

# A Constrained $\ell_1$ Minimization Approach to Sparse Precision Matrix Estimation

Presented by Xingguo Li

CSCI 8980, UMN

Authors : Tony Cai, Weidong Liu, and Xi Luo

April 14, 2014

# Outline

- CLIME
- Statistical Properties
  - Convergence Rates under Norms, *i.e.*  $\|\cdot\|_2$ ,  $\|\cdot\|_F^2$  and  $|\cdot|_\infty$
  - Convergence Rates of Expectation & Ordered Variables
  - Model Selection Consistency
  - Comparison with  $\ell_1$ -MLE
- Numerical Experiments

# Motivation

- $\ell_1$  regularized log-determinant (Banerjee *et al.*, 2008):

$$\hat{\Omega}_{\text{Glasso}} := \underset{\Omega \succ 0}{\operatorname{argmin}} \{ \langle \Omega, \Sigma_n \rangle - \log \det(\Omega) + \lambda_n \|\Omega\|_1 \} \quad (1)$$

Optimality condition:

$$\hat{\Omega}_{\text{Glasso}}^{-1} - \Sigma_n = \lambda_n \hat{\mathbf{Z}}, \quad \hat{\mathbf{Z}} \in \partial \|\hat{\Omega}_{\text{Glasso}}\|_1$$

Dantzig type problem:

$$\min \|\Omega\|_1 \text{ s.t. } |\Omega^{-1} - \Sigma_n|_\infty \leq \lambda_n, \Omega \in \mathbb{R}^{p \times p}$$

- CLIME

$$\min \|\Omega\|_1 \text{ s.t. } |\mathbf{I} - \Sigma_n \Omega|_\infty \leq \lambda_n, \Omega \in \mathbb{R}^{p \times p}, \quad (2)$$

# Properties

Compared with  $\ell_1$ -MLE (1),

- No requirement of positive definiteness of  $\Omega$
- Columnwise decomposibility: For all  $i = 1, \dots, p$ ,

$$\min \|\beta\|_1 \text{ s.t. } |\mathbf{e}_i - \Sigma_n \beta|_\infty \leq \lambda_n, \beta \in \mathbb{R}^p. \quad (3)$$

## Lemma 1

Let  $\{\widehat{\Omega}_1\}$  be the solution set of (2), and let  $\{\widehat{\mathbf{B}}\} := \{(\widehat{\beta}_1, \dots, \widehat{\beta}_p)\}$ , where  $\widehat{\beta}_i$  are solutions to (3) for  $i = 1, \dots, p$ . Then  $\{\widehat{\Omega}_1\} = \{\widehat{\mathbf{B}}\}$ .

- Improved convergence rate (polynomial-type tails)
- Improved model selection consistency (polynomial-type tails)

# Proof of Lemma 1

Remind:

$$\widehat{\Omega}_1 = (\widehat{\omega}_1^1, \dots, \widehat{\omega}_p^1) = \operatorname{argmin}_{\Omega \in \mathbb{R}^{p \times p}} \|\Omega\|_1 \text{ s.t. } |\mathbf{I} - \Sigma_n \Omega|_\infty \leq \lambda_n$$

$$\widehat{\mathbf{B}} = (\widehat{\beta}_1, \dots, \widehat{\beta}_p), \widehat{\beta}_i = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|\beta\|_1 \text{ s.t. } |\mathbf{e}_i - \Sigma_n \beta|_\infty \leq \lambda_n, \forall i$$

(1) We have

$$|\widehat{\omega}_i^1|_1 \geq |\widehat{\beta}_i|_1, \forall 1 \leq i \leq p \quad (4)$$

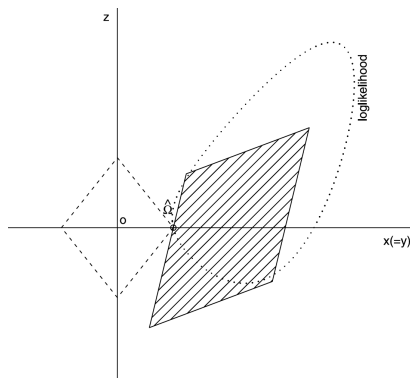
$$\|\widehat{\Omega}_1\|_1 \leq \|\widehat{\mathbf{B}}\|_1, \quad (5)$$

$$\Rightarrow \widehat{\mathbf{B}} \in \{\widehat{\Omega}_1\}$$

(2) If  $\widehat{\Omega}_1 \notin \{\widehat{\mathbf{B}}\}$ ,  $\exists i$  s.t.  $|\widehat{\omega}_i^1|_1 > |\widehat{\beta}_i|_1$ , then by (4)  $\|\widehat{\Omega}_1\|_1 > \|\widehat{\mathbf{B}}\|_1$   
 $\Rightarrow \Leftarrow$  (5)

