

Graphical models for high-dimensional data

Tianhao Wang, Liangchen He

Department of Statistics and Finance, USTC

November 20, 2022



- 1 Some basics
- 2 Estimation of Gaussian graphical models
- 3 Graphical models in exponential form
- 4 Graphs with corrupted or hidden variables

- 1 Some basics
- 2 Estimation of Gaussian graphical models
- 3 Graphical models in exponential form
- 4 Graphs with corrupted or hidden variables

Undirected graphical models

- An undirected graph $G = (V, E)$ consists of a set of vertices $V = 1, 2, \dots, d$ joined together by a collection of edges E . An edge (j, k) is an unordered pair of distinct vertices $j, k \in V$.
- We associate to each vertex $j \in V$ a random variable X_j , taking values in some space \mathcal{X}_j . We then consider the distribution \mathbb{P} of the d -dimensional random vector $X = (X_1, \dots, X_d)$.
- A clique C means that $(j, k) \in E$ for all distinct vertices $j, k \in C$. A maximal clique is a clique that is not a subset of any other clique. We use \mathfrak{C} to denote the set of all cliques in G , and for each clique $C \in \mathfrak{C}$, we use ψ_C to denote a function of the subvector $x_C := (x_j, j \in C)$. This clique compatibility function takes inputs from the Cartesian product space $\mathcal{X}^C := \bigotimes_{j \in C} \mathcal{X}_j$, and returns non-negative real numbers.

Factorization

Definition (11.1)

The random vector (X_1, \dots, X_d) factorizes according to the graph G if its density function p can be represented as

$$p(x_1, \dots, x_d) \propto \prod_{C \in \mathcal{C}} \psi_C(x_C) \quad (11.1)$$

for some collection of clique compatibility functions $\psi_C : \mathcal{X}^C \rightarrow [0, \infty)$.

Example 11.3

Any non-degenerate Gaussian distribution with zero mean can be parameterized in terms of its inverse covariance matrix $\Theta^* = \Sigma^{-1}$, also known as the precision matrix. In particular, its density can be written as

$$p(x_1, \dots, x_d; \Theta^*) = \frac{\sqrt{\det(\Theta^*)}}{(2\pi)^{d/2}} e^{-\frac{1}{2}x^T \Theta^* x} \quad (11.2)$$

By expanding the quadratic form, we see that

$$e^{-\frac{1}{2}x^T \Theta^* x} = \exp\left(-\frac{1}{2} \sum_{(j,k) \in E} \Theta_{jk}^* x_j x_k\right) = \prod_{(j,k) \in E} \underbrace{-\frac{1}{2} \Theta_{jk}^* x_j x_k}_{\psi_{jk}(x_j, x_k)},$$

showing that any zero-mean Gaussian distribution can be factorized in terms of functions on edges, or cliques of size two.

Conditional independence

- A vertex cutset S is a subset of vertices whose removal from the graph breaks it into two or more disjoint pieces.
- Removing S from the vertex set V leads to the vertex-induced subgraph $G(V \setminus S)$, consisting of the vertex set $V \setminus S$, and the residual edge set

$$E(V \setminus S) := \{(j, k) \in E \mid j, k \in V \setminus S\}. \quad (11.4)$$

- The set S is a vertex cutset if the residual graph $G(V \setminus S)$ consists of two or more disconnected non-empty components.

For any subset $A \subseteq V$, let $X_A := (X_j, j \in A)$ represent the subvector of random variables indexed by vertices in A . For any three disjoint subsets, say A , B and S , of the vertex set V , we use $X_A \perp\!\!\!\perp X_B \mid X_S$ to mean that the subvector X_A is conditionally independent of X_B given X_S .

Definition (11.5)

A random vector $X = (X_1, \dots, X_d)$ is Markov with respect to a graph G if, for all vertex cutsets S breaking the graph into disjoint pieces A and B , the conditional independence statement $X_A \perp\!\!\!\perp X_B \mid X_S$ holds.

Example 11.6

The Markov chain graph on vertex set $V = \{1, 2, \dots, d\}$ contains the edges $(j, j + 1)$ for $j = 1, 2, \dots, d - 1$. For such a chain graph, each vertex $j \in \{2, 3, \dots, d - 1\}$ is a non-trivial cutset, breaking the graph into the "past" $P = \{1, 2, \dots, j - 1\}$ and "future" $F = \{j + 1, \dots, d\}$. These singleton cutsets define the essential Markov property of a Markov time-series model—namely, that the past X_P and future X_F are conditionally independent given the present X_j .

Theorem (11.8 Hammersley-Clifford)

For a given undirected graph and any random vector $X = (X_1, \dots, X_d)$ with strictly positive density p , the following two properties are equivalent:

- (a) The random vector X factorizes according to the structure of the graph G , as in Definition 11.1.*
- (b) The random vector X is Markov with respect to the graph G , as in Definition 11.5.*

Proof

(a) \Rightarrow (b): Let S be an arbitrary vertex cutset of the graph such that subsets A and B are separated by S . We may assume without loss of generality that both A and B are non-empty, and we need to show that $X_A \perp\!\!\!\perp X_B \mid X_S$.

- $\mathfrak{C}_A := \{C \in \mathfrak{C} \mid C \cap A \neq \emptyset\}$, $\mathfrak{C}_B := \{C \in \mathfrak{C} \mid C \cap B \neq \emptyset\}$ and $\mathfrak{C}_S := \{C \in \mathfrak{C} \mid C \subseteq S\}$.
- $\mathfrak{C} = \mathfrak{C}_A \cup \mathfrak{C}_S \cup \mathfrak{C}_B$.

$$p(x_A, x_S, x_B) = \frac{1}{Z} \underbrace{\left[\prod_{C \in \mathfrak{C}_A} \psi_C(x_C) \right]}_{\Psi_A(x_A, x_S)} \underbrace{\left[\prod_{C \in \mathfrak{C}_S} \psi_C(x_C) \right]}_{\Psi_S(x_S)} \underbrace{\left[\prod_{C \in \mathfrak{C}_B} \psi_C(x_C) \right]}_{\Psi_B(x_B, x_S)}.$$

