

Bayesian structure learning in graphical models

Seung Bong Jung

Seoul National University

February 17, 2021

- Banerjee and Ghoshal, "Bayesian structure learning in graphical models", Journal of Multivariate Analysis, 2015, 136, 147-162

- 1 Notations and preliminaries
- 2 Prior and posterior concentration
- 3 Posterior computation
 - Approximating model posterior probabilities
 - Ignorability of non-regular models
 - Error in Laplace approximation
- 4 Simulation results

- 1 Notations and preliminaries
- 2 Prior and posterior concentration
- 3 Posterior computation
 - Approximating model posterior probabilities
 - Ignorability of non-regular models
 - Error in Laplace approximation
- 4 Simulation results

Notation 1.1

Let V be the set of indices of p vertices of p -dimensional random vector $X \sim N_p(0, \Omega^{-1})$. If induced undirect graph G has edges $E \subset \{(i, j) \in V \times V : i < j\}$, denote $G = (V, E)$.

Notation 1.2

- \mathcal{M} : the linear space of symmetric matrices of order p .
- \mathcal{M}^+ : the cone of positive definite matrices of order p .
- \mathcal{P}_G : the cone of positive definite matrices of p with zero entries corresponding to each missing edge in E .

Notations

Notation 1.3

Suppose $A = ((a_{ij}))$ is a $p \times p$ matrix and $x \in \mathbb{R}^p$.

- If $1 \leq r < \infty$, $\|A\|_r = (\sum_{i,j=1}^p |a_{ij}|^r)^{\frac{1}{r}}$, $\|x\|_r = (\sum_{i=1}^p |x_i|^r)^{\frac{1}{r}}$
- $\|A\|_\infty = \max_{i,j} |a_{ij}|$, $\|x\|_\infty = \max_i |x_i|$.
- If $A \in \mathcal{M}$, $\text{eig}_i(A)$ denotes i th biggest eigenvalue of A .

Notation 1.4

Suppose $r \times s$ matrix $A = ((a_{ij}))$ is given.

- $\|A\|_{(r,s)} = \sup(\|Ax\|_s : \|x\|_r = 1, x \in \mathbb{R}^p)$

Remark 1.1

Assume $A \in \mathcal{M}$. Clearly, $\|A\|_2 = \sqrt{\text{tr}(A^t A)}$, $\|A\|_{(2,2)} = \max_i |\text{eig}_i(A)|$

Notation 1.5

Suppose r_n, s_n are two numerical sequences.

- If $\frac{r_n}{s_n}$ is bounded as $n \rightarrow \infty$, we say $r_n = O(s_n)$ or equivalently $r_n \lesssim s_n$.
In case $r_n \lesssim s_n$ and $s_n \lesssim r_n$ both hold, denote this by $r_n \asymp s_n$
- If $\frac{r_n}{s_n} \rightarrow 0$ as $n \rightarrow \infty$. Then we denote this by $r_n = o(s_n)$ or equivalently $r_n \ll s_n$.
- If $\frac{r_n}{s_n} \rightarrow 1$ as $n \rightarrow \infty$, we say $r_n \sim s_n$.

Notation 1.6

Suppose random sequence X_n is given. We say $X_n = O_p(\delta_n)$, if $\exists M > 0$ such that $P(|X_n| > M\delta_n) \rightarrow 0$ as $n \rightarrow \infty$. Similarly, $X_n = o_p(\delta_n)$, if $\forall \epsilon > 0, P(|X_n| > \epsilon\delta_n) \rightarrow 0$ as $n \rightarrow \infty$.

Notation 1.7

- $\#T$ denotes the cardinality of a set T . Define the symmetric matrix $E_{(i,j)} = ((\mathbb{1}_{\{(i,j),(j,i)\}}(l, m)))$, where $\mathbb{1}$ is an indicator function.
- For a subset A of a metric space (S, d) , the minimum number of d -balls of size ϵ in S needed to cover A is denoted by $N(\epsilon, A, d)$

Notation 1.8

- The Hellinger distance between two probability densities q_1 and q_2 is defined as $h(q_1, q_2) = \|\sqrt{q_1} - \sqrt{q_2}\|_2 = (\int (\sqrt{q_1(x)} - \sqrt{q_2(x)})^2 dx)^{\frac{1}{2}}$

Theorem 1.1

Let $A, B \in \mathcal{M}$. Then followings hold :

$$\begin{aligned}\|A\|_{\infty} &\leq \|A\|_{(2,2)} \leq \|A\|_2 \leq p\|A\|_{\infty} \\ \|AB\|_2 &\leq \min\{\|A\|_{(2,2)}\|B\|_2, \|B\|_{(2,2)}\|A\|_2\}\end{aligned}$$

(Proof) Use spectral decomposition to A, B .

Theorem 1.2

Suppose $X \sim N_p(0, \Omega^{-1})$, where $\Omega = ((\omega_{ij}))$. Then $\omega_{ij} = \omega_{ji} = 0$ implies conditional independence of X_i and X_j given $(X_r : r \neq i, j)$ with $i \neq j$.

(Proof) w.l.o.g. assume $i = 1, j = 2$. Let $\Sigma = \Omega^{-1}$. Consider Σ as following :

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

where Σ_{11} is 2×2 matrix, Σ_{22} is $(p-2) \times (p-2)$ matrix, and Σ_{12} is $2 \times (p-2)$ matrix, $\Sigma_{21} = \Sigma_{12}^t$.

Let $\Sigma_{11.2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} = \Omega_{11}^{-1}$. Then, we have $X_1, X_2 | X_{r:r \neq 1,2} \sim N_2(\mathbf{0}, \Sigma_{11.2}^{-1})$. Hence, $\text{Cov}(X_1, X_2 | X_{r:r \neq 1,2}) = \omega_{12} = \omega_{21}$.

Under normality, X_1, X_2 given $(X_r : r \neq 1, 2)$ are independent if and only if their conditional covariance given other covariates is 0. Thus, we get the desired result.

- 1 Notations and preliminaries
- 2 Prior and posterior concentration
- 3 Posterior computation
 - Approximating model posterior probabilities
 - Ignorability of non-regular models
 - Error in Laplace approximation
- 4 Simulation results

Graphical lasso

Let $X_1, X_2, \dots, X_n \stackrel{i.i.d.}{\sim} N_p(0, \Sigma)$, where $\Omega = \Sigma^{-1}$. Then we clearly see that $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X_i X_i^t$ is MLE of Σ . Then the graphical lasso solves the following optimization problem with penalty parameter $\lambda \geq 0$ and constraint $\Omega \in \mathcal{M}^+$:

$$\text{minimize} \quad \underbrace{-\log \det(\Omega) + \text{tr}(\hat{\Sigma}\Omega)}_{\text{negative data log-likelihood}} + \underbrace{\frac{\lambda}{n} \|\Omega\|_1}_{\text{Lasso type penalty}} \quad (1)$$

Remark 2.1

Let $\mathcal{U}(\epsilon_0, s) = \{\Omega \in \mathcal{M}^+ : \#\{(i, j) : 1 \leq i < j \leq p, \omega_{ij} \neq 0\} \leq s, 0 < \epsilon_0 \leq \text{eig}_1(\Omega) \leq \text{eig}_p(\Omega) \leq \epsilon_0^{-1} < \infty\}$. Denote the minimizer of (1) by Ω^* with true precision matrix Ω_0 . Assuming $\Omega_0 \in \mathcal{U}(\epsilon_0, s)$, Rothman et al.(2008) has shown $\|\Omega^* - \Omega_0\|_2 = O_p(\sqrt{\frac{(p+s) \log p}{n}})$.

Hence $\|\Omega^* - \Omega_0\|_2 \xrightarrow{P} 0$ as $n^{-1}(p+s) \log p \rightarrow 0$.

