

Supplementary Information A: Proofs of theorems

Proof of Theorem 1

Proof. We first define the following quantities.

$$\mathbf{x}_s \mathbf{x}_s^T = \sum_i^m \mathbf{x}_i \mathbf{x}_i^T + \sum_{i<j}^m (\mathbf{x}_i \mathbf{x}_j^T + \mathbf{x}_j \mathbf{x}_i^T) \quad (4)$$

$$\mathbf{x}_{ij} \mathbf{x}_{ij}^T = \mathbf{x}_i \mathbf{x}_i^T + \mathbf{x}_j \mathbf{x}_j^T - \mathbf{x}_i \mathbf{x}_j^T - \mathbf{x}_j \mathbf{x}_i^T \quad (5)$$

$$\sum_{i<j}^m \mathbf{x}_{ij} \mathbf{x}_{ij}^T = (m-1) \sum_i^m \mathbf{x}_i \mathbf{x}_i^T - \sum_{i<j}^m (\mathbf{x}_i \mathbf{x}_j^T + \mathbf{x}_j \mathbf{x}_i^T) \quad (6)$$

$$\sum_i^m \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{m} \left(\mathbf{x}_s \mathbf{x}_s^T + \sum_{i<j}^m \mathbf{x}_{ij} \mathbf{x}_{ij}^T \right) \quad \text{From Eqs. (4) and (6)} \quad (7)$$

To prove the theorem, we visit each term in the model in Theorem 1, $\mathbf{y}^T \mathbf{\Lambda} \mathbf{y}$, $\mathbf{x}^T \mathbf{\Theta} \mathbf{y}$, $|\mathbf{\Lambda}|$, and $\exp(-\frac{1}{2} \mathbf{x}^T \mathbf{\Theta} \mathbf{\Lambda}^{-1} \mathbf{\Theta}^T \mathbf{x})$ and re-write it as follows:

$$\begin{aligned} \text{tr}(\mathbf{y} \mathbf{y}^T \mathbf{\Lambda}) &= \text{tr} \left(\left(\sum_i^m \mathbf{y}_i \mathbf{y}_i^T \right) (\mathbf{\Omega} + \mathbf{\Delta}) + \left(\sum_{i<j}^m (\mathbf{y}_i \mathbf{y}_j^T + \mathbf{y}_j \mathbf{y}_i^T) \right) \mathbf{\Omega} \right) \\ &= \text{tr} \left(\left(\sum_i^m \mathbf{y}_i \mathbf{y}_i^T + \sum_{i<j}^m (\mathbf{y}_i \mathbf{y}_j^T + \mathbf{y}_j \mathbf{y}_i^T) \right) \mathbf{\Omega} + \left(\sum_i^m \mathbf{y}_i \mathbf{y}_i^T \right) \mathbf{\Delta} \right) \\ &= \text{tr}(\mathbf{y}_s \mathbf{y}_s^T \mathbf{\Omega} + \left(\sum_i^m \mathbf{y}_i \mathbf{y}_i^T \right) \mathbf{\Delta}) \quad \text{From Eq. (4)} \\ &= \text{tr}(\mathbf{y}_s \mathbf{y}_s^T \mathbf{\Omega} + \frac{1}{m} \left(\mathbf{y}_s \mathbf{y}_s^T + \sum_{i<j}^m \mathbf{y}_{ij} \mathbf{y}_{ij}^T \right) \mathbf{\Delta}) \quad \text{From Eq. (7)} \\ &= \text{tr}(\mathbf{y}_s \mathbf{y}_s^T \left(\mathbf{\Omega} + \frac{1}{m} \mathbf{\Delta} \right) + \left(\sum_{i<j}^m \mathbf{y}_{ij} \mathbf{y}_{ij}^T \right) \frac{1}{m} \mathbf{\Delta}) \\ \text{tr}(\mathbf{y} \mathbf{x}^T \mathbf{\Theta}) &= \text{tr} \left(\left(\sum_i^m \mathbf{y}_i \mathbf{x}_i^T \right) \mathbf{\Pi} + \left(\sum_i^m \mathbf{y}_i \mathbf{z}^T \right) \mathbf{\Xi} \right) \\ &= \text{tr}(\mathbf{y}_s \mathbf{x}_s^T \left(\frac{1}{m} \mathbf{\Pi} \right) + \left(\sum_{i<j}^m \mathbf{y}_{ij} \mathbf{x}_{ij}^T \right) \frac{1}{m} \mathbf{\Pi} + \mathbf{y}_s \mathbf{z}^T \mathbf{\Xi}) \\ \text{tr}(\mathbf{x} \mathbf{x}^T \mathbf{\Theta} \mathbf{\Sigma} \mathbf{\Theta}^T) &= \text{tr} \left(\mathbf{x}_s \mathbf{x}_s^T \left(\frac{1}{m} \mathbf{\Pi} \right) \left(\mathbf{\Omega} + \frac{1}{m} \mathbf{\Delta} \right)^{-1} \left(\frac{1}{m} \mathbf{\Pi} \right)^T \right. \\ &\quad + \left(\sum_{i<j}^m \mathbf{x}_{ij} \mathbf{x}_{ij}^T \right) \frac{1}{m} \mathbf{\Pi} \left(\frac{1}{m} \mathbf{\Delta} \right)^{-1} \frac{1}{m} \mathbf{\Pi}^T \\ &\quad \left. + \mathbf{z} \mathbf{z}^T \mathbf{\Xi} \left(\mathbf{\Omega} + \frac{1}{m} \mathbf{\Delta} \right)^{-1} \mathbf{\Xi}^T \right) \\ \det(\mathbf{\Lambda}) &= \det(m \mathbf{\Omega} + \mathbf{\Delta}) \det(\mathbf{\Delta})^{m-1} \end{aligned}$$

The equation above on determinants follows from Theorem 3 in Sylvester (2000). \square

Proof of Theorem 2

We first show that the following lemma holds, which is then used to prove Theorem 2.

Lemma 1. Let $\mathbf{K} = \mathbf{C}^+\mathbf{C} = \mathbf{C}^T\mathbf{C}^{+T} = \mathbf{K}^T$. Then, for $\Sigma = \Lambda^{-1} = [\mathbf{I}_{m \times m} \otimes \Delta + \mathbf{1}_{m \times m} \otimes \Omega]^{-1}$, we have $\mathbf{K}\Sigma = \Sigma\mathbf{K}$.

