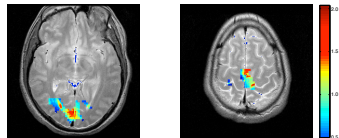


# Estimation and Inference in Large Graphical Models

Cun-Hui Zhang, Rutgers University

April 25, 2014

20th Applied Probability Day in Honor/Memory of Larry Shepp  
Columbia University

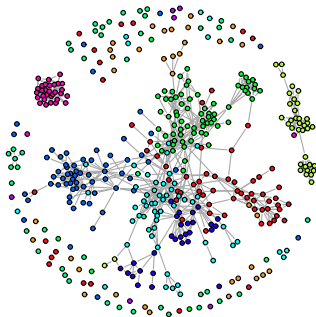


## Pick up sticks

*Larry Shepp, Doron Zeilberger, Cun-Hui Zhang*






October 23, 2012

## Large graphical models:

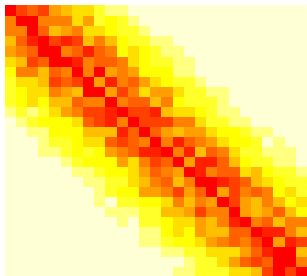
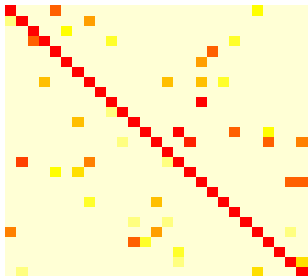


- Variables:  $X_1, \dots, X_p$
- Vertices represent variables
- Edges represent association
- Data:  $\mathbf{x}_1, \dots, \mathbf{x}_p$  in  $\mathbb{R}^n$
- Copulas:  $f_j(\mathbf{x}_j)$ , unknown  $f_j$
- $p \gg n$

## Acknowledgements:

-  Sun, T. and C. H. Zhang (2013). Sparse matrix inversion with scaled lasso. *Journal of Machine Learning Research* 14, 3385–3418.
-  Sun, T. and Zhang, C.-H. (2012). Comments on: Optimal rates of convergence for sparse covariance matrix estimation. *Statistica Sinica* 22, 1354–1358.
-  Ren, Z., Sun, T., Zhang, C.-H. and Zhou, H.H. (2013). Asymptotic normality and optimalities in estimation of large Gaussian graphical model, preprint.
-  Liu, H., Han, F. and Zhang, C.-H. (2012). Transelliptical graphical modeling under a hierarchical latent variable framework.
-  Mitra, R. and Zhang, C.-H. (2014). Multivariate analysis of nonparametric estimates of large correlation matrices. arXiv:1403.6195

## Examples of sparsity in $p \gg n$ models



## Measurements of association:

- Covariance/correlation
  - Parameters:  $\underline{\Sigma}^{jk} = \text{Cov}(X_j, X_k)$ ,  $\underline{\Sigma} = (\underline{\Sigma}^{jk})^{d \times d}$
  - Raw estimator:  $\underline{\Sigma} = \mathbf{x}_\top^j \mathbf{x}_k / n$  or  $\text{rank}(\underline{\Sigma}) \leq n \ll d$
  - Estimation of  $\underline{\Sigma}$  when  $\underline{\Sigma}$  is sparse: thresholding/tapering  $\underline{\Sigma}$
  - Matrix denoising problem
- Partial correlation (Markov graphical models)
  - Corr( $X_j, X_k | X_{\ell} \notin \{j, k\}$ ) or  $\underline{\Theta}^{jk} / \sqrt{\underline{\Theta}^{jj} \underline{\Theta}^{kk}}$ ,  $\underline{\Theta} = \underline{\Sigma}^{-1}$
  - Estimation of  $\underline{\Theta}$  when  $\underline{\Theta}$  is sparse
  - Inversion of an ill-conditioned, noisy matrix  $\underline{\Sigma}$
  - Confidence intervals/hypothesis testing for  $\underline{\Theta}^{jk}$  and  $\underline{\Theta}^{jk} / \sqrt{\underline{\Theta}^{jj} \underline{\Theta}^{kk}}$
- Gaussian and elliptical copulas
  - $Y_j = f_j(X_j)$  with unknown  $f_j$ ,  $j \in [d]$
  - $(X_1, \dots, X_d) \sim N(0, \underline{\Sigma})$  or  $\underline{\Sigma}^{1/2} \mathbf{u}$  with  $\xi \perp \mathbf{u} \& \mathbf{u} \sim \text{unif}(\mathbb{S}^{d-1})$
  - Data:  $\mathbf{y}_j = f_j(\mathbf{x}_j)$ , ranks are sufficient
  - Raw estimator: Spearman's rank correlation
  - Raw estimator: Kendall's tau, estimation of  $\text{sgn}(Y_j - Y_k) - \text{sgn}(Y_j - Y_k)$
  - Error bounds for matrices of nonparametric estimates

