



Bayesian Regression With Undirected Network Predictors With an Application to Brain Connectome Data

Sharmistha Guha & Abel Rodriguez

To cite this article: Sharmistha Guha & Abel Rodriguez (2020): Bayesian Regression With Undirected Network Predictors With an Application to Brain Connectome Data, Journal of the American Statistical Association, DOI: [10.1080/01621459.2020.1772079](https://doi.org/10.1080/01621459.2020.1772079)

To link to this article: <https://doi.org/10.1080/01621459.2020.1772079>



View supplementary material [↗](#)



Accepted author version posted online: 28 May 2020.
Published online: 07 Jul 2020.



Submit your article to this journal [↗](#)



Article views: 150



View related articles [↗](#)



View Crossmark data [↗](#)



Bayesian Regression With Undirected Network Predictors With an Application to Brain Connectome Data

Sharmistha Guha^a and Abel Rodriguez^b

^aDepartment of Statistical Science, Duke University, Durham, NC; ^bDepartment of Statistics, University of California, Santa Cruz, Santa Cruz, CA

ABSTRACT

This article focuses on the relationship between a measure of creativity and the human brain network for subjects in a brain connectome dataset obtained using a diffusion weighted magnetic resonance imaging procedure. We identify brain regions and interconnections that have a significant effect on creativity. Brain networks are often expressed in terms of symmetric adjacency matrices, with row and column indices of the matrix representing the regions of interest (ROI), and a cell entry signifying the estimated number of fiber bundles connecting the corresponding row and column ROIs. Current statistical practices for regression analysis with the brain network as the predictor and the measure of creativity as the response typically vectorize the network predictor matrices prior to any analysis, thus failing to account for the important structural information in the network. This results in poor inferential and predictive performance in presence of small sample sizes. To answer the scientific questions discussed above, we develop a flexible Bayesian framework that avoids reshaping the network predictor matrix, draws inference on brain ROIs and interconnections significantly related to creativity, and enables accurate prediction of creativity from a brain network. A novel class of *network shrinkage priors* for the coefficient corresponding to the network predictor is proposed to achieve these goals simultaneously. The Bayesian framework allows characterization of uncertainty in the findings. Empirical results in simulation studies illustrate substantial inferential and predictive gains of the proposed framework in comparison with the ordinary high-dimensional Bayesian shrinkage priors and penalized optimization schemes. Our framework yields new insights into the relationship of brain regions with creativity, also providing the uncertainty associated with the scientific findings. Supplementary materials for this article, including a standardized description of the materials available for reproducing the work, are available as an online supplement.

ARTICLE HISTORY

Received March 2018
Accepted May 2020

KEYWORDS

Brain connectome;
High-dimensional regression;
Influential edges; Influential
nodes; Network predictors;
Network shrinkage prior

1. Introduction

In recent years, there has been a growing interest in the study of functional and structural interconnection networks in the human brain (Fornito, Zalesky, and Breakspear 2013; Park and Friston 2013; Jbabdi et al. 2015), also referred to as *brain connectomics*. This article focuses on structural connections in the human brain, coinciding with fiber bundles that are estimated using diffusion weighted magnetic resonance imaging (DWI). The complex structural organization of the white matter of the brain can be depicted in vivo in great detail with advanced DWI (Hagmann et al. 2006). Using data from DWI, the brain can be segmented into different regions of interest (ROI), with the fiber bundles connecting the different regions estimated.

Information on fiber connections in the human brain can be used to generate network valued object data, with the network nodes corresponding to ROI and the estimated number of fiber bundles connecting any two ROIs signifying the strength of interconnection between them. There are several compelling possibilities in terms of relating these objects to brain related phenotypic traits of individuals (Zhang, Descoteaux, and Dunson 2019). For instance, the number of fibers connecting pairs of ROIs varies across individuals, and

perhaps features of these fiber connections as well as certain ROIs in the brain relate to traits of individuals, such as their creativity. In this article, we provide necessary Bayesian statistical tools to draw such inferences. In particular, we consider the problem of predicting the creativity of individuals (measured using the *composite creativity index* or CCI) based on neuroimaging data signifying the connectivity network among different brain regions. The objective of these studies is 2-fold. First, neuroscientists are often interested in identifying regions of the brain and interconnections between them which are involved in creative thinking. Second, it is important to determine how the strength of connection among these influential regions affects the level of creativity of the individual.

From a statistical point of view, our inferential problem can be best cast as a network regression problem. In a typical network regression framework, one is interested in the relationship between the structure of the network and one or more global attributes of the experimental unit on which the network data are collected. However, an overwhelming literature on network objects mainly relies on modeling the network as a random variable, rather than treating it as a predictor in a typical network regression. A number of classic models have

emerged in the literature viewing the network as a random variable with a goal of either predicting an unobserved link in a network or investigating *homophily*, that is, the process of formation of social ties due to matching individual traits. Examples include exponential random graph models (Frank and Strauss 1986), social space models (Hoff, Raftery, and Handcock 2002; Hoff 2005) including random dot product graph (RDPG) models (Young and Scheinerman 2007), and stochastic block models (Nowicki and Snijders 2001). Alternatively, models which explain covariates related to network nodes as a function of the network structure (e.g., Fowler and Christakis 2008; Shoham et al. 2015) and joint models for the co-evolution of the network structure and nodal attributes (e.g., De la Haye et al. 2010; Fosdick and Hoff 2015; Durante and Dunson 2018; Guhaniyogi and Rodriguez 2018) have also found prominence in the statistical literature of networks.

However, the inferential goals in this article in the context of the brain connectome data mentioned earlier do not quite fit in to the aforementioned literature. To address these, we employ a Bayesian model for network regression with the brain network as the predictor and the CCI as the response. In particular, we construct a Bayesian *network shrinkage prior* that combines ideas from spectral decomposition methods and spike-and-slab priors to generate a model that takes into account the topological structure of the network predictor. The model produces accurate predictions, allows us to identify both ROIs and interconnections among them that have influence on creativity, and yields well-calibrated interval estimates for the model parameters. For instance, a notable finding from the brain connectome dataset using our model is that a substantial number of ROIs influencing creativity are located in the *temporal* and *frontal* lobes of the human brain.

Earlier literature involving a network predictor and a scalar response has not amply exploited the relational nature of the network predictor. A common approach to network regression is to use a few summary measures from the network in the context of a flexible regression or classification approach (see, e.g., Bullmore and Sporns 2009 and references therein). Clearly, the success of this approach is highly dependent on selecting the right summaries to include. Furthermore, such an approach cannot identify the impact of specific nodes on the response, which is of clear interest in our setting. Alternatively, a number of authors have proceeded to vectorize the network predictor (originally obtained in the form of a symmetric *adjacency* matrix). Subsequently, the continuous response would be regressed on the high-dimensional collection of edge weights (see, e.g., Craddock et al. 2009; Richiardi et al. 2011). This approach can take advantage of the recent developments in high-dimensional regression, consisting of both penalized optimization (Tibshirani 1996) and Bayesian shrinkage (Park and Casella 2008; Carvalho, Polson, and Scott 2010). However, this approach treats the links of the network as fully exchangeable, ignoring the fact that coefficients involving common nodes can be expected to be correlated a priori.

Recently, Arroyo Reli3n et al. (2019) proposed a penalized optimization scheme that not only enables classification using network predictors, but also identifies important nodes and interconnections/edges between nodes. Although this model

seems to perform well for prediction problems, uncertainty quantification is difficult because standard bootstrap methods are not consistent for Lasso-type methods (see, e.g., Kyung et al. 2010). Modifications of the bootstrap that produce well-calibrated confidence intervals in the context of standard Lasso regression have been proposed (see, e.g., Chatterjee and Lahiri 2011), but it is not clear whether they extend to the kind of group Lasso penalties discussed in Arroyo Reli3n et al. (2019). Viewing the network predictor matrix as a two dimensional symmetric tensor, recent developments in tensor regression (see, e.g., Guhaniyogi, Qamar, and Dunson 2017; Fan, Gong, and Zhu 2019; Raskutti, Yuan, and Chen 2019) are also relevant to our work. However, these approaches do not generally take into account the symmetric constraint in the network predictor matrix, tend to focus mainly on prediction and identification of important network edges, and are not specifically designed to detect important nodes impacting the response. Moreover, their specification leads to a low-rank solution of the tensor predictor coefficient, whereas the estimated network predictor coefficient from our approach does not bear any guarantee of being low-rank.

The rest of the article evolves as follows. Section 2 provides a description of the brain connectome dataset we analyze in this article. Section 3 proposes the novel network shrinkage prior and discusses posterior computation for the proposed model. Empirical investigations with various simulation studies are presented in Section 4. Section 5 analyzes the brain connectome dataset, presenting scientific findings on influential ROIs and edges. Finally, Section 6 concludes the article with an eye toward future work.

2. Human Brain Network Data Description

Human creativity has been at the crux of the evolution of the human civilization, and has been the topic of research in several disciplines, including neuroscience. Though creativity can be defined in numerous ways, one could envision a creative idea as one that is unusual as well as effective in a given social context (Flaherty 2005). Neuroscientists generally concur that a coalescence of several cognitive processes determines the creative process, which often involves a *divergence of ideas* to conceivable solutions for a given problem. To measure the creativity of an individual, Jung et al. (2010) proposed the CCI, which is formulated by linking measures of divergent thinking and creative achievement to cortical thickness of young (23.7 ± 4.2 years), healthy subjects. Three independent judges grade the creative products of a subject from which the CCI is derived.