Proof. Using a block inversion of a matrix, we can show $\Sigma = \mathbf{I}_{m \times m} \otimes \mathbf{E} + \mathbf{1}_{m \times m} \otimes \Phi$, where $\Phi = -\Delta^{-1}\Omega(m\Omega + \delta)^{-1}$ and $\mathbf{E} = \Delta^{-1}$, and verify this by checking $\Sigma\Lambda = \mathbf{I}$. \mathbf{K} is an orthogonal projection matrix with $\mathbf{K}_{i,i} = \mathbf{K}_{i,2i} = \mathbf{K}_{2i,i} = \mathbf{K}_{2i,2i} = \frac{1}{2}$ for $i \in H$ for outputs with missing individual-level observations. For observed outputs $i \in V$, $\mathbf{K}_{i,i} = \mathbf{K}_{2i,2i} = 1$, and $\mathbf{K}_{i,2i} = \mathbf{K}_{2i,i} = 0$. Then, using the repeat block structure of Σ , we define all the elements of $\mathbf{K}\Sigma$ and $\Sigma\mathbf{K}$:

$$\begin{aligned} (\mathbf{K}\Sigma)_{i,i} &= (\mathbf{K}\Sigma)_{i,2i} = (\mathbf{K}\Sigma)_{2i,i} = (\mathbf{K}\Sigma)_{2i,2i} = \Phi_{i,i} + \frac{1}{2}\epsilon_i \text{ for } i \in H, \text{ and} \\ (\mathbf{K}\Sigma)_{i,i} &= (\mathbf{K}\Sigma)_{i,2i} = (\mathbf{K}\Sigma)_{2i,i} = (\mathbf{K}\Sigma)_{2i,2i} = \Phi_{i,i} + \epsilon_i \text{ for } i \in V. \end{aligned}$$

For $(i, j) \in \{i, j : i \neq j, i = 1, \dots, q, j = 1, \dots, q\}$, $(\mathbf{K}\Sigma)_{i,j} = (\mathbf{K}\Sigma)_{i,2j} = (\mathbf{K}\Sigma)_{2i,j} = (\mathbf{K}\Sigma)_{2i,2j} = \Phi_{i,j}$. The elements of $\Sigma\mathbf{K}$ are identical to those of $\mathbf{K}\Sigma$. \square

Now we prove Theorem 2.

Proof. From our model in Eq. (2), we derive the probability distribution for the observed outputs. We begin by re-writing Eq. (2) in the form of Gaussian distribution to make the marginal distribution explicitly represented:

$$p(\mathbf{y} | \mathbf{x}; \Lambda, \Theta) \sim \mathcal{N}(-\mathbf{C}_a\Lambda^{-1}\Theta^T\mathbf{x}, \mathbf{C}_a\Lambda^{-1}\mathbf{C}_a), \quad (8)$$

where $\mathbf{C}_a = \begin{bmatrix} \mathbf{I}_{q \times q} & \mathbf{0}_{q \times q} \\ \mathbf{0}_{q \times q} & \mathbf{I}_{q \times q} \end{bmatrix}$. Then, the output sum variables are also Gaussian distributed:

$$p(\mathbf{y}_s | \mathbf{x}; \Lambda, \Theta) \sim \mathcal{N}(-\mathbf{C}_s\Lambda^{-1}\Theta^T\mathbf{x}, \mathbf{C}_s\Lambda^{-1}\mathbf{C}_s^T), \quad (9)$$

where $\mathbf{C}_s = [\mathbf{I}_{q \times q} \ \mathbf{I}_{q \times q}]$.

We combine the marginal distributions for the observed output variables in V from Eq. (8) and for output variables in H with missing individual-level observations but only with group-level sum data from Eq. (9) to form a joint distribution of the observed variables \mathbf{y}_s^H and \mathbf{y}^V :

$$p(\mathbf{y}_s^H, \mathbf{y}^V | \mathbf{x}; \Lambda, \Theta) \sim \mathcal{N}(-\mathbf{C}\Lambda^{-1}\Theta^T\mathbf{x}, \mathbf{C}\Lambda^{-1}\mathbf{C}^T), \quad (10)$$

given $\mathbf{C} = \begin{bmatrix} \mathbf{C}_s^H \\ \mathbf{C}_a^V \end{bmatrix}$. \mathbf{C}_s^H and \mathbf{C}_a^V are the submatrices of \mathbf{C}_s and \mathbf{C}_a with only the rows corresponding to the variables in \mathbf{y}_s^H and \mathbf{y}^V . To show the distribution in Eq. (10) is a CGGM, we only need to show Eq. (10) can be written as

$$p(\mathbf{y}_s^H, \mathbf{y}^V | \mathbf{x}; \Lambda, \Theta) = \mathcal{N}(-\tilde{\Lambda}^{-1}\tilde{\Theta}^T\mathbf{x}, \tilde{\Lambda}^{-1}), \quad (11)$$

by obtaining the explicit forms of the mean and inverse covariance of the above, since expanding the quadratic term in this distribution above leads to the CGGM in Theorem 2.

We first show $\tilde{\Lambda} = [\mathbf{C}\Lambda^{-1}\mathbf{C}^T]^{-1} = \mathbf{C}^{+T}\Lambda\mathbf{C}^+$, where \mathbf{C}^+ is the Moore-Penrose inverse of \mathbf{C} , by showing $\mathbf{C}\Lambda^{-1}\mathbf{C}^T\mathbf{C}^{+T}\Lambda\mathbf{C}^+ = \mathbf{I}$:

$$\begin{aligned} \mathbf{C}\Sigma\mathbf{C}^T\mathbf{C}^{+T}\Lambda\mathbf{C}^+ &= \mathbf{C}\mathbf{C}^T\mathbf{C}^{+T}\Lambda^{-1}\Lambda\mathbf{C}^+ \quad \text{from Lemma 1} \\ &= \mathbf{C}\mathbf{C}^T\mathbf{C}^{+T}\mathbf{C}^+ \\ &= \mathbf{C}\mathbf{C}^+\mathbf{C}\mathbf{C}^+ \quad \text{since } \mathbf{C}^+\mathbf{C} \text{ is symmetric from the definition of Moore-Penrose inverse} \\ &= \mathbf{C}\mathbf{C}^+ \quad \text{from } \mathbf{C}\mathbf{C}^+ = \mathbf{I} \\ &= \mathbf{I} \end{aligned}$$

Then the mean of Eq. (10) can be written as

$$\begin{aligned}
 -\mathbf{C}\boldsymbol{\Lambda}^{-1}\boldsymbol{\Theta}^T\mathbf{x} &= -\mathbf{C}\mathbf{C}^+\mathbf{C}\boldsymbol{\Lambda}^{-1}\boldsymbol{\Theta}^T\mathbf{x} && \text{from } \mathbf{C} = \mathbf{C}\mathbf{C}^+\mathbf{C} \\
 &= -\mathbf{C}\boldsymbol{\Lambda}^{-1}\mathbf{C}^+\mathbf{C}\boldsymbol{\Theta}^T\mathbf{x} && \text{from Lemma 1} \\
 &= -\mathbf{C}\boldsymbol{\Lambda}^{-1}\mathbf{C}^T\mathbf{C}^{+T}\boldsymbol{\Theta}^T\mathbf{x} && \text{from } \mathbf{C}^+\mathbf{C} = \mathbf{C}^T\mathbf{C}^{+T} \text{ since } \mathbf{C}^+\mathbf{C} \text{ is symmetric} \\
 &= -\tilde{\boldsymbol{\Lambda}}^{-1}\tilde{\boldsymbol{\Theta}}^T\mathbf{x} && \tilde{\boldsymbol{\Theta}} = \boldsymbol{\Theta}\mathbf{C}^+
 \end{aligned}$$

Thus, we have Eq. (11) with $\tilde{\boldsymbol{\Theta}} = \boldsymbol{\Theta}\mathbf{C}^+$ and $\tilde{\boldsymbol{\Lambda}} = \mathbf{C}^{+T}\boldsymbol{\Lambda}\mathbf{C}^+$. □

The corollary below follows from the theorem above.

