

Appendix A: Proofs of the simple properties

Here, we provide the proofs for the three lemmas in Section 2.3 in the paper.

Lemma 1. *The joint distribution of X_1, \dots, X_p and the latent variables L_1, \dots, L_k is Gaussian: $(\mathbf{X}, \mathbf{L}) \sim N(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{XL}})$, where $\boldsymbol{\Sigma}_{\mathbf{XL}}$ is a $(p+k) \times (p+k)$ covariance matrix.*

Proof.

$$\begin{aligned}
 P(\mathbf{X}, \mathbf{L} | \boldsymbol{\Sigma}_{\mathbf{L}}) &= P(\mathbf{X} | \mathbf{L}) P(\mathbf{L} | \boldsymbol{\Sigma}_{\mathbf{L}}) \\
 &= \prod_{i=1}^p \frac{1}{\sqrt{2\pi\sigma_{Z_i}^2}} \exp\left\{-\frac{(X_i - L_{Z_i})^2}{2\sigma_{Z_i}^2}\right\} \cdot \frac{1}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}_{\mathbf{L}}|}} \exp\left\{-\frac{1}{2} \mathbf{L}^\top \boldsymbol{\Sigma}_{\mathbf{L}}^{-1} \mathbf{L}\right\} \\
 &= \frac{1}{\sqrt{(2\pi)^{p+k} (\sigma_{Z_i}^2)^p |\boldsymbol{\Sigma}_{\mathbf{L}}|}} \exp\left\{-\frac{1}{2} \left(\sum_{i=1}^p \frac{(X_i - L_{Z_i})^2}{\sigma_{Z_i}^2} + \mathbf{L}^\top \boldsymbol{\Sigma}_{\mathbf{L}}^{-1} \mathbf{L}\right)\right\}. \quad (1)
 \end{aligned}$$

The above distribution can be rewritten as:

$$P(\mathbf{X}, \mathbf{L} | \boldsymbol{\Sigma}_{\mathbf{L}}) = \frac{1}{\sqrt{(2\pi)^{p+k} |\boldsymbol{\Sigma}_{\mathbf{XL}}|}} \exp\left\{-\frac{1}{2} [\mathbf{X}^\top \mathbf{L}^\top] \boldsymbol{\Sigma}_{\mathbf{XL}}^{-1} [\mathbf{X}^\top \mathbf{L}^\top]^\top\right\}, \quad (2)$$

where

$$\boldsymbol{\Sigma}_{\mathbf{XL}} = \begin{bmatrix} \mathbf{A} & \mathbf{C}^\top \\ \mathbf{C} & \mathbf{B} \end{bmatrix}^{-1},$$

where \mathbf{A} is a $p \times p$ matrix whose element $A_{ij} = 1/\sigma_{Z_i}^2$ if $i = j$ and 0 otherwise; \mathbf{C} is a $k \times p$ matrix whose element $C_{ij} = -(1/\sigma_i^2)$ if $X_j \in \mathcal{M}_i$ and 0 otherwise; and

$$\mathbf{B} = \boldsymbol{\Sigma}_{\mathbf{L}}^{-1} + \begin{bmatrix} |\mathcal{M}_1|/\sigma_1^2 & \dots & 0 \\ \vdots & \vdots & \vdots \\ 0 & \dots & |\mathcal{M}_k|/\sigma_k^2 \end{bmatrix}.$$

□

Lemma 2. *The marginal probability distribution of the observed variables X_1, \dots, X_p is Gaussian: $\mathbf{X} \sim N(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{X}})$, where $\boldsymbol{\Sigma}_{\mathbf{X}}$ is a $p \times p$ covariance matrix.*

Proof. If joint distribution of $\mathbf{X} = \{X_1, \dots, X_p\}$ and the latent variables $\mathbf{L} = \{L_1, \dots, L_k\}$ is a multivariate Gaussian, then the marginal distribution over a subset of variables is also a multivariate Gaussian distribution. □