Along with CCI measurements, brain network information for $n = 79$ subjects is gathered using DWI. DWI is an imaging technique that enables measurement of the restricted diffusion of water in tissue to produce neural tract images. The brain imaging data we use has been preprocessed using the NDMG preprocessing pipeline (Kiar et al. 2016). In the context of DWI, the human brain is divided according to the Desikan atlas (Desikan et al. 2006) that identifies 34 cortical ROIs in each of the left and right hemispheres of the human brain, implying 68 cortical ROIs in all. These 68 ROIs are contained in 6 lobes

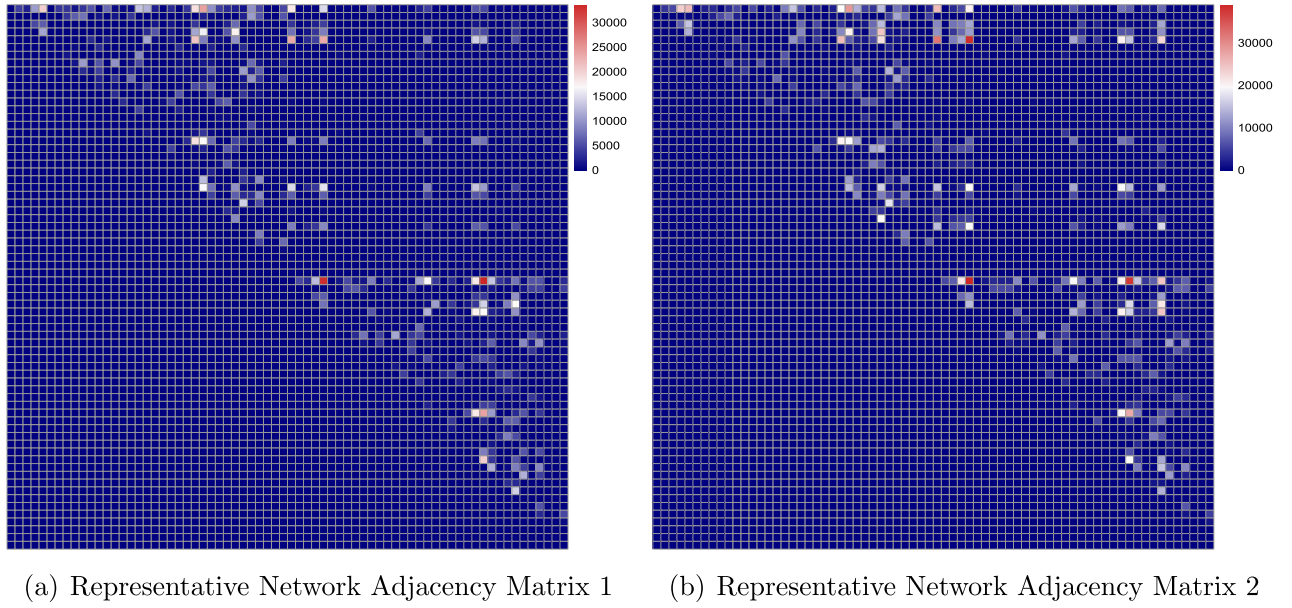


Figure 1. Maps of the brain network (weighted adjacency matrices) for two representative individuals in the sample. Since the (k, l) th off-diagonal entry in any adjacency matrix corresponds to the number of *fibers* connecting the k th and the l th ROIs, the adjacency matrices are symmetric. Hence, the figure only shows the upper triangular portion.

each in the left and the right hemispheres, namely the *temporal*, *frontal*, *occipital*, *parietal*, *insula*, and *cingulate* lobes. A “brain network” for each subject is represented by a 68×68 symmetric adjacency matrix whose rows and columns correspond to different ROIs and entries correspond to estimates of the number of “fibers” connecting pairs of brain regions. Figure 1 shows maps of the brain network for two representative individuals in the sample.

Our goal is to predict the CCI of a subject from his/her brain network, and to identify brain regions (nodes in the brain network) that are involved with creativity, as well as influential connections between different brain regions. Along with identifying a brain region as being influential, we also aim to provide the probability of the same, which serves as a measure of uncertainty. This necessitates the development of a flexible but parsimonious Bayesian model with a network predictor and a scalar response in a network regression context. The model and prior constructions follow in the next section.

3. Model Formulation

Let $y_i \in \mathbb{R}$ and $\mathbf{A}_i \in \mathbb{R}^{V \times V}$ represent the observed continuous scalar response and the weighted network predictor for the i th sample, $i = 1, \dots, n$, respectively. It is assumed that all network predictors are defined on a common set of nodes. For example, in our brain connectome application, y_i corresponds to a phenotype (CCI), while \mathbf{A}_i encodes the strength of the network connections between different regions of the brain for the i th individual. Mathematically, this amounts to \mathbf{A}_i being a $V \times V$ matrix, with the (k, l) th entry of \mathbf{A}_i denoted by $a_{i,k,l} \in \mathbb{R}$. In this article, we focus on networks that contain no self-relationship, that is, $a_{i,k,k} \equiv 0$, and are undirected ($a_{i,k,l} = a_{i,l,k}$). The brain connectome application that motivates our work naturally justifies these assumptions. However, as will become evident

later, extensions to directed networks with self-relations are straightforward.

3.1. Bayesian Network Regression Model

We propose the high-dimensional regression model of the response y_i for the i th individual on the undirected network predictor $\mathbf{A}_i = ((a_{i,k,l}))_{k,l=1}^V$ as

$$y_i = \mu + \langle \mathbf{A}_i, \mathbf{B} \rangle_F + \epsilon_i, \quad \epsilon_i \stackrel{\text{iid}}{\sim} N(0, \tau^2), \quad (1)$$

where \mathbf{B} is the symmetric network coefficient matrix of dimension $V \times V$ whose (k, l) th element is given by $\gamma_{k,l}/2$ and $\langle \mathbf{A}_i, \mathbf{B} \rangle_F = \text{Trace}(\mathbf{B}' \mathbf{A}_i)$ denotes the Frobenius inner product between \mathbf{A}_i and \mathbf{B} . The Frobenius inner product is the natural inner product in the space of matrices and is a generalization of the dot product from vector to matrix spaces. Similar to the network predictor, the network coefficient matrix \mathbf{B} is assumed to be symmetric with zero diagonal entries. The parameter τ^2 is the variance of the observational error.

Since self-relationship is absent and both \mathbf{A}_i and \mathbf{B} are symmetric, $\langle \mathbf{A}_i, \mathbf{B} \rangle_F = \sum_{1 \leq k < l \leq V} a_{i,k,l} \gamma_{k,l}$, and (1) can be rewritten as

$$y_i = \mu + \sum_{1 \leq k < l \leq V} a_{i,k,l} \gamma_{k,l} + \epsilon_i, \quad \epsilon_i \sim N(0, \tau^2). \quad (2)$$

Equation (2) connects the network regression model with the linear regression framework with $a_{i,k,l}$'s as predictors and $\gamma_{k,l}$'s as the corresponding coefficients. However, while in ordinary linear regression the predictor coefficients are indexed by the natural numbers \mathcal{N} , model (2) indexes the predictor coefficients by their positions in the matrix \mathbf{B} . This is done to keep tabs not only on the edge itself but also on the nodes connecting the edges.

3.2. Developing the Network Shrinkage Prior

3.2.1. Vector Shrinkage Prior

High-dimensional regression with vector predictors has recently been of interest in Bayesian statistics. Continuous shrinkage priors, which strongly shrink coefficients corresponding to unimportant variables to zero while minimizing the shrinkage of coefficients corresponding to influential variables, have become particularly popular. Many of these priors can be expressed as a scale mixture of normal distributions, commonly referred to as *global-local scale mixtures* (Polson and Scott 2010), that enable fast computation employing simple conjugate Gibbs sampling. More precisely, in the context of model (2), a global-local scale mixture prior would take the form $\gamma_{k,l} \sim N(0, s_{k,l}\tau^2)$, $s_{k,l} \sim g_1$, $\tau^2 \sim g_2$, $1 \leq k < l \leq V$.

Note that $s_{1,2}, \dots, s_{V-1,V}$ are local scale parameters controlling the shrinkage of the coefficients, while τ^2 is the global scale parameter. Different choices of g_1 and g_2 lead to different classes of Bayesian shrinkage priors. The Bayesian Lasso (Park and Casella 2008) shrinkage prior takes g_1 as exponential and g_2 as the Jeffreys' prior.

The direct application of this global-local prior in the context of (2) is unappealing. In practice, we expect the matrix of coefficients \mathbf{B} (which itself can be regarded as describing a weighted network) to exhibit transitivity effects, that is, we expect that if the interactions between regions k and l and between regions k and l' both influence the response, the interaction between regions l and l' will likely be influential (see, e.g., Li et al. 2013). Ordinary global-local shrinkage priors do not necessarily conform to such an important restriction.

3.2.2. Network Shrinkage Prior

We propose a shrinkage prior on the coefficients $\gamma_{k,l}$ and refer to it as the *Bayesian network shrinkage prior* (BNSP). The prior borrows ideas from low-order spectral representations of matrices, and aims to capture transitivity effects in the matrix of regression coefficients. Let $\mathbf{u}_1, \dots, \mathbf{u}_V \in \mathbb{R}^R$ be a collection of R -dimensional latent variables, one for each node, such that \mathbf{u}_k corresponds to node k . We draw each $\gamma_{k,l}$ conditionally independent from a density that can be represented as a location and scale mixture of normals. More precisely,

$$\begin{aligned} \gamma_{k,l} | s_{k,l}, \mathbf{u}_k, \mathbf{u}_l, \tau^2 &\sim N(\mathbf{u}_k' \mathbf{A} \mathbf{u}_l, \tau^2 s_{k,l}), \\ s_{k,l} &\sim \text{Exp}(\theta/2), \\ \theta &\sim \text{Gamma}(\zeta, \iota), \end{aligned} \quad (3)$$

where $s_{k,l}$ is the scale parameter corresponding to each $\gamma_{k,l}$, and $\mathbf{A} = \text{diag}(\lambda_1, \dots, \lambda_R)$ is an $R \times R$ diagonal matrix with $\lambda_r \in \{0, -1, 1\}$.

Conditional on the latent variables $\mathbf{u}_k, \mathbf{u}_l$ and \mathbf{A} , if $s_{k,l} = 0$ then (3) implies a reduced rank-decomposition $\mathbf{\Gamma} = 2\mathbf{B} = \mathbf{U}' \mathbf{A} \mathbf{U}$, where \mathbf{U} is an $R \times V$ matrix whose k th column corresponds to \mathbf{u}_k and $\mathbf{\Gamma} = ((\gamma_{k,l}))_{k,l=1}^V$. Drawing intuition from the RDPG models (Young and Scheinerman 2007), we can interpret the latent vectors $\mathbf{u}_1, \dots, \mathbf{u}_V$ as the positions of the nodes in a latent space, with the strength of the edge effect being controlled by the inner product or the angular distance between the vectors. In this interpretation, $\sum_{r=1}^R \delta_{\lambda_r \in \{-1,1\}} = R_{\text{eff}} \leq R$, represents the *effective dimensionality* of the latent space, where

$\delta_{\lambda_r \in \{-1,1\}}$ takes the value 1 if $\lambda_r \in \{-1, 1\}$, and is 0 otherwise. The effect of the interaction between the k th and l th nodes has a positive, negative or neutral impact on the response depending on whether $\mathbf{u}_k' \mathbf{A} \mathbf{u}_l > 0$, $\mathbf{u}_k' \mathbf{A} \mathbf{u}_l < 0$ or $\mathbf{u}_k' \mathbf{A} \mathbf{u}_l = 0$, respectively. This kind of bilinear structure is also commonly used to model social and biological networks because of its ability to capture the kind of transitive effects we discussed before (Hoff 2005).