Corollary 2. *Learning the collapsed multi-level CGGM in Eq. (3) with randomly missing individual-level output data is equivalent to learning the full multi-level model with imputed data $\tilde{\mathbf{y}} = \mathbf{C}^+ \begin{bmatrix} \mathbf{y}_s^H \\ \mathbf{y}_v \end{bmatrix}$.*

Proof. From Theorem 2, given $\tilde{\boldsymbol{\Lambda}} = \mathbf{C}^{+T}\boldsymbol{\Lambda}\mathbf{C}^+$ and $\tilde{\boldsymbol{\Theta}}^T = \mathbf{C}^{+T}\boldsymbol{\Theta}^T$, we can re-write the log probability as

$$\begin{aligned}
 \log p(\mathbf{y}_s^H, \mathbf{y}_v^V | \mathbf{x}; \boldsymbol{\Lambda}, \boldsymbol{\Theta}) &= -\frac{1}{2} \begin{bmatrix} \mathbf{y}_s^H \\ \mathbf{y}_v^V \end{bmatrix}^T \tilde{\boldsymbol{\Lambda}} \begin{bmatrix} \mathbf{y}_s^H \\ \mathbf{y}_v^V \end{bmatrix} - \mathbf{x}^T \tilde{\boldsymbol{\Theta}} \begin{bmatrix} \mathbf{y}_s^H \\ \mathbf{y}_v^V \end{bmatrix} - \log Z(\mathbf{x}) \\
 &= -\frac{1}{2} \begin{bmatrix} \mathbf{y}_s^H \\ \mathbf{y}_v \end{bmatrix}^T \mathbf{C}^{+T}\boldsymbol{\Lambda}\mathbf{C}^+ \begin{bmatrix} \mathbf{y}_s^H \\ \mathbf{y}_v \end{bmatrix} - \mathbf{x}^T \boldsymbol{\Theta}\mathbf{C}^+ \begin{bmatrix} \mathbf{y}_s^H \\ \mathbf{y}_v \end{bmatrix} - \log Z(\mathbf{x}) \\
 &= -\frac{1}{2} \left[\mathbf{C}^+ \begin{bmatrix} \mathbf{y}_s^H \\ \mathbf{y}_v \end{bmatrix} \right]^T \boldsymbol{\Lambda} \left[\mathbf{C}^+ \begin{bmatrix} \mathbf{y}_s^H \\ \mathbf{y}_v \end{bmatrix} \right] - \mathbf{x}^T \boldsymbol{\Theta} \left[\mathbf{C}^+ \begin{bmatrix} \mathbf{y}_s^H \\ \mathbf{y}_v \end{bmatrix} \right] - \log Z(\mathbf{x}) \\
 &= -\frac{1}{2} \tilde{\mathbf{y}}^T \tilde{\boldsymbol{\Lambda}} \tilde{\mathbf{y}} - \mathbf{x}^T \tilde{\boldsymbol{\Theta}} \tilde{\mathbf{y}} - \log Z(\mathbf{x}).
 \end{aligned}$$
□

Supplementary Information B: Details of the Optimization Method for Multi-level Conditional Gaussian Graphical Models

Alternating Newton coordinate descent

We provide the details of the learning algorithm for our multi-level model, using the representation of the model in Theorem 1. The L_1 -regularized negative log-likelihood is given as

$$\begin{aligned}
 \operatorname{argmin}_{\boldsymbol{\Omega} > 0, \boldsymbol{\Delta}, \boldsymbol{\Xi}, \boldsymbol{\Pi}} & -\log |\boldsymbol{\Omega} + \frac{1}{m}\boldsymbol{\Delta}| + \operatorname{tr} \left(\mathbf{S}_{\mathbf{y}_s} \left(\boldsymbol{\Omega} + \frac{1}{m}\boldsymbol{\Delta} \right) + 2\mathbf{S}_{\mathbf{y}_s\mathbf{x}_s} \frac{1}{m}\boldsymbol{\Pi} + 2\mathbf{S}_{\mathbf{y}_s\mathbf{z}} \boldsymbol{\Xi} \right. \\
 & \left. + \left(\boldsymbol{\Omega} + \frac{1}{m}\boldsymbol{\Delta} \right)^{-1} \left(\frac{1}{m^2}\boldsymbol{\Pi}^T \mathbf{S}_{\mathbf{x}_s} \boldsymbol{\Pi} + \frac{1}{m}\boldsymbol{\Pi}^T \mathbf{S}_{\mathbf{x}_s\mathbf{z}} \boldsymbol{\Xi} + \frac{1}{m}\boldsymbol{\Xi}^T \mathbf{S}_{\mathbf{x}_s\mathbf{z}}^T \boldsymbol{\Pi} + \boldsymbol{\Xi}^T \mathbf{S}_{\mathbf{z}} \boldsymbol{\Xi} \right) \right) \\
 & - (m-1) \log |\boldsymbol{\Delta}| + \operatorname{tr} \left(\mathbf{S}_{\mathbf{y}_d} \frac{1}{m}\boldsymbol{\Delta} + 2\mathbf{S}_{\mathbf{y}_d\mathbf{x}_d} \frac{1}{m}\boldsymbol{\Pi} + \frac{1}{m}\boldsymbol{\Delta}^{-1} \boldsymbol{\Pi}^T \mathbf{S}_{\mathbf{x}_d} \boldsymbol{\Pi} \right) \\
 & + \lambda_{\boldsymbol{\Omega}} \|\boldsymbol{\Omega}\|_1 + \lambda_{\boldsymbol{\Delta}} \|\boldsymbol{\Delta}\|_1 + \lambda_{\boldsymbol{\Pi}} \|\boldsymbol{\Pi}\|_1 + \lambda_{\boldsymbol{\Xi}} \|\boldsymbol{\Xi}\|_1,
 \end{aligned} \tag{12}$$

using components of covariances $\mathbf{S}_{\mathbf{y}_w\mathbf{y}_s} = [\mathbf{S}_{\mathbf{y}_s\mathbf{x}_s} \mathbf{S}_{\mathbf{y}_s\mathbf{z}}]$ and $\mathbf{S}_w = \begin{bmatrix} \mathbf{S}_{\mathbf{x}_s} & \mathbf{S}_{\mathbf{x}_s\mathbf{z}_s} \\ \mathbf{S}_{\mathbf{x}_s\mathbf{z}_s}^T & \mathbf{S}_{\mathbf{z}} \end{bmatrix}$.