Theorem 2.1

Suppose Ω^* is the solution of (1) with true precision matrix $\Omega_0 \in \mathcal{U}(\epsilon_0, s)$. Then $\|\Omega^*\|_{(2,2)} = O_p(1)$ and $\|\Omega^{*-1}\|_{(2,2)} = O_p(1)$

(proof)

- For the first equation, by triangle inequality and result of theorem 1.1,
$$\|\Omega^*\|_{(2,2)} \leq \|\Omega^* - \Omega_0\|_{(2,2)} + \|\Omega_0\|_2.$$
- By remark 2.1 and result of theorem 1.1, we see that
$$\|\Omega^* - \Omega_0\|_{(2,2)} = o_p(1). \text{ Since } \|\Omega_0\|_{(2,2)} \text{ is constant, we have } \|\Omega^*\|_{(2,2)} = O_p(1).$$
- Similarly, by triangle inequality and Cauchy-Schwartz inequality, one can see

$$\begin{aligned} \|\Omega^{*-1}\|_{(2,2)} &\leq \|\Omega_0^{-1}\|_{(2,2)} + \|\Omega^{*-1} - \Omega_0^{-1}\|_{(2,2)} \\ &\leq \|\Omega_0^{-1}\|_{(2,2)} + \|\Omega_0^{-1}\|_{(2,2)} \|\Omega^{*-1} - \Omega_0^{-1}\|_{(2,2)} \|\Omega^{*-1}\|_{(2,2)} \end{aligned}$$

$$\text{,which leads to } \|\Omega^{*-1}\|_{(2,2)} \leq \frac{\|\Omega_0^{-1}\|_{(2,2)}}{1 - \|\Omega_0^{-1}\|_{(2,2)} \|\Omega^{*-1} - \Omega_0^{-1}\|_{(2,2)}}.$$

- Again, for $\|\Omega^* - \Omega_0\|_{(2,2)} = o_p(1)$ and $\|\Omega_0^{-1}\|_{(2,2)}$ is constant, we have the second equation.

Wang's method(2012)

- Put Exponential prior $\lambda \exp(-\lambda x)$ on diagonal elements and Laplace prior $\frac{\lambda}{2} \exp(-\lambda|x|)$ on off-diagonal elements.
- But by absolute continuity of priors, the probability of event $\{\omega_{ij} = 0\}$ is 0 and so is the posterior probability of the event.
- Hard to introduce sparsity to the model.

Banerjee & Ghosal's method(2015)

- Put point-mass prior on the events $\{\omega_{ij} = 0\}$.
- Able to make posterior inference about the sparse structure of the model.

Prior settings

- Let $\Gamma = (\gamma_{ij} : 1 \leq i < j \leq p)$ be a $\binom{p}{2}$ -dimensional vector of edge-inclusion vector, where $\gamma_{ij} = \mathbb{1}\{(i, j) \in E\}$.
- $\#\Gamma$: the number of non-zero elements in Γ .
- Put the same priors on precision matrix suggested by Wang.
- Define \mathcal{M}_0^+ to be subset of \mathcal{M}^+ , consisting of positive definite matrices of order p , whose eigenvalues are bounded by two fixed positive numbers.
- Now assume true precision matrix belongs to \mathcal{M}_0^+ , unless otherwise stated in this section.
- Suppose $\gamma_{ij} \stackrel{i.i.d}{\sim} \text{Ber}(q)$, where q is fixed. To introduce sparsity to the precision matrix, put restriction on $\#\Gamma$, namely the size of model.

Prior 1

- For fixed positive number $\bar{\gamma}$, we have

$$p(\Gamma) \propto q^{\#\Gamma} (1 - q)^{\binom{p}{2} - \#\Gamma} \mathbb{1}(\#\Gamma \leq \bar{\gamma})$$

Prior 2

- Choose the prior distribution on \bar{R} such that $p(\bar{R} > a_1 m) \leq e^{-a_2 m \log m}$ for some $a_1, a_2 > 0$ and $\forall m \in \mathbb{N}$.
- Then, we have another prior $p(\Gamma | \bar{R}) \propto q^{\#\Gamma} (1 - q)^{\binom{p}{2} - \#\Gamma} \mathbb{1}(\#\Gamma \leq \bar{R})$ leading to $p(\Gamma) \propto q^{\#\Gamma} (1 - q)^{\binom{p}{2} - \#\Gamma} P(\bar{R} \geq \#\Gamma)$.

Posterior

- From Prior 1,2, we have following posterior :

$$\begin{aligned} p(\Omega, \Gamma | X^{(n)}) &\propto p(X^{(n)} | \Omega, \Gamma) p(\Omega | \Gamma) p(\Gamma) \\ &= (2\pi)^{-np/2} [\det(\Omega)]^{n/2} \exp\left(-\frac{n}{2} \text{tr}(\hat{\Sigma}\Omega)\right) \\ &\quad \times \prod_{\gamma_{ij}=1} \left\{ \frac{\lambda}{2} \exp(-\lambda |\omega_{ij}|) \right\} \prod_{i=1}^p \left\{ \frac{\lambda}{2} \exp\left(-\frac{\lambda}{2} \omega_{ii}\right) \right\} \\ &\quad \times p(\Gamma) \mathbb{1}_{\mathcal{M}_0^+}(\Omega) \end{aligned}$$

Posterior

- Consequently, $p(\Omega, \Gamma | X^{(n)}) \propto C_{\Gamma} Q(\Omega, \Gamma | X^{(n)})$, where

$$C_{\Gamma} = (2\pi)^{-np/2} q^{\#\Gamma} (1-q)^{\binom{p}{2}-\#\Gamma} (\lambda/2)^{p+\#\Gamma} \beta(\Gamma)$$

$$\beta(\Gamma) = \begin{cases} p(\bar{R} \geq \#\Gamma), & \text{for Prior 2} \\ \mathbb{1}(\#\Gamma \leq \bar{r}), & \text{for Prior 1} \end{cases}$$

$$Q(\Omega, \Gamma | X^{(n)}) = [\det(\Omega)]^{n/2} \exp\left(-\frac{n}{2} \text{tr}(\hat{\Sigma}\Omega)\right) \prod_{\gamma_{ij}=1} \exp(-\lambda |\omega_{ij}|) \\ \times \prod_{i=1}^p \exp\left(-\frac{\lambda}{2} \omega_{ii}\right) \mathbb{1}_{\mathcal{M}_0^+}(\Omega)$$

- The posterior is very intractable due to positive constraint on Ω . One possible method is RJMCMC, but the result may be unreliable due to too large number of models to search over. $(2^{\binom{p}{2}})$.

Claim 2.1

Suppose $A, B \in \mathcal{M}$. Then, $\text{tr}([AB]^2) \leq \text{tr}(A^2B^2)$

(proof) Use the equation $2(AB)^2 - 2A^2B^2 = (AB - BA)^2$ and the fact that $AB - BA$ is skew-hermitian if $A, B \in \mathcal{M}$.

Lemma 2.1

If p_{Ω_k} is the density of $N_p(0, \Omega_k^{-1})$, $k = 1, 2$, then for all $\Omega_k \in \mathcal{M}_0^+$, $k = 1, 2$, and $d_i, i = 1, \dots, p$, eigenvalues of $A = \Omega_1^{-1/2}\Omega_2\Omega_1^{-1/2}$, we have followings for sufficiently small $\delta > 0$ and some constant $c_0 > 0$.

$$(i) \quad c_0^{-1} \|\Omega_1 - \Omega_2\|_2^2 \leq \sum_{i=1}^p |d_i - 1|^2 \leq c_0 \|\Omega_1 - \Omega_2\|_2^2$$

$$(ii) \quad h(p_{\Omega_1}, p_{\Omega_2}) < \delta \text{ implies } \max_i |d_i - 1| < 1$$

$$\text{and } \|\Omega_1 - \Omega_2\|_2^2 \leq c_0 h^2(p_{\Omega_1}, p_{\Omega_2})$$

$$(iii) \quad h^2(p_{\Omega_1}, p_{\Omega_2}) \leq c_0 \|\Omega_1 - \Omega_2\|_2^2.$$

(Proof of (i))

- Since $\Omega_1 \in \mathcal{M}_+^+$, both $\|\Omega_1\|_2$ and $\|\Omega_1^{-1}\|_2$ are bounded by some constant.

•

$$\begin{aligned}\|\Omega_1 - \Omega\|_2^2 &= \|\Omega_1^{\frac{1}{2}}(\mathbf{I}_p - \mathbf{A})\Omega_1^{\frac{1}{2}}\|_2^2 \leq \|\Omega_1\|_{(2,2)}^2 \|\mathbf{I}_p - \mathbf{A}\|_2^2 \\ &= \|\Omega_1\|_{(2,2)}^2 \text{tr}(\mathbf{I}_p - \mathbf{A})^2 = \|\Omega_1\|_{(2,2)}^2 \sum_{i=1}^p (d_i - 1)^2\end{aligned}$$

•

$$\begin{aligned}\sum_{i=1}^p (d_i - 1)^2 &= \|\mathbf{I}_p - \mathbf{A}\|_2^2 = \|\Omega_1^{-\frac{1}{2}}(\Omega_1 - \Omega_2)\Omega_1^{-\frac{1}{2}}\|_2^2 \\ &\leq \|\Omega_1^{-1}\|_{(2,2)}^2 \|\Omega_1 - \Omega_2\|_2^2\end{aligned}$$

(Proof of (ii) & (iii))

- By direct calculations,

$$\frac{1}{2}h^2(p_{\Omega_1}, p_{\Omega_2}) = 1 - \left\{ \prod_{i=1}^p \frac{1}{2}(d_i^{\frac{1}{2}} + d_i^{-\frac{1}{2}}) \right\}^{-\frac{1}{2}} \cdots (*)$$

