatsopoulos, Z. Haga, L. E. Miller, and K. P. Körding. Bayesian inference of functional connectivity and network structure from spikes. *Neural Systems and Rehabilitation Engineering, IEEE Transactions on*, 17(3):203–213, 2009.
- [20] W. Truccolo, U. T. Eden, M. R. Fellows, J. P. Donoghue, and E. N. Brown. A point process framework for relating neural spiking activity to spiking history, neural ensemble, and extrinsic covariate effects. *Journal of Neurophysiology*, 93(2):1074–1089, 2005. doi: 10.1152/jn.00697.2004.

- [21] M. Vidne, Y. Ahmadian, J. Shlens, J. W. Pillow, J. Kulkarni, A. M. Litke, E. Chichilnisky, E. Simoncelli, and L. Paninski. Modeling the impact of common noise inputs on the network activity of retinal ganglion cells. *Journal of computational neuroscience*, 33(1):97–121, 2012.
- [22] J. Windle, N. G. Polson, and J. G. Scott. Sampling Pólya-gamma random variates: alternate and approximate techniques. *arXiv preprint arXiv:1405.0506*, 2014.
- [23] M. Zhou, L. Li, L. Carin, and D. B. Dunson. Lognormal and gamma mixed negative binomial regression. In *Proceedings of the 29th International Conference on Machine Learning*, pages 1343–1350, 2012.

A Model comparison via Predictive Log Likelihood

How can we compare different network models in a principled manner? Predictive log likelihood of held-out time bins is insufficient since it only depends directly on \mathbf{A} and \mathbf{W} . The network prior does aid in the estimation of \mathbf{A} and \mathbf{W} , but this is only an indirect effect, and it may be small relative to the effect of the data. Instead, we hold out neurons rather than time bins. Accurate network models play a crucial role in predicting held-out neurons’ activity, since the distribution of incident connections to the held-out neuron are informed solely by the network model.

Formally, we estimate the probability of a held-out neuron’s spike train $\mathbf{s}_{n^*} = [s_{1,n^*}, \dots, s_{T,n^*}]$, given the observed spike trains. We integrate over the latent variables and parameters underlying the observed spike train, as well as those underlying the new spike train, using Monte Carlo. Let $\mathbf{Z} = \{\{\mathbf{w}_n, \mathbf{a}_n, \nu_n, \mathbf{u}_n, \mathbf{v}_n\}_{n=1}^N, \boldsymbol{\theta}\}$, and let $\mathbf{z}_{n^*} = \{\nu_{n^*}, \mathbf{w}_{n^*}, \mathbf{a}_{n^*}, \mathbf{u}_{n^*}, \mathbf{v}_{n^*}\}$. Then,

$$p(\mathbf{s}_{n^*} | \mathbf{S}) \approx \int p(\mathbf{s}_{n^*} | \mathbf{z}_{n^*}, \mathbf{S}) p(\mathbf{z}_{n^*} | \mathbf{Z}) p(\mathbf{Z} | \mathbf{S}) d\mathbf{z}_{n^*} d\mathbf{Z} \approx \frac{1}{Q} \sum_{q=1}^Q p(\mathbf{s}_{n^*} | \mathbf{z}_{n^*}^{(q)}, \mathbf{S}),$$

$$\mathbf{z}_{n^*}^{(q)} \sim p(\mathbf{z}_{n^*} | \mathbf{Z}^{(q)}), \quad \mathbf{Z}^{(q)} \sim p(\mathbf{Z} | \mathbf{S}).$$

The samples $\{\mathbf{Z}^{(q)}\}_{q=1}^Q$ are posterior samples generated our MCMC algorithm given \mathbf{S} . While a proper Bayesian approach would impute \mathbf{s}_{n^*} , for large N this approximation suffices. For each sample, we draw a set of latent variables and connections for neuron n^* given the parameters $\mathbf{Z}^{(q)}$. These, combined with the spike train, enable us to compute the likelihood of \mathbf{s}_{n^*} .