In each iteration, we alternately optimize for each of the parameters $\boldsymbol{\Omega}$, $\boldsymbol{\Delta}$, $\boldsymbol{\Xi}$, and $\boldsymbol{\Pi}$ by solving the optimization problem above for each of the parameter given the previous estimates of all the other parameters. We maintain and update $\mathbf{V} = \boldsymbol{\Omega} + \frac{1}{m}\boldsymbol{\Delta}$ after each update of $\boldsymbol{\Omega}$ and $\boldsymbol{\Delta}$ in order to avoid repeatedly computing $\boldsymbol{\Omega} + \frac{1}{m}\boldsymbol{\Delta}$. In the rest of this section, we write $\mathbf{V} = \boldsymbol{\Omega} + \frac{1}{m}\boldsymbol{\Delta}$ whenever this term occurs.

Algorithm 1: Alternating Newton Coordinate Descent for Multi-level Conditional Gaussian Graphical Models

Input: $\mathbf{S}_{y_d}, \mathbf{S}_{y_{x_d}}, \mathbf{S}_{x_d}, \mathbf{Y}_s, \mathbf{X}_s, \mathbf{X}_d, \mathbf{Z}$, regularization parameters $\lambda_\Omega, \lambda_\Pi, \lambda_\Xi$, and α

Initialize: $\Pi, \Xi \leftarrow 0, \Omega, \Delta \leftarrow \mathbf{I}_{q \times q}$

for $t = 0, 1, \dots$ **do**

Determine active sets $\mathcal{S}_\Omega, \mathcal{S}_\Delta, \mathcal{S}_\Pi, \mathcal{S}_\Xi$

Find Newton directions via coordinate descent

$$D_\Omega = \underset{\delta_\Omega}{\operatorname{argmin}} \bar{g}(\Omega + \delta_\Omega, \Delta, \Pi, \Xi) + h(\Omega + \delta_\Omega, \Delta, \Pi, \Xi)$$

Update $\Omega^+ = \Omega + \alpha D_\Omega$, where α is found via line search

$$D_\Delta = \underset{\delta_\Delta}{\operatorname{argmin}} \bar{g}(\Omega, \Delta + \delta_\Delta, \Pi, \Xi) + h(\Omega, \Delta + \delta_\Delta, \Pi, \Xi)$$

Update $\Delta^+ = \Delta + \alpha D_\Delta$, where α is found via line search

Solve via coordinate descent:

$$\Pi^+ = \underset{\Pi}{\operatorname{argmin}} g(\Pi) + h(\Pi)$$

$$\Xi^+ = \underset{\Xi}{\operatorname{argmin}} g(\Xi) + h(\Xi)$$

Update for Ξ : Optimizing for Ξ given all the other parameters corresponds to solving the following optimization problem:

$$\underset{\Xi}{\operatorname{argmin}} g_{\Omega, \Delta, \Pi}(\Xi) + \lambda_\Xi \|\Xi\|_1,$$

where

$$g_{\Omega, \Delta, \Pi}(\Xi) = \operatorname{tr} \left(2\mathbf{S}_{yz}\Xi + \mathbf{V}^{-1} \left(\frac{1}{m} \Pi^T \mathbf{S}_{xz} \Xi + \frac{1}{m} \Xi^T \mathbf{S}_{xz}^T \Pi + \Xi^T \mathbf{S}_z \Xi \right) \right).$$

Update for Π : Optimizing for Π given all the other parameters corresponds to solving the following optimization problem:

$$\underset{\Pi}{\operatorname{argmin}} g_{\Omega, \Delta, \Xi}(\Pi) + \lambda_\Pi \|\Pi\|_1,$$

where

$$g_{\Omega, \Delta, \Xi}(\Pi) = \operatorname{tr} \left(2\mathbf{S}_{y_{x_s}} \frac{1}{m} \Pi + 2\mathbf{S}_{y_{x_d}} \frac{1}{m} \Pi + \mathbf{V}^{-1} \left(\frac{1}{m^2} \Pi^T \mathbf{S}_{x_s} \Pi + \frac{1}{m} \Pi^T \mathbf{S}_{xz} \Xi + \frac{1}{m} \Xi^T \mathbf{S}_{xz}^T \Pi \right) + \frac{1}{m} \Delta^{-1} \Pi^T \mathbf{S}_{x_d} \Pi \right).$$

The two problems above for Ξ and Π are Lasso, which can be solved efficiently using a coordinate descent algorithm.

Update for Ω : To optimize Eq. (12) for Ω given all the other parameters, we use the Newton's method. We first find the Newton descent direction by minimizing the second-order approximation of the objective in Eq. (12) with respect to Ω and then update Ω using this Newton direction and step size found by line search. The Newton direction D_Ω is found by solving the following optimization problem:

$$D_\Omega = \underset{\delta_\Omega}{\operatorname{argmin}} \bar{g}_{\Delta, \Xi, \Pi}(\delta_\Omega) + \lambda_\Omega \|\Omega + \delta_\Omega\|_1,$$

where $\bar{g}_{\Delta, \Xi, \Pi}(\delta_\Omega)$ is the second-order Taylor expansion of the data log-likelihood in Eq. (12) with respect to Ω and is given as

$$\bar{g}_{\Delta, \Xi, \Pi}(\delta_\Omega) = \operatorname{vec}(\nabla_\Omega g_{\Delta, \Xi, \Pi}(\Omega))^T \operatorname{vec}(\delta_\Omega) + \frac{1}{2} \operatorname{vec}(\delta_\Omega)^T \nabla_\Omega^2 g_{\Delta, \Xi, \Pi}(\Omega) \operatorname{vec}(\delta_\Omega),$$

with the gradient and Hessian

$$\nabla_\Omega g_{\Delta, \Xi, \Pi}(\Omega) = \mathbf{S}_{y_s} - \mathbf{V}^{-1} - \mathbf{V}^{-1} \left(\frac{1}{m^2} \Pi^T \mathbf{S}_{x_s} \Pi + \frac{1}{m} \Pi^T \mathbf{S}_{xz} \Xi + \frac{1}{m} \Xi^T \mathbf{S}_{xz}^T \Pi + \Xi^T \mathbf{S}_z \Xi \right) \mathbf{V}^{-1}$$

$$\nabla_\Omega^2 g_{\Delta, \Xi, \Pi}(\Omega) = \mathbf{V}^{-1} \otimes \left(\mathbf{V}^{-1} + 2\mathbf{V}^{-1} \left(\frac{1}{m^2} \Pi^T \mathbf{S}_{x_s} \Pi + \frac{1}{m} \Pi^T \mathbf{S}_{xz} \Xi + \frac{1}{m} \Xi^T \mathbf{S}_{xz}^T \Pi + \Xi^T \mathbf{S}_z \Xi \right) \mathbf{V}^{-1} \right).$$

The above optimization problem is a Lasso problem and can be solved using a coordinate descent method.