Therefore,  $\{\widehat{\mathbf{B}}\} = \{\widehat{\Omega}_1\}$ . □

# Properties



**Figure 1 :** Plot of the elementwise  $l_\infty$  constrained feasible set (shaded polygon) and the elementwise  $l_1$  norm objective (dashed diamond near the origin) from CLIME. The log-likelihood function as in Glasso is represented by the dotted line. (Cai *et al.*, 2011)

# Parameter Class

$\Omega_0 \in \mathcal{U}$ : Uniformity class of matrices,

$$\mathcal{U} := \mathcal{U}(q, s_0(p))$$

$$= \left\{ \Omega : \Omega \succ 0, \|\Omega\|_{L_1} \leq M, \max_{1 \leq i \leq p} \sum_{j=1}^p |w_{ij}|^q \leq s_0(p), 0 \leq q < 1 \right\}$$

- $q = 0$ ,  $\mathcal{U}(0, s_0(p))$  is a class of  $s_0(p)$ -sparse matrices
- Wider class of precision matrix than truly sparse matrices, *i.e.*  $s_0(p)$  is small when many entries are small.

# Tail Class

Two types of tails:

(C1) Exponential-type tails:  $\exists$  some constant  $0 < \eta < 1/4$  such that  $\log p/n \leq \eta$  and for bounded constant  $K$

$$\mathbb{E}e^{t(X_i - \mu_i)^2} \leq K < \infty \text{ for all } |t| \leq \eta, \text{ for all } i$$

(C2) Polynomial-type tails: For some  $\gamma, c_1, \delta > 0$  and  $p \leq c_1 n^\gamma$ ,

$$\mathbb{E}|X_i - \mu_i|^{4\gamma+4+\delta} \leq K \text{ for all } i$$

- Bounded  $\theta := \max_{ij} \theta_{ij} = \max_{i,j} \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j) - \sigma_{ij}^0]^2$



# Convergence Rate under Norms

- Symmetrizing operation:  $\widehat{\Omega} = (\widehat{\omega}_{ij})$ , where  $\widehat{\Omega}_1 = (\widehat{\omega}_{ij}^1)$  and  $\widehat{\omega}_{ij} = \widehat{\omega}_{ji} = \widehat{\omega}_{ij}^1 I\{|\widehat{\omega}_{ij}^1| \leq |\widehat{\omega}_{ji}^1|\} + \widehat{\omega}_{ji}^1 I\{|\widehat{\omega}_{ij}^1| > |\widehat{\omega}_{ji}^1|\}$

## Theorem 1

Suppose  $\Omega_0 \in \mathcal{U}(q, s_0(p))$ . Assume (C1) or (C2) holds. Let  $\lambda_n = C_{0i} M \sqrt{\log p/n}$  and  $\tau > 0$ , then w.h.p.

$$\|\widehat{\Omega} - \Omega_0\|_2 \leq C_{1i} M^{2-2q} s_0(p) (\log p/n)^{(1-q)/2} \quad (6)$$

$i = 1, 2$  for (C1) and (C2) respectively.

- $\ell_1$ -MLE estimator for polynomial-type of tails when  $q = 0$ :

$$\|\widehat{\Omega} - \Omega_0\|_2 = \mathcal{O}(s_0(p) \sqrt{\frac{p^{\tau/(\gamma+1+\delta/4)}}{n}})$$

- Other norms:  $\frac{1}{p} \|\widehat{\Omega} - \Omega_0\|_F^2 = \mathcal{O}(s_0(p) (\frac{\log p}{n})^{1-q/2})$

$$\|\widehat{\Omega} - \Omega_0\|_\infty = \mathcal{O}(\sqrt{\frac{\log p}{n}})$$

# Convergence Rate of $\sup_{\Omega_0 \in \mathcal{U}} \mathbb{E} \|\widehat{\Omega} - \Omega_0\|_2^2$

- Replace  $\Sigma_n$  with  $\Sigma_{n,\rho} = \Sigma_n + \rho \mathbf{I}$  to (1) ensure the existence of  $\mathbb{E} \|\widehat{\Omega}_\rho - \Omega_0\|_2^2$  and (2) get a feasible initial value of  $\widehat{\Omega}$ .

