

---

## Supplement Materials: Bayesian Joint Spike-and-Slab Graphical Lasso

---

Zehang Richard Li<sup>1</sup> Tyler H. McCormick<sup>2,3</sup> Samuel J. Clark<sup>4</sup>

### A. Proof of Proposition 1 and 2

Here we provide a proof of Proposition 2 in the main paper. The same arguments generalize to Proposition 1 directly by fixing all binary indicators to 1 and thus are not repeated.

**Proof of DSS-GGL prior** We first consider the GGL and DSS-GGL penalties. The joint distribution of all the parameters under the parameterization of the scale mixture of Normal distributions is

$$\begin{aligned}
 p(\{\Omega\}, \tau, \rho, \delta, \xi, \pi_\delta, \pi_\xi) &= \prod_g \prod_{j < k} \exp\left(-\frac{1}{2}(\omega_{jk}^{(g)})^2 \left(\frac{v_{\delta_{jk}}}{\tau_{jkg}} + \frac{v_{\xi_{jk}^*}}{\rho_{jk}}\right)\right) \prod_g \prod_j \exp\left(-\frac{\lambda_0}{2} \omega_{jj}^{(g)}\right) \\
 &\times \prod_g \prod_{j < k} \tau_{jkg}^{-\frac{1}{2}} \exp\left(-\frac{\lambda_1^2}{2} \tau_{jkg}\right) \prod_{j < k} \rho_{jk}^{-\frac{1}{2}} \exp\left(-\frac{\lambda_2^2}{2} \rho_{jk}\right) \\
 &\times \prod_{j < k} \pi_\delta^{\delta_{jk}} (1 - \pi_\delta)^{1 - \delta_{jk}} \prod_{j < k} \pi_\xi^{\xi_{jk}} (1 - \pi_\xi)^{1 - \xi_{jk}} \\
 &\times \prod_{j < k} \pi_\delta^{\alpha_1 - 1} (1 - \pi_\delta)^{b_1 - 1} \prod_{j < k} \pi_\xi^{\alpha_2 - 1} (1 - \pi_\xi)^{b_2 - 1}
 \end{aligned}$$

The following two identities provide the key steps to connect the scale mixture of normal distributions to the Laplace representation in the penalty function:

$$\int_0^\infty \frac{1}{\sqrt{2\pi s}} \exp\left(-\frac{z^2}{2s}\right) \frac{\lambda^2}{2} \exp\left(-\frac{\lambda^2 s}{2}\right) ds = \frac{\lambda}{2} \exp(-\lambda|z|) \quad (1)$$

$$\int_0^\infty \frac{1}{\sqrt{2\pi s}} \exp\left(-\frac{\|z\|_2^2}{2s}\right) \left(\frac{\lambda^2}{2}\right)^{\frac{p+1}{2}} \exp\left(-\frac{\lambda^2 s}{2}\right) s^{\frac{p+1}{2} - 1} \frac{1}{\Gamma\left(\frac{p+1}{2}\right)} ds = \frac{\lambda}{2} \exp(-\lambda\|z\|_2) \quad (2)$$

where  $z$  is a vector of length  $p$  and  $\|z\|_2 = \sqrt{\sum_{i=1}^p z_i^2}$ . By rearranging the terms and plugging in the two identities above, it can be seen that

$$p(\{\Omega\} | \delta, \xi) \propto \exp\left(-\sum_{j < k} \frac{\lambda_1}{v_{\delta_{jk}}} |\omega_{jk}^{(g)}| - \sum_{j < k} \frac{\lambda_2}{v_{\xi_{jk}^*}} \|\omega_{jk}\|_2 - \sum_g \sum_j \frac{\lambda_0}{2} \omega_{jj}^{(g)}\right).$$

The conditional distribution of  $\{\Omega\} | \delta, \xi$  takes the form of  $\exp(-pen(\{\Omega\} | \delta, \xi))$  for DSS-GGL, and thus the mode of the posterior is equivalent to the DSS-GGL solution. It still remains to be seen that the intractable constant terms are all finite,

---

<sup>\*</sup>Equal contribution <sup>1</sup>Department of Biostatistics, Yale School of Public Health, New Haven, Connecticut, USA <sup>2</sup>Department of Statistics, University of Washington, Seattle, Washington, USA <sup>3</sup>Department of Sociology, University of Washington, Seattle, Washington, USA <sup>4</sup>Department of Sociology, Ohio State University, Columbus, Ohio, USA. Correspondence to: Zehang Richard Li <zehang.li@yale.edu>.

so that each of the three conditional distributions are proper. This can be seen as follows:

$$\begin{aligned}
 C_{\tau, \rho} C_{\delta, \xi} &= \int \prod_{j < k} \text{Normal}(\omega_{jk}; 0, \Theta) \prod_g \prod_j \text{Exp}(\omega_{jj}^{(g)}; \frac{\lambda_0}{2}) \mathbf{1}_{\{\Omega\} \in M^+} d\{\Omega\} \\
 &< \int \prod_{j < k} \text{Normal}(\omega_{jk}; 0, \Theta) \prod_g \prod_j \text{Exp}(\omega_{jj}^{(g)}; \frac{\lambda_0}{2}) d\{\Omega\} = 1 \\
 C_{\tau, \rho}^{-1} &\propto \int \prod_{j < k} \left( \exp\left(-\frac{\lambda_1^2}{2} \sum_g \tau_{jkg} - \frac{\lambda_2^2}{2} \rho_{jk}\right) \rho_{jk}^{-\frac{1}{2}} \prod_g \left(\tau_{jkg} \left(\frac{1}{\tau_{jkg}} + \frac{1}{\rho_{jk}}\right)\right)^{-\frac{1}{2}} \right) d\{\tau_{jk}, \rho_{jk}\} \\
 &< \int \prod_{j < k} \left( \exp\left(-\frac{\lambda_1^2}{2} \sum_g \tau_{jkg} - \frac{\lambda_2^2}{2} \rho_{jk}\right) \rho_{jk}^{-\frac{1}{2}} \right) d\rho_{jk} \\
 &< \int \prod_{j < k} \left( \exp\left(-\frac{\lambda_2^2}{2} \rho_{jk}\right) \rho_{jk}^{-\frac{1}{2}} \right) d\rho_{jk} = (2\pi/\lambda_2^2)^{\frac{p(p-1)}{4}}
 \end{aligned}$$

The above inequalities completes the proof that the conditional prior distributions are all proper for DSS-GGL. For GGL prior, the proof is essentially the same by fixing  $\delta$  and  $\xi$  to be 1.