Update for Δ : To optimize Eq. (12) for Δ given all the other parameters, we again use the Newton's method. We find the Newton descent direction D_Δ by minimizing a second-order approximation of the objective in Eq. (12) with respect to Δ and update Δ with this Newton direction and a step size found by line search. The Newton descent direction D_Δ can be found by solving the following:

$$D_\Delta = \underset{\delta_\Delta}{\operatorname{argmin}} \bar{g}_{\Omega, \Xi, \Pi}(\delta_\Delta),$$

where

$$\bar{g}_{\Omega, \Xi, \Pi}(\delta_\Delta) = \operatorname{vec}(\nabla_\Delta g_{\Omega, \Xi, \Pi}(\Delta))^T \operatorname{vec}(\delta_\Delta) + \frac{1}{2} \operatorname{vec}(\delta_\Delta)^T \nabla_\Delta^2 g_{\Omega, \Xi, \Pi}(\Delta) \operatorname{vec}(\delta_\Delta),$$

with the gradient and Hessian

$$\begin{aligned} \nabla_\Delta g_{\Omega, \Xi, \Pi}(\Delta) &= \frac{1}{m} \left(\mathbf{S}_{\mathbf{y}_s} + \mathbf{S}_{\mathbf{y}_d} - \mathbf{V}^{-1} - (m-1)\Gamma^{-1} \right. \\ &\quad - \mathbf{V}^{-1} \left(\frac{1}{m^2} \Pi^T \mathbf{S}_{\mathbf{x}_s} \Pi + \frac{1}{m} \Pi^T \mathbf{S}_{\mathbf{x}_z} \Xi + \frac{1}{m} \Xi^T \mathbf{S}_{\mathbf{x}_z}^T \Pi + \Xi^T \mathbf{S}_z \Xi \right) \mathbf{V}^{-1} \\ &\quad \left. - \frac{1}{m^2} \Delta^{-1} \Pi^T \mathbf{S}_{\mathbf{x}_d} \Pi \Delta^{-1} \right) \\ \nabla_\Delta^2 g_{\Omega, \Xi, \Pi}(\Delta) &= \frac{1}{m^2} \left[\mathbf{V}^{-1} \otimes (\mathbf{V}^{-1} + 2\mathbf{V}^{-1} \left(\frac{1}{m^2} \Pi^T \mathbf{S}_{\mathbf{x}_s} \Pi + \frac{1}{m} \Pi^T \mathbf{S}_{\mathbf{x}_z} \Xi + \frac{1}{m} \Xi^T \mathbf{S}_{\mathbf{x}_z}^T \Pi + \Xi^T \mathbf{S}_z \Xi \right) \mathbf{V}^{-1}) \right. \\ &\quad \left. + \Gamma^{-1} \otimes \left((m-1)\Gamma^{-1} + \frac{2}{m^2} \Delta^{-1} \Pi^T \mathbf{S}_{\mathbf{x}_d} \Pi \Delta^{-1} \right) \right]. \end{aligned}$$

This optimization problem for finding a Newton direction is again a Lasso problem, which can be solved using a coordinate descent algorithm.

In order to improve the efficiency of the coordinate descent algorithm, we extend the strategies used in McCarter and Kim (2016). First, we restrict the coordinate descent updates to the active set of variables given as

$$\begin{aligned} \mathcal{S}_\Omega &= \{(\delta_\Omega)_{ij} : |(\nabla_\Omega g(\Omega))_{ij}| > \lambda_\Omega \vee \Omega_{ij} \neq 0\} \\ \mathcal{S}_\Pi &= \{\Pi_{ij} : |(\nabla_\Pi g(\Pi))_{ij}| > \lambda_\Pi \vee \Pi_{ij} \neq 0\} \\ \mathcal{S}_\Xi &= \{\Xi_{ij} : |(\nabla_\Xi g(\Xi))_{ij}| > \lambda_\Xi \vee \Xi_{ij} \neq 0\} \\ \mathcal{S}_\Delta &= \{(\delta_\Delta)_{ij} : |(\nabla_\Delta g(\Delta))_{ij}| > \lambda_\Delta \vee \Delta_{ij} \neq 0\}. \end{aligned}$$

Second, we compute and store the intermediate results at the beginning of the coordinate descent updates for Ξ and Π and for δ_Ω and δ_Δ in each Newton iteration. We store the intermediate results $U_\Xi := \Xi \mathbf{V}^{-1}$ for Ξ , $U_{\Pi_1} := \Pi \mathbf{V}^{-1}$ and $U_{\Pi_2} := \Pi \Delta^{-1}$ for Π , $U_\Omega := \delta_\Omega \mathbf{V}^{-1}$ for δ_Ω , and $U_{\delta_1} := \delta_\Delta \mathbf{V}^{-1}$ and $U_{\delta_2} := \delta_\Delta \Delta^{-1}$ for δ_Δ . When Ξ_{ij} and Π_{ij} are updated, the i th row of the corresponding intermediate results is updated. After $(\delta_\Omega)_{ij}$ update, the i th and j th rows of U_Ω are updated and after $(\delta_\Delta)_{ij}$ update, the i th and j th rows of U_Δ is updated.

Alternating Newton block coordinate descent

We further adopt the block-wise optimization method in McCarter and Kim (2016) to remove the memory requirement for storing large dense matrices during coordinate descent optimization. The coordinate descent updates require precomputing and storing the large dense matrices, such as $\mathbf{S}_{\mathbf{x}_s}$, $\mathbf{S}_{\mathbf{x}_d}$, and \mathbf{V}^{-1} . Instead, we perform a block-wise update of Ξ , Π , δ_Ω , and δ_Δ , precomputing and reusing the portions of the large matrices that are required for updating the current block of the parameters.