- Assuming $h(p_{\Omega_1}, p_{\Omega_2}) < \delta$ and rearranging terms, we have

$$\prod_{i=1}^p \frac{1}{2}(d_i^{\frac{1}{2}} + d_i^{-\frac{1}{2}}) \leq (1 - \frac{\delta^2}{2})^{-2} = 1 + \eta$$

- From GM-AM inequality, $\max_i \frac{1}{2}(d_i^{\frac{1}{2}} + d_i^{-\frac{1}{2}}) \leq 1 + \eta$.
- With sufficiently small $\delta > 0$, η can be made sufficiently small, so that solving the inequality $\frac{1}{2}(d_i^{\frac{1}{2}} + d_i^{-\frac{1}{2}})$, we have $|d_i - 1| < 1, \forall i$.
- Since $|d_i - 1| < 1$, by a Taylor's expansion, for some constant $c_1, c_2 > 0$,

$$\begin{aligned} 1 + c_1 \sum_{i=1}^p (d_i - 1)^2 &\leq 2^{-p} \prod_{i=1}^p (d_i^{\frac{1}{2}} + d_i^{-\frac{1}{2}}) \approx (1 - h^2)^{-2} \\ &\leq 1 + c_2 \sum_{i=1}^p (d_i - 1)^2 \end{aligned}$$

- Since h is made sufficiently small, using Taylor's expansion, $(1 - h^2)^{-2} \sim 1 + h^2$. By the lower estimate in the above inequality, $c_1 \sum_{i=1}^p (d_i - 1)^2 \leq h^2$. With the result of (i), this establishes (ii).
- To prove (iii), since Hellinger distance is bounded, we may assume $\|\Omega_1 - \Omega_2\|_2$ is sufficiently small. Again by (i), $\max_i |d_i - 1| < 1$ and so with the upper estimate in the above inequality and the result of (i), we obtain (iii).

Claim 2.2

Suppose p_0, p_1 are two probability densities. Then $K(p_0, p_1) \geq h^2(p_0, p_1)$

(note) $K(p_0, p_1) = \int p_0 \log(\frac{p_0}{p_1})$, which is KL divergence.

(proof) If $x < 1$, we have $-\log(1 - x) \geq x$. So,

$$\begin{aligned} K(p_0, p_1) &= -2 \int p_0 \log\left(\sqrt{\frac{p_1}{p_0}}\right) \geq 2 \int p_0 \left(1 - \sqrt{\frac{p_1}{p_0}}\right) = 2 - 2 \int \sqrt{p_0 p_1} \\ &= h^2(p_0, p_1) \end{aligned}$$

Claim 2.3

Let p_{Ω_i} be density of $N_p(0, \Omega_i^{-1})$, $i = 0, 1$ and $d_i, i = 1, 2, \dots, p$ be eigenvalues of $\Omega_0^{-\frac{1}{2}} \Omega_1 \Omega_0^{-\frac{1}{2}}$. Then followings hold:

$$K(p_{\Omega_0}, p_{\Omega_1}) = -\frac{1}{2} \sum_{i=1}^p (\log d_i - 1 + d_i)$$

$$V(p_{\Omega_0}, p_{\Omega_1}) = \frac{1}{4} \left(\sum_{i=1}^p (\log d_i - 1 + d_i) \right)^2 + \frac{1}{2} \sum_{i=1}^p (1 - d_i)^2$$

where $V(p_0, p_1) = \int p_0 \log^2 \left(\frac{p_0}{p_1} \right)$

(proof) Use the fact $E(Z^t A Z) = \text{tr}(A \Sigma)$, $\text{Var}(Z^t A Z) = 2 \text{tr}(A \Sigma A \Sigma)$, where $Z \sim N_p(0, \Sigma)$ and $A \in \mathcal{M}$. The result is direct from some computations.

Posterior concentration

Claim 2.4

Let $p \in \mathbb{N}$ and choose $\bar{r} \in \mathbb{N}$ such that $\bar{r} < \binom{p}{2}/2$. If $j \leq \bar{r}$,

$$\binom{\binom{p}{2}}{j} \leq \binom{p + \binom{p}{2}}{\bar{r}}$$

Theorem 2.2

Let $\mathbf{X}^{(n)} = (\mathbf{X}_1, \dots, \mathbf{X}_n) \sim N_p(0, \Omega_0^{-1})$, where $\Omega_0 \in \mathcal{U}(\epsilon_0, s)$ for some $0 < \epsilon_0 < \infty$ and $0 \leq s \leq p(p-1)/2$. Also assume that priors are either Prior1 or Prior2 with $q < 1/2$ and the range of eigenvalues of matrices in \mathcal{M}_0^+ is sufficiently broad to contain $[\epsilon_0, \epsilon_0^{-1}]$. Then the posterior distribution of Ω satisfies

$$E_0[P\{\|\Omega - \Omega_0\|_2 > M\epsilon_n | \mathbf{X}^{(n)}\}] \rightarrow 0,$$

for $\epsilon_n = n^{-1/2}(p+s)^{1/2}(\log n)^{1/2}$ and a sufficiently large $M > 0$.

(Proof)

- Let $B(p_{\Omega_0}, \epsilon_n) = \{p_{\Omega} : K(p_{\Omega_0}, p_{\Omega}) \leq \epsilon_n^2, V(p_{\Omega_0}, p_{\Omega}) \leq \epsilon_n^2\}$. Denote the prior of Ω by Π .
- Also let \mathcal{P}_n be the set of all densities p_{Ω} such that the graph induced by Ω has maximum number of edges $\bar{r} < \binom{p}{2}/2$ and each entry of Ω is at most L in absolute value, where \bar{r} and L depend only on n and to be determined.
- The proof is by verifying each condition in the remark 2.2, which is the result of theorem 2.1. in "Convergence Rates of Posterior Distributions"

Remark 2.2

Suppose that for a sequence ϵ_n with $\epsilon_n \rightarrow 0$ and $n\epsilon_n^2 \rightarrow \infty$, a constant $C > 0$ and sets $\mathcal{P}_n \subset \mathcal{P}$, where \mathcal{P} is a space of densities, we have

- (Condition 1) $\log N(\epsilon_n, \mathcal{P}_n, d) \leq n\epsilon_n^2$
- (Condition 2) $\Pi(\mathcal{P}_n^c) \leq \exp(-n\epsilon_n^2(C + 4))$
- (Condition 3) $\Pi(B(p_{\Omega_0}, \epsilon_n)) \geq \exp(-n\epsilon_n^2 C)$

Then for sufficiently large M , $E_0[\Pi(d(p, p_0) \geq M\epsilon_n) | X^{(n)}] \rightarrow 0$.

- We shall verify (Condition 1) first. We set metric d as Frobenius distance. Then $N(\epsilon_n, \mathcal{P}_n, d) = \sum_{j=1}^{\bar{r}} \left(\frac{L}{\epsilon_n}\right)^j \binom{\bar{p}}{j}$. Taking sufficiently large n and fixed constants $b_1, b_2 > 0$,

$$\begin{aligned} \log \left\{ \sum_{i=1}^{\bar{r}} \left(\frac{L}{\epsilon_n} \right)^j \binom{\binom{p}{2}}{j} \right\} &\leq \log \left\{ \bar{r} \left(\frac{L}{\epsilon_n} \right)^{\bar{r}} \binom{p + \binom{p}{2}}{\bar{r}} \right\} \\ &\lesssim \log \bar{r} + \bar{r} \log L + \bar{r} \log \epsilon_n^{-1} + \bar{r} \log p \end{aligned} \quad (2)$$

- Choose $\bar{r} \sim \frac{b_1 n \epsilon_n^2}{\log n}$ and $L \sim \frac{b_2 n \epsilon_n^2}{\log n}$. Then we have

$$\log \bar{r} + \bar{r} \log L + \bar{r} \log \epsilon_n^{-1} + \bar{r} \log p \approx n \epsilon_n^2$$

- Next, we verify (Condition2).