**Proof of DSS-FGL prior** For DSS-FGL, the joint distribution of all the parameters using the scale Normal mixture representation is

$$\begin{aligned}
 p(\{\Omega\}, \tau, \rho, \delta, \xi, \pi_\delta, \pi_\xi) &= \prod_g \prod_{j < k} \exp\left(-\frac{1}{2} (\omega_{jk}^{(g)})^T \Theta_{jk} \omega_{jk}^{(g)}\right) \prod_g \prod_j \exp\left(-\frac{\lambda_0}{2} \omega_{jj}^{(g)}\right) \\
 &\times \prod_g \prod_{j < k} \tau_{jkg}^{-\frac{1}{2}} \exp\left(-\frac{\lambda_1^2}{2} \tau_{jkg}\right) \prod_{j < k} \rho_{jk}^{-\frac{1}{2}} \exp\left(-\frac{\lambda_2^2}{2} \rho_{jk}\right) \\
 &\times \prod_{j < k} \pi_\delta^{\delta_{jk}} (1 - \pi_\delta)^{1 - \delta_{jk}} \prod_{j < k} \pi_\xi^{\xi_{jk}} (1 - \pi_\xi)^{1 - \xi_{jk}} \\
 &\times \prod_{j < k} \pi_\delta^{a_1 - 1} (1 - \pi_\delta)^{b_1 - 1} \prod_{j < k} \pi_\xi^{a_2 - 1} (1 - \pi_\xi)^{b_2 - 1}
 \end{aligned}$$

where the first term can be rewritten as the same form as  $\exp(-pen(\{\Omega\}|\delta, \xi))$  for DSS-FGL:

$$(\omega_{jk}^{(g)})^T \Theta_{jk} \omega_{jk}^{(g)} = \sum_g \frac{v_{\delta_{jk}}}{\tau_{jkg}} (\omega_{jk}^{(g)})^2 + \sum_{g < g'} \frac{v_{\xi_{jk}^*}}{\phi_{jkgg'}} (\omega_{jk}^{(g)} - \omega_{jk}^{(g')})^2.$$

Using the same identity as before, we can rewrite the conditional distribution below into the form of the DSS-FGL penalty:

$$p(\{\Omega\}|\delta, \xi) \propto \exp\left(-\sum_{j < k} \frac{\lambda_1}{v_{\delta_{jk}}} |\omega_{jk}^{(g)}| - \sum_{j < k} \sum_{g < g'} \frac{\lambda_2}{v_{\xi_{jk}^*}} |\omega_{jk}^{(g)} - \omega_{jk}^{(g')}| - \sum_g \sum_j \frac{\lambda_0}{2} \omega_{jj}^{(g)}\right).$$

The proof of the DSS-FGL conditional distributions being proper is similar to the previous case. We first note that

$$\Theta_{jk} = \text{diag}\left(\left\{\frac{1}{\tau_{jkg}}\right\}_{g=1, \dots, G}\right) + \mathbf{L}_{jk}$$

where  $\mathbf{L}_{jk}$  is a graph Laplacian matrix and is positive semi-definite. Then by Minkowski inequality,  $\det(\Theta_{jk}) \geq \det(\text{diag}\left(\left\{\frac{1}{\tau_{jkg}}\right\}_{g=1, \dots, G}\right))$ . Then we have

$$\begin{aligned}
 C_{\tau, \phi}^{-1} &\propto \int \prod_{j < k} (\det(\Theta_{jk})^{-\frac{1}{2}} \exp(-\frac{\lambda_1^2}{2} \sum_g \tau_{jkg} - \frac{\lambda_2^2}{2} \sum_{g < g'} \phi_{jkgg'})) \prod_g \tau_{jkg}^{-\frac{1}{2}} \prod_{g < g'} \phi_{jkgg'}^{-\frac{1}{2}} d\{\tau, \phi\} \\
 &\leq \int \prod_{j < k} (\exp(-\frac{\lambda_1^2}{2} \sum_g \tau_{jkg} - \frac{\lambda_2^2}{2} \sum_{g < g'} \phi_{jkgg'}) \prod_{g < g'} \phi_{jkgg'}^{-\frac{1}{2}}) d\{\tau, \phi\} \\
 &\leq \int \prod_{j < k} (\exp(-\frac{\lambda_2^2}{2} \sum_{g < g'} \phi_{jkgg'}) \prod_{g < g'} \phi_{jkgg'}^{-\frac{1}{2}}) d\phi,
 \end{aligned}$$

which is again finite since the integral consists of products of Gamma densities. The rest of the argument follows in the same way as the DSS-GGL case.

## B. Details of the EM algorithm implementation

Assuming no missing data, the full objective function in the  $t$ -th iteration of the EM algorithm described in 5 is the expectation of the complete data log likelihood, i.e.,

$$\begin{aligned}
 Q(\{\Omega\}, \pi_\delta, \pi_\xi | \{\Omega\}^{(t)}, \pi_\delta^{(t)}, \pi_\xi^{(t)}) &= E_{\delta, \xi | \{\Omega\}^{(t)}, \pi_\delta^{(t)}, \pi_\xi^{(t)}, \mathbf{X}} (\log p(\{\Omega\}, \pi_\delta, \pi_\xi | \mathbf{X}) | \{\Omega\}^{(t)}, \pi_\delta^{(t)}, \pi_\xi^{(t)}, \mathbf{X}) \\
 &= \text{constant} + \sum_g \frac{n_g}{2} \log |\Omega_g| - \frac{1}{2} \sum_g \text{tr}(\mathbf{S}_g \Omega_g) - \frac{\lambda_0}{2} \sum_j \sum_g |\omega_{jj}^{(g)}| \\
 &\quad - \lambda_1 \sum_{j < k} \sum_g |\omega_{jk}^{(g)}| E_{\cdot | \cdot} \left[ \frac{1}{v_0(1 - \delta_{jk}) + v_1 \delta_{jk}} \right] + \sum_{j < k} \log \left( \frac{\pi_\delta}{1 - \pi_\delta} \right) E_{\cdot | \cdot}(\delta_{jk}) \\
 &\quad - \lambda_2 \sum_{j < k} \widetilde{pen}(\omega_{jk}) E_{\cdot | \cdot} \left[ \frac{1}{v_0(1 - \delta_{jk} \xi_{jk}) + v_1 \delta_{jk} \xi_{jk}} \right] + \sum_{j < k} \log \left( \frac{\pi_\xi}{1 - \pi_\xi} \right) E_{\cdot | \cdot}(\xi_{jk}) \\
 &\quad + (a_1 - 1) \log(\pi_\delta) + (b_1 + \frac{p(p-1)}{2} - 1) \log(1 - \pi_\delta) \\
 &\quad + (a_2 - 1) \log(\pi_\xi) + (b_2 + \frac{p(p-1)}{2} - 1) \log(1 - \pi_\xi),
 \end{aligned}$$

where  $E_{\cdot | \cdot}$  denotes conditional expectation  $E_{\delta, \xi | \{\Omega\}^{(t)}, \pi_\delta^{(t)}, \pi_\xi^{(t)}, \mathbf{X}}$ , and  $\widetilde{pen}(\omega_{jk}) = \|\omega_{jk}\|_2$  for DSS-GGL and  $\widetilde{pen}(\omega_{jk}) = \sum_{g < g'} |\omega_{jk}^{(g)} - \omega_{jk}^{(g')}|$  for DSS-FGL.