Lemma 3. *Let $\boldsymbol{\Sigma}_{\mathbf{L}}$ be a $k \times k$ covariance matrix of \mathbf{L} . The relationship between $\boldsymbol{\Sigma}_{\mathbf{X}}$ and $\boldsymbol{\Sigma}_{\mathbf{L}}$ is as*

$$\boldsymbol{\Sigma}_{\mathbf{X}} = \{\mathbf{A} - \mathbf{C}^\top \mathbf{B}^{-1} \mathbf{C}\}^{-1}, \quad (3)$$

where \mathbf{A} is a $p \times p$ matrix whose element $A_{ij} = 1/\sigma_{Z_i}^2$ if $i = j$ and 0 otherwise; \mathbf{C} is a $k \times p$ matrix whose element $C_{ij} = -(1/\sigma_i^2)$ if $X_j \in \mathcal{M}_i$ and 0 otherwise; and

$$\mathbf{B} = \boldsymbol{\Sigma}_{\mathbf{L}}^{-1} + \begin{bmatrix} |\mathcal{M}_1|/\sigma_1^2 & \dots & 0 \\ \vdots & \vdots & \vdots \\ 0 & \dots & |\mathcal{M}_k|/\sigma_k^2 \end{bmatrix}$$

$|\mathcal{M}_k|$ meaning the number of X variables in the module k .

Proof. From Lemma 1, the distribution described in (1) is equivalent to $N(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{XL}})$, which leads to the following:

$$\boldsymbol{\Sigma}_{\mathbf{XL}}^{-1} = \begin{bmatrix} \mathbf{A} & \mathbf{C}^\top \\ \mathbf{C} & \mathbf{B} \end{bmatrix},$$

where \mathbf{B} and \mathbf{C} are defined in (3). Making use of the Schur complement, we obtain the following:

$$\boldsymbol{\Sigma}_{\mathbf{XL}} = \begin{bmatrix} (\mathbf{A} - \mathbf{C}^\top \mathbf{B}^{-1} \mathbf{C})^{-1} & -(\mathbf{A} - \mathbf{C}^\top \mathbf{B}^{-1} \mathbf{C})^{-1} \mathbf{C} \mathbf{B}^{-1} \\ -\mathbf{B}^{-1} \mathbf{C} (\mathbf{A} - \mathbf{C}^\top \mathbf{B}^{-1} \mathbf{C})^{-1} & \mathbf{B}^{-1} \mathbf{C} (\mathbf{A} - \mathbf{C}^\top \mathbf{B}^{-1} \mathbf{C})^{-1} \mathbf{C} \mathbf{B}^{-1} + \mathbf{B}^{-1} \end{bmatrix}.$$

Therefore, the covariance matrix $\boldsymbol{\Sigma}_{\mathbf{X}}$ is the following:

$$\boldsymbol{\Sigma}_{\mathbf{X}} = \{\mathbf{A} - \mathbf{C}^\top \mathbf{B}^{-1} \mathbf{C}\}^{-1}.$$

□

Appendix B: Derivation of MGL algorithm

Here, we present our learning algorithm that optimizes the likelihood function based on the joint distribution in the following form:

$$\begin{aligned} P(\mathbf{X}, \mathbf{L}, \mathbf{Z}, \boldsymbol{\Sigma}_{\mathbf{L}}) &= \prod_{i=1}^p P(X_i | L_{Z_i}) P(\mathbf{L} | \boldsymbol{\Sigma}_{\mathbf{L}}) P(\boldsymbol{\Sigma}_{\mathbf{L}}^{-1}) P(\mathbf{Z}) \\ &= \prod_{i=1}^p \frac{1}{\sqrt{2\pi\sigma_{Z_i}^2}} \exp\left\{-\frac{(X_i - L_{Z_i})^2}{2\sigma_{Z_i}^2}\right\} \cdot \prod_{j \neq j'} \frac{\lambda}{2} \exp\{-\lambda |(\boldsymbol{\Sigma}_{\mathbf{L}}^{-1})_{jj'}|\} \\ &\quad \cdot \frac{1}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}_{\mathbf{L}}|}} \exp\left\{-\frac{1}{2} \mathbf{L}^\top \boldsymbol{\Sigma}_{\mathbf{L}}^{-1} \mathbf{L}\right\} P(\mathbf{Z}). \end{aligned} \quad (4)$$

We begin by summarizing our model's variables:

- $\mathbf{L} = \{L_1, \dots, L_k\}$: A set of latent variables $\mathbf{L} \sim N(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{L}})$.
- $\boldsymbol{\Sigma}_{\mathbf{L}}$: A $k \times k$ covariance matrix of \mathbf{L} . The non-zero pattern of $\boldsymbol{\Sigma}_{\mathbf{L}}^{-1}$ corresponds to the graph structure.
- $\mathbf{Z} = \{Z_1, \dots, Z_p\}$: Each Z_i indicates the module which X_i is assigned to. ($1 \leq Z_i \leq k$)
- $\mathbf{X} = \{X_1, \dots, X_p\}$: A set of observed variables each of which is modeled as $X_i \sim N(L_{Z_i}, \sigma_{Z_i}^2)$.
- $\boldsymbol{\Sigma}_{\mathbf{X}}$: A $p \times p$ covariance matrix of \mathbf{X} . We can compute $\boldsymbol{\Sigma}_{\mathbf{X}}$ based on $\boldsymbol{\Sigma}_{\mathbf{L}}$ and \mathbf{Z} (see Lemma 3).