## Theorem 2

Suppose  $\Omega_0 \in \mathcal{U}(q, s_0(p))$  and (C1) holds. Let  $\rho = \sqrt{\log p/n}$ ,  $\lambda_n = C_0 M \sqrt{\log p/n}$  and  $\tau$  sufficiently large. If  $p = n^\xi$  for some  $\xi > 0$ , then

$$\sup_{\Omega_0 \in \mathcal{U}} \mathbb{E} \|\widehat{\Omega}_\rho - \Omega_0\|_2^2 = \mathcal{O} \left( M^{4-4q} s_0(p)^2 \left( \frac{\log p}{n} \right)^{1-q} \right). \quad (7)$$

- Hold for  $\min \left( \sqrt{\frac{\log p}{n}}, p^{-\alpha} \right) \leq \rho \leq \sqrt{\frac{\log p}{n}}$  with any  $\alpha > 0$ .
- Same order of rate for  $\|\cdot\|_\infty^2$  and  $\|\cdot\|_F^2$  with the rates under norms.

# Convergence Rate of Ordered Variables

- $\mathcal{U}_o(\alpha, B) = \{\Omega : \max_j \sum_i \{|\omega_{ij}| : |i - j| \geq k\} \leq B(k + 1)^{-\alpha}, \Omega \succ 0, \forall k \geq 0\}$  for some  $\alpha > 0$
- Better rates can be obtained

## Theorem 3

Let  $\Omega_0 \in \mathcal{U}_o(\alpha, B)$  and  $\lambda_n = CB\sqrt{\log p/n}$  with sufficiently large  $C$ .

(a) If (C1) or (C2) holds, then w.h.p.,

$$\|\widehat{\Omega} - \Omega_0\|_2 = \mathcal{O}\left(B^2(\log p/n)^{\alpha/(2\alpha+2)}\right) \quad (8)$$

(b) Suppose  $p \geq n^\xi, \xi > 0$ . If (C1) holds and  $\rho = \sqrt{\log p/n}$ , then

$$\sup_{\Omega_0 \in \mathcal{U}_o(\alpha, B)} \mathbb{E}\|\widehat{\Omega}_\rho - \Omega_0\|_2^2 = \mathcal{O}\left(B^4(\log p/n)^{\alpha/(\alpha+1)}\right) \quad (9)$$

- $s_0(p)$  term disappears from the bounds.

# A General Result

## Theorem 6

Let  $\Omega_0 \in \mathcal{U}(q, s_0(p))$  and  $\rho > 0$ . If  $\lambda_n \geq \|\Omega_0\|_{L_1} (\max_{ij} |\hat{\sigma}_{ij} - \sigma_{ij}^0| + \rho)$ , then

$$|\hat{\Omega}_\rho - \Omega_0|_\infty = \mathcal{O}(\|\Omega_0\|_{L_1} \lambda_n), \quad (10)$$

$$\|\hat{\Omega}_\rho - \Omega_0\|_2 = \mathcal{O}(\|\Omega_0\|_{L_1}^{1-q} s_0(p) \lambda_n^{1-q}), \quad (11)$$

$$\frac{1}{p} \|\hat{\Omega}_\rho - \Omega_0\|_F^2 = \mathcal{O}(\|\Omega_0\|_{L_1}^{2-q} s_0(p) \lambda_n^{2-q}). \quad (12)$$

- Need to show  $\max_{ij} |\hat{\sigma}_{ij} - \sigma_{ij}^0| = \mathcal{O}(\sqrt{\log p/n})$  w.h.p. with corresponding constant for each result.