- Suppose p_Ω does not belong to \mathcal{P}_n . By definition of \mathcal{P}_n , there exists an entry of Ω that exceeds L in magnitude or the number of edges in the graph induced by Ω exceeds \bar{r} . The probability of latter event is $P(\bar{R} > \bar{r}) \leq \exp(-a'_2 b_1 n \epsilon_n^2)$. Note b_1 is chosen to be sufficiently large. To bound the probability of the former event, if we naively say that prior on all entries of Ω is given by Exponential prior, the probability of each entry exceeds L is $\int_L^\infty \lambda \exp(-\lambda \omega_{ij}) d\omega_{ij} = \exp(-\lambda L)$. For there are total $p + \binom{p}{2}$ distinct entries in Ω , assuming independence on entries for convenience, we have the probability of former event $(p + \binom{p}{2}) \exp(-\lambda L) \lesssim \binom{p}{2} \exp(-L)$. Hence to sum up, there exists a constant $C > 0$ such that

$$\Pi(\mathcal{P}_n^c) \leq \exp(-n \epsilon_n^2 (C + 4))$$

Posterior concentration

- Assume $p_\Omega \in B(p_{\Omega_0}, \epsilon_n)$. By claim 2.2., $K(p_{\Omega_0}, p_\Omega) \geq h^2(p_{\Omega_0}, p_\Omega)$.
- By Taylor's expansion to $K(p_{\Omega_0}, p_\Omega)$, provided that it is valid, and using the result of theorem 1.1 and lemma 2.1,

$$K(p_{\Omega_0}, p_\Omega) \sim \frac{1}{4} \sum_{i=1}^p (1 - d_i)^2 \leq c^{-2} p^2 \|\Omega - \Omega_0\|_\infty^2$$

- Thus, if $c^{-2} p^2 \|\Omega - \Omega_0\|_\infty^2 \leq \epsilon_n^2$, then the above approximation is valid and so $\mathcal{P} = \{p_\Omega : \|\Omega - \Omega_0\|_\infty \leq c\epsilon_n/p\} \subset B(p_{\Omega_0}, \epsilon_n)$.
- Hence it suffices to get a lower estimate of the prior probability of \mathcal{P} .
- Because of constraint \mathcal{M}_0^+ , the components of Ω are not independently distributed. However, a small neighborhood of $\Omega_0 \in \mathcal{U}(\epsilon_0, s)$ lies within \mathcal{M}_0^+ , which was the assumption. Hence, the constraint can only increase prior concentration of $B(p_{\Omega_0}, \epsilon_n)$, as the prior probability of \mathcal{P} tends to increase. So in order to obtain lower bound for the prior probability of \mathcal{P} , we may assume independence on components of Ω .

- So we have the estimate

$$\Pi(\|\Omega_0 - \Omega\|_\infty \leq c\epsilon_n/p) \gtrsim \left(\frac{c\epsilon_n}{p}\right)^{p+s} \quad (3)$$

- Because $\log\left(\left(\frac{c\epsilon_n}{p}\right)^{p+s}\right) \approx n\epsilon_n^2$, we have (Condition 3).

Remark 2.3

Similar to frequentist graphical lasso, we have $\|\Omega - \Omega^\|_2 = O(\epsilon_n)$ with posterior probability tending to one in probability. This gives,*

$$\frac{\int_{\|\Omega - \Omega^*\|_2 \leq \epsilon_n} f(\Omega) \prod_{(i,j) \in \bar{E}} d\omega_{ij}}{\int_{\Omega \in \mathcal{M}_0^+} f(\Omega) \prod_{(i,j) \in \bar{E}} d\omega_{ij}} \rightarrow 1$$

, for a bounded and positive measurable function $f(\cdot)$ of Ω

- 1 Notations and preliminaries
- 2 Prior and posterior concentration
- 3 Posterior computation
 - Approximating model posterior probabilities
 - Ignorability of non-regular models
 - Error in Laplace approximation
- 4 Simulation results

Approximation of Posterior

- Let $\Omega_\Gamma = ((\omega_{\Gamma,ij}))$ be the precision matrix in the model Γ
- Reformulate $Q(\Omega, \Gamma|X^{(n)})$ as

$$Q(\Omega, \Gamma|X^{(n)}) = \exp\left(-\frac{n}{2}h_\Gamma(\Omega_\Gamma)\right) \mathbb{1}_{\mathcal{M}_0^+}(\Omega_\Gamma)$$

where

$$h_\Gamma(\Omega_\Gamma) = -\log \det(\Omega_\Gamma) + \text{tr}(\hat{\Sigma}\Omega_\Gamma) + \frac{2\lambda}{n} \sum_{\gamma_{ij}=1} |\omega_{\Gamma,ij}| + \frac{\lambda}{n} \sum_{i=1}^p \omega_{\Gamma,ii}.$$

- Let $\bar{E}_\Gamma = \{(i,j) \in V \times V : i = j, \text{ or } \gamma_{ij} = 1, i \neq j\}$
- By integrating $p(\Gamma, \Omega|X^{(n)})$ with respect to $\omega_{\Gamma,ij}$,

$$p(\Gamma|X^{(n)}) \propto C_\Gamma \int_{\Omega_\Gamma \in \mathcal{M}_0^+} \exp\left(-\frac{n}{2}h_\Gamma(\Omega_\Gamma)\right) \prod_{(i,j) \in \bar{E}_\Gamma} d\omega_{\Gamma,ij}$$

Approximation of Posterior

- Suppose $h_{\Gamma}(\Omega_{\Gamma})$ is uniquely minimized at Ω_{Γ}^* , which is graphical lasso. Define $\Delta_{\Gamma} = \Omega_{\Gamma} - \Omega_{\Gamma}^* = ((u_{\Gamma,ij}))$
- $p(\Gamma|X^{(n)}) \propto C_{\Gamma} \int_{\Omega_{\Gamma}^* + \Delta_{\Gamma} \in \mathcal{M}_0^+} \exp(-\frac{n}{2} h_{\Gamma}(\Delta_{\Gamma} + \Omega_{\Gamma}^*)) \prod_{(i,j) \in \bar{E}_{\Gamma}} du_{\Gamma,ij}$
- $g_{\Gamma}(\Delta_{\Gamma}) = -\log \det(\Delta_{\Gamma} + \Omega_{\Gamma}^*) + \text{tr}(\hat{\Sigma} \Delta_{\Gamma}) + \frac{2\lambda}{n} \sum_{\gamma_{ij}=1} (|u_{\Gamma,ij} + \omega_{\Gamma,ij}^*| - |\omega_{\Gamma,ij}^*|) + \frac{\lambda}{n} \sum_{i=1}^p u_{\Gamma,ii}$
- $h_{\Gamma}(\Delta_{\Gamma} + \Omega_{\Gamma}^*) = g_{\Gamma}(\Delta_{\Gamma}) + h_{\Gamma}(\Omega_{\Gamma}^*) + \log \det(\Omega_{\Gamma}^*)$
-

$$p(\Gamma|X^{(n)}) \propto C_{\Gamma} \exp(-\frac{n}{2} h_{\Gamma}(\Omega_{\Gamma}^*)) [\det(\Omega_{\Gamma}^*)]^{-\frac{n}{2}} \\ \times \int_{\Omega_{\Gamma}^* + \Delta_{\Gamma} \in \mathcal{M}_0^+} \exp(-\frac{n}{2} g_{\Gamma}(\Delta_{\Gamma})) \prod_{(i,j) \in \bar{E}_{\Gamma}} du_{\Gamma,ij}$$

Approximation of Posterior

- Recall that Laplace approximation for twice differentiable function h on \mathbb{R}^d with unique minimizer \hat{x} and $M > 0$,

$$\int e^{-Mh(x)} dx \approx e^{-Mh(\hat{x})} (2\pi)^{\frac{d}{2}} M^{-\frac{d}{2}} [\det(D^2 h(x))|_{x=\hat{x}}]^{-\frac{1}{2}} \quad (4)$$

- Define the matrix $H_B = ((h_B\{(i, j), (l, m)\}))$, where $h_B\{(i, j), (l, m)\} = \text{tr}(B^{-1}E_{(i,j)}B^{-1}E_{(l,m)})$.
- The Hessian matrix of g_Γ is the $\#\bar{E}_\Gamma \times \#\bar{E}_\Gamma$ matrix $H_{\Delta_\Gamma + \Omega_\Gamma^*}$, whose $\{(i, j), (l, m)\}$ th entry for $(i, j), (l, m) \in \bar{E}_\Gamma$ is given by

$$\frac{\partial^2 g_\Gamma(\Delta_\Gamma)}{\partial u_{\Gamma, ij} \partial u_{\Gamma, lm}} = \text{tr}((\Delta_\Gamma + \Omega_\Gamma^*)^{-1} E_{(i,j)} (\Delta_\Gamma + \Omega_\Gamma^*)^{-1} E_{(l,m)})$$

- Since Ω_Γ^* is the unique minimizer of h_Γ , g_Γ is minimized only at 0.

- Assuming that $\omega_{\Gamma,ij}^*$ is not 0, by (2),

$$p^*(\Gamma|X^{(n)}) \propto C_{\Gamma} \exp\left(-\frac{n}{2}h_{\Gamma}(\Omega_{\Gamma}^*)\right)\left(\frac{4\pi}{n}\right)^{\frac{\#\bar{E}}{2}}[\det(H_{\Omega_{\Gamma}^*})]^{-\frac{1}{2}} \quad (5)$$

- If $\omega_{\Gamma,ij}^*$ is 0 for at least one (i,j) , the approximation is not valid because g_{Γ} is not differentiable. We may call the model making Laplace approximation as "regular model".
- By showing that we can ignore non-regular model, we can consider only regular models, which leads to significantly reduced number of models to search over.