Blockwise optimization for Ω : A coordinate-descent optimization for $(\delta_\Omega)_{ij}$ requires computing the i th and j th columns of large dense $q \times q$ matrices \mathbf{V}^{-1} and $\mathbf{V}^{-1} \left(\frac{1}{m^2} \Pi^T \mathbf{S}_{\mathbf{x}_s} \Pi + \frac{1}{m} \Pi^T \mathbf{S}_{\mathbf{x}_z} \Xi + \frac{1}{m} \Xi^T \mathbf{S}_{\mathbf{x}_z}^T \Pi + \Xi^T \mathbf{S}_z \Xi \right) \mathbf{V}^{-1}$ found in the gradient and Hessian of the objective with respect to Ω . We represent $\mathbf{V}^{-1} \left(\frac{1}{m^2} \Pi^T \mathbf{S}_{\mathbf{x}_s} \Pi + \frac{1}{m} \Pi^T \mathbf{S}_{\mathbf{x}_z} \Xi + \frac{1}{m} \Xi^T \mathbf{S}_{\mathbf{x}_z}^T \Pi + \Xi^T \mathbf{S}_z \Xi \right) \mathbf{V}^{-1}$ as $(R_1 + R_2 + R_3)^T (R_1 + R_2 + R_3)$, where

$R_1 = Q_1 \mathbf{V}^{-1}$, $R_2 = Q_2 \mathbf{V}^{-1}$, and $R_3 = Q_3 \mathbf{V}^{-1}$. We precompute and store $n \times q$ matrices, $Q_1 = \frac{1}{m} \mathbf{X}_s \mathbf{\Pi}$, $Q_2 = \mathbf{X}_s \mathbf{\Xi}$, and $Q_3 = \mathbf{Z} \mathbf{\Xi}$, which are used to update R_1 , R_2 , and R_3 after each update of \mathbf{V} . Instead of storing \mathbf{V}^{-1} and $(R_1 + R_2 + R_3)^T (R_1 + R_2 + R_3)$ persistently, we perform blockwise coordinate descent, whereby the active set \mathcal{S}_Ω of δ_Ω is clustered into smaller blocks, and we perform coordinate-descent updates for elements in the active set in each block. We cluster the rows and columns of δ_Ω by partitioning $\{1, \dots, q\}$ into k_Ω sets, C_1, \dots, C_{k_Ω} . For each (C_r, C_z) block of δ_Ω , we do blockwise computation for $(\mathbf{V}^{-1})_{(:,C_r)}$ and $(\mathbf{V}^{-1})_{(:,C_z)}$, and compute the i th column in each block as $\mathbf{V}(\mathbf{V}^{-1})_{(:,i)} = \mathbf{e}_i$ using conjugate gradient method, where \mathbf{e}_i is a vector of q 0's except for 1 in the i th element. Similarly for $[(R_1 + R_2 + R_3)^T (R_1 + R_2 + R_3)]_{(:,C_r)}$, and $[(R_1 + R_2)^T (R_1 + R_2)]_{(:,C_z)}$, we obtain the i th column in each block using $(R_1 + R_2 + R_3)^T (R_1 + R_2 + R_3)_{(:,i)}$. After updating each $(\delta_\Omega)_{ij}$, we update the corresponding i th and j th rows of the intermediate results $(U_\Omega)_{(:,C_z)}$ and $(U_\Omega)_{(:,C_r)}$.

Blockwise optimization for Δ : A coordinate descent optimization for $(\delta_\Delta)_{ij}$ requires computing the i th and j th columns of large dense $q \times q$ matrices \mathbf{V}^{-1} , $(R_1 + R_2 + R_3)^T (R_1 + R_2 + R_3)$, Δ^{-1} , and $\Delta^{-1} \mathbf{\Pi}^T \mathbf{S}_{\mathbf{x}_d} \mathbf{\Pi} \Delta^{-1}$ found in the gradient and Hessian of the objective with respect to Δ . The optimization method for Δ is similar to the one for Ω . We represent $\frac{1}{m^2} \Delta^{-1} \mathbf{\Pi}^T \mathbf{S}_{\mathbf{x}_d} \mathbf{\Pi} \Delta^{-1}$ as $R_4^T R_4$, where $R_4 = Q_4 \Delta^{-1}$. We precompute and store a $n \times q$ matrix, $Q_4 = \frac{1}{m} \mathbf{X}_d \mathbf{\Pi}$, which is used to update R_4 after each update of Δ .

Blockwise optimization for Ξ : The coordinate descent update for $(\Xi)_{ij}$ requires the i th column of $\mathbf{S}_{\mathbf{x}_s}$ and \mathbf{S}_z and the j th column of \mathbf{V}^{-1} . Instead of storing $\mathbf{S}_{\mathbf{x}_s}$, \mathbf{S}_z , and \mathbf{V} persistently, we perform blockwise coordinate descent. We cluster the columns of Ξ by partitioning $\{1, \dots, q\}$ into k_Ξ sets, C_1, \dots, C_{k_Ξ} . For each block (i, C_r) , where $i \in \{1, \dots, p+r\}$, we compute $(\mathbf{S}_{\mathbf{x}_s})_{(:,i)}$, $(\mathbf{S}_z)_{(:,i)}$, and $(\mathbf{V}^{-1})_{(:,C_r)}$. After updating each $(\Xi)_{ij}$, we update the corresponding i th row of the intermediate result $(U_\Xi)_{(:,C_r)}$.

Blockwise optimization for Π : The coordinate descent update for $(\Pi)_{ij}$ requires the i th column of $\mathbf{S}_{\mathbf{x}_s}$ and $\mathbf{S}_{\mathbf{x}_d}$ and the j th column of \mathbf{V}^{-1} and Δ^{-1} . Instead of storing $\mathbf{S}_{\mathbf{x}_s}$, $\mathbf{S}_{\mathbf{x}_d}$, \mathbf{V} , and Δ persistently, we perform blockwise coordinate descent. We cluster the columns of Π by partitioning $\{1, \dots, q\}$ into k_Π sets, C_1, \dots, C_{k_Π} . For each block (i, C_r) , where $i \in \{1, \dots, p\}$, we compute $(\mathbf{S}_{\mathbf{x}_s})_{(:,i)}$, $(\mathbf{S}_{\mathbf{x}_d})_{(:,i)}$, $(\mathbf{V}^{-1})_{(:,C_r)}$, and $(\Delta^{-1})_{(:,C_r)}$. After updating each $(\Pi)_{ij}$, we update the corresponding i th row of the intermediate results $(U_{\Pi_1})_{(:,C_r)}$ and $(U_{\Pi_2})_{(:,C_r)}$.