# Graphical Model Selection Consistency

- Threshold  $\tilde{\Omega} = (\tilde{\omega}_{ij})$  with  $\tilde{\omega}_{ij} = \hat{\omega}_{ij} I\{|\hat{\omega}_{ij}| \geq \tau_n\}$ , for  $\tau_n \geq 4M\lambda_n$
- Define:  $\mathcal{M}(\Omega) = \{\text{sign}(\omega_{ij}), 1 \leq i, j \leq p\}$ ,  
 $S(\Omega) = \{(i, j) : \omega_{ij} \neq 0\}$ ,  $\theta_{\min} = \min_{(i,j) \in S(\Omega_0)} |\omega_{ij}^0|$

## Theorem 7

Suppose (C1) or (C2) holds and  $\Omega_0 \in \mathcal{U}(0, s_0(p))$ . If  $\theta_{\min} > 2\tau_n$ , then  $\mathcal{M}(\tilde{\Omega}) = \mathcal{M}(\Omega_0)$  w.h.p.

- Sign consistency: Recover both sparsity pattern and signs of nonzero elements
- $\theta_{\min} > 2\tau_n$ : Ensure nonzero elements are correctly retained
- If  $M \ll n, p$ , then  $\tau_n = \mathcal{O}(\sqrt{\log p/n})$

# Comparison CLIME with $\ell_1$ -MLE (Ravikumar *et al.*, 2008)

- CLIME:  $\min \|\Omega\|_1$  s.t.  $|\mathbf{I} - \Sigma_n \Omega|_\infty \leq \lambda_n, \Omega \in \mathbb{R}^{p \times p}$
- $\ell_1$ -MLE:  $\min_{\Theta_{\gamma>0}} \{ \langle \Omega, \Sigma_n \rangle - \log \det(\Omega) + \lambda_n \|\Omega\|_{1,\text{off}} \}$

	CLIME	$\ell_1$ -MLE
Irrepresentability <sup>1</sup>	No	Yes
$(n, p)$ Scale	$\log p = o(n)$	$n > Cs_0^2(p) \log p$
Sparsity	Allow small values	Only truly sparse
Conv. Rate (poly)	$\mathcal{O} \left( s_0(p) \sqrt{\log p/n} \right)$	$\mathcal{O} \left( s_0(p) \sqrt{p^{\tau/(\gamma+1+\delta/4)}/n} \right)$
Model Selection	$\theta_{\min} \geq C \sqrt{\log p/n}$	$\theta_{\min} \geq C \sqrt{p^{\tau/\gamma+1+\delta/4}/n}$

<sup>1</sup>  $\|\Gamma_{SS}(\Gamma_{SS})^{-1}\|_{L_1} \leq 1 - \alpha, \alpha \in (0, 1], \Gamma = \Omega_0^{-1} \otimes \Omega_0^{-1}$

# Numerical Experiments

- CLIME:  $\min \|\beta\|_1$  s.t.  $|\mathbf{e}_i - \Sigma_n \beta|_\infty \leq \lambda_n, \beta \in \mathbb{R}^p, i = 1, \dots, p$
- LP:  $\min \sum_{j=1}^p u_j$  s.t.  $\forall 1 \leq j \leq p, \forall 1 \leq k \leq p$   
 $-\beta_j \leq u_j, -\hat{\sigma}_k^T + I\{k=i\} \leq \lambda_n$   
 $+\beta_j \leq u_j, +\hat{\sigma}_k^T - I\{k=i\} \leq \lambda_n$

- Refit (correct bias): Let  $\hat{S} = S(\tilde{\Omega}), \hat{S}^c = \{\omega_{ij}, (i, j) \in \hat{S}^c\},$   
 $\check{\Omega} = \operatorname{argmin}_{\Omega_{\hat{S}^c}=0} \langle \Omega, \Sigma \rangle - \log \det(\Omega)$

- Compare with Glasso and SCAD (Fan *et al.*, 2001)

$$\hat{\Omega}_{\text{Glasso}} := \operatorname{argmin}_{\Theta_{>0}} \langle \Omega, \Sigma_n \rangle - \log \det(\Omega) + \lambda_n \|\Omega\|_1$$

$$\hat{\Omega}_{\text{SCAD}} := \operatorname{argmin}_{\Theta_{>0}} \langle \Omega, \Sigma_n \rangle - \log \det(\Omega) + \sum_{i=1}^p \sum_{j=1}^p P_{\lambda_n, a}^{\text{SCAD}}(\omega_{ij})$$

$$\text{where } P_{\lambda, a}^{\text{SCAD}}(x) = \begin{cases} \lambda|x| & \text{if } |x| \leq \lambda; \\ -\left(\frac{|x|^2 - 2a\lambda|x| + \lambda^2}{2(a-1)}\right) & \text{if } \lambda \leq |x| \leq a\lambda; \\ \frac{(a+1)\lambda^2}{2} & \text{if } |x| \geq a\lambda \end{cases}$$