Estimation of L

In order to estimate L given Z and Θ_L based on the joint log-likelihood function (4), we solve the following optimization problem:

$$\text{maximize}_{L_1, \dots, L_k} \left\{ -\text{trace}(LL^\top \Theta_L) - \sum_{i=1}^p \frac{\|X_i - L_{Z_i}\|_2^2}{\sigma_{Z_i}^2} \right\}. \quad (5)$$

where X_i denotes the i th row of the matrix X , L_i denotes the i th row of the matrix L , and the trace term can be rewritten as:

$$\begin{aligned} \text{trace}(LL^\top \Theta_L) &= L_1 L_1^\top (\Theta_L)_{11} + L_1 L_2^\top (\Theta_L)_{12} + \dots + L_1 L_k^\top (\Theta_L)_{1k} \\ &+ L_2 L_1^\top (\Theta_L)_{21} + L_2 L_2^\top (\Theta_L)_{22} + \dots + L_2 L_k^\top (\Theta_L)_{2k} \\ &+ \dots \\ &+ L_k L_1^\top (\Theta_L)_{k1} + L_k L_2^\top (\Theta_L)_{k2} + \dots + L_k L_k^\top (\Theta_L)_{kk}. \end{aligned}$$

Setting the derivative of the objective function in Eq (5) to zero with respect to L_m , we obtain:

$$((\Theta_L)_{1m} L_1 + (\Theta_L)_{2m} L_2 + \dots + (\Theta_L)_{km} L_k) - \frac{1}{\sigma_m^2} \sum_{X_i \in \mathcal{M}_m} (X_i - L_m) = 0. \quad (6)$$

Here, \mathcal{M}_m means a set of X_i that belongs to the m th module: $\mathcal{M}_m = \{X_i | Z_i = m\}$, and $|\mathcal{M}_m|$ means the number of variables that belong to \mathcal{M}_m . Solving the Eq (6) with respect to L_m leads to:

$$L_m = \frac{\sum_{X_i \in \mathcal{M}_m} X_i - \sigma_m^2 \sum_{i \neq m} (\Theta_L)_{im} L_i}{|\mathcal{M}_m| + \sigma_m^2 (\Theta_L)_{mm}}. \quad (7)$$

We update L_m for each m ($1 \leq m \leq k$), based on the current values on the other latent variables $L_1, \dots, L_{(m-1)}, L_{(m+1)}, \dots, L_k$.

If all entries in Θ_L equal to zero, L_m would be updated to be the average of X_i 's that belong to the m th module. In other words, L_m would be set to be the centroid of the m th module. This leads to a nice interpretation of the MGL learning algorithm with respect to the k-means clustering algorithm.

The k-means clustering algorithm is the special case of the MGL when the latent variables are assumed to be independent; this means that the elements of Θ_L are forced to be zero.

Estimation of Z

In order to estimate Z given L and Θ_L , we solve the following optimization problem:

$$\text{maximize}_{Z_1, \dots, Z_p} \left\{ - \sum_{i=1}^p \frac{\|X_i - L_{Z_i}\|_2^2}{\sigma_{Z_i}^2} \right\}, \quad (8)$$

where X_i denotes the i th row of the matrix X and L_i denotes the i th row of the matrix L . When $\sigma_1, \dots, \sigma_k = 1$, this amounts to finding the assignment for each X_i that minimizes the Euclidean distance with the latent variable.

$$Z_i = \arg \min_{Z_i \in \{1, \dots, k\}} \{ \| X_i - L_{Z_i} \|_2^2 \}, \quad (9)$$

Estimation of Θ_L

To estimate Θ_L given L and Z , we solve the following optimization problem:

$$\text{maximize}_{\Theta_L > 0} \left\{ \log \det \Theta_L - \text{trace} (S_L \Theta_L) - \lambda \sum_{j \neq j'} |(\Theta_L)_{jj'}| \right\}, \quad (10)$$

where S_L is the empirical estimate of the covariance matrix of L : $S_L = \frac{1}{n} L L^\top$. Since L is given, the optimization problem (10) can be solved by the standard graphical lasso algorithm applied to L .