Approximation of Posterior

- In fact, the validity of approximation also depends on $[\det(H_{\Omega_r^*})]^{-\frac{1}{2}}$. But this term is always bounded away from 0, regardless of regularity, since the minimum eigenvalue of the Hessian $H_{\Omega_r^*}^*$ is bounded away from 0.
- The idea to find posterior mode of $p^*(\Gamma|X^{(n)})$ is bad. The posterior heavily decays, even if we consider only regular model. Also, there are too many models to search over.

- 1 Notations and preliminaries
- 2 Prior and posterior concentration
- 3 Posterior computation
 - Approximating model posterior probabilities
 - Ignorability of non-regular models
 - Error in Laplace approximation
- 4 Simulation results

Ignorability of non-regular models

- For convenience, suppose that the vector Γ has t 1s on its first components and the rest of them are 0. The model Γ is non-regular when the corresponding graphical lasso to the model Γ has at least one 0 entry on the edge which is presented in the model Γ . So among those t 1s, say the last r of them have corresponding graphical solution equal to 0. In this context, we obtain submodel Γ' of Γ with first $(t - r)$ 1s and rest 0s

Lemma 3.1

For the corresponding regular submodel Γ' of Γ , the graphical lasso solutions corresponding to models Γ and Γ' are identical.

(proof)

- Our convex optimization problem is to minimize

$$f(\Omega_{\Gamma}) = -\log \det(\Omega_{\Gamma}) + \text{tr}(\hat{\Sigma}\Omega_{\Gamma}) + \frac{2\lambda}{n} \sum_{\gamma_{ij}=1} |\omega_{\Gamma,ij}| + \frac{\lambda}{n} \sum_{i=1}^p \omega_{\Gamma,ii}$$

- Define a function $a : \mathcal{M} \rightarrow \mathcal{M}$ by $a(X) = 2X - \text{diag}(X)$. Then, it is clear that a is injection and linear map. Hence $a(X) = 0$ only when $X = 0$.

Ignorability of non-regular models

- By KKT condition,

$$\frac{\partial f(\Omega_{\Gamma})}{\partial \Omega_{\Gamma}} = -a(\Omega_{\Gamma}^{-1} - \hat{\Sigma} - \frac{\lambda}{n}G) = 0 \Leftrightarrow \Omega_{\Gamma}^{-1} - \hat{\Sigma} - \frac{\lambda}{n}G = 0$$

where $G = ((g_{ij}))$ is a matrix with if $i = j$, $g_{ij} = 1$ and otherwise $g_{ij} = \frac{\omega_{\Gamma,ij}}{|\omega_{\Gamma,ij}|}$ for $\omega_{\Gamma,ij} \neq 0$ and $|g_{ij}| \leq 1$ for $\omega_{\Gamma,ij} = 0$.

- Since Ω_{Γ}^* is graphical lasso, by (30), $\Omega_{\Gamma}^{*-1} - \hat{\Sigma} - \frac{\lambda}{n}G = 0$. Let (i, j) be a pair such that $\gamma_{ij} = 1$ but $\omega_{\Gamma^*,ij}^* = 0$. By the definition of regular submodel Γ' , $\Omega_{\Gamma'}^*$ also satisfies $\Omega_{\Gamma'}^{*-1} - \hat{\Sigma} - \frac{\lambda}{n}G = 0$ because $\omega_{\Gamma',ij}^* = 0$ for any $\gamma'_{ij} = 0$. $\Omega_{\Gamma'}^* = \Omega_{\Gamma}^*$ follows from the uniqueness of solution of the given optimization problem.

Theorem 3.1

Consider the prior on Γ either Prior1 or Prior2 with $q < \frac{1}{2}$. The posterior probability of a non-regular model Γ is always less than that of the corresponding regular submodel Γ' .

(proof)

- It is clear that $\{u_{\Gamma,ij} : \|\Delta_{\Gamma}\|_2 \leq \epsilon_n\} \subset \{u_{\Gamma',ij} : \|\Delta_{\Gamma'}\|_2 \leq \epsilon_n\} \cdots (*)$.

- Recall that

$$p(\Gamma|X^{(n)}) \propto C_{\Gamma} \int_{\Omega_{\Gamma}^* + \Delta_{\Gamma} \in \mathcal{M}_0^+} \exp\left(-\frac{n}{2} h_{\Gamma}(\Delta_{\Gamma} + \Omega_{\Gamma}^*)\right) \prod_{(i,j) \in \bar{E}_{\Gamma}} du_{\Gamma,ij} \text{ and}$$

$$h_{\Gamma}(\Delta_{\Gamma} + \Omega_{\Gamma}^*) = g_{\Gamma}(\Delta_{\Gamma}) + h_{\Gamma}(\Omega_{\Gamma}^*) + \log \det(\Omega_{\Gamma}^*)$$

Ignorability of non-regular models

- As a consequence of Lemma 3.1, the term $h_{\Gamma}(\Omega_{\Gamma}^*) + \log \det(\Omega_{\Gamma}^*)$ coincide, because $\Omega_{\Gamma}^* = \Omega_{\Gamma'}^*$.
- Thus
$$\frac{p(\Gamma|X^{(n)})}{p(\Gamma'|X^{(n)})} = \frac{C_{\Gamma} \int_{\Delta_{\Gamma} + \Omega_{\Gamma}^* \in \mathcal{M}_0^+} \exp(-\frac{n}{2} g_{\Gamma}(\Delta_{\Gamma})) \prod_{(i,j) \in \bar{E}_{\Gamma}} du_{\Gamma,ij}}{C_{\Gamma'} \int_{\Delta_{\Gamma'} + \Omega_{\Gamma'}^* \in \mathcal{M}_0^+} \exp(-\frac{n}{2} g_{\Gamma'}(\Delta_{\Gamma'})) \prod_{(i,j) \in \bar{E}_{\Gamma'}} du_{\Gamma',ij}}$$
- Recalling the remark 2.2,

$$\begin{aligned} \frac{p(\Gamma|X^{(n)})}{p(\Gamma'|X^{(n)})} &\approx \frac{C_{\Gamma} \int_{\|\Delta_{\Gamma}\|_2 \leq \epsilon_n} \exp(-\frac{n}{2} g_{\Gamma}(\Delta_{\Gamma})) \prod_{(i,j) \in \bar{E}_{\Gamma}} du_{\Gamma,ij}}{C_{\Gamma'} \int_{\|\Delta_{\Gamma'}\|_2 \leq \epsilon_n} \exp(-\frac{n}{2} g_{\Gamma'}(\Delta_{\Gamma'})) \prod_{(i,j) \in \bar{E}_{\Gamma'}} du_{\Gamma',ij}} \\ &\leq \frac{C_{\Gamma} \int_{\|\Delta_{\Gamma}\|_2 \leq \epsilon_n} \exp(-\frac{n}{2} g_{\Gamma}(\Delta_{\Gamma})) \prod_{(i,j) \in \bar{E}_{\Gamma}} du_{\Gamma,ij}}{C_{\Gamma'} \int_{\|\Delta_{\Gamma}\|_2 \leq \epsilon_n} \exp(-\frac{n}{2} g_{\Gamma'}(\Delta_{\Gamma'})) \prod_{(i,j) \in \bar{E}_{\Gamma'}} du_{\Gamma',ij}} \end{aligned} \quad (6)$$

- The term $-\log \det(\Delta_{\Gamma} + \Omega_{\Gamma}^*) + \text{tr}(\hat{\Sigma} \Delta_{\Gamma})$ in g_{Γ} is larger than that in $g_{\Gamma'}$.