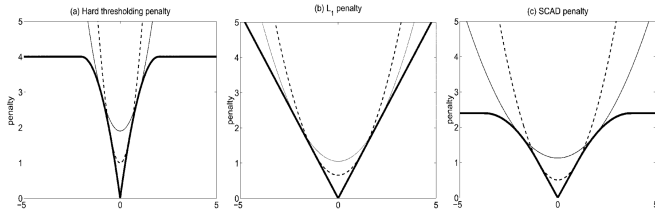


Figure 2 : Penalty functions and their quadratic approximations. (Fan *et al.*, 2001)

- Model 1.  $\omega_{ij}^0 = 0.6^{|i-j|}$

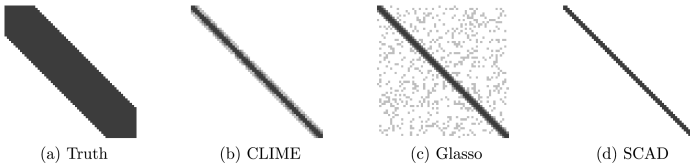


Figure 3 : Heatmaps of the frequency of the zeros identified for each entry of the precision matrix (when  $p = 60$ ) out of 100 replications. (Cai *et al.*, 2011)



# Breast Cancer Dataset (Hass *et al.*, 2006)

Classification performance criterion:

- Specificity:  $\frac{TN}{TN+FP}$
- Sensitivity:  $\frac{TP}{TP+FN}$
- MCC:  $\frac{TP \times TN - FP \times FN}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$

Method	Specificity	Sensitivity	MCC	Nonzero entries in $\hat{\Omega}$
Glasso	0.768 (0.009)	0.630 (0.021)	0.366 (0.018)	3923 (2)
Adaptive lasso	0.787 (0.009)	0.622 (0.022)	0.381 (0.018)	1233 (1)
SCAD	0.794 (0.009)	0.634 (0.022)	0.402 (0.020)	674 (1)
CLIME	0.749 (0.005)	0.806 (0.017)	0.506 (0.020)	492 (7)

**Figure 4 :** Comparison of classification performance. Glasso, Adaptive lasso, and SCAD results are taken from Fan *et al.*, 2009.

## Proof of Theorem 6

Let  $\rho = 0$ . Same proof for  $\rho > 0$ .

Assumption:  $\lambda_n \geq \|\Omega_0\|_{L_1} (\max_{ij} |\hat{\sigma}_{ij} - \sigma_{ij}^0|) \Leftrightarrow |\Sigma_0 - \Sigma_n|_\infty \leq \lambda_n / \|\Omega_0\|_{L_1}$

$$(1) \frac{|\hat{\Omega} - \Omega_0|_\infty}{\uparrow} \leq 4 \|\Omega_0\|_{L_1} \lambda_n$$

$$|\hat{\Omega} - \Omega_0|_\infty \leq |\hat{\Omega}_1 - \Omega_0|_\infty \leq \|\Omega_0\|_{L_1} \underbrace{|\Sigma_0(\hat{\Omega}_1 - \Omega_0)|_\infty}_{(i)} \leq 4 \|\Omega_0\|_{L_1} \lambda_n$$

$$(*) \|\mathbf{AB}\|_\infty \leq \|\mathbf{A}\|_{L_1} \|\mathbf{B}\|_\infty$$

$$(i) \leq \underbrace{|\Sigma_n(\hat{\Omega}_1 - \Omega_0)|_\infty}_{(ii)} + \underbrace{|(\Sigma_n - \Sigma_0)(\hat{\Omega}_1 - \Omega_0)|_\infty}_{(iii)} \leq 4\lambda_n$$

$$(ii) \leq |\Sigma_n \hat{\Omega}_1 - \mathbf{I}|_\infty + |\mathbf{I} - \Sigma_n \Omega_0|_\infty \leq \lambda_n + \|\Omega_0\|_{L_1} |\Sigma_0 - \Sigma_n|_\infty \leq 2\lambda_n$$

$$(iii) \leq \|\hat{\Omega}_1 - \Omega_0\|_{L_1} |\Sigma_n - \Sigma_0|_\infty$$

$$\leq \|\hat{\Omega}_1\|_{L_1} |\Sigma_n - \Sigma_0|_\infty + \|\Omega_0\|_{L_1} |\Sigma_n - \Sigma_0|_\infty \leq 2\lambda_n$$