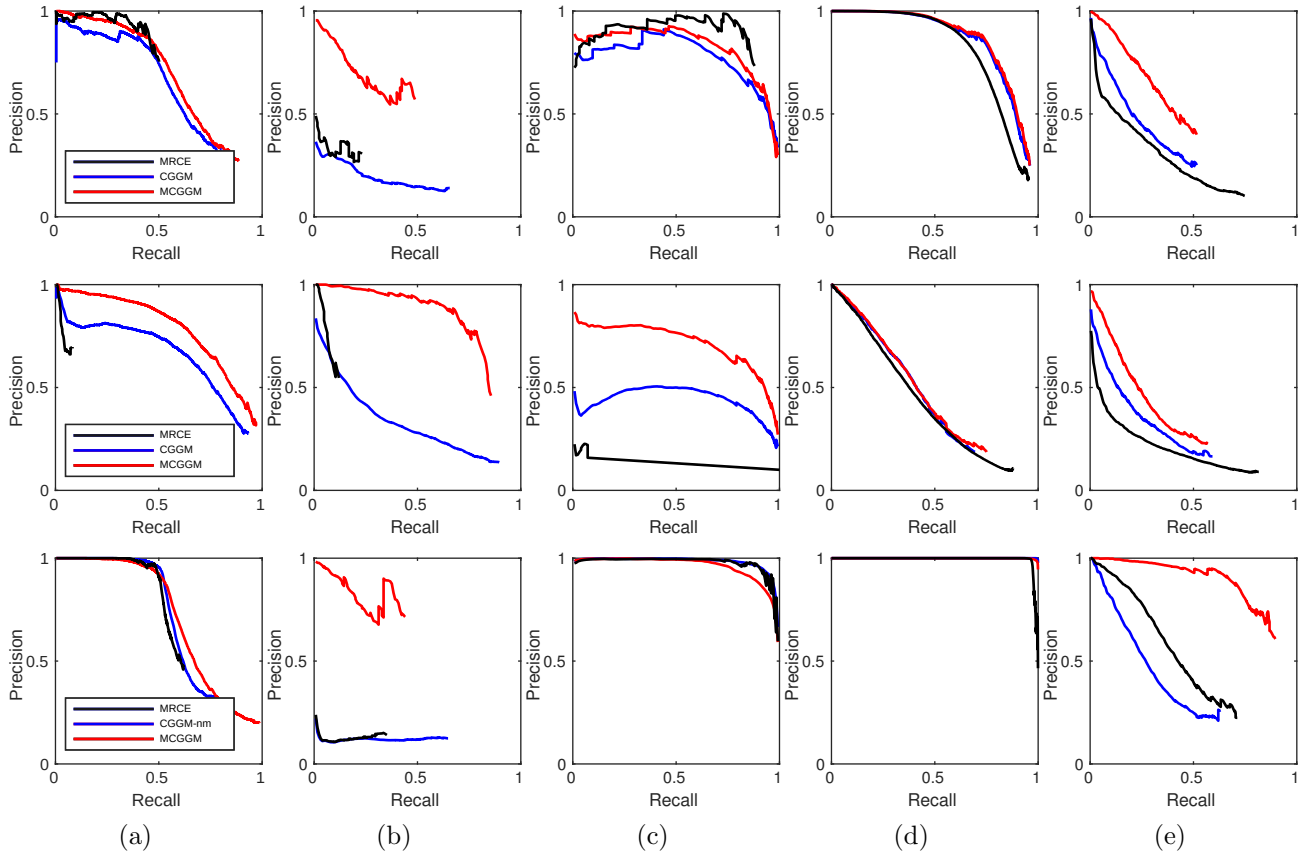


Figure 4: Comparison of methods on simulated data. Precision-recall curves for the recovery of (a) $\Omega + \Delta$, (b) Ω , and (c) Δ , (d) Π , and (e) Ξ . Top row for Case 1, middle row for Case 2, and bottom row for Case 3 with group size $m = 6$. Ω is set as a clustered network.

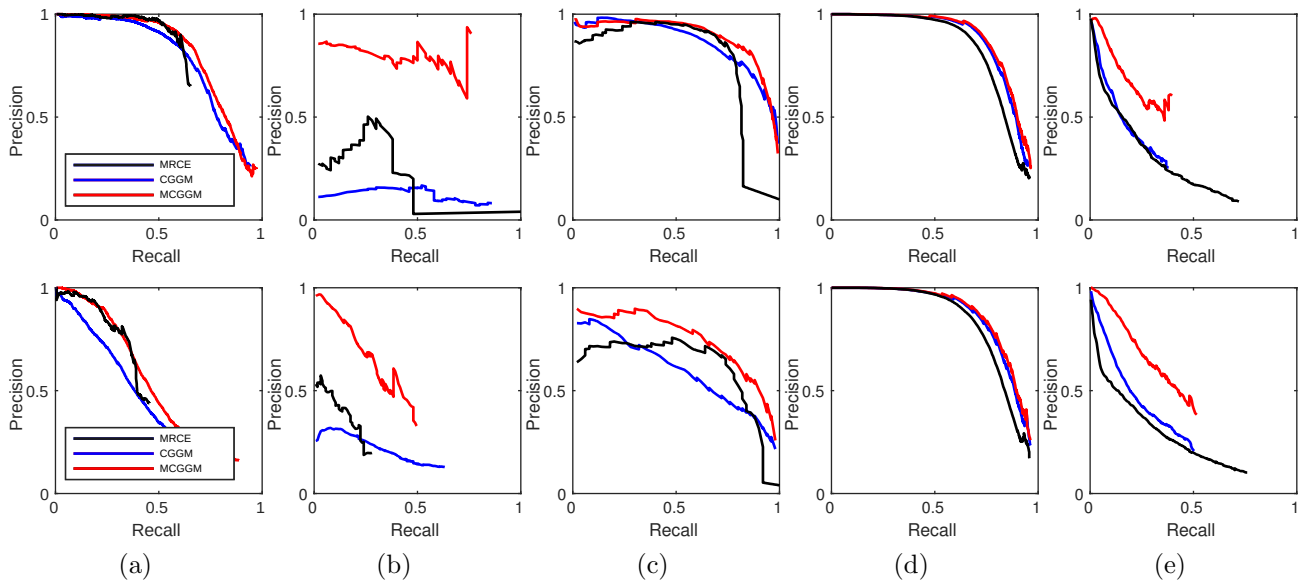


Figure 5: Comparison of methods on simulated data. Precision-recall curves for the recovery of (a) $\Omega + \Delta$, (b) Ω , and (c) Δ , (d) Π , and (e) Ξ . The results from Case 4 (top) and Case 5 (bottom) are shown for group size $m = 6$ and for scale-free networks for Ω and Δ .