Defining the quantities

$$Z_A(x_S) := \sum_{x_A} \psi_A(x_A, x_S) \text{ and } Z_B(x_S) := \sum_{x_B} \psi_B(x_B, x_S),$$

we then obtain the following expressions for the marginal distributions of interest:

$$p(x_S) = \frac{Z_A(x_S) Z_B(x_S)}{Z} \psi_S(x_S),$$

$$p(x_A, x_S) = \frac{Z_B(x_S)}{Z} \psi_A(x_A, x_S) \psi_S(x_S),$$

$$p(x_B, x_S) = \frac{Z_A(x_S)}{Z} \psi_B(x_B, x_S) \psi_S(x_S).$$

for any x_S for which $p(x_S) > 0$

$$\frac{p(x_A, x_S, x_B)}{p(x_S)} = \frac{\psi_A(x_A, x_S) \psi_B(x_B, x_S)}{Z_A(x_S) Z_B(x_S)},$$

$$\frac{p(x_A, x_S)}{p(x_S)} = \frac{\psi_A(x_A, x_S)}{Z_A(x_S)},$$

$$\frac{p(x_B, x_S)}{p(x_S)} = \frac{\psi_B(x_B, x_S)}{Z_B(x_S)}.$$

$$\begin{aligned} \Rightarrow p(x_A, x_B | x_S) &= \frac{p(x_A, x_B, x_S)}{p(x_S)} = \frac{p(x_A, x_S)}{p(x_S)} \frac{p(x_B, x_S)}{p(x_S)} \\ &= p(x_A | x_S) p(x_B | x_S). \end{aligned}$$

- 1 Some basics
- 2 Estimation of Gaussian graphical models
- 3 Graphical models in exponential form
- 4 Graphs with corrupted or hidden variables

The graphical Lasso estimator

$$\widehat{\Theta} \in \arg \min_{\Theta \in \mathcal{S}^{d \times d}} \underbrace{\{\langle \Theta, \widehat{\Sigma} \rangle - \log \det \Theta\}}_{\mathcal{L}_n(\Theta)} + \lambda_n \|\Theta\|_{1, \text{off}},$$

where $\|\Theta\|_{1, \text{off}} := \sum_{j \neq k} |\Theta_{jk}|$ corresponds to the ℓ_1 -norm applied to the off-diagonal entries of Θ .

The following result is based on a sample covariance matrix $\widehat{\Sigma}$ formed from n i.i.d. samples $\{x_i\}_{i=1}^n$ of a zero-mean random vector in which each coordinate has σ -sub-Gaussian tails.

Theorem (11.9 Frobenius norm bounds for graphical Lasso)

Suppose that the inverse covariance matrix Θ^* has at most m non-zero entries per row, and we solve the graphical Lasso (11.10) with regularization parameter $\lambda_n = 8\sigma^2 \left(\sqrt{\frac{\log d}{n}} + \delta \right)$ for some $\delta \in (0, 1]$. Then as long as $(\|\Theta^*\|_2 + 1)^2 \lambda_n \sqrt{md} < 1$, the graphical Lasso estimate $\widehat{\Theta}$ satisfies

$$\left\| \widehat{\Theta} - \Theta^* \right\|_F^2 \leq \frac{9}{(\|\Theta^*\|_2 + 1)^4} md \lambda_n^2 \quad (11.11)$$

with probability at least $1 - 8e^{-\frac{1}{16}n\delta^2}$.

Proof

- $\mathbb{B}_F(1) = \{\Delta \in \mathcal{S}^{d \times d} \mid \|\Delta\|_F \leq 1\}$
- $\nabla \mathcal{L}_n(\Theta) = \widehat{\Sigma} - \Theta^{-1}$ and $\nabla^2 \mathcal{L}_n(\Theta) = \Theta^{-1} \otimes \Theta^{-1}$
- $\underbrace{\mathcal{L}_n(\Theta^* + \Delta) - \mathcal{L}_n(\Theta^*) - \langle \nabla \mathcal{L}_n(\Theta^*), \Delta \rangle}_{\mathcal{E}_n(\Delta)} =$
 $\frac{1}{2} \text{vec}(\Delta)^T \nabla^2 \mathcal{L}_n(\Theta^* + t\Delta) \text{vec}(\Delta)$ for some $t \in [0, 1]$
- $\|\mathbf{A}^{-1} \otimes \mathbf{A}^{-1}\|_2 = \frac{1}{\|\mathbf{A}\|_2^2}$ for any symmetric invertible matrix

Verifying restricted strong convexity

For any $\Delta \in \mathbb{B}_F(1)$,

$$\mathcal{E}_n(\Delta) \geq \frac{1}{2} \gamma_{\min} \left(\nabla^2 \mathcal{L}_n(\Theta^* + t\Delta) \right) \|\text{vec}(\Delta)\|_2^2 = \frac{1}{2} \frac{\|\Delta\|_F^2}{\|\Theta^* + t\Delta\|_2^2}.$$

$t\|\Delta\|_2 \leq t\|\Delta\|_F \leq 1$ implies that $\|\Theta^* + t\Delta\|_2^2 \leq (\|\Theta^*\|_2 + 1)^2$. Then

$$\mathcal{E}_n(\Delta) \geq \frac{\kappa}{2} \|\Delta\|_F^2 \quad \text{where } \kappa := \left(\|\Theta^*\|_2 + 1 \right)^{-2}$$

showing that the RSC condition from Definition 9.15 holds over $\mathbb{B}_F(1)$ with tolerance $\tau_n^2 = 0$.

Computing the subspace Lipschitz constant

Letting S denote the support set of Θ^* , we define the subspace $\mathbb{M}(S) = \{\Theta \in \mathbb{R}^{d \times d} \mid \Theta_{jk} = 0 \text{ for all } (j, k) \notin S\}$. Then we have

$$\psi^2(\mathbb{M}(S)) = \sup_{\Theta \in \mathbb{M}(S)} \frac{(\sum_{j \neq k} |\Theta_{jk}|)^2}{\|\Theta\|_F^2} \leq |S| \stackrel{(i)}{\leq} md.$$

where inequality (i) follows since Θ^* has at most m non-zero entries per row.

Verifying event $\mathbb{G}(\lambda_n)$

Using Lemma 6.26, we have

$$\mathbb{P} \left[\|\widehat{\Sigma} - \Sigma\|_{\max, \text{off}} \geq \sigma^2 t \right] \leq 8e^{-\frac{\pi}{16} \min\{t, t^2\} + 2 \log d} \quad \text{for all } t > 0.$$

Setting $t = \lambda_n / \sigma^2$ shows that the event $\mathbb{G}(\lambda_n)$ from Corollary 9.20 holds with the claimed probability. Consequently, Proposition 9.13 implies that the error matrix $\widehat{\Delta}$ satisfies the bound

$$\|\widehat{\Delta}_{S^c}\|_1 \leq 3 \|\widehat{\Delta}_S\|_1, \text{ and hence } \|\widehat{\Delta}\|_1 \leq 4 \|\widehat{\Delta}_S\|_1 \leq 4 \sqrt{md} \|\widehat{\Delta}\|_F,$$

where the final inequality again uses the fact that $|S| \leq md$. In order to apply Corollary 9.20, the only remaining detail to verify is that $\widehat{\Delta}$ belongs to the Frobenius ball $\mathbb{B}_F(1)$.