When missing data exists, we need to also calculate the expectation of  $\mathbf{S}_g$  given the observed data. That is, We need to replace the  $\text{tr}(\mathbf{S}_g \Omega_g)$  term in the above objective function by

$$E_{\cdot | \cdot}(\mathbf{S}_g \Omega_g) = E_{\cdot | \cdot} \left( \left( \frac{1}{n_g} \sum_i \mathbf{x}_i^{(g)} (\mathbf{x}_i^{(g)})^T \right) \Omega_g \right) = \frac{1}{n_g} \left( \sum_i E_{\mathbf{x}_{i,m}^{(g)} | \mathbf{x}_{i,o}^{(g)}} (\mathbf{x}_i^{(g)} (\mathbf{x}_i^{(g)})^T) \right) \Omega_g.$$

where  $\mathbf{x}_{i,o}^{(g)}$  and  $\mathbf{x}_{i,m}^{(g)}$  denote the observed and missing cells in  $\mathbf{x}_i^{(g)}$  respectively.  $\mathbf{x}_i^{(g)}$  follows a multivariate Gaussian distribution. Without loss of generality, if we let  $\mathbf{x}_i^{(g)} = \begin{pmatrix} \mathbf{x}_{i,o}^{(g)} \\ \mathbf{x}_{i,m}^{(g)} \end{pmatrix}$ , we know

$$\begin{aligned}
 E_{\mathbf{x}_{i,m}^{(g)} | \mathbf{x}_{i,o}^{(g)}} (\mathbf{x}_{i,m}^{(g)}) &= \Sigma_{mo} \Sigma_{oo}^{-1} \mathbf{x}_{i,o}^{(g)} \\
 E_{\mathbf{x}_{i,m}^{(g)} | \mathbf{x}_{i,o}^{(g)}} (\mathbf{x}_i^{(g)} (\mathbf{x}_i^{(g)})^T) &= E_{\cdot | \cdot} (\mathbf{x}_i^{(g)}) E_{\cdot | \cdot} (\mathbf{x}_i^{(g)})^T + \begin{pmatrix} \mathbf{0}_{oo} & \mathbf{0}_{om} \\ \mathbf{0}_{mo} & \Sigma_{mm} - \Sigma_{mo} \Sigma_{oo}^{-1} \Sigma_{om} \end{pmatrix}
 \end{aligned}$$

where  $\Sigma_{oo}$ ,  $\Sigma_{om}$ ,  $\Sigma_{mo}$  and  $\Sigma_{mm}$  are the corresponding submatrices of  $\Sigma_g$ .

### C. Gibbs sampler of the proposed models

The EM algorithm introduced in the main paper maximizes the complete data likelihood by looking at the Laplace representation after integrating out all the latent parameters. In this section, we show that these latent parameters,  $\tau$ ,  $\phi$  and  $\rho$ , facilitates efficient block Gibbs sampling algorithms for fully Bayesian inference.

We start by describing the posterior sampling of  $\{\Omega\}$ . The basic idea is to sample each column and row for all the precision matrices jointly. To simplify notation, we separate out the last column and row in  $\Omega_g$  and  $S_g$  and define

$$\Omega_g = \begin{pmatrix} \Omega_{11}^{(g)} & \omega_{12}^{(g)} \\ \omega_{21}^{(g)} & \omega_{22}^{(g)} \end{pmatrix}, \quad S_g = \begin{pmatrix} S_{11}^{(g)} & s_{12}^{(g)} \\ s_{21}^{(g)} & s_{22}^{(g)} \end{pmatrix}.$$

We further let  $\omega_{12} = [(\omega_{12}^{(1)})^T, (\omega_{12}^{(2)})^T, \dots, (\omega_{12}^{(G)})^T]^T$ , and  $s_{12} = [(s_{12}^{(1)})^T, (s_{12}^{(2)})^T, \dots, (s_{12}^{(G)})^T]^T$ , each denoting a vector of length  $(p-1)G$ , and  $\omega_{22} = [\omega_{22}^{(1)}, \omega_{22}^{(2)}, \dots, \omega_{22}^{(G)}]$ .

The conditional distribution of  $(\omega_{12}, \omega_{22})$  given the rest of the elements in  $\{\Omega\}$  does not seem to take any standard form. However, if we perform a change of variables and let  $\theta_g = \omega_{22}^{(g)} - \omega_{21}^{(g)}(\Omega_{11}^{(g)})^{-1}\omega_{12}^{(g)}$ , the conditional distribution of  $(\omega_{12}, \theta)$  becomes

$$\begin{aligned} p(\omega_{12}, \theta) &\propto \sum_g \theta_g^{\frac{n_g}{2}} \exp\left(-\frac{s_{22}^{(g)} + \lambda_0}{2} \theta_g - s_{12}^T \omega_{12} - \frac{1}{2} \omega_{12}^T \mathbf{A} \omega_{12}\right) \\ &= \prod_g \text{Gamma}(\theta_g; \frac{n_g}{2} + 1, \frac{s_{22}^{(g)} + \lambda_0}{2}) \times \text{Normal}(\omega_{12}; -\mathbf{A}^{-1} s_{12}, \mathbf{A}^{-1}). \end{aligned}$$

The  $\mathbf{A}$  matrix can be calculated by  $\mathbf{A} = \mathbf{U} + \mathbf{V}$ , where  $\mathbf{U}$  is a matrix by rearranging the precision matrices so that its  $((g-1)(p-1) + k, (g'-1)(p-1) + k)$ -th element is the  $(g, g')$ -element in  $\Theta_{jk}$  defined in (12) and (13) of the main paper, and

$$\mathbf{V} = \begin{pmatrix} (\lambda_0 + s_{22}^{(1)})(\Omega_{11}^{(1)})^{-1} & & & \\ & \ddots & & \\ & & & (\lambda_0 + s_{22}^{(G)})(\Omega_{11}^{(G)})^{-1} \end{pmatrix},$$

For DSS-GGL, we notice that  $\mathbf{A}$  is block diagonal, thus we can alternatively sample  $\omega_{12}^{(g)}$  independently by

$$\omega_{12}^{(g)} | \cdot \sim \text{Normal}(-\mathbf{A}_g^{-1} s_{12}^{(g)}, \mathbf{A}_g^{-1})$$

where  $\mathbf{A}_g = (\lambda_0 + s_{22}^{(g)})(\Omega_{11}^{(g)})^{-1} + \tilde{\Theta}_{11}^{(g)}$ .

Given  $\{\Omega\}$ , the latent parameters in DSS-GGL have simple conditional distribution as follows:

$$\begin{aligned} \tau_{jkg}^{-1} | \cdot &\sim \text{InvGaussian}\left(\frac{\lambda_1}{v_{\delta_{jk}}^{\frac{1}{2}} |\omega_{jk}^{(g)}|}, \lambda_1^2\right), \quad j, k = 1, \dots, p, g = 1, \dots, G \\ \rho_{jk}^{-1} | \cdot &\sim \text{InvGaussian}\left(\frac{\lambda_2}{v_{\xi_{jk}^*} \sum_g (\omega_{jk}^{(g)})^2}, \lambda_2^2\right), \quad j, k = 1, \dots, p \\ \delta_{jk}, \xi_{jk} | \cdot &\sim p^*(\delta_{jk}, \xi_{jk}), \quad j, k = 1, \dots, p \\ \pi_{\delta} &\sim \text{Beta}\left(\sum \delta_{jk} + a_1, \sum (1 - \delta_{jk}) + b_1\right) \\ \pi_{\xi} &\sim \text{Beta}\left(\sum \xi_{jk} + a_2, \sum (1 - \xi_{jk}) + b_2\right) \end{aligned}$$

where  $p^*(\delta, \xi)$  is defined in (5) in the main paper.