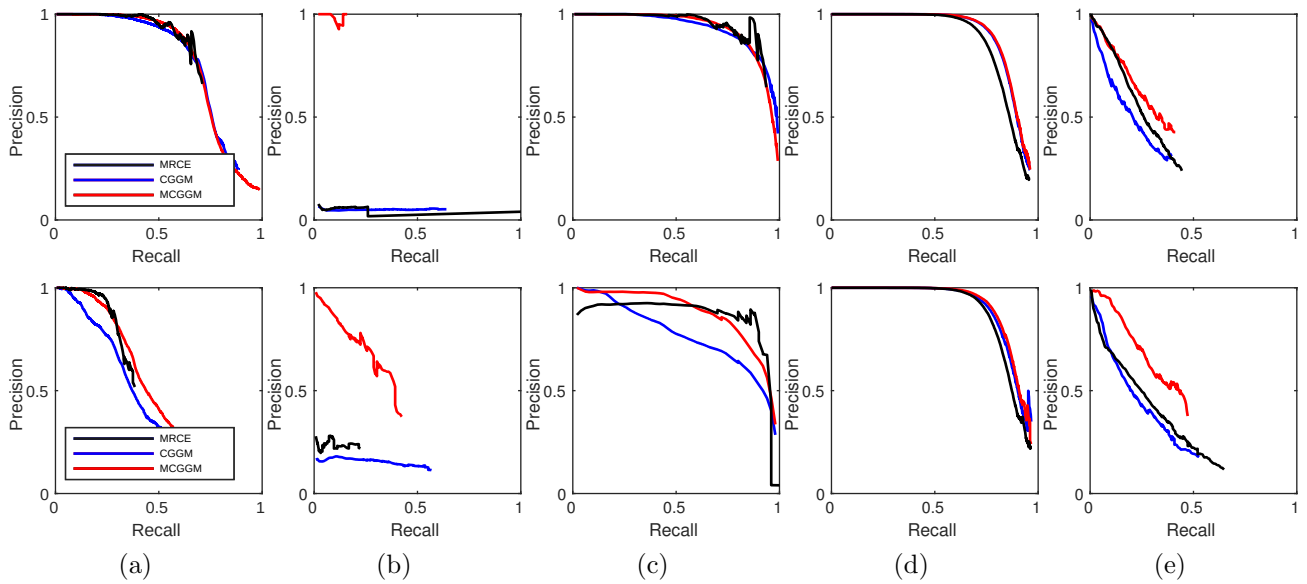


Figure 6: Comparison of methods on simulated data. Precision-recall curves for the recovery of (a) $\Omega + \Delta$, (b) Ω , and (c) Δ , (d) Π , and (e) Ξ . Results from Case 4 (top) and Case 5 (bottom) are shown for group size $m = 6$ and for Ω with a clustered network.

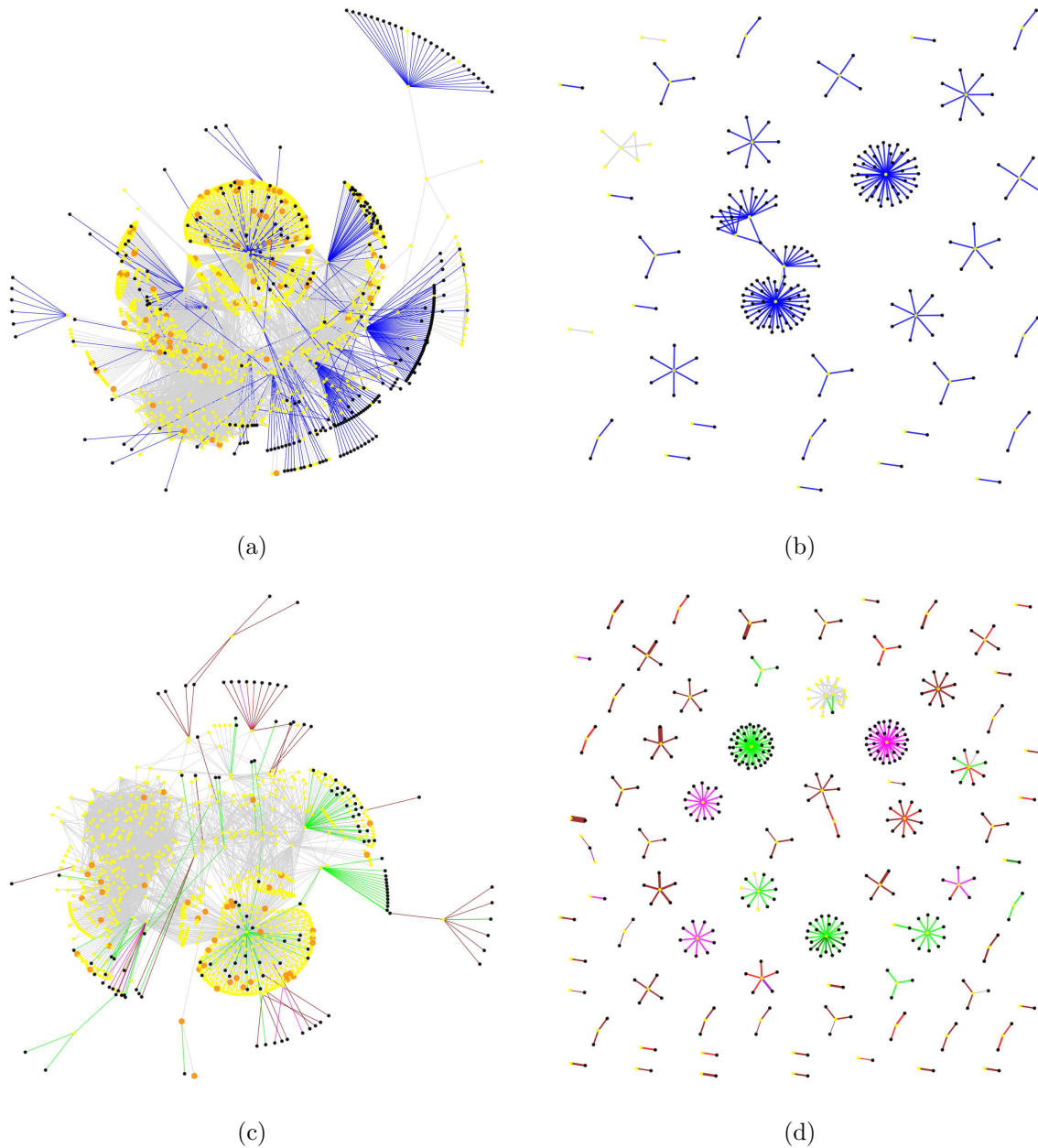


Figure 7: Single-level and multi-level conditional Gaussian graphical models estimated from GTEx data. The estimated single-level model is shown for (a) the largest connected component and (b) the rest of the gene networks, where gene network edges are shown as gray and the effects of variants on gene expressions are shown as blue. The estimated multi-level model is shown for (c) the largest connected component and (d) the rest of the gene networks. In the multi-level model, the influence of *cis*-acting variants is shown as brown if they have been previously annotated as having impact on gene expression, red if they have not been annotated, and pink if the distance to the gene it is influencing is greater than 1 megabase. The influence of *trans*-acting edges is shown as green. Known transcription factors are colored as big orange dots.