Ignorability of non-regular models

- Hence, combining this fact with the inequality (6),

$$\begin{aligned}\frac{p(\Gamma|X^{(n)})}{p(\Gamma'|X^{(n)})} &\leq \frac{C_\Gamma}{C_{\Gamma'}} \int_{\|\Delta_\Gamma\|_2 \leq \epsilon_n} \exp\left(-\frac{n}{2} \frac{2\lambda}{n} \sum_{\gamma_{ij}=1, \gamma'_{ij}=0} |u_{\Gamma,ij}|\right) \prod_{(i,j) \in \bar{E}_\Gamma \cap \bar{E}_{\Gamma'}^c} du_{\Gamma,ij} \\ &\leq \frac{C_\Gamma}{C_{\Gamma'}} \int \exp\left(-\lambda \sum_{\gamma_{ij}=1, \gamma'_{ij}=0} |u_{\Gamma,ij}|\right) \prod_{(i,j) \in \bar{E}_\Gamma \cap \bar{E}_{\Gamma'}} du_{\Gamma,ij} \\ &= \frac{C_\Gamma}{C_{\Gamma'}} \left(\frac{2}{\lambda}\right)^{\#\Gamma - \#\Gamma'} = \left(\frac{q}{1-q}\right)^r \frac{\beta(\Gamma)}{\beta(\Gamma')} \leq \left(\frac{q}{1-q}\right)^r\end{aligned}$$

- The last inequality follows from the fact if the Prior2 is used, then $P(\bar{R} \geq \#\Gamma) \leq P(\bar{R} \geq \#\Gamma')$ since $\#\Gamma > \#\Gamma'$. Similarly, one can see $\frac{\beta(\Gamma)}{\beta(\Gamma')} \leq 1$ for Prior 1.
- For $q < \frac{1}{2}$, $\frac{q}{1-q} < 1$. So we conclude that posterior of Γ' is larger than that of Γ . This leads us to concentrate only on regular model with $q < \frac{1}{2}$.

- 1 Notations and preliminaries
- 2 Prior and posterior concentration
- 3 Posterior computation
 - Approximating model posterior probabilities
 - Ignorability of non-regular models
 - Error in Laplace approximation
- 4 Simulation results

Lemma 3.2

For any regular model Γ , with probability tending to one, the remainder term of function $h_\Gamma(\Omega_\Gamma)$, defined in the beginning of section (3), is bounded by $(p + \#\Gamma)\|\Delta_\Gamma\|_2^2(C_1\|\Delta_\Gamma\|_2 + C_2\|\Delta_\Gamma\|_2^2)/2$ for some constants $C_1, C_2 > 0$.

- By Taylor's expansion,

$$h_\Gamma(\Omega_\Gamma) = h_\Gamma(\Omega_\Gamma^*) + \frac{1}{2}\text{vec}(\Delta_\Gamma)^t H_{\Omega_\Gamma^*} \text{vec}(\Delta_\Gamma) + R_n \quad (7)$$

- Using the integral form of the remainder,

$$h_\Gamma(\Omega_\Gamma) = h_\Gamma(\Omega_\Gamma^*) + \text{vec}(\Delta_\Gamma)^t \left[\int_0^1 (1-v) H_{\Omega_\Gamma^* + v\Delta_\Gamma} dv \right] \text{vec}(\Delta_\Gamma) \quad (8)$$

Error in Laplace approximation

- Subtracting (8) from (7)

$$R_n = \text{vec}(\Delta_\Gamma)^t \left[\int_0^1 (1 - \nu) (H_{\Omega_\Gamma^* + \nu \Delta_\Gamma} - H_{\Omega_\Gamma^*}) d\nu \right] \text{vec}(\Delta_\Gamma) \quad (9)$$

- Applying Cauchy-Schwartz inequality and theorem 1.1 to (9),

$$\begin{aligned} |R_n| &\leq \|\Delta_\Gamma\|_2^2 \left\| \int_0^1 (1 - \nu) (H_{\Omega_\Gamma^* + \nu \Delta_\Gamma} - H_{\Omega_\Gamma^*}) d\nu \right\|_{(2,2)} \\ &\leq \|\Delta_\Gamma\|_2^2 \int_0^1 (1 - \nu) \|H_{\Omega_\Gamma^* + \nu \Delta_\Gamma} - H_{\Omega_\Gamma^*}\|_{(2,2)} d\nu \\ &\leq \frac{1}{2} \|\Delta_\Gamma\|_2^2 \max_{0 \leq \nu \leq 1} \|H_{\Omega_\Gamma^* + \nu \Delta_\Gamma} - H_{\Omega_\Gamma^*}\|_{(2,2)} \\ &\leq \frac{1}{2} \|\Delta_\Gamma\|_2^2 (p + \#\Gamma) \max_{0 \leq \nu \leq 1} \|H_{\Omega_\Gamma^* + \nu \Delta_\Gamma} - H_{\Omega_\Gamma^*}\|_\infty \end{aligned} \quad (10)$$

Error in Laplace approximation

- It can be shown that each entry of $H_{\Omega_\Gamma^* + \nu \Delta_\Gamma} - H_{\Omega_\Gamma^*} \leq C_1 \|\Delta_\Gamma\|_2 + C_2 \|\Delta_\Gamma\|_2^2$ for all $0 \leq \nu \leq 1$ and some constant $C_1, C_2 > 0$.
- Therefore, $\|H_{\Omega_\Gamma^* + \nu \Delta_\Gamma} - H_{\Omega_\Gamma^*}\|_\infty \leq C_1 \|\Delta_\Gamma\|_2 + C_2 \|\Delta_\Gamma\|_2^2$ for all $0 \leq \nu \leq 1$ and so $\max_{0 \leq \nu \leq 1} \|H_{\Omega_\Gamma^* + \nu \Delta_\Gamma} - H_{\Omega_\Gamma^*}\|_\infty \leq C_1 \|\Delta_\Gamma\|_2 + C_2 \|\Delta_\Gamma\|_2^2$
- From inequality (10), we get the desired result

Error in Laplace approximation

Theorem 3.2

The error in the Laplace approximation tends to 0 in probability if $(p + \#\Gamma)^2 \epsilon_n = n^{-\frac{1}{2}}(p + \#\Gamma)^{\frac{5}{2}}(\log p)^{\frac{1}{2}} \rightarrow 0$ in probability, hence asymptotically negligible, where $\epsilon_n = n^{-\frac{1}{2}}(p + \#\Gamma)^{\frac{1}{2}}(\log p)^{\frac{1}{2}}$, which is posterior convergence rate established in theorem 2.2.

(proof)

- With Taylor's expansion,

$$p(\Gamma|X^{(n)}) \propto C_\Gamma \int_{\Omega_\Gamma^* + \Delta_\Gamma \in \mathcal{M}_0^+} \exp\left(-\frac{n}{2}(h_\Gamma(\Omega_\Gamma^*) + \frac{1}{2}\text{vec}(\Delta_\Gamma)^t H_{\Omega_\Gamma^*} \text{vec}(\Delta_\Gamma) + R_n)\right) \prod_{(i,j) \in \bar{E}_\Gamma} du_{\Gamma,ij} \quad (11)$$

- For notational simplicity, denote $\prod_{(i,j) \in \bar{E}} du_{\Gamma,ij}$ by $d\Delta_\Gamma$.

Error in Laplace approximation

- By the remark 2.2,

$$\frac{\int_{\|\Delta_\Gamma\|_2 \leq \epsilon_n} \exp(-\frac{n}{4} \text{vec}(\Delta_\Gamma)^t H_{\Omega_\Gamma^*} \text{vec}(\Delta_\Gamma) - \frac{n}{2} R_n) dU_{\Gamma,ij}}{\int_{\Omega_\Gamma^* + \Delta_\Gamma \in \mathcal{M}_0^+} \exp(-\frac{n}{4} \text{vec}(\Delta_\Gamma)^t H_{\Omega_\Gamma^*} \text{vec}(\Delta_\Gamma) - \frac{n}{2} R_n) dU_{\Gamma,ij}} \rightarrow 1 \quad (12)$$

- With sufficiently large n , for $\|\Delta_\Gamma\|_2 \leq \epsilon_n$, by the result of lemma 3.2, $|R_n| \leq \frac{1}{2} C(p + \#\Gamma) \|\Delta_\Gamma\|_2^2 \epsilon_n$ for some constant $C > 0$
- Thus, the upper and lower bounds of the integral

$\int_{\|\Delta_\Gamma\|_2 \leq \epsilon_n} \exp(-\frac{n}{4} \text{vec}(\Delta_\Gamma)^t H_{\Omega_\Gamma^*} \text{vec}(\Delta_\Gamma) - \frac{n}{2} R_n) dU_{\Gamma,ij}$ are given by

$$\int_{\|\Delta_\Gamma\|_2 \leq \epsilon_n} \exp(-\frac{n}{4} \text{vec}(\Delta_\Gamma)^t (H_{\Omega_\Gamma^*} \mp C(p + \#\Gamma) \epsilon_n I) \text{vec}(\Delta_\Gamma)) dU_{\Gamma,ij}$$

- If $(p + \#\Gamma) \rightarrow 0$, since the minimum eigenvalue of $H_{\Omega_\Gamma^*}$ is bounded away from 0,

$$\int_{\|\Delta_\Gamma\|_2 > \epsilon_n} \exp(-\frac{n}{4} \text{vec}(\Delta_\Gamma)^t (H_{\Omega_\Gamma^*} \mp C(p + \#\Gamma) \epsilon_n I) \text{vec}(\Delta_\Gamma)) dU_{\Gamma,ij} \rightarrow 0.$$