For DSS-FGL, the conditional distribution of  $\tau$ ,  $\delta$ ,  $\xi$ ,  $\pi_{\delta}$ , and  $\pi_{\xi}$  are the same as DSS-GGL. The conditional distribution of  $\phi$  is

$$\phi_{jkgg}^{-1} \sim \text{InvGaussian}\left(\frac{\lambda_2}{v_{\xi_{jk}^*} |\omega_{jk}^{(g)} - \omega_{jk}^{(g')}|}, \lambda_2^2\right), \quad j, k = 1, \dots, p, g, g' = 1, \dots, G$$

The Gibbs sampler is then complete by circling through and sampling each blocks of  $\{\Omega\}$  and the latent parameters with the above posterior conditional distributions.

## D. Additional illustrating example

In this section, we provide a more detailed description to the small 2-class simulated example described in Section 6 of the main paper. We let  $n_g = 150$  for  $g = 1, 2$ , and  $p = 100$ . The first 10 variables in the first class form a 10-node block with an AR(1) precision matrix, i.e.,  $(\Omega^{-1})_{jk} = 0.7^{|j-k|}$ . The rest of the 90 variables are independent noises from a standard Gaussian distribution. The second class shares the first 4 edges of the first class. The first 5 variables form another AR(1) block with different strength of correlations so that  $(\Omega^{-1})_{jk} = 0.9^{|j-k|}$ . The rest of the 95 variables all follow independent standard Gaussian distribution.

The best performance of FGL in our experiments was achieved with  $\lambda_2 = 0.1$ . We then obtained the regularization path along different values of  $\lambda_1$ . We also fit DSS-FGL and SS-FGL with  $\lambda_1 = 1, \lambda_2 = 1$ . The latter assumes  $\xi_{jk} = 1$  for all edges, i.e., the penalization of similarities is always proportional to the penalization of sparsity.

Figure 1 shows both the two solution paths presented in the main paper, and the solution path from SS-FGL. It can be seen that although SS-FGL achieves similar bias as DSS-FGL, it also estimates several more false positive edges. This can be seen from the formulation of the doubly spike-and-slab selection: with only one spike-and-slab mixture of the penalties, the selected edges from the slab distributions receive also only weak penalization for between-class similarities. Thus it is more likely to pick up spurious edges due to noises that happen to exist in one class. This full illustration of the simple example shows the advantage of having the doubly spike-and-slab setup.

## E. Additional simulation evidence

Here we describe our procedure in simulating the dataset described in Section 7 of the main paper. We first generate three networks with  $p$  features with 10 equal sized unconnected subnetwork. Each of the subnetwork follow a power law degree distribution, which is generally harder to estimate than simpler structures (Peng et al., 2009). The first class contains all ten subnetworks, and the second and third classes each has one and two subnetworks removed. Given the network structure, we generate  $\Omega_g$  from the  $G$ -Wishart distribution  $W_G(3, I_p)$ , and rescale them so that  $\Omega_g^{-1}$  have unit variances. Finally, we generate  $n = 150$  independent and identically distributed samples from  $\text{Normal}(\mathbf{0}, \Omega_g)$  in each class. The resulted graph for  $p = 500$  is shown in Figure 2. We fit GGL and FGL with various choice of fixed  $\lambda_2$  and a sequence of  $\lambda_1$ . We fit DSS-GGL and DSS-FGL with  $\lambda_1 = 1, \lambda_2 = 2/30$  in this case. We explore more choices of  $\lambda_2$  in the moderate dimensional experiments below and found no substantial changes in the performance of the final models.

Additional results for  $p = 100$  and  $p = 200$  are shown in Figure 3 and 4. The dots correspond to DSS-GGL and DSS-FGL with  $\lambda_1 = 1, \lambda_2 = 0.1$ . We also examined different choices of  $\lambda_2$  are found no substantial differences in performance. An exploratory sensitivity analysis is presented in the next subsection.

### E.1. Sensitivity to hyperparameters

Figure 5 and 6 shows the converged regions over 100 replications on space of the true positive against false positive discoveries for edges and differential edges respectively when  $p = 200$ . It can be seen that the performance is relatively stable under different choices of  $\lambda_2$ .