To learn which components of \mathbf{u}_k are informative for (3), we assign a hierarchical prior

$$\lambda_r \sim \begin{cases} 0, & \text{w.p. } \pi_{1,r}, \\ 1, & \text{w.p. } \pi_{2,r}, \\ -1, & \text{w.p. } \pi_{3,r}, \end{cases}$$

$$(\pi_{1,r}, \pi_{2,r}, \pi_{3,r}) \sim \text{Dirichlet}(r^\eta, 1, 1), \quad \eta > 1.$$

The choice of hyper-parameters of the beta distribution is crucial. In particular, note that $E[\delta_{\lambda_r \in \{-1,1\}}] = 2/(2 + r^\eta) \rightarrow 0$ as $r \rightarrow \infty$ and that $\sum_{r=1}^R \text{var}(\delta_{\lambda_r \in \{-1,1\}}) = \sum_{r=1}^R \left[\frac{2(r^\eta+1)}{(r^\eta+2)^2(r^\eta+3)} + \frac{2(r^\eta+1)}{(r^\eta+3)(r^\eta+4)} \right] < \infty$ as $R \rightarrow \infty$. The first property provides (weak) identifiability of the different latent dimensions, while the second ensures that $\lim_{R \rightarrow \infty} \text{var}(R_{\text{eff}}) < \infty$.

In fact, we can think of our model as a level- R truncation of an infinite dimensional model, similar in spirit to the stick-breaking construction of the Indian buffet process (Teh, Grür, and Ghahramani 2007). Therefore, as long as R is chosen to be “large enough,” the inferences will be roughly invariant to this choice. In our illustrations, we perform sensitivity analyses to determine an optimal value of R that maintains computational efficiency, and at the same time ensures the robustness of the results.

To determine which nodes are most influential in explaining the response, we assign a *spike-and-slab* mixture prior (Ishwaran and Rao 2005) to the latent factor \mathbf{u}_k ,

$$\mathbf{u}_k \sim \begin{cases} N(\mathbf{0}, \mathbf{M}), & \text{if } \xi_k = 1, \\ \delta_0, & \text{if } \xi_k = 0, \end{cases} \quad \xi_k \sim \text{Ber}(\Delta), \quad (4)$$

where δ_0 is the Dirac-delta function at $\mathbf{0}$ and \mathbf{M} is a covariance matrix of order $R \times R$. The parameter Δ corresponds to the probability of the nonzero mixture component. Note that if the k th node of the network predictor is not influential in predicting the response then, a-posteriori, ξ_k should provide high probability to 0. Thus, based on the posterior probability of ξ_k , it will be possible to identify unimportant nodes, which we loosely refer to as “uninfluential nodes,” in the network regression.

The rest of the hierarchy is accomplished by assigning prior distributions on $\Delta \sim \text{Beta}(a_\Delta, b_\Delta)$ and $\mathbf{M} \sim \text{IW}(\nu, \mathbf{I})$, where $\text{IW}(\nu, \mathbf{I})$ denotes an Inverse-Wishart distribution with identity scale matrix \mathbf{I} and degrees of freedom ν . Finally, we choose a non-informative prior on (μ, τ^2) such that $p(\mu, \tau^2) \propto \frac{1}{\tau^2}$. Appendix A in the supplementary materials shows the propriety of the posterior distribution under this prior. Note that integrating over the $s_{k,l}$'s alone leads to double exponential priors that are reminiscent of the Bayesian Lasso. On the other hand, while no closed form expression exists for the marginal prior after integrating over $\mathbf{u}_1, \dots, \mathbf{u}_V$, it is easy to see that, marginally, the edge coefficients have mean zero and are not independent. Hence,

from this point of view, the latent positions $\mathbf{u}_1, \dots, \mathbf{u}_V$ simply provide a mechanism to sparsely model the prior dependence among coefficients.

Before concluding this section, we note that the BNSP prior in (3) can be alternatively written as $\gamma_{k,l} = \mathbf{u}'_k \mathbf{\Lambda} \mathbf{u}_l + \tilde{\gamma}_{k,l}$, $\tilde{\gamma}_{k,l} \sim N(0, \tau^2 s_{k,l})$, $s_{k,l} \sim \text{Exp}(\theta^2/2)$. Integrating over $s_{k,l}$, the distribution of $\tilde{\gamma}_{k,l}$ becomes a double exponential distribution. Notably, the distributions of $\tilde{\gamma}_{k,l}$'s are independent over all (k, l) pairs ($k < l$) and are independent of the distribution of \mathbf{u}_k and \mathbf{u}_l , and $\tilde{\gamma}_{k,l} \in (-\infty, \infty)$. Since there is no restriction imposed on $\tilde{\mathbf{\Gamma}} = ((\tilde{\gamma}_{k,l}/2))_{k,l=1}^V$, $\mathbf{B} = \mathbf{U}' \mathbf{\Lambda} \mathbf{U}/2 + \tilde{\mathbf{\Gamma}}$ does not assume low-rankness. Therefore, the proposed prior should ideally be able to accurately estimate the true predictor coefficient matrix even when it does not assume any specific low-rank structure.

3.3. Posterior Computation

Although summaries of the posterior distribution cannot be computed in closed form, full conditional distributions for all the parameters are available and correspond, in most cases, to standard families. Thus, posterior computation can proceed through a Markov chain Monte Carlo algorithm. We note, however, that a naive implementation of such an algorithm to update $\boldsymbol{\gamma}$, the vector composed of the upper triangular elements of $\mathbf{\Gamma}$, would have complexity $O(q^3)$, where $q = V(V-1)/2$. The resulting algorithm would therefore be computationally too expensive for situations such as our real data application, where $V = 68$ and $q = 2278$. To address this issue, we develop a procedure that proposes the use of the Woodbury matrix identity (Harville 1998) to instead compute the inverse of an $n \times n$ matrix. Since in the type of applications with which this article is concerned, n is typically much smaller than q , this approach leads to substantial computational savings that make real-life applications viable. Details of the Markov chain Monte Carlo algorithm and the efficient sampling procedure for $\boldsymbol{\gamma}$ are presented in Appendix B of the online supplementary materials.

While inferences on the latent positions $\mathbf{u}_1, \dots, \mathbf{u}_V$ is not our main focus, being able to visualize these positions can be helpful in terms of interpreting the model results. However, note that vectors $\mathbf{u}_1, \dots, \mathbf{u}_V$ are not identifiable because the model is invariant to rotations of the latent space. Hence, before we can use the posterior samples generated by our algorithm to conduct inferences on these latent positions we must first rotate them to a common orientation. This is done using a ‘‘Procrustean’’ transformation (Sibson 1978; Hoff, Raftery, and Handcock 2002; Hoff 2005). For each posterior sample $\mathbf{U}^{(\ell)}$ we find the rotation $\tilde{\mathbf{U}}^{(\ell)}$ that has the smallest sum of squared deviations from an arbitrary fixed reference matrix \mathbf{U}_0 . This rotation is given by $\tilde{\mathbf{U}}^{(\ell)} = \mathbf{U}_0 (\mathbf{U}^{(\ell)})' \left\{ \mathbf{U}^{(\ell)} \mathbf{U}_0' \mathbf{U}_0 (\mathbf{U}^{(\ell)})' \right\}^{-1/2} \mathbf{U}^{(\ell)}$. In our analysis, we use the first iterate after burn-in, $\mathbf{U}^{(1)}$, as the reference matrix \mathbf{U}_0 . Appendix E of the online supplementary materials provides an illustration of the inference on latent positions using $\tilde{\mathbf{U}}^{(\ell)}$ in one of the simulation cases.

To identify whether the k th node is important in terms of predicting the response, we rely on the post burn-in L

samples $\xi_k^{(1)}, \dots, \xi_k^{(L)}$ of ξ_k . Node k is said to be influential if $\frac{1}{L} \sum_{l=1}^L \xi_k^{(l)} > 0.5$. To identify influential edges we utilize a modification of the algorithm proposed in Li and Pati (2017) that allows us to estimate the false discovery rate of the procedure as a function of the number of discoveries. Details are provided in Appendix C of the online supplementary materials. Finally, an estimate of $P(R_{\text{eff}} = r | \text{Data})$ is given by $\frac{1}{L} \sum_{l=1}^L I(\sum_{m=1}^R \lambda_m^{(l)} = r)$, where $I(A)$ for an event A is 1 if the event A happens and 0 otherwise, and $\lambda_m^{(1)}, \dots, \lambda_m^{(L)}$ are the L post burn-in MCMC samples of λ_m .

4. Simulation Studies

This section comprehensively contrasts both the inferential and predictive performances of our proposed approach with a number of competitors in various simulation settings. As competitors, we consider both penalized likelihood methods as well as Bayesian shrinkage priors for high-dimensional regression. Our first set of competitors use generic variable selection and shrinkage methods that treat edges between nodes as ‘‘bags of predictors’’ and rely on high-dimensional regression, thereby ignoring the relational nature of the predictor. More specifically, we use Lasso (Tibshirani 1996), which is a popular penalized optimization scheme, and the Bayesian Lasso (Park and Casella 2008) and Horseshoe priors (Carvalho, Polson, and Scott 2010), which are popular Bayesian shrinkage regression methods. The Horseshoe in particular is considered to be a state-of-the-art Bayesian shrinkage prior and is known to perform well, both in sparse and not-so-sparse regression settings. We use the `glmnet` package in R (Friedman, Hastie, and Tibshirani 2010) to implement Lasso regression, and the `monomvn` package in R (Gramacy 2013) to implement the Bayesian Lasso (BLasso for short) and the Horseshoe. A thorough comparison with these methods will indicate the relative advantage of exploiting the structure of the network predictor.