Localizing the error matrix

The result of Exercise 9.10 then implies that

$$\langle \langle \nabla \mathcal{L}_n(\Theta^* + \Delta) - \nabla \mathcal{L}_n(\Theta^*), \Delta \rangle \rangle \geq \kappa \|\Delta\|_F \quad \text{for all } \Delta \in \mathcal{S}^{d \times d} \setminus \mathbb{B}_F(1).$$

By the optimality of $\widehat{\Theta}$, we have $0 = \langle \langle \nabla \mathcal{L}_n(\Theta^* + \widehat{\Delta}) + \lambda_n \widehat{\mathbf{Z}}, \widehat{\Delta} \rangle \rangle$, where $\widehat{\mathbf{Z}} \in \partial \|\widehat{\Theta}\|_{1, \text{off}}$. By adding and subtracting terms, we find that

$$\begin{aligned} \langle \langle \nabla \mathcal{L}_n(\Theta^* + \widehat{\Delta}) - \nabla \mathcal{L}_n(\Theta^*), \widehat{\Delta} \rangle \rangle &\leq \lambda_n |\langle \widehat{\mathbf{Z}}, \widehat{\Delta} \rangle| + \langle \langle \nabla \mathcal{L}_n(\Theta^*), \widehat{\Delta} \rangle \rangle \\ &\leq \left\{ \lambda_n + \|\nabla \mathcal{L}_n(\Theta^*)\|_{\max} \right\} \|\widehat{\Delta}\|_1. \end{aligned}$$

the right-hand side is at most $\frac{3\lambda_n}{2} \|\widehat{\Delta}\|_1 \leq 6\lambda_n \sqrt{md} \|\widehat{\Delta}\|_F$. If $\|\widehat{\Delta}\|_F > 1$, then we obtain $\kappa \|\widehat{\Delta}\|_F \leq \frac{3\lambda_n}{2} \|\widehat{\Delta}\|_1 \leq 6\lambda_n \sqrt{md} \|\widehat{\Delta}\|_F$. This inequality leads to a contradiction whenever $\frac{6\lambda_n \sqrt{md}}{\kappa} < 1$, which completes the proof.

- $S := E \cup \{(j, j) \mid j \in V\}$, $S^c = (V \times V) \setminus S$.
- the matrix $\Gamma^* := \nabla^2 \mathcal{L}_n(\Theta^*)$ is α -incoherent if

$$\max_{e \in S^c} \left\| \Gamma_{eS}^* (\Gamma_{SS}^*)^{-1} \right\|_1 \leq 1 - \alpha \quad \text{for some } \alpha \in (0, 1].$$
- $\widehat{E} := \{(j, k) \in [d] \times [d] \mid j < k \text{ and } \widehat{\Theta}_{jk} \neq 0\}$.
- $\|\mathbf{A}\|_2 \leq m + 1$ for any graph of degree at most m .

Theorem (11.10)

Consider a zero-mean d -dimensional Gaussian distribution based on an α -incoherent inverse covariance matrix Θ^* . Given a sample size lower bounded as $n > c_0 \left(1 + 8\alpha^{-1}\right)^2 m^2 \log d$, suppose that we solve the graphical Lasso (11.10) with a regularization

parameter $\lambda_n = \frac{c_1}{\alpha} \sqrt{\frac{\log d}{n}} + \delta$ for some $\delta \in (0, 1]$. Then with

probability at least $1 - c_2 e^{-c_3 n \delta^2}$, we have the following:

(a) The graphical Lasso solution leads to no false inclusions-that is, $\widehat{\Theta}_{jk} = 0$ for all $(j, k) \notin E$.

(b) It satisfies the sup-norm bound

$$\left\| \widehat{\Theta} - \Theta^* \right\|_{\max} \leq \underbrace{c_4 \left\{ \left(1 + 8\alpha^{-1}\right) \sqrt{\frac{\log d}{n}} \right\}}_{\tau(n, d, \alpha)} + \lambda_n. \quad (11.16)$$

Corollary (11.11)

Under the conditions of Proposition 11.10, consider the graphical Lasso estimate $\widehat{\Theta}$ with regularization parameter $\lambda_n = \frac{c_1}{\alpha} \sqrt{\frac{\log d}{n}} + \delta$ for some $\delta \in (0, 1]$. Then with probability at least $1 - c_2 e^{-c_3 n \delta^2}$, we have

$$\|\widehat{\Theta} - \Theta^*\|_2 \leq c_4 \|\mathbf{A}\|_2 \left\{ (1 + 8\alpha^{-1}) \sqrt{\frac{\log d}{n}} + \lambda_n \right\}, \quad (11.17a)$$

where \mathbf{A} denotes the adjacency matrix of the graph G (including ones on the diagonal). In particular, if the graph has maximum degree m , then

$$\|\widehat{\Theta} - \Theta^*\|_2 \leq c_4 (m + 1) \left\{ (1 + 8\alpha^{-1}) \sqrt{\frac{\log d}{n}} + \lambda_n \right\}. \quad (11.17b)$$

Neighborhood-based methods

- $\mathcal{N}^+(j) := \{j\} \cup \mathcal{N}(j)$, $X_j \perp\!\!\!\perp X_{V \setminus \mathcal{N}^+(j)} \mid X_{\mathcal{N}(j)}$.

Lasso-based neighborhood regression:

1 For each node $j \in V$:

(a) Extract the column vector $X_j \in \mathbb{R}^n$ and the submatrix

$\mathbf{X}_{\setminus\{j\}} \in \mathbb{R}^{n \times (d-1)}$.

(b) Solve the Lasso problem:

$$\widehat{\theta} = \arg \min_{\theta \in \mathbb{R}^{d-1}} \left\{ \frac{1}{2n} \|X_j - \mathbf{X}_{\setminus\{j\}} \theta\|_2^2 + \lambda_n \|\theta\|_1 \right\}.$$

(c) Return the neighborhood estimate $\widehat{\mathcal{N}}(j) = \{k \in V \setminus \{j\} \mid \widehat{\theta}_k \neq 0\}$.