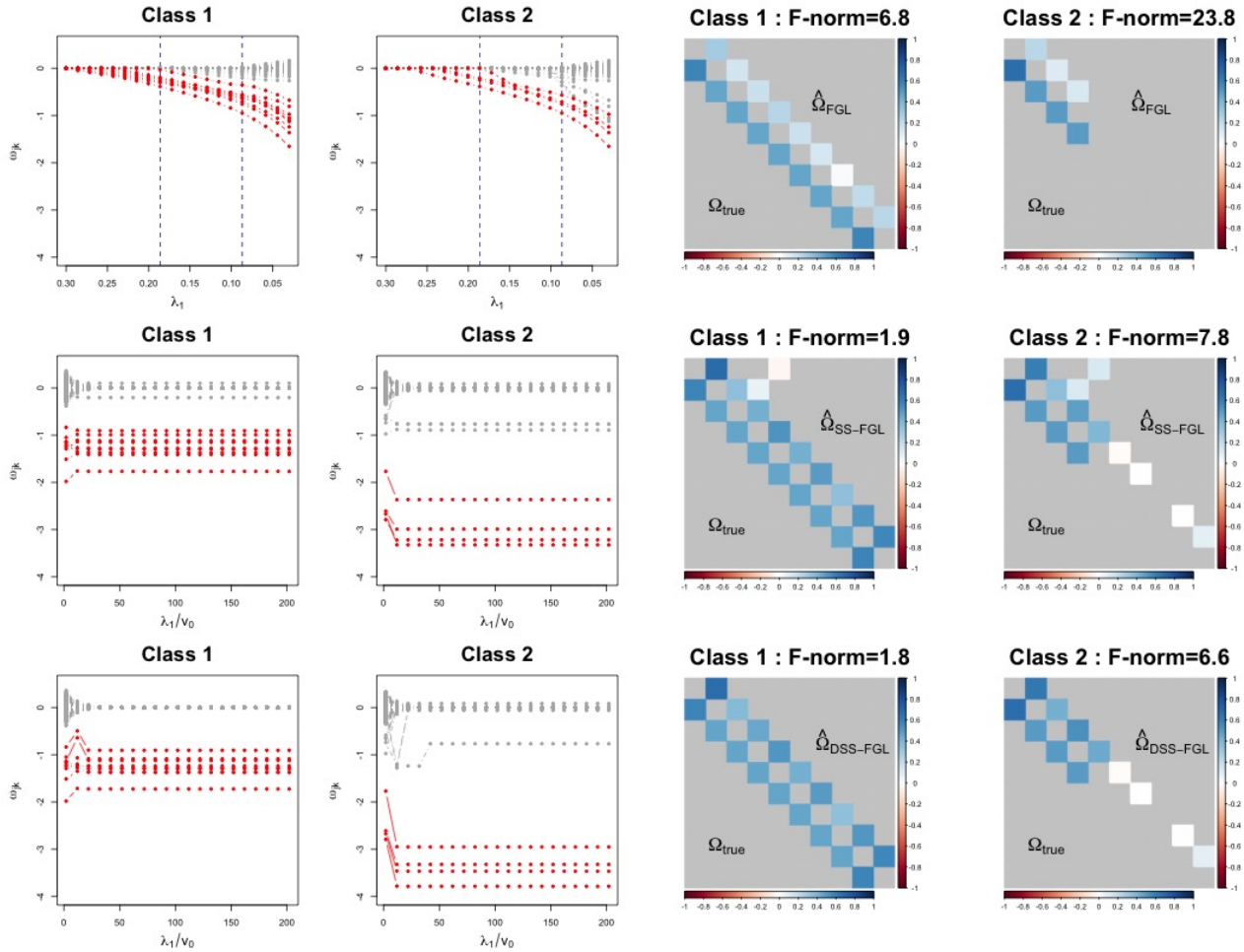


Figure 1. The solution paths and estimated precision matrices of FGL (upper row), SS-FGL (middle row) and DSS-FGL (lower row). The red nodes correspond to true edges and the gray nodes correspond to 0's. The two vertical lines in the FGL solution path indicate the model that best matches the true sparsity (left) and the model with the lowest AIC (right). The block containing the edges is plotted for the estimated values (upper triangular) against the truth (lower triangular). The model that best matches the true graphs is plotted for FGL. The off-diagonal values are rescaled and negated to partial correlations, and 0's are colored with light gray background for easier visual comparison. The bias of the estimated precision matrix measured by the Frobenius norm,  $\|\hat{\Omega}_g - \Omega_g\|_F$ , is also printed in the captions.

### Graph structure

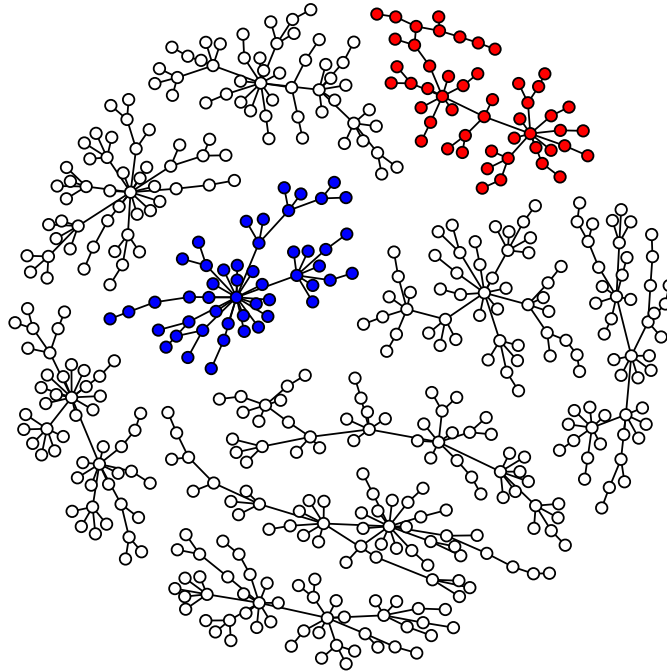


Figure 2. Graph structure of the simulated dataset. The edges between the red nodes are removed from the second class, and edges between both the red and blue nodes are removed from the third class.

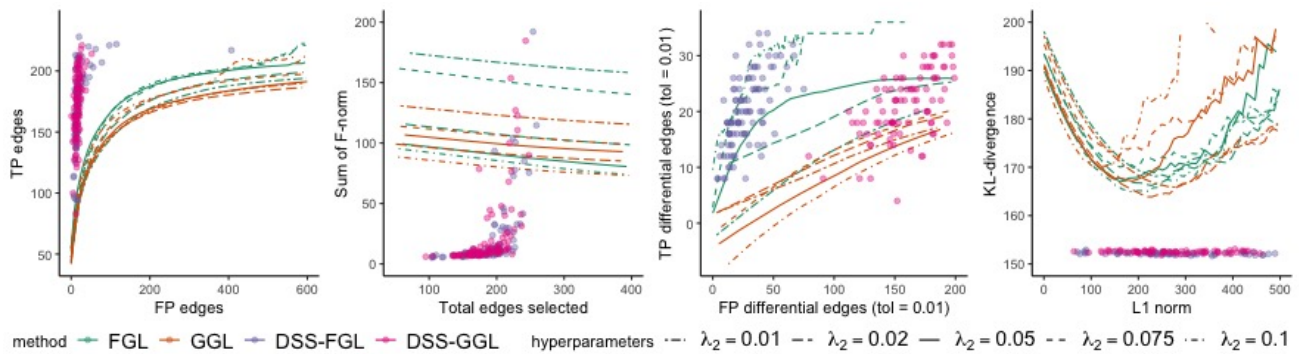


Figure 3. Performance of FGL, GGL, DSS-FGL, and DSS-GGL over 100 replications,  $p = 100$ . The dots represent the metrics for the 100 selected models under DSS-FGL and DSS-GGL, and the lines represent the average performance of FGL and GGL over 100 replications under different tuning parameters.

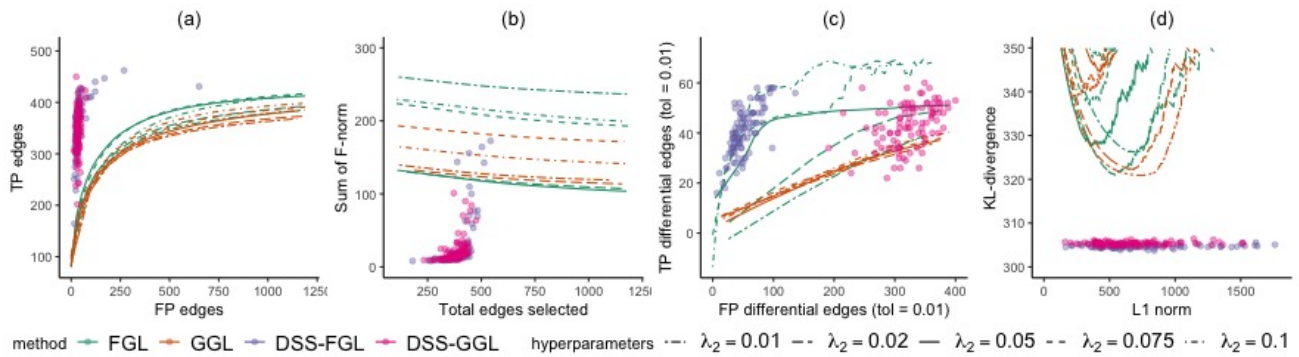


Figure 4. Performance of FGL, GGL, DSS-FGL, and DSS-GGL over 100 replications,  $p = 200$ . The dots represent the metrics for the 100 selected models under DSS-FGL and DSS-GGL, and the lines represent the average performance of FGL and GGL over 100 replications under different tuning parameters.