Additionally, we compare our method to a frequentist approach that develops network regression in the presence of a network predictor and scalar response (Arroyo Reli3n et al. 2019). To be precise, we adapt Arroyo Reli3n et al. (2019) to a *continuous response* context and propose to estimate the network regression coefficient matrix \mathbf{B} by solving

$$\hat{\mathbf{B}} = \arg \min_{\mathbf{B} \in \mathbb{R}, \mathbf{B} = \mathbf{B}', \text{diag}(\mathbf{B}) = \mathbf{0}} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \mu - \langle \mathbf{A}_i, \mathbf{B} \rangle_F)^2 + \frac{\varphi}{2} \|\mathbf{B}\|_F^2 + \varsigma \left(\sum_{k=1}^V \|\mathbf{B}_{(k)}\|_2 + \rho \|\mathbf{B}\|_1 \right) \right\}, \quad (5)$$

where $\|\mathbf{B}\|_F = \sqrt{\langle \mathbf{B}, \mathbf{B} \rangle_F}$ denotes the Frobenius norm, $\|\mathbf{B}\|_1$ is the sum of the absolute values of all the elements of matrix \mathbf{B} , $\|\cdot\|_2$ is the l_2 norm of a vector, $\mathbf{B}_{(k)}$ is the k th row of \mathbf{B} and φ, ρ, ς are tuning parameters. The best possible choice of the tuning parameter triplet $(\varphi, \rho, \varsigma)$ is made using cross-validation over a grid of possible values. Comparison with (5) will highlight the advantages of a carefully structured Bayesian network shrinkage prior over the penalized optimization scheme incorporating network information. In the absence of open source code, we

implemented the algorithm in Arroyo Reli3n et al. (2019) ourselves. While a better and more efficient implementation of (5) is perhaps possible, our implementation provides an approximate understanding of its performance. All Bayesian competitors are allowed to draw 50,000 MCMC samples, with posterior inference carried out on the last 20,000 MCMC samples after suitable thinning.

4.1. Predictor and Response Data Generation

In all simulation studies, y_i is generated according to the network regression model

$$y_i = \mu_0 + \langle \mathbf{A}_i, \mathbf{B}_0 \rangle_F + \epsilon_i, \quad \epsilon_i \sim N(0, \tau_0^2), \quad (6)$$

with τ_0^2 as the true noise variance. All simulations use $V = 20$ nodes and $n = 70$ samples.

4.1.1. Simulation 1

In this group of simulations, the (k, l) th entry of \mathbf{B}_0 is given by $\frac{\mathbf{w}_k' \mathbf{w}_l}{2}$, where the vectors $\mathbf{w}_1, \dots, \mathbf{w}_V$, each of dimension R_{gen} , are generated from a mixture

$$\mathbf{w}_k \sim \pi_w N_{R_{\text{gen}}}(\mathbf{w}_{\text{mean}}, \mathbf{w}_{\text{SD}}^2) + (1 - \pi_w) \delta_0, \quad k \in \{1, \dots, V\}, \quad (7)$$

where δ_0 is the Dirac-delta function and π_w is the probability of any \mathbf{w}_k being nonzero. Since $(1 - \pi_w)$ is the probability of a node not being influential, it is referred to as the *node sparsity parameter*. This data generation mechanism is quite similar (although not identical) to our hierarchical prior. Hence, the goal of this first simulation is to evaluate the ability of the model to recover the true data-generation mechanism and, in particular, its ability to identify the true dimension of the latent space, as well as the sensitivity of the results to the choice of the maximum latent dimension R . The ability of the model to identify the true dimension of the latent space is assessed by looking at the posterior mode of the effective dimensionality R_{eff} (discussed in Sections 3.2.2 and 3.3). As we illustrate in Appendix D of the online supplementary materials, the posterior mode of R_{eff} coincides with the true dimension R_{gen} in all cases under *Simulation 1*.

For a comprehensive picture of *Simulation 1*, we consider 10 different cases as summarized in Table 1. In each of these cases, the network predictor coefficient and the response are generated by changing the sparsity $(1 - \pi_w)$ and the true dimension R_{gen} of the latent variables \mathbf{w}_k 's. The table also presents the maximum dimension R , used to fit the model, of the latent variables \mathbf{u}_k 's for the network regression model (2). Note that we include various cases of model mis-specification in which $R > R_{\text{gen}}$.

Table 1. The different cases for *Simulation 1*.

Quantity	Cases									
	Case 1	Case 2	Case 3	Case 4	Case 5	Case 6	Case 7	Case 8	Case 9	Case 10
R_{gen}	2	2	2	2	2	2	3	2	3	3
R	2	3	5	4	5	3	4	4	5	5
Sparsity	0.5	0.6	0.3	0.4	0.5	0.4	0.5	0.7	0.5	0.6

NOTE: The true dimension R_{gen} is the dimension of vector object \mathbf{w}_k using which data has been generated. The maximum dimension R is the dimension of vector object \mathbf{u}_k using which the model has been fit. *Sparsity* refers to the fraction of generated $\mathbf{w}_k = \mathbf{0}$, that is, $(1 - \pi_w)$.

For all simulations, \mathbf{w}_{mean} and \mathbf{w}_{SD}^2 are set as $0.8 \times \mathbf{1}_{R_{\text{gen}}}$ and $\mathbf{I}_{R_{\text{gen}} \times R_{\text{gen}}}$, respectively, and the variance τ_0^2 is fixed at 1. In Cases 1–8, the entries of the network predictor \mathbf{A}_i for the i th sample are simulated from a standard normal distribution. In Cases 9 and 10, the network predictor \mathbf{A}_i for the i th sample follows a stochastic blockmodel. In Case 9, we assume that each brain network has three local clusters with high within-cluster and low between-cluster connectivity. More specifically, the matrices \mathbf{A}_i 's consist of three symmetric block diagonal matrices of dimensions 6×6 , 7×7 , and 7×7 , respectively. Elements in these matrices are drawn from $N(j, j^2)$ where $j \in \{1, 2, 3\}$, for the j th block diagonal. The off-diagonal blocks are highly sparse, with very few non-sparse elements denoting connections between nodes in different clusters, randomly chosen from $N(0, 1)$. In Case 10, \mathbf{A}_i 's consist of 3 block diagonal matrices of dimensions 5×5 , 8×8 , and 7×7 . As before, the elements in these matrices have been drawn from $N(j, j^2)$ where $j \in \{1, 2, 3\}$, for the j th block diagonal. However, in this case the elements in the off-diagonal matrices have been drawn from $N(4, 1)$, $N(5, 1)$, and $N(6, 1)$.

4.1.2. Simulation 2

In this case, \mathbf{B}_0 is constructed by first generating V binary indicators ξ_1^0, \dots, ξ_V^0 independently from a $\text{Ber}(\pi_{2,w})$, one for each node in the network. If both $\xi_k^0 = 1$ and $\xi_l^0 = 1$, the edge coefficient connecting the k th and the l th nodes ($k < l$) is simulated from $N(0.8, 1)$. Otherwise, we set the (k, l) th edge coefficient to be 0. Similar to *Simulation 1*, we refer to $(1 - \pi_{2,w})$ as the *node sparsity parameter*. While this simulation scenario has some similarities to our proposed model, the mean effect for influential edges is constant. Therefore, the goal of this simulation is to evaluate the performance of the model in situations where there are weak network effects in the matrix of coefficients. The network predictor \mathbf{A}_i for the i th sample is simulated by drawing $a_{i,k,l}$ independently from a $N(0, 1)$ distribution for $k < l$ and setting $a_{i,k,l} = a_{i,l,k}$ and $a_{i,k,k} = 0$ for all $k, l \in \{1, \dots, V\}$. Finally, the variance τ_0^2 is fixed at 1 as in *Simulation 1*. Table 2 presents the two cases we consider for *Simulation 2*, which are obtained by varying the node sparsity

Table 2. The different cases for *Simulation 2*.

Quantity	Cases	
	Case 1	Case 2
R	5	5
Sparsity	0.7	0.2

NOTE: The maximum dimension R is the dimension of vector object \mathbf{u}_k using which the model has been fit. *Simulation 2* only has one sparsity parameter $\pi_{2,w}$.

parameter. Unlike in *Simulation 1*, the true network predictor coefficient matrix B_0 is not guaranteed to be low-rank in *Simulation 2*.

4.1.3. Simulation 3

In this case, we draw V indicator variables ξ_1^0, \dots, ξ_V^0 from a $\text{Ber}(\pi_{2,w})$ corresponding to the V nodes of the network. If both $\xi_k^0 = 1$ and $\xi_l^0 = 1$, then the edge coefficient connecting the k th and the l th nodes ($1 \leq k < l \leq V$) is simulated from a mixture distribution given by $\pi_{3,w}N(0.8, 1) + (1 - \pi_{3,w})\delta_0$. Otherwise, if $\xi_k^0 = 0$ for any k , we set the (k, l) th edge coefficient to be 0 for all l . Contrary to *Simulation 2*, *Simulation 3* allows the possibility of an edge between the k th and the l th nodes having no impact on the response even when both ξ_k^0 and ξ_l^0 are nonzero. In the context of *Simulation 3*, $(1 - \pi_{2,w})$ and $(1 - \pi_{3,w})$ are referred to as the *node sparsity* and the *edge sparsity* parameters, respectively. Hence, the goal of this simulation is to evaluate the impact of edge sparsity and its interaction with node sparsity on model performance in situations where there are weak network effects in the matrix of coefficients. Network predictors are randomly generated using the same mechanism as in *Simulation 2* and the true variance τ_0^2 is again fixed at 1 for all cases. Table 3 presents the four cases we consider in this evaluation, which are generated by changing the *node sparsity* and *edge sparsity*. Similar to *Simulation 2*, the true network predictor coefficient matrix B_0 is not guaranteed to be low-rank in *Simulation 3*.

Table 3. The different cases for *Simulation 3*.

Quantity	Cases			
	Case 1	Case 2	Case 3	Case 4
R	5	5	5	5
Node sparsity	0.7	0.2	0.7	0.2
Edge sparsity	0.5	0.5	0.3	0.7

NOTE: The maximum dimension R is the dimension of vector object u_k using which the model has been fit. While *Simulation 2* only has a sparsity parameter, *Simulation 3* has a node sparsity ($\pi_{2,w}$) and an edge sparsity ($\pi_{3,w}$) parameter, respectively.

4.2. Results

In all simulation results shown in this section, our BNSP model is fitted with the choices of the hyper-parameters given by $v = 10$, $a_\Delta = 1$, $b_\Delta = 1$, $\zeta = 1$, and $\iota = 1$. Our extensive simulation studies reveal that both inference and prediction are robust to various choices of the hyper-parameters.