Error in Laplace approximation

- Thus the bounds can be simplified into

$$\int_{\Delta_{\Gamma} + \Omega_{\Gamma}^* \in \mathcal{M}_0^+} \exp\left(-\frac{n}{4} \text{vec}(\Delta_{\Gamma})^t (H_{\Omega_{\Gamma}^*} \mp C(p + \#\Gamma)\epsilon_n I) \text{vec}(\Delta_{\Gamma})\right) d u_{\Gamma, ij}$$

- Using the above bound, the ratio of the actual integral to the approximate integral has upper bound and lower bounds given by

$$\frac{\int_{\Delta_{\Gamma} + \Omega_{\Gamma}^* \in \mathcal{M}_0^+} \exp\left(-\frac{n}{4} \text{vec}(\Delta_{\Gamma})^t (H_{\Omega_{\Gamma}^*} \mp C(p + \#\Gamma)\epsilon_n I) \text{vec}(\Delta_{\Gamma})\right) d u_{\Gamma, ij}}{\int_{\Omega_{\Gamma}^* + \Delta_{\Gamma} \in \mathcal{M}_0^+} \exp\left(-\frac{n}{4} \text{vec}(\Delta_{\Gamma})^t H_{\Omega_{\Gamma}^*} \text{vec}(\Delta_{\Gamma})\right) d u_{\Gamma, ij}} \\ = [\det(I \mp C(p + \#\Gamma)\epsilon_n (H_{\Omega_{\Gamma}^*})^{-1})]^{-\frac{1}{2}} \quad (13)$$

- The above expression must lie between

$[1 \mp C(p + \#\Gamma)\epsilon_n \{\text{eig}_1(H_{\Omega_{\Gamma}^*})\}^{-1}]^{-\frac{p + \#\Gamma}{2}}$, as the dimension of $H_{\Omega_{\Gamma}^*}$ is $p + \#\Gamma$. Again the minimum eigenvalue of $H_{\Omega_{\Gamma}^*}$ is bounded away from 0. So if $(p + \#\Gamma)^2 \epsilon_n \rightarrow 0$, the above bound on the ratio goes to 1.

Outline

- 1 Notations and preliminaries
- 2 Prior and posterior concentration
- 3 Posterior computation
 - Approximating model posterior probabilities
 - Ignorability of non-regular models
 - Error in Laplace approximation
- 4 Simulation results

- Simulation with models with sparse precision matrix
- The mean of models used is 0. The models with precision matrix $\Omega = ((\omega_{ij}))$ or covariance matrix $\Sigma = ((\sigma_{ij}))$ used in the simulation are followings:

AR(1): $\sigma_{ij} = 0.7^{|i-j|}$

AR(2): $\omega_{ii} = 1, \omega_{i-1,i} = \omega_{i,i-1} = 0.5, \omega_{i,i-2} = \omega_{i-2,i} = 0.25$

Star: $\omega_{ii} = 1, \omega_{1,i} = \omega_{i,1} = 0.1$ and $\omega_{ij} = 0$, otherwise.

Circle: $\omega_{ii} = 2, \omega_{i-1,i} = \omega_{i,i-1} = 1, \omega_{1,p} = \omega_{p,1} = 0.9$

- We use specificity(SP), sensitivity(SE), and Matthews Correlation Coefficient(MCC) to measure the performance of the algorithm which will be described soon. Here, specificity and sensitivity are defined in the context of edge.

$$\begin{aligned} SP &= \frac{TN}{TN + FP}, \quad SE = \frac{TP}{TP + FN}, \\ MCC &= \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \end{aligned} \quad (14)$$

- where TP,TN,FP, and FN respectively denote the true positives(edges included which are present in the true model), true negatives(edges excluded which are absent in the true model), false positives(edges included which are absent in the true model), and false negatives(edges excluded which are present in the true model)

Simulation

- Note that for estimated $\hat{\Omega}$, we say there is no edge between node i and j if $|\hat{\omega}_{ij}| \leq 0.1^3$.
- To compare Bayesian method with frequentist method, we also do simulation for frequentist method with the same models.
- Corresponding to each model, generate samples of size $n = 100, 200$ and dimension $p = 30, 50, 100$.
- The penalty parameter λ for the graphical lasso algorithm is chosen such that $\frac{\lambda}{n} = 0.5$ for model AR(1), AR(2), and Circle, and $\frac{\lambda}{n} = 0.2$ for Star. Set $q = 0.4$.
- Run 50 replications of each model and calculate SP, SE, and MCC in each replication. Then, average each measure over the replications and calculate standard deviation of each measure also.
- Provide ROC curves corresponding to the various models with values of the penalty parameter ranging between 0 1, plotting Sensitivity against False Positive Rate in case $n = 100$ and $p = 30, 50$

- For Bayesian method, we set the restriction of model size \bar{r} by that of frequentist method with the same penalty parameter.
- The idea to find posterior mode is not good. For example, in AR(1) model with $n = 100, p = 30$, we have $\bar{r} = 29$. We have to search total $\sum_{k=0}^{29} \binom{\binom{p}{2}}{k}$ models. Considering $\binom{\binom{p}{2}}{k} \geq 10^{45}$ when $k = 29$, this is impossible.
- Also, the posterior tends to decay.
- Hence we need to resort to other method. One possibility was RJMCMC, which turned out to be inappropriate for simulation.

Simulation

- We may use Metropolis-Hastings algorithm to obtain MCMC samples of model indicator Γ . In each replication of model, we choose final model by Median Probability Model (denoted by "MPP"). We may suggest symmetric proposal distribution. (500 MCMC sample with 100 burn-in)
- Suppose there are two p -dimension models $G = (V, E)$ and $G' = (V, E')$ with model indicators Γ, Γ' respectively. Considering Γ as vector, we denote $\Gamma(i, j) = 1$ if there is edge between node i and j and $\Gamma(i, j) = 0$ otherwise. Let $\pi(\Gamma|\Gamma')$ symmetric proposal distribution, which samples Γ' uniformly from the graphs that differ from Γ in one position. To be more specific, we set $\Gamma'(i, j) = 1$ for one pair (i, j) uniformly sampled from $V \times V \setminus E$, i.e., $\Gamma(i, j) = 0$ and $\Gamma'(i, j) = 0$ for one pair (i, j) uniformly sampled from E , i.e., $\Gamma(i, j) = 1$. So we have equal probability to set $\Gamma'(i, j) = 1$ or $\Gamma'(i, j) = 0$. Set $\Gamma'(-i, -j) = \Gamma(i, j)$, i.e. set Γ' as the same as Γ except for pair (i, j) . Using this symmetric proposal distribution, we have following algorithm for MCMC sampling of Γ .

Algorithm 1 Metropolis-Hastings algorithm for sampling Γ

- 1: Initial model indicator : $\Gamma^{(0)}$, Given data : $\mathbf{X}^{(n)}$
 - 2: **for** $i=0,1,\dots,k$ **do**
 - 3: $\Gamma^{\text{temp}} \sim \pi(\Gamma|\Gamma^{(i)})$
 - 4: Check regularity of Γ^{temp}
 - 5: If Γ^{temp} is regular, $\Gamma^{\text{cand}} = \Gamma^{\text{temp}}$. Else, repeat 3 \sim 4.
 - 6: $\alpha_i = \min\{1, \frac{p^*(\Gamma^{\text{cand}}|\mathbf{X}^{(n)})}{p^*(\Gamma^{(i)}|\mathbf{X}^{(n)})}\}$
 - 7: $U_i \sim U(0, 1)$
 - 8: If $U_i \leq \alpha_i$, $\Gamma^{(i+1)} = \Gamma^{\text{cand}}$. Else, $\Gamma^{(i+1)} = \Gamma^{(i)}$.
 - 9: **end for**
-

Simulation

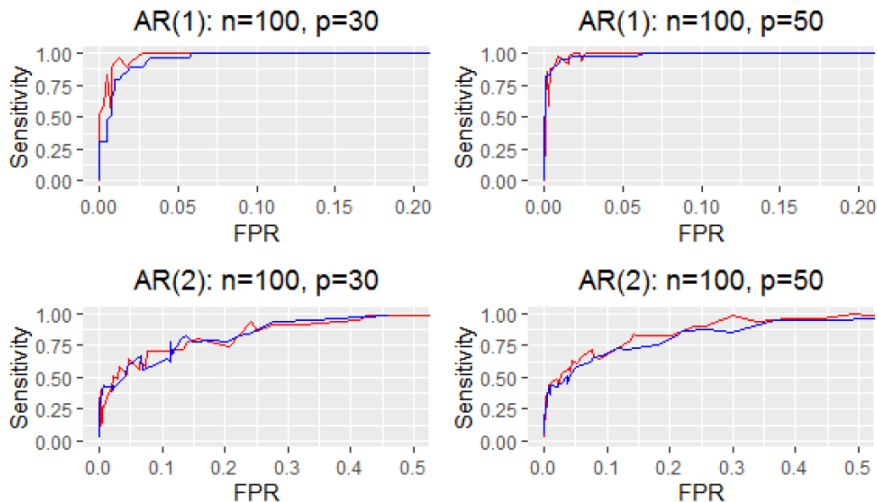


Figure 1: ROC curves for AR(1) and AR(2) structures. Red line : GL, Blue line: MPP