## Estimation of $\Theta = \Sigma^{-1}$

- Graphical Lasso/SCAD (Yuan-Lin 07, Friedman et al 08, Ravikuma et al 08, Rothman et al 08, Lam-Fan 09)

- Minimize  $\langle \Theta, \underline{\Sigma} \rangle + \log \det(\Theta) + \lambda \sum_{j \neq k} \text{pen}(|\Theta_{jk}|)$
- $\|\Theta\|_2^F \lesssim (\#\{j, k : \Theta_{jk} \neq 0\}) (\log p)/n$ ;  $\|\Theta\|_S^2$

CLIME (Cai et al 10) or Dantzig Selector-by-row (Yuan 10)

- Minimize  $\|\Theta_{j^*}\|_1$  subject to  $\|\underline{\Sigma} \Theta - I^{p \times p}\|_{\max} \leq \lambda$ ,  $j \in [p]$
- $\|\Theta\|_S \lesssim d^* \max_j \|\Theta_{j^*}\|_1 \sqrt{(\log p)/n}$ ;  $\Theta_{j^*}$ ?
- $d^* = \max_j \#\{k : \Theta_{jk} \neq 0\} = \text{max degree}$

- Lasso-by-row (Meinshausen-Bühlmann 06; Sun-Z 13)

- Minimize  $(\Theta \underline{\Sigma} \Theta)_{jj} \Theta_{jj}^{-3/2} + \Theta_{jj}^{-1/2} + 4 \sqrt{(\log p)/n} \|\Theta_{j \setminus j}^c\|_1$ ,  $j \in [p]$
- $\|\Theta\|_S \lesssim d^* \sqrt{(\log p)/n}$ ;  $\Theta_{j^*}$ ?

Analysis

- $\|\underline{\Sigma}\|_S + \|\Theta\|_S = O(1)$
- Large deviation of  $\max_{j,k} |\underline{\Sigma}_{jk} - \Sigma_{jk}|$
- Large deviation of  $\max_{|A^c| \leq d^*} \|\underline{\Sigma} - \Sigma\|_{A^c \times A^c}$

## Confidence intervals and hypothesis testing for $\Theta_{jk}$

- Confidence intervals/p-values after model selection
- In the presence of small signals, it is impossible to construct approximate confidence intervals or  $p$  value for  $\Theta_{jk}$  based on a selected model (LeCam 72; Leeb-Pötscher 06)
- Conservative methods (Berk et al 09, 11; Laber-Murphy 11)
- Semiparametric approach (Z 11)
- Bias correction and Low-dimensional projection methods (Z-Zhang 11), desparsification (van de Geer et al 13)
- Thresholded-LDPE with multiplicity adjustment (Z-Zhang 11), resparsification (van de Geer et al 13)
- Estimation of variance in linear regression (Sun-Z 12)
- Equivalent to the estimation of  $1/\Theta_{jk}^{\text{eff}}$
- Minimize  $(\Theta_{jk}^{\text{eff}})^2 \Theta_{jk}^{-3/2} + \Theta_{jk}^{-1/2} + 4\sqrt{(\log p)/n} \|\Theta_{jk}^{\text{eff}}\|_{\text{F}, (j,c)}/\Theta_{jk}^{\text{eff}}$
- $\sqrt{n}(\widehat{\Theta}_{jk} - \Theta_{jk}^{\text{eff}})/\Theta_{jk}^{\text{eff}} \rightarrow N(0, 2)$
- Analysis: large deviation of  $\max_{|A| \leq p^*} \|\underline{\Sigma} - \Sigma\|_{A \times A}$