4.2.1. Identification of Influential Nodes

Figures 2 and 3 show the posterior probability of the k th node being detected as influential, that is, $P(\xi_k = 1 | \text{Data})$, for each node and each case within *Simulation 1*, *Simulation 2*, and *Simulation 3*. In the case of *Simulation 1*, the model is able to accurately identify nodes influencing the response for any reasonable cutoff threshold. Indeed, the receiver operating characteristic (ROC) curves associated with all these simulations have areas under the curve (AUC) very close to 1. For both cases in *Simulation 2*, using our default threshold of 0.5, the model identifies all nodes to be inactive, though the performance is a bit better in Case 1 with a higher node sparsity. In particular, in Case 1, the model tends to assign lower posterior probabilities to truly non-influential nodes, and hence, the associated AUC is quite high (1.00). In *Simulation 3*, the model performs well when the node sparsity is high (AUC 1.00 for both Cases 1 and 3), and somewhat unsatisfactorily when the node sparsity is low. Furthermore, under lower node sparsity, the AUC is higher with lower edge sparsity (AUC for Case 2 and Case 4 are 0.81 and 0.57, respectively). Among the competitors, only Arroyo Relión et al. (2019) allow identification of influential nodes, and when applied to these simulations, selects all nodes as significant in every case.

4.2.2. Parameter Estimation

Tables 4 and 5 present the mean squared error (MSE) of all the competitors in *Simulations 1*, 2, and 3, respectively. Given that both the fitted network regression coefficient B and the true coefficient B_0 are symmetric, the MSE is calculated as $\frac{2}{V(V-1)} \sum_{k < l} (\hat{\gamma}_{k,l} - \gamma_{k,l,0})^2$, where $\hat{\gamma}_{k,l}$ is the point estimate of

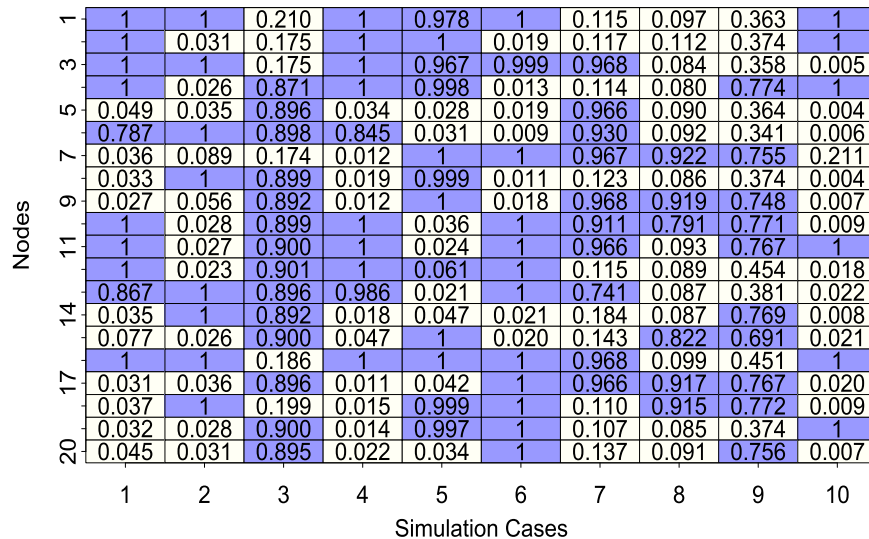


Figure 2. Posterior probability that a node is influential, $P(\xi_k = 1 | \text{Data})$, for each node and each of the 10 cases associated with *Simulation 1*. Dark cells correspond to the truly influential nodes.

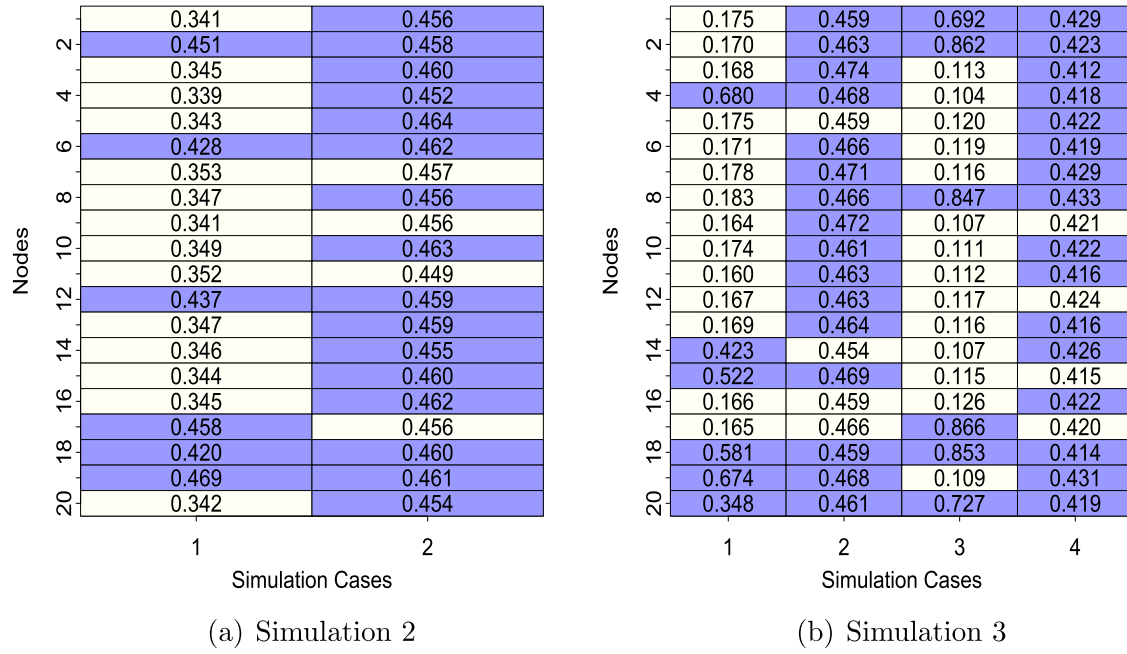


Figure 3. Posterior probability that a node is influential, $P(\xi_k = 1 | \text{Data})$, for each node and all cases associated with *Simulation 2* and *Simulation 3*. Dark cells correspond to the truly influential nodes.

Table 4. Performance of BNSP vis-a-vis competitors for cases in *Simulation 1*.

Method	MSE									
	Case 1	Case 2	Case 3	Case 4	Case 5	Case 6	Case 7	Case 8	Case 9	Case 10
BNSP	0.007	0.008	0.058	0.006	0.007	0.005	0.050	0.006	0.396	0.004
Lasso	0.438	0.660	1.295	0.455	0.371	0.412	1.344	0.010	2.231	0.025
Relión(2019)	0.524	0.929	1.117	0.552	0.493	0.655	1.629	0.069	2.207	0.047
BLasso	0.472	0.863	1.060	0.465	0.699	0.464	1.638	0.008	0.751	0.018
Horseshoe	0.395	0.012	1.070	0.393	0.299	0.493	1.381	0.007	0.706	0.012

NOTE: Point estimation of edge coefficients has been captured through the mean squared error (MSE). The minimum MSE among competitors for any case is made bold.

Table 5. Performance of BNSP vis-a-vis competitors for cases in *Simulation 2* and *Simulation 3*, respectively.

Method	MSE					
	Simulation 2		Simulation 3			
	Case 1	Case 2	Case 1	Case 2	Case 3	Case 4
BNSP	0.053	0.732	0.006	0.485	0.010	0.128
Lasso	0.012	0.843	0.006	0.636	0.004	0.178
Relión(2019)	0.036	0.859	0.017	0.617	0.036	0.145
BLasso	0.008	0.836	0.004	0.669	0.005	0.182
Horseshoe	0.006	0.948	0.002	0.629	0.004	0.145

NOTE: Point estimation of edge coefficients has been captured through the mean squared error (MSE). The minimum MSE among competitors for any case is made bold.

$\gamma_{k,l}$. For Bayesian models (including our proposed model), $\hat{\gamma}_{k,l}$ is taken to be the posterior mean of $\gamma_{k,l}$.

Table 4 shows that BNSP outperforms all its competitors in all cases of *Simulation 1*. In Cases 1–7, where the sparsity parameter is low to moderate, we perform overwhelmingly better than all the competitors. When the sparsity parameter in *Simulation 1* is high (Case 8), our simulation scheme sets a very large proportion of $\gamma_{k,l,0}$'s to zero. As a result, BNSP only slightly outperforms BLasso and Horseshoe. BNSP also shows superior performance when the network predictor has modular structure (Cases 9 and 10). While BNSP is expected to perform much better than

BLasso, Horseshoe, and Lasso due to incorporation of network information, it is important to note that the carefully chosen global-local shrinkage prior with a well formulated hierarchical mean structure seems also to outperform Arroyo Relión et al. (2019), which is explicitly designed to account for the network structure.

For *Simulations 2* and *3*, Table 5 demonstrates that, when node or edge sparsity are high, BNSP is marginally outperformed by Horseshoe. This might be due to the fact that a high degree of sparsity in the edge coefficients in the truth favors ordinary high-dimensional regression. As node sparsity decreases, so that more edge coefficients are nonzero in the truth and the network structure in the predictors dominates, BNSP tends to show increasing advantage in terms of estimating the network coefficient \mathbf{B} .

4.2.3. Identifying Influential Edges

Tables 6 and 7 show the true positive rates (TPR) and false positive rates (FPR) associated with the detection of important edges for *Simulation 1* and *Simulation 3* using BNSP, Lasso, and Arroyo Relión et al. (2019). Since *Simulation 3* is a more general case of *Simulation 2*, a corresponding table for *Simulation 2* is omitted in the interest of space. The results for our method are based on controlling the FDR at 0.05 using

Table 6. True positive rates (TPR) and false positive rates (FPR) for edges for cases in *Simulation 1*.

Method		Case 1	Case 2	Case 3	Case 4	Case 5	Case 6	Case 7	Case 8	Case 9	Case 10
BNSP	TPR	0.98	1	0.90	1	0.82	0.98	0.93	0.87	0.83	0.85
	FPR	0.05	0.02	0.02	0.08	0.06	0.08	0.06	0.01	0.02	0.06
Lasso	TPR	0.60	0.86	0.14	0.53	0.47	0.59	0.60	0.73	0.58	0.61
	FPR	0.29	0.25	0.05	0.23	0.27	0.29	0.27	0.22	0.18	0.17
Relión	TPR	1	1	1	1	1	1	1	1	1	1
	FPR	1	1	1	1	1	1	1	1	1	1

Table 7. True positive rates (TPR) and false positive rates (FPR) for edges for cases in *Simulation 3*.