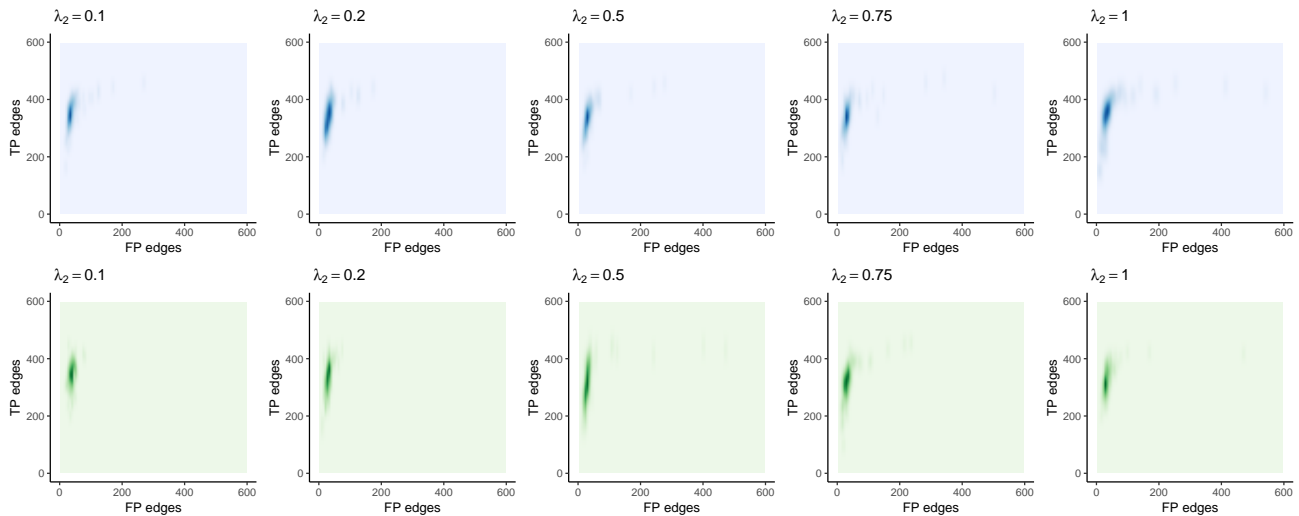


Figure 5. The density plot of true positive edges against false positive edges for DSS-FGL (top row), and DSS-GGL (bottom row) under different choices of  $\lambda_2$ .  $\lambda_1$  is set to 1.



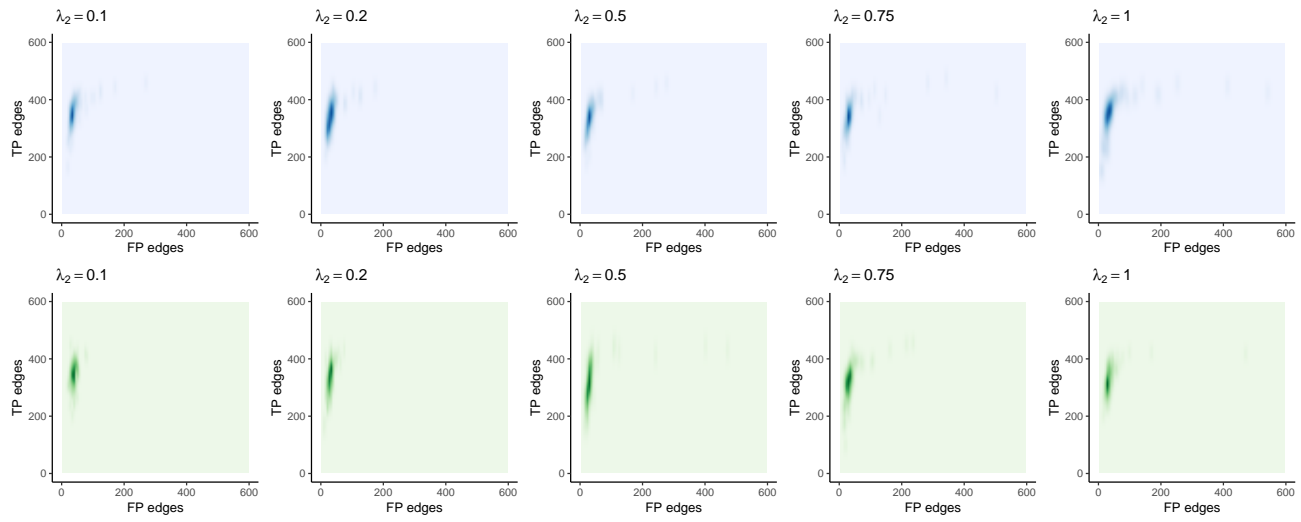


Figure 6. The density plot of true positive differential edges against false positive positive edges for DSS-FGL (top row), and DSS-GGL (bottom row) under different choices of  $\lambda_2$ .  $\lambda_1$  is set to 1.

## F. Details on the verbal autopsy data analysis

### F.1. List of symptoms

Table 1 shows the questions with continuous responses used in the analysis in the main paper.

Table 1. List of symptoms considered in this analysis.

Abbreviation	Questionnaire item
ill	For how long was [name] ill before s/he died? [days]
fever	How many days did the fever last? [days]
rash	How many days did [name] have the rash? [days]
ulcer	For how many days did the ulcer ooze pus? [days]
yellow discoloration	For how long did [name] have the yellow discoloration? [days]
ankle swelling	For how long did [name] have ankle swelling? [days]
puffiness face	For how long did [name] have puffiness of the face? [days]
puffiness body	For how long did [name] have puffiness all over his/her body? [days]
cough	For how long did [name] have a cough? [days]
difficulty breathing	For how long did [name] have difficulty breathing? [days]
fast breathing	For how long did [name] have fast breathing? [days]
liquid stool	For how long before death did [name] have loose or liquid stools? [days]
vomit	For how long before death did [name] vomit? [days]
difficulty swallowing	For how long before death did [name] have difficulty swallowing? [days]
belly pain	For how long before death did [name] have belly pain? [days]
protruding belly	For how long before death did [name] have a protruding belly? [days]
mass belly	For how long before death did [name] have a mass in the belly [days]
headaches	For how long before death did [name] have headaches? [days]
stiff neck	For how long before death did [name] have stiff neck? [days]
unconsciousness	For how long did the period of loss of consciousness last? [days]
confusion	For how long did the period of confusion last? [days]
convulsion	For how long before death did the convulsions last? [days]
paralysis	For how long before death did [name] have paralysis? [days]
period overdue	For how many weeks was her period overdue? [days]
tobacco	How much pipe/chewing tobacco did [name] use daily?
cigarettes	How many cigarettes did [name] smoke daily?
age	Age [years]

### F.2. Comparing with JGL

The estimated symptom network from the DSS-FGL and DSS-GGL are summarized in Figure 7.