Simulation

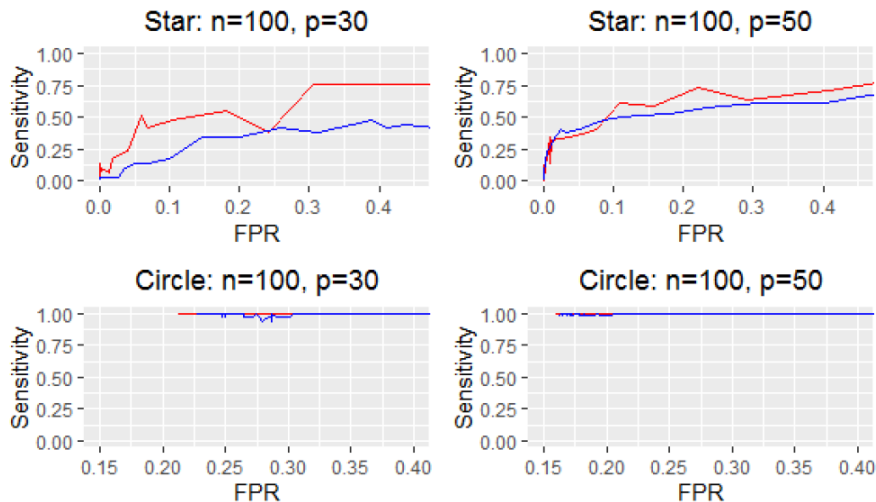


Figure 2: ROC curves for Star and Circle. Red line : GL, Blue line: MPP

Simulation

Model	p	n=100						n=200					
		MPP			GL			MPP			GL		
		SP	SE	MCC	SP	SE	MCC	SP	SE	MCC	SP	SE	MCC
AR(1)	30	0.977	0.919	0.813	0.977	0.970	0.844	0.980	0.972	0.861	0.982	0.992	0.886
		(0.009)	(0.054)	(0.048)	(0.011)	(0.045)	(0.055)	(0.007)	(0.038)	(0.041)	(0.007)	(0.012)	(0.035)
	50	0.986	0.893	0.802	0.985	0.969	0.839	0.990	0.962	0.869	0.990	0.992	0.888
		(0.004)	(0.044)	(0.038)	(0.005)	(0.038)	(0.041)	(0.003)	(0.032)	(0.032)	(0.004)	(0.011)	(0.035)
	100	0.994	0.872	0.804	0.993	0.957	0.837	0.995	0.959	0.878	0.995	0.994	0.888
		(0.001)	(0.038)	(0.030)	(0.002)	(0.029)	(0.029)	(0.001)	(0.027)	(0.028)	(0.001)	(0.010)	(0.022)
AR(2)	30	0.983	0.434	0.546	0.978	0.475	0.562	0.988	0.485	0.613	0.988	0.488	0.614
		(0.007)	(0.050)	(0.046)	(0.010)	(0.044)	(0.038)	(0.007)	(0.031)	(0.029)	(0.006)	(0.029)	(0.029)
	50	0.987	0.409	0.518	0.983	0.483	0.559	0.993	0.464	0.609	0.993	0.485	0.621
		(0.004)	(0.034)	(0.032)	(0.005)	(0.032)	(0.026)	(0.003)	(0.024)	(0.026)	(0.003)	(0.023)	(0.027)
	100	0.991	0.400	0.499	0.989	0.475	0.536	0.997	0.442	0.607	0.996	0.486	0.630
		(0.002)	(0.023)	(0.022)	(0.002)	(0.026)	(0.022)	(0.001)	(0.019)	(0.017)	(0.001)	(0.017)	(0.013)
Star	30	0.965	0.229	0.225	0.954	0.297	0.256	0.990	0.167	0.263	0.994	0.237	0.384
		(0.006)	(0.077)	(0.080)	(0.013)	(0.104)	(0.098)	(0.006)	(0.093)	(0.122)	(0.004)	(0.084)	(0.109)
	50	0.965	0.319	0.264	0.948	0.479	0.331	0.994	0.387	0.516	0.994	0.498	0.606
		(0.007)	(0.074)	(0.061)	(0.009)	(0.104)	(0.067)	(0.002)	(0.078)	(0.074)	(0.003)	(0.088)	(0.061)
	100	0.940	0.992	0.484	0.940	1.000	0.489	0.988	1.000	0.793	0.988	1.000	0.787
		(0.005)	(0.010)	(0.018)	(0.004)	(0.000)	(0.013)	(0.002)	(0.000)	(0.019)	(0.002)	(0.000)	(0.023)
Circle	30	0.784	1.000	0.448	0.746	1.000	0.411	0.754	1.000	0.418	0.724	1.000	0.392
		(0.021)	(0.000)	(0.022)	(0.015)	(0.000)	(0.013)	(0.013)	(0.000)	(0.012)	(0.018)	(0.000)	(0.015)
	50	0.828	0.976	0.395	0.830	1.000	0.408	0.836	0.996	0.413	0.834	1.000	0.412
		(0.010)	(0.019)	(0.016)	(0.009)	(0.000)	(0.011)	(0.007)	(0.009)	(0.010)	(0.006)	(0.000)	(0.008)
	100	0.892	0.984	0.371	0.891	1.000	0.376	0.905	0.997	0.402	0.904	1.000	0.400
		(0.003)	(0.012)	(0.008)	(0.004)	(0.000)	(0.006)	(0.003)	(0.006)	(0.007)	(0.003)	(0.000)	(0.006)

- Bayesian method performs slightly better than the frequentist method in terms of specificity, but suffers in sensitivity.
- This is because we chose $q < \frac{1}{2}$ to focus on regular models, allowing less edges to the model.
- The performance of Bayesian model tends to get lower as the dimension grows.
- Sensitivity was not good in AR(2) and Star models for both methods.
- In ROC curves, as the penalty parameter λ tends to get larger, sensitivity grows at the cost of higher false positive rate.

Illustration with real data

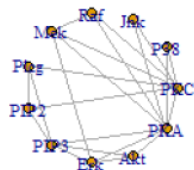
- Check tendency of performance of Bayesian method with real data
- Protein data from R package sparsebn consisting of $n = 7466$ observations of $p = 11$ continuous variables corresponding to different proteins in human immune system cells.
- Penalty parameter λ is chosen such that $\frac{\lambda}{n} = 0.4$ and $q = 0.4$.
- For Bayesian method, we obtain 10000 MCMC samples with 2000 burn-in.

Illustration with real data

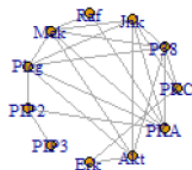
Method	SP	SE	MCC
MPP	0.636	0.364	0.000
GL	0.606	0.455	0.0602

Table 2: Simulation result for protein data

True Network



Frequentist Network



Bayesian Network



Illustration with real data

- The model estimated by Bayesian method shows less edges, hence suffers in sensitivity.
- But frequentist method performs is not good in specificity compared to Bayesian method and more edges are connected to Jnk, Akt in the model estimated by frequentist method than true network. Along with table 2, one can see what we wanted to verify with real data.

- Bayesian method has mainly two drawbacks. Efficiency and Sensitivity.
- To explain poor efficiency, we need to check regularity condition and sample Γ again until we have regular model.
- Also, in estimating precision matrix Ω_{Γ} corresponding to regular model Γ , we use frequentist method. The algorithm of it requires very long time as the dimension grows especially in sparse models.
- However, the method is brilliant in that it introduced the sparsity to the model. One may hope this method to be improved in efficiency and sensitivity. Also, one can consider the problem of choosing penalty parameter λ .

- Ghosal, Ghosh and Van Der Vaart, "Convergence Rates of Posterior Distributions", Annals of Statistics, 2000, 28(2), 500-531
- Wang, "Bayesian Graphical Lasso Models and Efficient Posterior Computation", Bayesian Analysis, 2012, 7(4), 867-886
- A.J.Rothman et al, "Sparse permutation invariant covariance estimation", Electronic Journal of Statistics, 2008, 2, 494-515
- Liu and Martin, "An empirical G-Wishart prior for sparse high-dimensional Gaussian graphical models", arXiv, 2019, 1912.03807, 1-36

The End