Method		Case 1	Case 2	Case 3	Case 4
BNSP	TPR	0.72	0.82	1	0.91
	FPR	0.14	0.62	0.51	0.84
Lasso	TPR	0.86	0.36	0.91	0.23
	FPR	0.20	0.21	0.15	0.07
Relión	TPR	1	1	1	1
	FPR	1	1	1	1

the algorithm described in Appendix C of the supplementary materials.

In *Simulation 1*, BNSP outperforms Lasso and Arroyo Relión et al. (2019). Lasso is competitive with BNSP in *Simulation 3*, although in this case all models tend to perform unsatisfactorily when node sparsity is low. Arroyo Relión et al. (2019) identified all edges as important in all the simulation scenarios, resulting in high FPRs, although their performance could perhaps be improved with a more appropriate choice of tuning parameters rather than the choice made by a 10-fold cross-validation scheme in this article (Leng, Lin, and Wahba 2006; Meinshausen and Bühlmann 2010). However, in the absence of any available package for our context, we refrain from a more

detailed investigation of this issue as it does not constitute the main Bayesian focus of this article.

4.2.4. Predictive Inference

We compare the out-of-sample predictive ability of the different models based on the point prediction and characterization of predictive uncertainties using test samples of size $n_{\text{pred}} = 30$. To assess point prediction, we employ the mean squared prediction error (MSPE). As measures of predictive uncertainty, we provide coverage and length of 95% predictive intervals. For frequentist competitors, 95% predictive intervals are obtained by using predictive point estimates plus and minus 1.96 times standard errors.

Tables 8 and 9 show results for *Simulation 1*, *Simulation 2*, and *Simulation 3*, respectively. For *Simulation 1*, BNSP clearly outperforms other competitors in terms of point prediction. Horseshoe becomes competitive in cases with a higher degree of sparsity (Cases 2 and 8). Lasso and BLasso are competitive only in Case 8, while our approach seems to dominate the method of Arroyo Relión et al. (2019) in all cases. In terms of prediction uncertainty, BNSP tends to generate the shortest intervals among the competitors, showing slight under-coverage only in Cases 2 and 8. As in the case of point prediction, Horseshoe seems to yield results very similar to our model in Cases 2 and 8.

Table 8. MSPE, coverage and length of 95% predictive intervals (PIs) under the BNSP vis-a-vis competitors for cases in *Simulation 1*.

MSPE										
Method	Case 1	Case 2	Case 3	Case 4	Case 5	Case 6	Case 7	Case 8	Case 9	Case 10
BNSP	0.021	0.003	0.052	0.011	0.025	0.006	0.025	0.041	0.186	0.037
Lasso	0.425	0.311	1.022	0.294	0.555	0.447	0.533	0.075	0.823	4.887
Relión	0.587	0.461	0.498	0.369	0.561	0.539	0.605	0.365	0.821	0.461
BLasso	0.451	0.414	0.815	0.249	0.759	0.821	0.572	0.060	0.381	0.252
Horseshoe	0.434	0.005	0.772	0.243	0.361	0.745	0.563	0.046	0.383	0.182
Coverage of 95% PI										
Method	Case 1	Case 2	Case 3	Case 4	Case 5	Case 6	Case 7	Case 8	Case 9	Case 10
BNSP	1.000	0.933	0.990	0.967	0.990	1.000	0.967	0.933	0.990	0.990
Lasso	1.000	0.967	0.900	1.000	0.933	0.967	0.967	1.000	0.933	0.965
Relión	0.667	0.767	0.767	0.900	0.633	0.667	0.633	0.900	0.867	0.833
BLasso	1.000	0.967	1.000	1.000	0.967	0.900	1.000	1.000	0.990	0.967
Horseshoe	0.967	1.000	1.000	1.000	1.000	0.967	0.967	0.967	0.990	0.900
Length of 95% PI										
Method	Case 1	Case 2	Case 3	Case 4	Case 5	Case 6	Case 7	Case 8	Case 9	Case 10
BNSP	10.923	6.118	24.650	7.894	8.697	8.008	27.491	4.753	34.726	5.638
Lasso	39.732	59.778	38.259	52.177	27.43	56.922	70.655	23.964	75.309	13.216
Relión	18.348	34.138	30.126	26.848	15.115	32.434	31.746	12.122	61.221	10.419
BLasso	40.297	60.158	67.251	39.728	43.027	75.322	83.132	8.578	55.603	11.485
Horseshoe	33.489	9.366	61.534	33.529	30.132	76.089	68.103	5.846	69.886	6.618

NOTE: Lowest MSPE for any case is made bold.

Table 9. MSPE, coverage and length of 95% predictive intervals (PIs) under the BNSP vis-a-vis competitors for cases in *Simulation 2* and *Simulation 3*, respectively.

Method	MSPE					
	Simulation 2		Simulation 3			
	Case 1	Case 2	Case 1	Case 2	Case 3	Case 4
BNSP	0.519	0.576	0.243	0.661	0.148	0.438
Lasso	0.252	0.671	0.252	0.815	0.134	0.696
Relión	0.453	0.739	0.452	0.867	0.407	0.576
BLasso	0.153	0.722	0.220	0.875	0.124	0.781
Horseshoe	0.122	0.816	0.210	0.873	0.119	0.595

Method	Coverage of 95% PI					
	Simulation 2		Simulation 3			
	Case 1	Case 2	Case 1	Case 2	Case 3	Case 4
BNSP	0.990	0.933	1.000	1.000	0.900	0.967
Lasso	1.000	0.700	1.000	0.900	1.000	0.433
Relión	0.833	0.600	0.967	0.633	0.900	0.633
BLasso	0.933	0.833	0.933	1.000	0.933	0.867
Horseshoe	0.967	0.867	0.967	1.000	0.833	0.967

Method	Length of 95% PI					
	Simulation 2		Simulation 3			
	Case 1	Case 2	Case 1	Case 2	Case 3	Case 4
BNSP	14.053	32.219	6.889	33.094	10.069	18.704
Lasso	15.940	34.413	13.080	27.121	22.774	9.397
Relión	9.544	20.117	7.877	17.913	11.760	11.675
BLasso	8.957	45.959	6.508	52.087	5.980	22.049
Horseshoe	6.879	43.834	5.268	50.072	5.654	20.227

NOTE: Lowest MSPE for any case is made bold.

In the case of *Simulations 2* and *3*, BNSP seems to outperform all other methods in terms of point prediction in situations where the node sparsity is low. Similar observations can be made with respect to the coverage of the intervals. BNSP seems to have shorter intervals and good coverage among all Bayesian competitors in Case 2 of *Simulation 2* and in Cases 2 and 4 of *Simulation 3*, making it the top performer among the Bayesian competitors. In Case 2 of *Simulation 2*, BNSP and Horseshoe are the two best performers in terms of characterizing uncertainties. For the remaining cases in *Simulations 2* and *3*, Horseshoe and BLasso also are competitive with our method. Notably under *Simulations 2* and *3*, the true coefficient matrix does not necessarily assume any low-rank structure. BNSP is a competitive performer in terms of parametric and predictive inference even under such settings.

4.2.5. Sensitivity to the Choice of R

To examine the behavior of the model with increasing R , we rerun our model for each simulation scenario with $R = 10, 15$, and 20 (in addition to our original choice of R). For the sake of brevity, we only provide results for the data corresponding to Case 5 in *Simulation 1* (see Table 10). The behavior of all metrics is quite stable. The only values that seem to be slightly affected are the posterior means of R_{eff} and the length of the 95% credible intervals, which show little increase when we go from $R = 5$ to $R = 20$.

4.2.6. Scalability and Computation Time

Computation times (in hours) for competing methods are provided in Table 11. It may be noted that computation times for frequentist methods correspond to the time required to

Table 10. Model behavior in terms of model performance metrics with changing values of R for data corresponding to *Simulation 1*, Case 5.

R	MSE	MSPE	Coverage	Length of 95% PI	Posterior Mean of R_{eff}
5	0.0066	0.025	0.990	8.697	2.12
10	0.0063	0.022	0.967	7.167	2.27
15	0.0062	0.021	0.967	7.251	2.86
20	0.0063	0.019	0.990	8.014	2.79

NOTE: We report MSE, MSPE, length and coverage of 95% predictive intervals, and the posterior mean of effective dimensionality R_{eff} .**Table 11.** Computation times of competing methods for different values of sample size (N) and number of nodes (V).

V	N	BNR	Lasso	Relión(2019)
20	70	0.312	0.0001	0.0005
20	100	0.363	0.0002	0.0006
20	150	0.484	0.0001	0.0016
40	70	1.595	0.0001	0.0008
40	100	2.069	0.0001	0.0008
40	150	2.876	0.0002	0.0008
60	70	7.051	0.0003	0.0016
60	100	9.412	0.0002	0.0025
60	150	13.458	0.0002	0.0026
80	70	16.541	0.0005	0.0026
80	100	25.605	0.0004	0.0011
80	150	39.68	0.0006	0.0017
100	70	47.23	0.0002	0.0011
100	100	63.41	0.0004	0.0011
100	150	75.90	0.0004	0.0023

NOTE: For the Bayesian method BNR, the table records the run time in hours for all 50,000 MCMC iterations. The last two columns record total run time for frequentist methods with a single combination of tuning parameter values.

compute the method for each combination of tuning parameters and each fold in the cross-validation procedure. For BNSP, we provide the total time taken to run 50,000 MCMC iterations. Since frequentist methods yield results just after a few iterations, they require substantially smaller computation times compared to BNSP.

5. Findings From Brain Connectome Data Using BNSP

This section presents the analysis of the brain connectome data using BNSP and highlights the important findings. Note that the interest lies in predicting the CCI of a subject from his/her brain network, and in identifying brain regions (nodes in the brain network) that are involved with creativity, as well as influential connections between different brain regions. Before carrying out our analysis, each cell of the adjacency matrix is standardized by subtracting the mean and dividing by the standard deviation with respect to all $n = 79$ samples. CCI is also standardized in a similar fashion. The MCMC chain for our model is run for 50,000 iterations, with the first 30,000 iterations discarded as burn-in. Convergence is assessed by comparing different simulated sequences of representative parameters started at different initial values (Gelman et al. 2014). We monitor the autocorrelation plots and effective sample sizes. Prior specification is identical as in the simulation studies. For the purpose of this data analysis, BNSP is fitted with $R = 9$. Later, we show that the results are robust to moderate increases in the value of R .