## Confidence intervals and hypothesis testing for $\Theta_{jk}$

- Estimation of covariance in linear regression (Sun-Z 12; Ren et al 2013)
- Multivariate linear regression
- Subject to  $\mathbf{B}^{A \times A} = \mathbf{I}_A$  with  $\mathbf{B} \in \mathbb{R}^{[p] \times A}$ , minimize

$$\text{diag}^{-1}(\sigma_A)(\mathbf{B}_T \underline{\Sigma} \mathbf{B})^{A \times A} + \|\sigma_A\|_1 + 4 \sqrt{(\log p)/n} \sum_{j \in A} \|\mathbf{B}_{*j}\|_1$$

- Estimate the  $A \times A$  diagonal block of  $\Theta$  by

$$\Theta^{A \times A} = (\mathbf{B}_T \underline{\Sigma} \mathbf{B})^{-1}_{A \times A}$$

- $\Theta_{jk}$  with  $A = \{j, k\}$

- Validity of statistical inference:

$$\mathbb{P} \max_{1 \leq j < k \leq p} \sup_t \left\{ \frac{\sqrt{n}(\Theta_{jk} - \Theta_{jk}^*)}{\sqrt{\Theta_{jj} + \Theta_{kk}}} \leq t \right\} \leq \Phi(t) \rightarrow 0$$

provided that  $\|\Sigma\|_S + \|\Theta\|_S = O(1)$  and  $d^*(\log p) = o(n^{1/2})$

- The Fisher information bound attained

## Confidence intervals and hypothesis testing for $\Theta_{jk}$

- Optimality of the proposed method (Ren et al 2013)

- For  $\mathcal{E}_{M_0, d_0} = \{\Theta : \|\Sigma\|_S + \|\Theta\|_S \leq M_0, d^* \leq d_0\}$  and  $c_0 > 0$ ,

$$\begin{aligned} \min_{1 \leq j \leq k \leq p} \inf_{\Theta \in \mathcal{E}_{M_0, d_0}} \mathbb{P} \left\{ \left| \Theta_{jk} - \Theta_{jk} \right| \geq c_0 \right\} &\geq c_0 \\ \inf_{\Theta \in \mathcal{E}_{M_0, d_0}} \mathbb{P} \left\{ \left| \Theta_{jk} - \Theta_{jk} \right|_{\max} \geq c_0 \right\} &\geq c_0 \end{aligned}$$

- When  $d_0(\log p)/n \leq a_0$  for a certain small  $a_0 > 0$ ,

$$\begin{aligned} \lim_{n \rightarrow \infty} \max_{1 \leq j \leq k \leq p} \sup_{\Theta \in \mathcal{E}_{M_0, d_0}} \mathbb{P} \left\{ \left| \Theta_{jk} - \Theta_{jk} \right| \geq t \right\} &= 0 \\ \sup_{\Theta \in \mathcal{E}_{M_0, d_0}} \mathbb{P} \left\{ \left| \Theta_{jk} - \Theta_{jk} \right|_{\max} \geq t_0 \right\} &\rightarrow 0, \mathbb{E} t_0 \end{aligned}$$

- Lower bound:  $\mathbb{E} I_{\frac{d^*}{d_0}} - 1 \rightarrow 0$  with suitable random  $\Theta \in \mathcal{E}_{M_0, d_0}$

## Copulas

- Data:  $\mathbf{Y} = (Y_{ij}^{(n \times p)}) = (Y_1^p, \dots, Y_p^p)$  with  $\uparrow\uparrow f_j$
- Kendall's (1938) tau between vectors  $\mathbf{Y}_j$  and  $\mathbf{Y}_k$  is
 
$$\tau_{jk} = 2\{n(n-1)\}^{-1} \sum_{1 \leq i_1 < i_2 \leq n} \text{sgn}(Y_{i_1 j} - Y_{i_2 j}) \text{sgn}(Y_{i_1 k} - Y_{i_2 k})$$
- Spearman's (1904) rank correlation between  $\mathbf{Y}_j$  and  $\mathbf{Y}_k$  is
 