2 Combine the neighborhood estimates to form an edge estimate \widehat{E} , using either the OR rule or the AND rule.

Assume $\text{diag}(\Sigma^*) \leq 1$ and $n \gtrsim m \log d$.

Theorem (11.12)

Consider a zero-mean Gaussian random vector with covariance Σ^* such that for each $j \in V$, the submatrix $\Sigma_{\setminus\{j\}}^* := \text{cov}(\mathbf{X}_{\setminus\{j\}})$ is α -incoherent with respect to $\mathcal{N}(j)$, and $\left\| \left(\Sigma_{\mathcal{N}(j), \mathcal{N}(j)}^* \right)^{-1} \right\|_{\infty} \leq b$ for some $b \geq 1$. Suppose that the neighborhood Lasso selection method is implemented with $\lambda_n = c_0 \left\{ \frac{1}{\alpha} \sqrt{\frac{\log d}{n}} + \delta \right\}$ for some $\delta \in (0, 1]$. Then with probability greater than $1 - c_2 e^{-c_3 n \min\{\delta^2, \frac{1}{m}\}}$, the estimated edge set \widehat{E} , based on either the AND or OR rules, has the following properties: (a) No false inclusions: it includes no false edges, so that $\widehat{E} \subseteq E$. (b) All significant edges are captured: it includes all edges (j, k) for which $|\Theta_{jk}^*| \geq 7b\lambda_n$.

- $\Gamma^* = \text{cov}(X_{\setminus \{j\}}), \widehat{\Gamma} = \frac{1}{n} \mathbf{X}_{\setminus \{j\}}^T \mathbf{X}_{\setminus \{j\}}, S = \mathcal{N}(j), S^c = V \setminus \mathcal{N}^+(j).$
- $\widehat{\Gamma}_{SS} :=$ the submatrix indexed by the subset S .

Proof of part (a)

We follow the proof of Theorem 7.21 until equation (7.53), namely

$$\widehat{\mathbf{z}}_{S^c} = \underbrace{\widehat{\mathbf{\Gamma}}_{S^c S} \left(\widehat{\mathbf{\Gamma}}_{SS} \right)^{-1}}_{\mu \in \mathbb{R}^{d-s}} \widehat{\mathbf{z}}_S + \underbrace{\mathbf{x}_{S^c}^T \left[\mathbf{I}_n - \mathbf{x}_S \left(\mathbf{x}_S^T \mathbf{x}_S \right)^{-1} \mathbf{x}_S^T \right]}_{V_{S^c} \in \mathbb{R}^{d-s}} \left(\frac{W_j}{\lambda_n n} \right)$$

As argued in Chapter 7, in order to establish that the Lasso support is included within S , it suffices to establish the strict dual feasibility condition $\|\widehat{\mathbf{z}}_S\|_\infty < 1$. We do so by establishing that

$$\mathbb{P} \left[\|\mu\|_\infty \geq 1 - \frac{3}{4}\alpha \right] \leq c_1 e^{-c_2 n \alpha^2 - \log d}, \quad (11.25a)$$

$$\mathbb{P} \left[\|V_{S^c}\|_\infty \geq \frac{\alpha}{4} \right] \leq c_1 e^{-c_2 n \delta^2 \alpha^2 - \log d}. \quad (11.25b)$$

Taken together, these bounds ensure that $\|\bar{z}_{S^c}\|_\infty \leq 1 - \frac{\alpha}{2} < 1$, and hence that the Lasso support is contained within $S = \mathcal{N}(j)$, with probability at least $1 - c_1 e^{-c_2 n \delta^2 \alpha^2 - \log d}$, where the values of the universal constants may change from line to line. Taking the union bound over all d vertices, we conclude that $\widehat{E} \subseteq E$ with probability at least $1 - c_1 e^{-c_2 n \delta^2 \alpha^2}$.

Proof of (11.25a): $\mathbf{X}_{S^c}^T = \mathbf{\Gamma}_{S^c S}^* (\mathbf{\Gamma}_{SS}^*)^{-1} \mathbf{X}_S^T + \tilde{\mathbf{W}}_{S^c}^T$, where

$\tilde{\mathbf{W}}_{S^c} \in \mathbb{R}^{n \times |S^c|}$ is a zero-mean Gaussian random matrix that is independent of \mathbf{X}_S . Since

$\text{cov}(\tilde{\mathbf{W}}_{S^c}) = \mathbf{\Gamma}_{S^c S^c}^* - \mathbf{\Gamma}_{S^c S}^* (\mathbf{\Gamma}_{SS}^*)^{-1} \mathbf{\Gamma}_{SS^c}^* \leq \mathbf{\Gamma}^*$ and $\text{diag}(\mathbf{\Gamma}^*) \leq 1$, we see that the elements of $\tilde{\mathbf{W}}_{S^c}$ have variance at most 1.

$$\begin{aligned}
\|\mu\|_\infty &= \left\| \mathbf{\Gamma}_{S^c S}^* (\mathbf{\Gamma}_{SS}^*)^{-1} \widehat{\mathbf{z}}_S + \frac{\tilde{\mathbf{W}}_{S^c}^T \mathbf{X}_S}{\sqrt{n}} (\widehat{\mathbf{\Gamma}}_{SS})^{-1} \widehat{\mathbf{z}}_S \right\|_\infty \\
&\stackrel{(i)}{\leq} (1 - \alpha) + \underbrace{\left\| \frac{\tilde{\mathbf{W}}_{S^c}^T \mathbf{X}_S}{\sqrt{n}} (\widehat{\mathbf{\Gamma}}_{SS})^{-1} \widehat{\mathbf{z}}_S \right\|_\infty}_{\tilde{V} \in \mathbb{R}^{|S^c|}} \quad (11.27)
\end{aligned}$$

$$\begin{aligned}
\frac{1}{\sqrt{n}} \left\| \frac{\mathbf{X}_S}{\sqrt{n}} (\widehat{\mathbf{\Gamma}}_{SS})^{-1} \widehat{\mathbf{z}}_S \right\|_2 &\leq \frac{1}{\sqrt{n}} \left\| \frac{\mathbf{X}_S}{\sqrt{n}} (\widehat{\mathbf{\Gamma}}_{SS})^{-1} \right\|_2 \|\widehat{\mathbf{z}}_S\|_2 \\
&\leq \frac{1}{\sqrt{n}} \sqrt{\left\| (\widehat{\mathbf{\Gamma}}_{SS})^{-1} \right\|_2} \sqrt{m} \\
&\leq 2 \sqrt{\frac{bm}{n}},
\end{aligned}$$

where inequality (i) follows with probability at least $1 - 4e^{-c_1 n}$, using standard bounds on Gaussian random matrices (see Theorem 6.1). Using this upper bound to control the conditional variance of \widetilde{V} , standard Gaussian tail bounds and the union bound then ensure that

$$\mathbb{P} \left[\|\widetilde{V}\|_{\infty} \geq t \right] \leq 2 |S^c| e^{-\frac{nn^2}{8bm}} \leq 2e^{-\frac{nn^2}{8bm} + \log d}.$$