B Further MCMC details

Gibbs sampling the parameters of the network model

- *Independent Model* Under an independent model, the neurons do not have latent variables so all we have to sample are the global parameters, $\boldsymbol{\theta}$. If the independent model applies to the adjacency matrix, then $\boldsymbol{\theta} = \bar{\rho}$. The model is conjugate with a beta prior. If the independent model applies to the weights, then $\boldsymbol{\theta} = \{\bar{\boldsymbol{\mu}}, \bar{\boldsymbol{\Sigma}}\}$, and the model is conjugate with a normal inverse-Wishart prior.
- *Stochastic Block Model (SBM) updates:* If a stochastic block model is used for either the adjacency matrix or the weights, then it is necessary to sample the class assignments from their conditional distribution. We iterate over each neuron and update its assignment given the rest by sampling from the conditional distribution. For example, if \mathbf{u}_n governs a stochastic block model for the adjacency matrix, the conditional distribution of the label for neuron n is given by,

$$p(\mathbf{u}_n = c | \{\mathbf{u}_{m \neq n}\}, \mathbf{A}, \boldsymbol{\theta}) \propto \pi_c \prod_{m=1}^N p(a_{m \rightarrow n} | \rho_{c_m \rightarrow}) p(a_{n \rightarrow m} | \rho_{c \rightarrow c_m}),$$

where $\boldsymbol{\theta} = \{\boldsymbol{\pi}, \{\rho_{c \leftarrow c'}\}\}$. For stochastic block models of the weight matrix, \mathbf{W} , the conditional distribution depends on $\mathbf{w}_{n \rightarrow m}$ and $\mathbf{w}_{m \rightarrow n}$ instead.

Given the class assignments and the network, the parameters $\rho_{c \leftarrow c'}$, $\boldsymbol{\mu}_{c \leftarrow c'}$, $\boldsymbol{\Sigma}_{c \leftarrow c'}$, and $\boldsymbol{\pi}$ are easily updating according to their conditional distributions, assuming $\boldsymbol{\pi}$ and $\rho_{c \rightarrow c'}$ are given conjugate Dirichlet and beta priors, respectively.

- *Latent location updates:* We resample the locations using hybrid Monte Carlo (HMC) [9]. Since the latent variables are continuous and unconstrained, this method is quite effective.

In addition to the locations, the latent distance model is parameterized by a location scale, η . Given the locations and an inverse gamma prior, the inverse gamma conditional distribution can be computed in closed form.

The remaining parameters include the log-odds, γ_0 , if the distance model applies to the adjacency matrix. This can be sampled alongside the locations with HMC. For a latent distance model of weights, the baseline mean and variance, (μ_0, σ^2) , are conjugate with a normal inverse-gamma prior.

Observation parameter updates The observation parameter updates depend on the particular distribution. Bernoulli observations have no parameters. In the binomial model, ν_n corresponds to the maximum number of possible spikes — this can often be set a priori, but it must upper bound the maximum observed spike count. For negative binomial spike counts, the shape parameter ν_n can be resampled as in Zhou et al. [23]. One possible extension is to introduce a transition operator that switches between binomial and negative binomial observations in order to truly capture both over- and under-dispersion.

We implemented our code in Python using Cython and OMP to parallelize the Pólya-gamma updates. This is available at <https://github.com/slinderman/pyglm>.

Number of MCMC iterations For both the synthetic and retinal ganglion cell results presented in the main paper, we ran our MCMC algorithm for 1000 iterations and used the last 500 samples to approximate posterior expectations. These limits were set based on monitoring the convergence of the joint log probability. Even after convergence, the weighted adjacency matrix continues to vary from sample to sample, reflecting genuine posterior uncertainty. In some cases, the samples of the stochastic block model seem to get stuck in local modes that are difficult to escape. This is a challenge with coclustering models like these, and more sophisticated transition operators could be considered, such as collapsing over block parameters in order to update block assignments.