We first focus on identifying influential ROIs in the brain network. To add robustness to this process in the absence of any ground truth, BNSP has been fitted 10 times on the real data. In each fit, the k th node is identified as *influential* if $P(\xi_k = 1 | \text{Data})$ exceeds 0.5. Finally, we report those nodes to be *influential* that are found as so in more than 50% of the 10 runs. This criteria identifies 19 out of 68 ROIs as *influential*.

Of the influential ROIs, 7 belong to the left hemisphere and 12 belong to the right hemisphere (see Table 12). Our findings coincide with results that have been previously presented in the literature. A large number of the 19 influential nodes detected by our method are part of the *frontal* (6) and *temporal* (6) cortices in both hemispheres. The frontal cortex has been scientifically associated with divergent thinking and problem solving ability, in addition to motor function, spontaneity, memory, language, initiation, judgment, impulse control, and social behavior (Stuss et al. 1985). Some of the other functions directly related to the frontal cortex seem to be behavioral spontaneity, interpreting environmental feedback and risk taking (Kolb and Milner 1981; Miller and Milner 1985; Razumnikova 2007). Similarly, Finkelstein, Vardi, and Hod (1991) report de novo artistic expression to be associated with the frontal and temporal regions. Our results also show substantial overlap with those of Jung et al. (2010), in which a regression model is used to understand the relationship between CCI and ROI-specific measures to account for the relationship between creativity and different brain regions. In particular, both approaches identify the *left superior temporal sulcus*, the *pars triangularis*, the *right superiorparietal*, the *orbitofrontal* and the *right posterior cingulate* regions as influencing CCI. However, although there is significant intersection between the findings of Jung et al. (2010) and our method, there are a few regions that we detect as influential and they do not, and vice versa. For example, our model detects the *right paracentral* and the *precentral* regions in both the hemispheres to be significantly related to CCI, while Jung et al. (2010) do not. On the other hand, they identify the *fusiform* region to be significant while we do not. Applying the method of Arroyo Relión et al. (2019) to our dataset leads to the identification of 65 out of 68 ROIs as influential. The three regions that are found

Table 12. Brain regions (ROIs) detected as influential for CCI by BNSP.

Region of interest (ROI)	Hemisphere	Lobe
Bank of the superior temporal sulcus	Left	Temporal
Rostral middle frontal gyrus	Left	Frontal
Pericalcarine	Left	Parietal
Medial orbitofrontal	Left	Temporal
Precentral	Left	Parietal
Frontal pole	Left	Temporal
Temporal pole	Left	Temporal
Medial orbitofrontal	Right	Temporal
Superior parietal lobule	Right	Temporal
Caudal anterior cingulate	Right	Cingulate
Paracentral	Right	Frontal
Isthmus cingulate cortex	Right	Occipital
Pars opercularis	Right	Frontal
Pars orbitalis	Right	Frontal
Pars triangularis	Right	Occipital
Posterior cingulate cortex	Right	Frontal
Precentral	Right	Parietal
Superior frontal gyrus	Right	Parietal
Rostral anterior cingulate cortex	Right	Frontal

Table 13. Predictive performance of competitors in terms of mean squared prediction error (MSPE), coverage and length of 95% predictive intervals, obtained through 10-fold cross-validation in the context of real data.

	BNSP	Lasso	Relión(2019)	BLasso	Horseshoe
MSPE	0.84	0.98	0.98	1.84	1.78
Coverage of 95% PI	0.91	0.97	0.97	0.97	0.93
Length of 95% PI	3.52	3.88	3.89	3.40	4.99

NOTE: Since the response has been standardized, an MSPE value greater than or around 1 will denote an inconsequential analysis.

Table 14. Predictive performance of BNSP with $R = 9, 10, 12, 15$ to assess the sensitivity of predictive inference with the choice of R .

	BNSP ($R = 9$)	BNSP ($R = 10$)	BNSP ($R = 12$)	BNSP ($R = 15$)
MSPE	0.84	0.85	0.85	0.85
Coverage of 95% PI	0.91	0.91	0.91	0.91
Length of 95% PI	3.52	3.61	3.65	3.64
Posterior mean of R_{eff}	3.37	3.57	3.62	3.59

to be uninfluential are the *frontalpole*, *temporalpole*, and the *transversetemporal* regions in the right hemisphere.

Along with influential ROIs, we are interested in identifying the statistically significant edges or connections between the 68 ROIs. Figure 4 plots the 98 interconnections that appear to be influential, controlling for a 0.05 FDR. A posteriori, the average of the mean effective dimensionality over 10 runs is 3.37.

Our interest turns now to the predictive ability of the Bayesian network regression model. Table 13 reports the average mean squared prediction error (MSPE) between observed and predicted responses, length and coverage of 95% predictive intervals for a 10-fold cross-validation exercise, over 10 independent model fits. As reference, we also present MSPE, length and coverage values for Lasso, BLasso, and Arroyo Relión et al. (2019). BNSP outperforms all other methods in terms of point prediction. In terms of prediction intervals, all methods perform similarly. The coverage of BNSP and Horseshoe are slightly under our target, while those of the other methods are slightly above target.

Finally, we assess the sensitivity of the model to the choice of R . Table 14 shows nearly identical results by choosing $R = 9, 10, 12$, and 15, suggesting that our original choice of R is sufficiently large for this application.

6. Conclusion and Future Work

This article finds ROI and interconnections among ROIs in the human brain predictive of the CCI using a novel Bayesian regression framework. This framework involves CCI as the response and the brain network as the predictor. To account for the correlation in the regression coefficients that is expected from the relational nature of the brain network predictor, we carefully construct a novel class of network shrinkage priors. A notable finding of this article is that the majority of brain regions significantly related to CCI belong to the *frontal* and *temporal* lobes in both hemispheres. Further, the Bayesian framework allows careful quantification of predictive uncertainty, as well as the uncertainty related to the identification of each ROI to be influential. Empirical results from simulation studies show that our method is superior to popular alternatives in situations

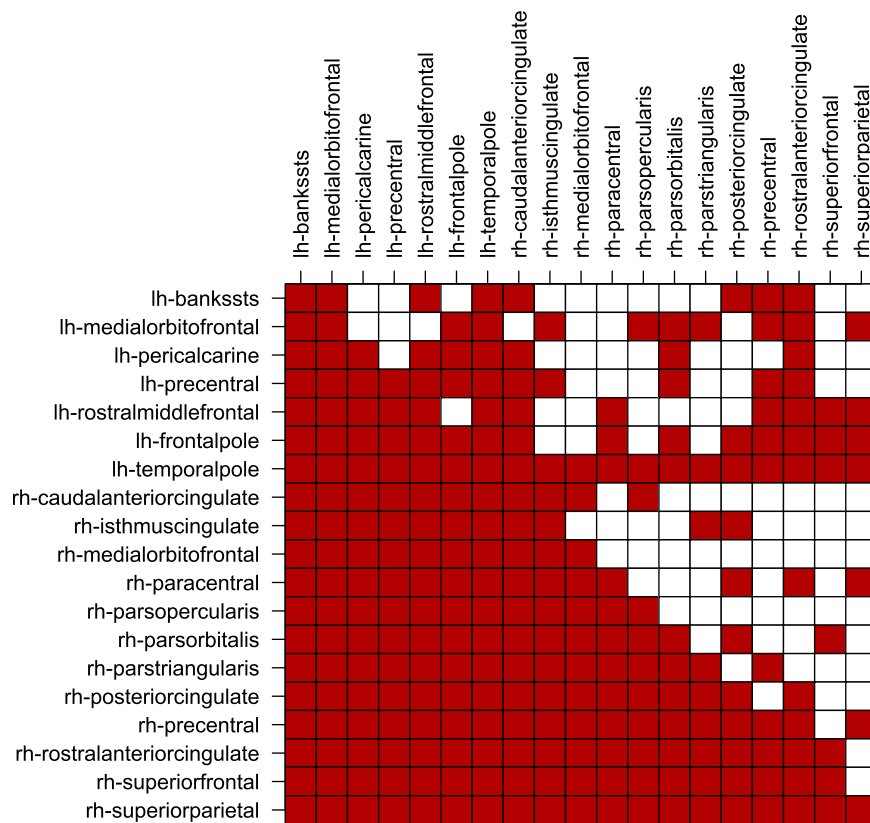


Figure 4. Significant inter-connections detected among influential brain regions of interest (ROIs) in the Desikan atlas. White cells show influential nodal associations among ROIs. Prefix “lh-” and “rh-” in the ROI names denote their positions in the left and right hemispheres of the brain, respectively.

where the level of network node sparsity is at least moderate, and mostly competitive in other circumstances.

A number of future directions emerge from this work. There exist brain connectome applications where there is a binary response phenotype (say intelligence quotient or IQ), and the goal is to simultaneously classify a sample of brain networks into “high” and “low” IQ groups and identify influential ROIs and the interconnections between ROIs predictive of IQ. The methodological problem in such applications requires extending our current approach to the context of a binary response. We will explore embedding our proposed hierarchical prior into a generalized linear model and extend the computational algorithms to deal with binary regression schemes using standard latent data augmentation (Albert and Chib 1993).

The efficient MCMC algorithm proposed here provides substantial benefits if N is small to moderate, but faces longer computation time when N is large. We plan to suitably adapt existing Bayesian methods for high-dimensional regression in the presence of a large sample size and a large number of predictors in the context of our proposed network regression, with a large sample size N and a larger number of nodes V . Finally, we could replace the global-local priors we use in this article with nonlocal priors (Johnson and Rossell 2012).

Supplementary Materials

Appendix A: This section shows posterior propriety of parameters in the BNSP model.

Appendix B: This section provides details of posterior computation for all the parameters.

Appendix C: This section describes the procedure for edge selection in BNSP.

Appendix D: This section presents inference on the *effective dimensionality* in BNSP.

Appendix E: This section presents inference on the *latent positions* in the BNSP model.

Funding

Sharmistha Guha was partially supported by funds from the Trinity College of Arts & Sciences at Duke University. Abel Rodriguez was partially supported by award NSF-DMS 1738053 and 1740850.