We now set $t = \left[\frac{64bm \log d}{n} + \frac{1}{64} \alpha^2 \right]^{1/2}$, a quantity which is less than $\frac{\alpha}{4}$ as long as $n \geq c \frac{bm \log d}{\alpha}$ for a sufficiently large universal constant. Thus, we have established that $\|\widetilde{V}\|_{\infty} \leq \frac{\alpha}{4}$ with probability at least $1 - c_1 e^{-c_2 n \alpha^2 - \log d}$. Combined with the earlier bound (11.27), the claim (11.25a) follows.

Proof of (11.25b): note that the matrix $\mathbf{\Pi} := \mathbf{I}_n - \mathbf{X}_S (\mathbf{X}_S^T \mathbf{X}_S)^{-1} \mathbf{X}_S^T$ has the range of \mathbf{X}_S as its nullspace. Thus, using the decomposition (11.26), we have

$$\mathbf{V}_{S^c} = \tilde{\mathbf{W}}_{S^c}^T \mathbf{\Pi} \left(\frac{\mathbf{W}_j}{\lambda_n n} \right),$$

where $\tilde{\mathbf{W}}_{S^c} \in \mathbb{R}^{|S^c|}$ is independent of $\mathbf{\Pi}$ and \mathbf{W}_j . Since $\mathbf{\Pi}$ is a projection matrix, we have $\|\mathbf{\Pi} \mathbf{W}_j\|_2 \leq \|\mathbf{W}_j\|_2$. The vector $\mathbf{W}_j \in \mathbb{R}^n$ has i.i.d. Gaussian entries with variance at most 1, and hence the event $\mathcal{E} = \left\{ \frac{\|\mathbf{W}_j\|_2}{\sqrt{n}} \leq 2 \right\}$ holds with probability at least $1 - 2e^{-n}$. Conditioning on this event and its complement, we find that

$$\mathbb{P} [\|\mathbf{V}_{S^c}\|_\infty \geq t] \leq \mathbb{P} [\|\mathbf{V}_{S^c}\|_\infty \geq t \mid \mathcal{E}] + 2e^{-c_3 n}.$$

Conditioned on \mathcal{E} , each element of V_{S^c} has variance at most $\frac{4}{\lambda_n^2 n}$, and hence

$$\mathbb{P}\left[\|V_{S^c}\|_\infty \geq \frac{\alpha}{4}\right] \leq 2e^{-\frac{\lambda_n^2 n^2}{256} + \log|S^c|} + 2e^{-n},$$

where we have combined the union bound with standard Gaussian tail bounds. Since $\lambda_n = c_0 \left\{ \frac{1}{\alpha} \sqrt{\frac{\log d}{n}} + \delta \right\}$ for a universal constant c_0 that may be chosen, we can ensure that $\frac{\lambda_n^2 n \alpha^2}{256} \geq c_2 n \delta^2 \alpha^2 + 2 \log d$ for some constant c_2 , for which it follows that

$$\mathbb{P}\left[\|V_{S^c}\|_\infty \geq \frac{\alpha}{4}\right] \leq c_1 e^{-c_2 n \delta^2 \alpha^2 - \log d} + 2e^{-n}.$$

Proof of part (b)

It suffices to establish ℓ_∞ -bounds on the error in the Lasso solution. Here we provide a proof in the case $m \leq \log d$. Again returning to the proof of Theorem 7.21, equation (7.54) guarantees that

$$\begin{aligned} \|\widehat{\theta}_S - \theta_S^*\|_\infty &\leq \left\| (\widehat{\Gamma}_{SS})^{-1} \mathbf{X}_S^T \frac{W_j}{n} \right\|_\infty + \lambda_n \left\| (\widehat{\Gamma}_{SS})^{-1} \right\|_\infty \\ &\leq \left\| (\widehat{\Gamma}_{SS})^{-1} \mathbf{X}_S^T \frac{W_j}{n} \right\|_\infty + \lambda_n \left\{ \left\| (\widehat{\Gamma}_{SS})^{-1} - (\Gamma_{SS}^*)^{-1} \right\|_\infty + \left\| (\Gamma_{SS}^*)^{-1} \right\|_\infty \right\} \end{aligned} \quad (11.28)$$

Now for any symmetric $m \times m$ matrix, we have

$$\|\mathbf{A}\|_\infty = \max_{i=1,\dots,m} \sum_{\ell=1}^m |A_{i\ell}| \leq \sqrt{m} \max_{i=1,\dots,m} \sqrt{\sum_{\ell=1}^m |A_{i\ell}|^2} \leq \sqrt{m} \|\mathbf{A}\|_2$$

Applying this bound to the matrix $\mathbf{A} = (\widehat{\boldsymbol{\Gamma}}_{SS})^{-1} - (\boldsymbol{\Gamma}_{SS}^*)^{-1}$, we find

$$\left\| (\widehat{\boldsymbol{\Gamma}}_{SS})^{-1} - (\boldsymbol{\Gamma}_{SS}^*)^{-1} \right\|_{\infty} \leq \sqrt{m} \left\| (\widehat{\boldsymbol{\Gamma}}_{SS})^{-1} - (\boldsymbol{\Gamma}_{SS}^*)^{-1} \right\|_2. \quad (11.29)$$