The estimated symptom network from the FGL and GGL are summarized in Figure 8. We fit both models under a 2-dimensional grids over  $\lambda_1$  and  $\lambda_2$ . As expected, AIC selects very dense graphs (first two rows of Figure 8) and are difficult to interpret. We also compare the FGL and GGL graph with the closest number of edges as those from DSS-FGL and DSS-GGL in the third and fourth row of Figure 8. The number of differential edges is typically smaller compared to DSS-FGL and DSS-GGL, which is likely due to over penalization of similarities, i.e., edges become too similar using FGL, and too sparse among half of the nodes using GGL.

## G. Details on prediction of missing mortality rates

The data we consider in this example consist of log mortality rates over  $n = 51$  years for  $p = 101$  age groups, and 2 classes representing female and male series respectively. The estimated graph structure from one of the cross-validation dataset using FGL and DSS-FGL are shown in Figure 9. The Lee-Carter model are estimated using the R package `ilc` (Butt et al., 2014) for each gender separately.

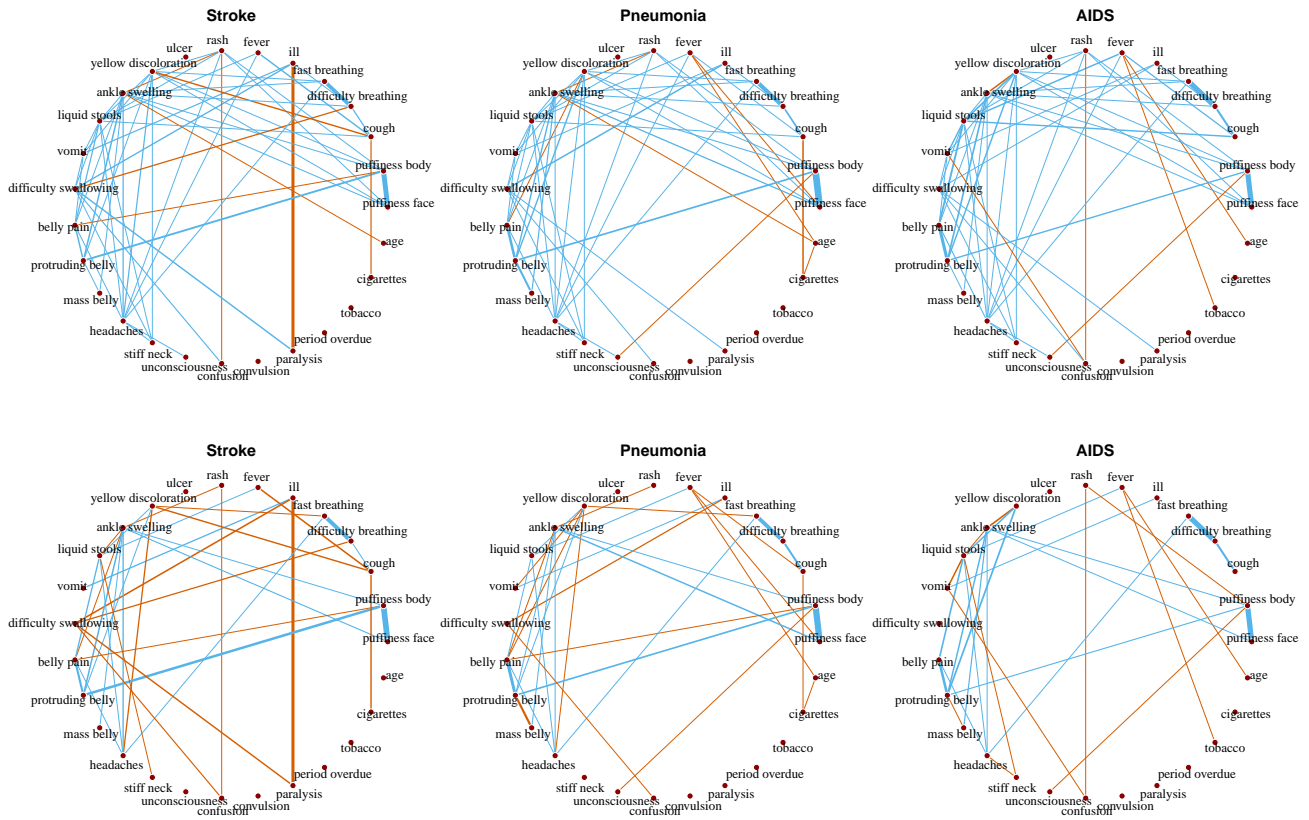


Figure 7. Estimated edges between the symptoms under the three causes using DSS-FGL (top row) and DSS-GGL (bottom row). The width of the edges are proportional to the size of  $|\omega_{jk}^{(g)}|$ . Common edges across all groups are colored in blue, and the differential edges are colored in red.

The DSS-FGL is able to pick up more conditional dependence structures along the diagonal among several age groups, while the FGL estimates mostly within adults only. It is interesting that both approaches identifies positive partial correlations between age 14 – 17 and 30 – 40 between male and female. This is likely due to the fact that male mortality around age 20 typically shows a hump of increase due to young adult accident mortality, which leads to the mean model more likely to underestimations for mortality during age 18 – 30 and overestimations both before and after that period. This relationship of the age curve, however, is not seen in female mortality.