## Proof of Theorem 6

$$(2) \frac{\|\widehat{\Omega} - \Omega_0\|_2}{\uparrow} \leq C_4 s_0(p) (4\|\Omega_0\|_{L_1} \lambda_n)^{1-q}, \quad C_4 \leq 2(1 + 2^{1-q} + 3^{1-q})$$

$$\|\widehat{\Omega} - \Omega_0\|_2 \leq \|\widehat{\Omega} - \Omega_0\|_{L_1} = \max_j |\widehat{\omega}_j - \omega_j^0|_1$$

$$(*) \|\mathbf{A}\|_2 \leq \sqrt{\|\mathbf{A}\|_{L_1} \|\mathbf{A}\|_\infty}, \quad \mathbf{A} = \mathbf{A}^T \Rightarrow \|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_{L_1}$$

$$\text{Let } \mathbf{h}_j = \widehat{\omega}_j - \omega_j^0, \quad \mathbf{h}_j^1 = (\widehat{\omega}_{ij} \mathbf{1}\{|\widehat{\omega}_{ij}| \geq 2t_n\}; 1 \leq i \leq p)^T - \omega_j^0,$$

$$\mathbf{h}_j^2 = \mathbf{h}_j - \mathbf{h}_j^1, \quad t_n = |\widehat{\Omega} - \Omega_0|_\infty \leq 4\|\Omega_0\|_{L_1} \lambda_n$$

$$\Rightarrow |\omega_j^0|_1 - |\mathbf{h}_j^1|_1 + |\mathbf{h}_j^2|_1 \leq |\omega_j^0 + \mathbf{h}_j^1|_1 + |\mathbf{h}_j^2|_1 = |\widehat{\omega}_j|_1 \leq |\widehat{\omega}_j^1|_1 \leq |\omega_j^0|_1$$

$$\Rightarrow |\mathbf{h}_j^2|_1 \leq |\mathbf{h}_j^1|_1 \Rightarrow |\mathbf{h}_j|_1 \leq |\mathbf{h}_j^1|_1 + |\mathbf{h}_j^2|_1 \leq 2|\mathbf{h}_j^1|_1$$

$$\dots \Rightarrow |\mathbf{h}_j|_1 \leq 2|\mathbf{h}_j^1|_1 \leq 2(1 + 2^{1-q} + 3^{1-q}) t_n^{1-q} s_0(p)$$

$$(3) \frac{1}{p} \|\widehat{\Omega} - \Omega_0\|_F^2 \leq C_5 s_0(p) (4\|\Omega_0\|_{L_1} \lambda_n)^{2-q}, \quad C_5 \leq C_4$$

$$(*) \|\mathbf{A}\|_F^2 \leq p \|\mathbf{A}\|_{L_1} \|\mathbf{A}\|_\infty$$

# Proof of Other Theorems

Based on Theorem 6, bound  $\max_{ij} |\hat{\sigma}_{ij} - \sigma_{ij}^0|$  w.h.p.

- Theorem 1 (a) and 4 (a), *i.e.* exponential-type tails,

$$\max_{ij} |\hat{\sigma}_{ij} - \sigma_{ij}^0| \leq 2\eta^{-2}(2 + \tau + \eta^{-1}e^2K^2)^2 \sqrt{\log p/n},$$

w.p.  $\geq 1 - 4p^{-\tau}$ .

- Theorem 1 (b) and 4 (b), *i.e.* polynomial-type tails,

$$\max_{ij} |\hat{\sigma}_{ij} - \sigma_{ij}^0| \leq \sqrt{(\theta + 1)(5 + \tau) \log p/n},$$

w.p.  $\geq 1 - \mathcal{O}(n^{-\delta/8} + p^{-\tau/2})$ .

- Theorem 2,3,5 are direct results of Theorem 6,1,4.

# Thank you!