Since $\|\boldsymbol{\Gamma}_{SS}^*\|_2 \leq \|\boldsymbol{\Gamma}_{SS}^*\|_{\infty} \leq b$, applying the random matrix bound from Theorem 6.1 allows us to conclude that

$$\left\| (\widehat{\boldsymbol{\Gamma}}_{SS})^{-1} - (\boldsymbol{\Gamma}_{SS}^*)^{-1} \right\|_2 \leq 2b \left(\sqrt{\frac{m}{n}} + \frac{1}{\sqrt{m}} + 10 \sqrt{\frac{\log d}{n}} \right),$$

with probability at least $1 - c_1 e^{-c_2 \frac{n}{m} - \log d}$. Combined with the earlier bound (11.29), we find that

$$\left\| (\widehat{\boldsymbol{\Gamma}}_{SS})^{-1} - (\boldsymbol{\Gamma}_{SS}^*)^{-1} \right\|_{\infty} \leq 2b \left(\sqrt{\frac{m^2}{n}} + 1 + 10 \sqrt{\frac{m \log d}{n}} \right) \stackrel{(i)}{\leq} 6b, \quad (11.30)$$

where inequality (i) uses the assumed lower bound $n \gtrsim m \log d \geq m^2$. Putting together the pieces in the bound (11.28) leads to

$$\left\| \widehat{\theta}_S - \theta_S^* \right\|_{\infty} \leq \underbrace{\left\| \left(\widehat{\mathbf{\Gamma}}_{SS} \right)^{-1} \mathbf{X}_S^T \frac{W_j}{n} \right\|_{\infty}}_{U_S} + 7b\lambda_n. \quad (11.31)$$

Now the vector $W_j \in \mathbb{R}^n$ has i.i.d. Gaussian entries, each zero-mean with variance at most $\text{var}(X_j) \leq 1$, and is independent of \mathbf{X}_S . Consequently, conditioned on \mathbf{X}_S , the quantity U_S is a zero-mean Gaussian m -vector, with maximal variance

$$\frac{1}{n} \left\| \text{diag} \left(\widehat{\mathbf{\Gamma}}_{SS} \right)^{-1} \right\|_{\infty} \leq \frac{1}{n} \left\{ \left\| \left(\widehat{\mathbf{\Gamma}}_{SS} \right)^{-1} - \left(\mathbf{\Gamma}_{SS}^* \right)^{-1} \right\|_{\infty} + \left\| \left(\mathbf{\Gamma}_{SS}^* \right)^{-1} \right\|_{\infty} \right\} \leq \frac{7b}{n},$$

where we have combined the assumed bound $\left\| \left(\mathbf{\Gamma}_{SS}^* \right)^{-1} \right\|_{\infty} \leq b$ with the inequality (11.30).

Therefore, the union bound combined with Gaussian tail bounds implies that

$$\mathbb{P} [\|U_S\|_\infty \geq b\lambda_n] \leq 2|S|e^{-\frac{n\lambda_n^2}{14}} \stackrel{(i)}{\leq} c_1 e^{-c_2 nb\delta^2 - \log d},$$

where, as in our earlier argument, inequality (i) can be guaranteed by a sufficiently large choice of the pre-factor c_0 in the definition of λ_n . Substituting back into the earlier bound (11.31), we find that

$\|\widehat{\theta}_S - \theta_S^*\|_\infty \leq 7b\lambda_n$ with probability at least $1 - c_1 e^{-c_2 n\{\delta^2 \wedge \frac{1}{m}\} - \log d}$. Finally, taking the union bound over all vertices $j \in V$ causes a loss of at most a factor $\log d$ in the exponent.

- 1 Some basics
- 2 Estimation of Gaussian graphical models
- 3 Graphical models in exponential form**
- 4 Graphs with corrupted or hidden variables

Let us now move beyond the Gaussian case, and consider the graph estimation problem for a more general class of graphical models that can be written in an exponential form. In particular, for a given graph $G = (V, E)$, consider probability densities that have a pairwise factorization of the form

$$p_{\Theta^*}(x_1, \dots, x_d) \propto \exp \left\{ \sum_{j \in V} \phi_j(x_j; \Theta_j^*) + \sum_{(j,k) \in E} \phi_{jk}(x_j, x_k; \Theta_{jk}^*) \right\},$$

Example 11.4(Ising model)

$$p(x_1, \dots, x_d; \theta^*) \propto \exp \left\{ \sum_{j \in V} \theta_j^* x_j + \sum_{(j,k) \in E} \theta_{jk}^* x_j x_k \right\}$$

A general form of neighborhood regression

- We can form a matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ with x_i^T as the i th row. For $j = 1, \dots, d$, we let $X_j \in \mathbb{R}^n$ denote the j th column of \mathbf{X} . Neighborhood regression is based on predicting the column $X_j \in \mathbb{R}^n$ using the columns of the submatrix $\mathbf{X}_{\setminus(j)} \in \mathbb{R}^{n \times (d-1)}$.
- Consider the conditional likelihood of $X_j \in \mathbb{R}^n$ given $\mathbf{X}_{\setminus(j)} \in \mathbb{R}^{n \times (d-1)}$, this conditional likelihood depends only on the vector of parameters

$$\Theta_{j+} := \{\Theta_j, \Theta_{jk}, k \in V \setminus \{j\}\}$$

Moreover, in the true model Θ^* , we are guaranteed that $\Theta_{jk}^* = 0$ whenever $(j, k) \notin E$, so that it is natural to impose some type of block-based sparsity penalty on Θ_{j+} .

$$\widehat{\Theta}_{j+} = \arg \min_{\Theta_{j+}} \underbrace{\left\{ -\frac{1}{n} \sum_{i=1}^n \log p_{\Theta_{j+}}(x_{ij} \mid x_{i \setminus \{j\}}) \right\}}_{\mathcal{L}_n(\Theta_{j+}; x_j, x_{\setminus \{j\}})} + \lambda_n \sum_{k \in V \setminus \{j\}} \|\Theta_{jk}\}.$$

Especially, or the Ising model, the neighborhood regression estimate:

$$\widehat{\theta}_{j+} = \arg \min_{\theta_{j+} \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{i=1}^n f \left(\theta_j x_{ij} + \sum_{k \in V \setminus \{j\}} \theta_{jk} x_{ij} x_{ik} \right) + \lambda_n \sum_{k \in V \setminus \{j\}} |\theta_{jk}| \right\}$$

where $f(t) = \log(1 + e^t)$ is the logistic function.

Let θ_{j+}^* denote the minimizer of the population objective function $\bar{\mathcal{L}}(\theta_{j+}) = \mathbb{E}[\mathcal{L}_n(\theta_{j+}; X_j, \mathbf{X}_{\setminus\{j\}})]$. We then consider the Hessian of the cost function $\bar{\mathcal{L}}$ evaluated at the "true parameter" θ_{j+}^* -namely, the d -dimensional matrix $\mathbf{J} := \nabla^2 \bar{\mathcal{L}}(\theta_{j+}^*)$.