References

- Albert, J. H., and Chib, S. (1993), “Bayesian Analysis of Binary and Polychotomous Response Data,” *Journal of the American Statistical Association*, 88, 669–679. [12]
- Arroyo Reli3n, J. D., Kessler, D., Levina, E., and Taylor, S. F. (2019), “Network Classification With Applications to Brain Connectomics,” *The Annals of Applied Statistics*, 13, 1648–1677. [2,5,6,7,8,9,11]
- Bullmore, E., and Sporns, O. (2009), “Complex Brain Networks: Graph Theoretical Analysis of Structural and Functional Systems,” *Nature Reviews Neuroscience*, 10, 186–198. [2]
- Carvalho, C. M., Polson, N. G., and Scott, J. G. (2010), “The Horseshoe Estimator for Sparse Signals,” *Biometrika*, 97, 465–480. [2,5]
- Chatterjee, A., and Lahiri, S. N. (2011), “Bootstrapping Lasso Estimators,” *Journal of the American Statistical Association*, 106, 608–625. [2]
- Craddock, R. C., Holtzheimer, P. E., Hu, X. P., and Mayberg, H. S. (2009), “Disease State Prediction From Resting State Functional Connectivity,” *Magnetic Resonance in Medicine*, 62, 1619–1628. [2]

- De la Haye, K., Robins, G., Mohr, P., and Wilson, C. (2010), "Obesity-Related Behaviors in Adolescent Friendship Networks," *Social Networks*, 32, 161–167. [2]
- Desikan, R. S., Ségonne, F., Fischl, B., Quinn, B. T., Dickerson, B. C., Blacker, D., Buckner, R. L., Dale, A. M., Maguire, R. P., Hyman, B. T., and Albert, M. S. (2006), "An Automated Labeling System for Subdividing the Human Cerebral Cortex on MRI Scans Into Gyral Based Regions of Interest," *Neuroimage*, 31, 968–980. [2]
- Durante, D., and Dunson, D. B. (2017), "Bayesian Inference and Testing of Group Differences in Brain Networks," *Bayesian Analysis*, 13, 29–58. [2]
- Fan, J., Gong, W., and Zhu, Z. (2019), "Generalized High-Dimensional Trace Regression via Nuclear Norm Regularization," *Journal of Econometrics*, 212, 177–202. [2]
- Finkelstein, Y., Vardi, J., and Hod, I. (1991), "Impulsive Artistic Creativity as a Presentation of Transient Cognitive Alterations," *Behavioral Medicine*, 17, 91–94. [11]
- Flaherty, A. W. (2005), "Frontotemporal and Dopaminergic Control of Idea Generation and Creative Drive," *Journal of Comparative Neurology*, 493, 147–153. [2]
- Fornito, A., Zalesky, A., and Breakspear, M. (2013), "Graph Analysis of the Human Connectome: Promise, Progress, and Pitfalls," *Neuroimage*, 80, 426–444. [1]
- Fosdick, B. K., and Hoff, P. D. (2015), "Testing and Modeling Dependencies Between a Network and Nodal Attributes," *Journal of the American Statistical Association*, 110, 1047–1056. [2]
- Fowler, J. H., and Christakis, N. A. (2008), "Dynamic Spread of Happiness in a Large Social Network: Longitudinal Analysis Over 20 Years in the Framingham Heart Study," *British Medical Journal*, 337, a2338. [2]
- Frank, O., and Strauss, D. (1986), "Markov Graphs," *Journal of the American Statistical Association*, 81, 832–842. [2]
- Friedman, J., Hastie, T., and Tibshirani, R. (2010), "Regularization Paths for Generalized Linear Models via Coordinate Descent," *Journal of Statistical Software*, 33, 1–22. [5]
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., and Rubin, D. B. (2014), *Bayesian Data Analysis* (Vol. 2), Boca Raton, FL: CRC Press. [10]
- Gramacy, R. B. (2013), "R Package monomvn." [5]
- Guhaniyogi, R., Qamar, S., and Dunson, D. B. (2017), "Bayesian Tensor Regression," *Journal of Machine Learning Research*, 18, 2733–2763. [2]
- Guhaniyogi, R., and Rodriguez, A. (2018), "Joint Modeling of Longitudinal Relational Data and Exogenous Variables," *Bayesian Analysis*, 15, 477–503. [2]
- Hagmann, P., Jonasson, L., Maeder, P., Thiran, J.-P., Wedeen, V. J., and Meuli, R. (2006), "Understanding Diffusion MR Imaging Techniques: From Scalar Diffusion-Weighted Imaging to Diffusion Tensor Imaging and Beyond," *Radiographics*, 26, S205–S223. [1]
- Harville, D. A. (1998), "Matrix Algebra From a Statistician's Perspective," *Technometrics*, 40, 164. [5]
- Hoff, P. D. (2005), "Bilinear Mixed-Effects Models for Dyadic Data," *Journal of the American Statistical Association*, 100, 286–295. [2,4,5]
- Hoff, P. D., Raftery, A. E., and Handcock, M. S. (2002), "Latent Space Approaches to Social Network Analysis," *Journal of the American Statistical Association*, 97, 1090–1098. [2,5]
- Ishwaran, H., and Rao, J. S. (2005), "Spike and Slab Variable Selection: Frequentist and Bayesian Strategies," *Annals of Statistics*, 33, 730–773. [4]
- Jbabdi, S., Sotiropoulos, S. N., Haber, S. N., Van Essen, D. C., and Behrens, T. E. (2015), "Measuring Macroscopic Brain Connections In Vivo," *Nature Neuroscience*, 18, 1546. [1]
- Johnson, V. E., and Rossell, D. (2012), "Bayesian Model Selection in High-Dimensional Settings," *Journal of the American Statistical Association*, 107, 649–660. [12]
- Jung, R. E., Segall, J. M., Jeremy Bockholt, H., Flores, R. A., Smith, S. M., Chavez, R. S., and Haier, R. J. (2010), "Neuroanatomy of Creativity," *Human Brain Mapping*, 31, 398–409. [2,11]
- Kiar, G., Gray Roncal, W., Mhembere, D., Bridgeford, E., Burns, R., and Vogelstein, J. (2016), "ndmg: Neurodata's MRI Graphs Pipeline." [2]
- Kolb, B., and Milner, B. (1981), "Performance of Complex Arm and Facial Movements After Focal Brain Lesions," *Neuropsychologia*, 19, 491–503. [11]
- Kyung, M., Gill, J., Ghosh, M., and Casella, G. (2010), "Penalized Regression, Standard Errors, and Bayesian Lassos," *Bayesian Analysis*, 5, 369–411. [2]
- Leng, C., Lin, Y., and Wahba, G. (2006), "A Note on the Lasso and Related Procedures in Model Selection," *Statistica Sinica*, 16, 1273. [9]
- Li, H., and Pati, D. (2017), "Variable Selection Using Shrinkage Priors," *Computational Statistics & Data Analysis*, 107, 107–119. [5]
- Li, Y., Qin, Y., Chen, X., and Li, W. (2013), "Exploring the Functional Brain Network of Alzheimer's Disease: Based on the Computational Experiment," *PLoS One*, 8, e73186. [4]
- Meinshausen, N., and Bühlmann, P. (2010), "Stability Selection," *Journal of the Royal Statistical Society, Series B*, 72, 417–473. [9]
- Miller, L., and Milner, B. (1985), "Cognitive Risk-Taking After Frontal or Temporal Lobectomy-II. The Synthesis of Phonemic and Semantic Information," *Neuropsychologia*, 23, 371–379. [11]
- Nowicki, K., and Snijders, T. A. B. (2001), "Estimation and Prediction for Stochastic Block Structures," *Journal of the American Statistical Association*, 96, 1077–1087. [2]
- Park, H.-J., and Friston, K. (2013), "Structural and Functional Brain Networks: From Connections to Cognition," *Science*, 342, 1238411. [1]
- Park, T., and Casella, G. (2008), "The Bayesian Lasso," *Journal of the American Statistical Association*, 103, 681–686. [2,4,5]
- Polson, N. G., and Scott, J. G. (2010), "Shrink Globally, Act Locally: Sparse Bayesian Regularization and Prediction," *Bayesian Statistics*, 9, 501–538. [4]
- Raskutti, G., Yuan, M., and Chen, H. (2019), "Convex Regularization for High-Dimensional Multiresponse Tensor Regression," *The Annals of Statistics*, 47, 1554–1584. [2]
- Razumnikova, O. M. (2007), "Creativity Related Cortex Activity in the Remote Associates Task," *Brain Research Bulletin*, 73, 96–102. [11]
- Richiardi, J., Eryilmaz, H., Schwartz, S., Vuilleumier, P., and Van De Ville, D. (2011), "Decoding Brain States From fMRI Connectivity Graphs," *Neuroimage*, 56, 616–626. [2]
- Shoham, D. A., Hammond, R., Rahmandad, H., Wang, Y., and Hovmand, P. (2015), "Modeling Social Norms and Social Influence in Obesity," *Current Epidemiology Reports*, 2, 71–79. [2]
- Sibson, R. (1978), "Studies in the Robustness of Multidimensional Scaling: Procrustes Statistics," *Journal of the Royal Statistical Society, Series B*, 40, 234–238. [5]
- Stuss, D., Ely, P., Hugenholtz, H., Richard, M., LaRochelle, S., Poirier, C., and Bell, I. (1985), "Subtle Neuropsychological Deficits in Patients With Good Recovery After Closed Head Injury," *Neurosurgery*, 17, 41–47. [11]
- Teh, Y. W., Grür, D., and Ghahramani, Z. (2007), "Stick-Breaking Construction for the Indian Buffet Process," in *Artificial Intelligence and Statistics*, pp. 556–563. [4]
- Tibshirani, R. (1996), "Regression Shrinkage and Selection via the Lasso," *Journal of the Royal Statistical Society, Series B*, 58, 267–288. [2,5]
- Young, S. J., and Scheinerman, E. R. (2007), "Random Dot Product Graph Models for Social Networks," in *International Workshop on Algorithms and Models for the Web-Graph*, Springer, pp. 138–149. [2,4]
- Zhang, Z., Descoteaux, M., and Dunson, D. B. (2019), "Nonparametric Bayes Models of Fiber Curves Connecting Brain Regions," *Journal of the American Statistical Association*, 114, 1505–1517. [1]