Supplement: Bayesian Joint Spike-and-Slab Graphical Lasso

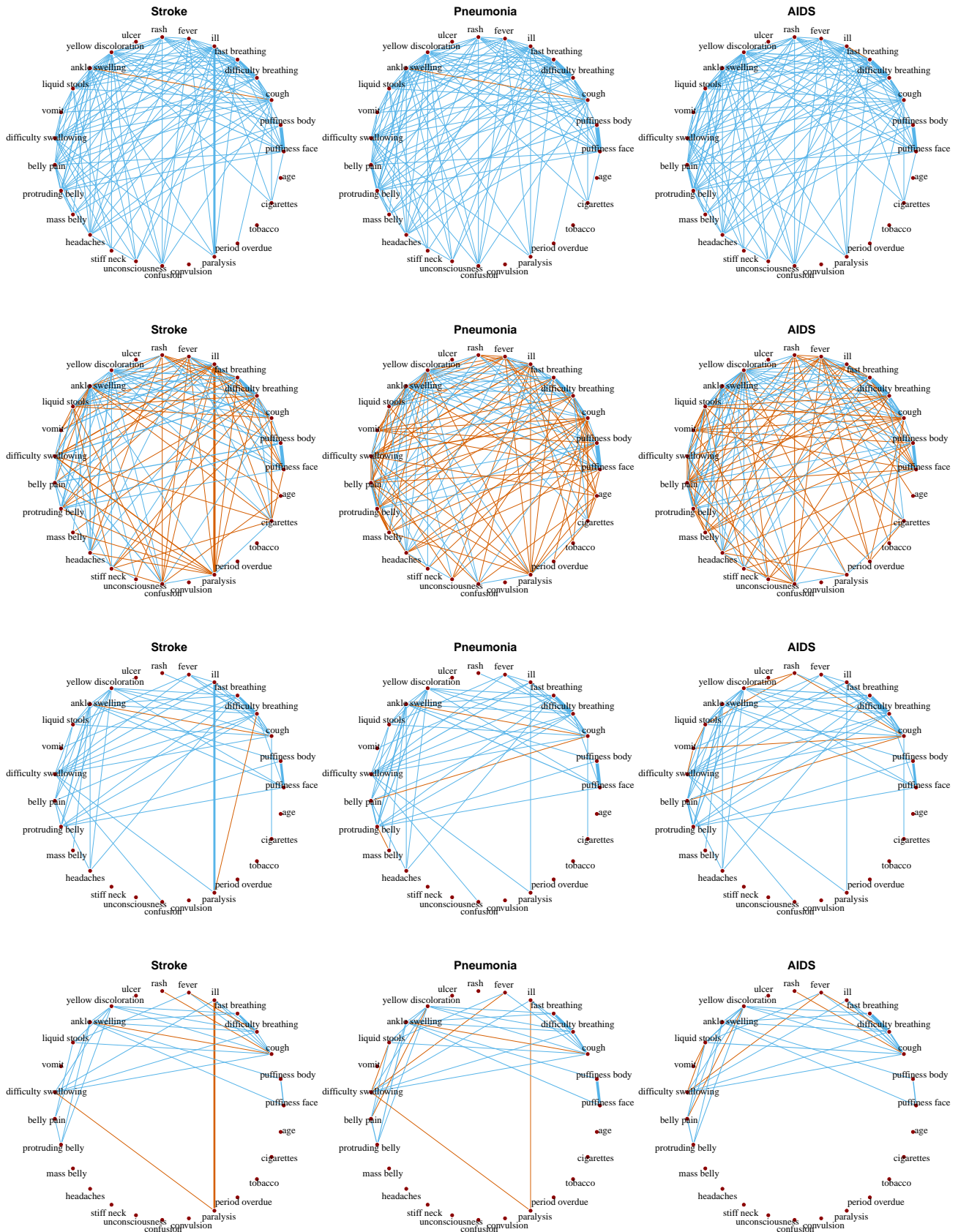


Figure 8. Estimated edges between the symptoms under the three causes using FGL using AIC (first row), GGL using AIC (second row), FGL with the same number of edges as selected by DSS-FGL (third row), and GGL with the same number of edges as selected by DSS-GGL (last row). The width of the edges are proportional to the size of  $|\omega_{jk}^{(g)}|$ . Common edges across all groups are colored in blue, and the differential edges are colored in red.

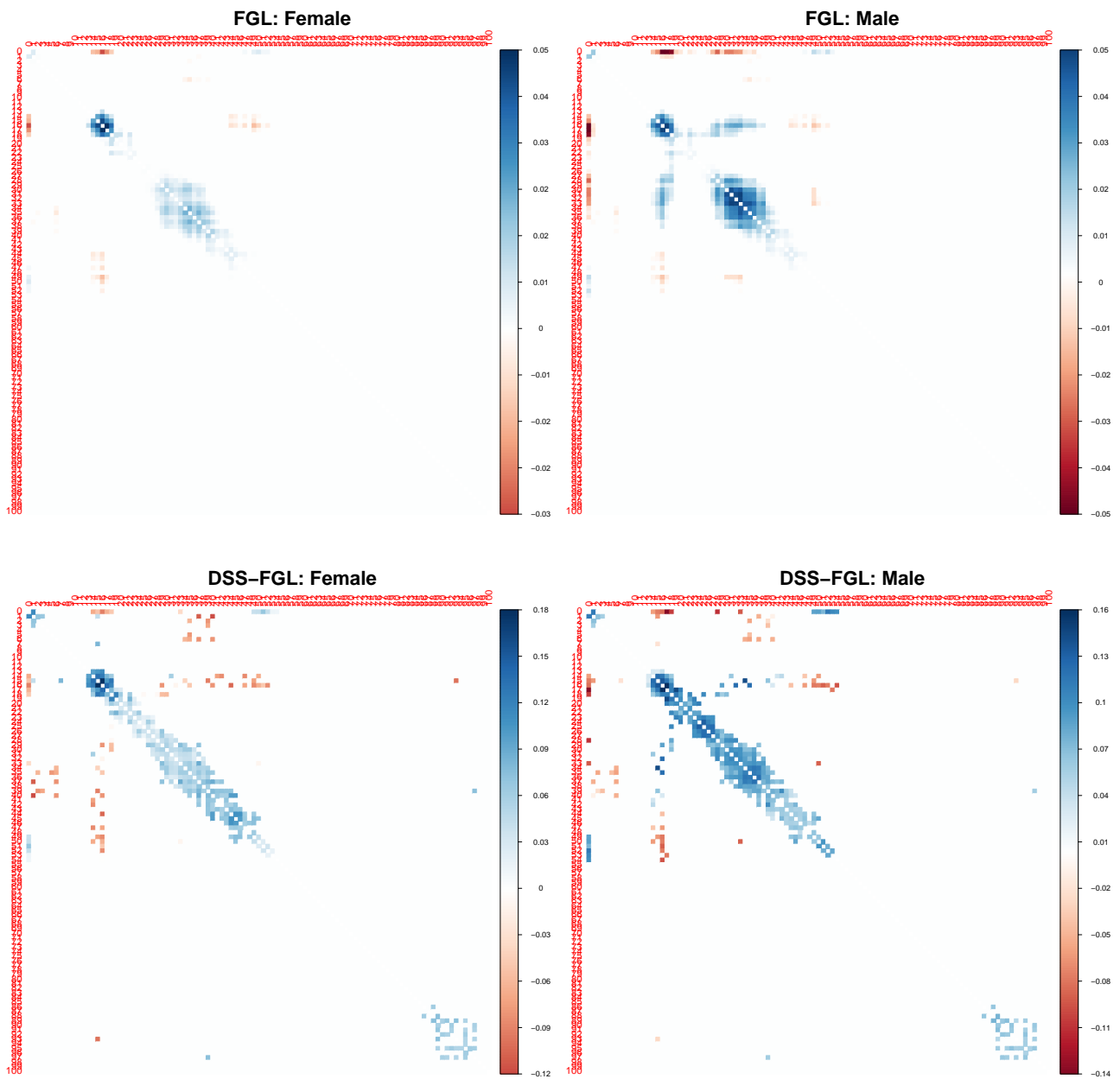


Figure 9. Estimated partial correlation matrix using one cross-validation dataset. The partial correlations among the 101 age groups are estimated using FGL with the same number of edges as selected by DSS-FGL (top row), and DSS-FGL (bottom row). DSS-FGL estimates 197 and 199 edges respectively for female and male. The closet configuration of FGL estimates 157 and 241 edges respectively. The precision matrices are rescaled and negated to partial correlations for easier interpretation.

## References

- Butt, Z., Haberman, S., and Shang, H. *ilc: Lee-Carter Mortality Models using Iterative Fitting Algorithms*, 2014. R package version 1.0.
- Peng, J., Wang, P., Zhou, N., and Zhu, J. Partial correlation estimation by joint sparse regression models. *Journal of the American Statistical Association*, 104(486):735–746, 2009.