- For a given $\alpha \in (0, 1]$, we say that \mathbf{J} satisfies an α -incoherence condition at node $j \in V$ if

$$\max_{k \notin S} \left\| J_{kS} (\mathbf{J}_{SS})^{-1} \right\|_1 \leq 1 - \alpha,$$

- We assume the submatrix \mathbf{J}_{SS} has its smallest eigenvalue lower bounded by some $c_{\min} > 0$.
- The following result applies to an Ising model defined on a graph G with d vertices and maximum degree at most m .

Theorem 11.15 Given n i.i.d. samples with $n > c_0 m^2 \log d$, consider the estimator with $\lambda_n = \frac{32}{\alpha} \sqrt{\frac{\log d}{n}} + \delta$ for some $\delta \in [0, 1]$. Then with probability at least $1 - c_1 e^{-c_2(n\delta^2 + \log d)}$, the estimate $\widehat{\theta}_{j+}$ has the following properties:

- (a) It has a support $\widehat{S} = \text{supp}(\widehat{\theta})$ that is contained within the neighborhood set $\mathcal{N}(j)$.
- (b) It satisfies the ℓ_∞ -bound $\left\| \widehat{\theta}_{j+} - \theta_{j+}^* \right\|_\infty \leq \frac{c_3}{c_{\min}} \sqrt{m} \lambda_n$.

Part (a) guarantees that the method leads to no false inclusions. On the other hand, the ℓ_∞ -bound in part (b) ensures that the method picks up all significant variables.

- 1 Some basics
- 2 Estimation of Gaussian graphical models
- 3 Graphical models in exponential form
- 4 Graphs with corrupted or hidden variables

Let us begin our exploration with the case of corrupted data. Letting $\mathbf{X} \in \mathbb{R}^{n \times d}$ denote the data matrix corresponding to the original samples, suppose that we instead observe a corrupted version \mathbf{Z} . In the simplest case, we might observe $\mathbf{Z} = \mathbf{X} + \mathbf{V}$, where the matrix \mathbf{V} represents some type of measurement error.

$$\widehat{\Theta}_{\text{NAI}} = \arg \min_{\Theta \in \mathcal{S}^{d \times d}} \left\{ \langle \Theta, \widehat{\Sigma}_z \rangle - \log \det \Theta + \lambda_n \|\Theta\|_{1, \text{off}} \right\}$$

where $\widehat{\Sigma}_z = \frac{1}{n} \mathbf{Z}^T \mathbf{Z} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^T$.

More generally, any unbiased estimate $\widehat{\Gamma}$ of Σ_X :

$$\widetilde{\Theta} = \arg \min_{\Theta \in \mathcal{S}_+^{d \times d}} \left\{ \langle \Theta, \widehat{\Gamma} \rangle - \log \det \Theta + \lambda_n \|\Theta\|_{1, \text{off}} \right\}.$$

Unbiased covariance estimate

In order to obtain a consistent estimator, we need to replace $\widehat{\Sigma}_Z$ with an unbiased estimator of $\text{cov}(x)$ based on the observed data matrix \mathbf{Z} . In order to develop intuition, let us explore an example. Suppose that each row v_i of the noise matrix \mathbf{V} is drawn i.i.d. from a zero-mean distribution, say with covariance Σ_v .

$$\widehat{\Gamma} := \frac{1}{n} \mathbf{Z}^T \mathbf{Z} - \Sigma_v.$$

As long as the noise matrix \mathbf{V} is independent of \mathbf{X} , then $\widehat{\Gamma}$ is an unbiased estimate of Σ_x .

Missing data

Example 11.17 (Missing data) In other settings, some entries of the data matrix \mathbf{X} might be missing, with the remaining entries observed. In the simplest model of missing data known as missing completely at random-entry (i, j) of the data matrix is missing with some probability $v \in [0, 1]$.

$$\tilde{Z}_{ij} = \begin{cases} \frac{Z_{ij}}{1-v} & \text{if entry } (i, j) \text{ is observed,} \\ 0 & \text{otherwise.} \end{cases}$$

With this choice, it can be verified that

$$\hat{\mathbf{\Gamma}} = \frac{1}{n} \tilde{\mathbf{Z}}^T \tilde{\mathbf{Z}} - v \operatorname{diag} \left(\frac{\tilde{\mathbf{Z}}^T \tilde{\mathbf{Z}}}{n} \right)$$

linear regression

In order to obtain a consistent form of linear regression, consider the following population-level objective function

$$\overline{\mathcal{L}}(\theta) = \frac{1}{2} \theta^T \mathbf{\Gamma} \theta - \langle \theta, \gamma \rangle,$$

where $\mathbf{\Gamma} := \text{cov}(x)$ and $\gamma := \text{cov}(x, y)$. By construction, the true regression vector is the unique global minimizer of $\overline{\mathcal{L}}$. Thus

$$\mathcal{L}_n(\theta) = \frac{1}{2} \theta^T \widehat{\mathbf{\Gamma}} \theta - \langle \theta, \widehat{\gamma} \rangle.$$

As in our analysis of the ordinary Lasso from Chapter 7, we impose a restricted eigenvalue (RE) condition on the covariance estimate $\widehat{\mathbf{\Gamma}}$: more precisely, we assume that there exists a constant $\kappa > 0$ such that

$$\langle \Delta, \widehat{\mathbf{\Gamma}} \Delta \rangle \geq \kappa \|\Delta\|_2^2 - c_0 \frac{\log d}{n} \|\Delta\|_1^2 \quad \text{for all } \Delta \in \mathbb{R}^d.$$

Proposition 11.18 Under the RE condition, suppose that the pair $(\widehat{\gamma}, \widehat{\Gamma})$ satisfy the deviation condition

$$\left\| \widehat{\Gamma} \theta^* - \widehat{\gamma} \right\|_{\max} \leq \varphi(\mathbb{Q}, \sigma_w) \sqrt{\frac{\log d}{n}},$$

for a pre-factor $\varphi(\mathbb{Q}, \sigma_w)$ depending on the conditional distribution \mathbb{Q} and noise standard deviation σ_w . Then for any regularization parameter $\lambda_n \geq 2(2c_0 + \varphi(\mathbb{Q}, \sigma_w)) \sqrt{\frac{\log d}{n}}$, any local optimum $\widetilde{\theta}$ to the program (11.48) satisfies the bound

$$\left\| \widetilde{\theta} - \theta^* \right\|_2 \leq \frac{2}{\kappa} \sqrt{s} \lambda_n.$$

Proof:



$$\begin{aligned} \langle \widehat{\Delta}, \nabla \mathcal{L}_n(\theta^* + \widehat{\Delta}) - \nabla \mathcal{L}_n(\theta^*) \rangle &\leq \left| \langle \widehat{\Delta}, \nabla \mathcal{L}_n(\theta^*) \rangle \right| - \lambda_n \langle \widehat{z}, \widehat{\Delta} \rangle \\ &\leq \|\widehat{\Delta}\|_1 \|\nabla \mathcal{L}_n(\theta^*)\|_\infty + \lambda_n \{\|\theta^*\|_1 - \|\widetilde{\theta}\|_1\} \end{aligned}$$

• $\|\theta^*\|_1 - \|\widetilde{\theta}\|_1 \leq \left\| \widehat{\Delta}_S \right\|_1 - \left\| \widehat{\Delta}_S \right\|_1. (\text{Theorem 7.8})$

• Since θ^* is s -sparse, we have $\|\theta^*\|_1 \leq \sqrt{s} \|\theta^*\|_2 \leq \sqrt{\frac{n}{\log d}}$, where the final inequality follows from the assumption that $n \geq s \log d$. Consequently, we have

$$\|\widehat{\Delta}\|_1 \leq \|\widetilde{\theta}\|_1 + \|\theta^*\|_1 \leq 2 \sqrt{\frac{n}{\log d}}.$$

• $\langle \widehat{\Delta}, \widehat{\Gamma} \widehat{\Delta} \rangle \geq \kappa \|\widehat{\Delta}\|_2^2 - c_0 \frac{\log d}{n} \|\widehat{\Delta}\|_1^2 \geq \kappa \|\widehat{\Delta}\|_2^2 - 2c_0 \sqrt{\frac{\log d}{n}} \|\widehat{\Delta}\|_1.$

Gaussian graph selection with hidden variables

Consider a family of $d + r$ random variables-say written as $X := (X_1, \dots, X_d, X_{d+1}, \dots, X_{d+r})$ and suppose that this full vector can be modeled by a sparse graphical model with $d + r$ vertices. Now suppose that we observe only the subvector $X_O := (X_1, \dots, X_d)$, with the other components $X_H := (X_{d+1}, \dots, X_{d+r})$ staying hidden.

$$\Theta^\diamond = \begin{bmatrix} \Theta_{OO}^\diamond & \Theta_{OH}^\diamond \\ \Theta_{HO}^\diamond & \Theta_{HH}^\diamond \end{bmatrix}.$$

$$(\Sigma_{OO}^*)^{-1} = \underbrace{\Theta_{OO}^\diamond}_{\Gamma^*} - \underbrace{\Theta_{OH}^\diamond (\Theta_{HH}^\diamond)^{-1} \Theta_{HO}^\diamond}_{\Lambda^*}.$$

The Hammersley-Clifford theorem implies that the inverse covariance matrix Θ^\diamond of the full vector $X = (X_O, X_H)$ is sparse. By our modeling assumptions, the matrix $\Gamma^* := \Theta_{OO}^\diamond$ is sparse.

When $n > d$, then the sample covariance matrix $\widehat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$ will be invertible with high probability, and hence setting $\mathbf{Y} := (\widehat{\Sigma})^{-1}$, we can consider an observation model of the form

$$\mathbf{Y} = \mathbf{\Gamma}^* - \mathbf{\Lambda}^* + \mathbf{W}.$$

For a threshold $v_n > 0$ to be chosen, we define the estimates

$$\widehat{\mathbf{\Gamma}} := T_{v_n}((\widehat{\Sigma})^{-1}) \text{ and } \widehat{\mathbf{\Lambda}} := \widehat{\mathbf{\Gamma}} - (\widehat{\Sigma})^{-1}.$$

Here the hard-thresholding operator is given by

$$T_{v_n}(v) = vI[|v| > v_n].$$

- (A1) $\|\Lambda^*\|_{\max} \leq \frac{\alpha}{d}$
- (A2) $\left\| \sqrt{\Theta^*} \right\|_{\infty} = \max_{j=1,\dots,d} \sum_{k=1}^d \left| \sqrt{\Theta^*} \right|_{jk} \leq \sqrt{M}$
- (A3) $v_n := M \left(4 \sqrt{\frac{\log d}{n}} + \delta \right) + \frac{\alpha}{d}$ for some $\delta \in [0, 1]$

Proposition 11.19 Consider a precision matrix Θ^* that can be decomposed as the difference $\Gamma^* - \Lambda^*$. Given $n > d$ i.i.d. samples from the $\mathcal{N}(0, (\Theta^*)^{-1})$ distribution and any $\delta \in (0, 1]$

$$\left\| \widehat{\Gamma} - \Gamma^* \right\|_{\max} \leq 2M \left(4 \sqrt{\frac{\log d}{n}} + \delta \right) + \frac{2\alpha}{d}$$

and

$$\left\| \widehat{\Lambda} - \Lambda^* \right\|_2 \leq M \left(2 \sqrt{\frac{d}{n}} + \delta \right) + s \left\| \widehat{\Gamma} - \Gamma^* \right\|_{\max}$$

with probability at least $1 - c_1 e^{-c_2 n \delta^2}$.

Proof

- $$(\widehat{\Sigma})^{-1} - \Theta^* = \sqrt{\Theta^*} \{n^{-1} \mathbf{V}^T \mathbf{V} - \mathbf{I}_d\} \sqrt{\Theta^*}$$

$$\|(\widehat{\Sigma})^{-1} - \Theta^*\|_{\max} = \max_{j,k=1,\dots,d} |e_j^T \sqrt{\Theta^*} \widetilde{\Sigma} \sqrt{\Theta^*} e_k|$$
- $$\leq \max_{j,k=1,\dots,d} \|\sqrt{\Theta^*} e_j\|_1 \|\widetilde{\Sigma} \sqrt{\Theta^*} e_k\|_{\infty}$$

$$\leq \|\widetilde{\Sigma}\|_{\max} \max_{j=1,\dots,d} \|\sqrt{\Theta^*} e_j\|_1^2.$$
- $$\|\widehat{\Gamma} - \Gamma^*\|_{\max} \leq \|\mathbf{Y} - \Theta^*\|_{\max} + \|\mathbf{Y} - T_{v_n}(\mathbf{Y})\|_{\max} + \|\Lambda^*\|_{\max}$$

$$\leq M \left(4 \sqrt{\frac{\log d}{n}} + \delta \right) + v_n + \frac{\alpha}{d}$$

$$\leq 2M \left(4 \sqrt{\frac{\log d}{n}} + \delta \right) + \frac{2\alpha}{d}$$