$$\widehat{\rho}_{jk} = \text{Corr}(\mathbf{r}_j, \mathbf{r}_k)$$
 where  $\mathbf{r}_j$  is the rank vector of  $\mathbf{Y}_j$ ,  $r_{i_0 j} = \sum_{i=1}^n \mathbb{1}\{Y_{ij} \leq Y_{i_0 j}\}$
- If  $\mathbf{X} = (X_{ij}^{(n \times p)})$  has iid rows and continuous elements, then
 
$$\begin{aligned} \widehat{\tau}_{jk} &= \tau_{jk} = \mathbb{E} \left\{ \text{sgn}(X_{1j} - X_{2j}) \text{sgn}(X_{1k} - X_{2k}) \right\} \\ \widehat{\rho}_{jk} &= \{(n-2)\rho_{jk} + 3\tau_{jk}\} / (n+1) \\ \rho_{jk} &= 3 \mathbb{E} \left\{ \text{sgn}(X_{1j} - X_{2j}) \text{sgn}(X_{1k} - X_{3k}) \right\} \end{aligned}$$
- Moreover,  $X_j \perp\!\!\!\perp X_k$  implies that  $\widehat{\tau}_{jk} = \widehat{\rho}_{jk} = 0$

## Copulas

- Kendall (1948), Kruskal (1958)
- Gaussian copula/nonparanormal: If the rows of  $\mathbf{X} = (x_{ij})_{n \times p}$  are iid Gaussian with *correlation* matrix  $\Sigma$ , then

$$\Sigma_{jk} = \sin(\pi\tau_{jk}/2) = 2 \sin(\pi\rho_{jk}/6)$$

- Elliptical copula/trans-elliptical: If the rows of  $\mathbf{X}$  are iid from

$$\xi \mathbf{A} \mathbf{U}$$

with uniform  $\mathbf{U}$  in the sphere  $\mathbb{S}^{p-1}$ ,  $\xi \in \mathbb{R}$  independent of  $\mathbf{U}$  and a deterministic  $\mathbf{A}$ . Then,

$$\Sigma_{jk} = \sin(\pi\tau_{jk}/2)$$

with  $\Sigma$  being the correlation matrix corresponding to  $\mathbf{A} \mathbf{A}^T$

- The Kendall-Kruskal formula motivates  $\widehat{\Sigma}^\tau$  and  $\widehat{\Sigma}^\rho$  with elements

$$\widehat{\Sigma}_{jk}^\tau = \sin\left(\frac{\pi}{2} \widehat{\tau}_{jk}\right) \quad \widehat{\Sigma}_{jk}^\rho = 2 \sin\left(\frac{\pi}{6} \widehat{\rho}_{jk}\right).$$

What happens if we carry out high-dimensional multivariate analysis/learning using these nonparametric statistics?

Will the “best” results match those of the sample correlation matrix  $\widehat{\Sigma}^s$  based on Gaussian data  $\mathbf{X}$ ? Here

$$\widehat{\Sigma}^s = \mathbf{D}^{-1/2} \mathbf{X}^T \mathbf{X} \mathbf{D}^{-1/2}, \quad \mathbf{D} = \text{diag}(\mathbf{X}^T \mathbf{X}).$$

Examples:

- Estimation of sparse  $\Sigma$
- PCA
- Estimation of sparse  $\Sigma^{-1}$
- Statistical inference about  $(\Sigma^{-1})_{jk}$  in graphical models

- In Liu-Lafferty-Wasserman (2009), Xue-Zou (2012a,b), Liu et al (2012), Han-Liu (2012), Liu-Han-Z (2012) and more, high-dimensional methodologies were proposed and analyzed based on error bounds on

$$\|\underline{\Sigma} - \underline{\Sigma}_{jk}\|_{\max} = \max_{j,k} \underline{\Sigma}_{jk} - \underline{\Sigma}_{jk}$$

- Since  $\underline{\Sigma}_{\tau}^{jk}$  and  $\underline{\Sigma}_{\rho}^{jk}$  are Lipschitz in U-statistics with bounded kernels, its large deviation property is comparable to that of the sample correlation:

$$\mathbb{P}\left\{ \underline{\Sigma} - \underline{\Sigma}_{\max} > Ct \mid \leq p^2 \max_{j,k} \mathbb{P}\left\{ \underline{\Sigma}_{jk} - \underline{\Sigma}_{jk} > Ct \right\} \leq 2p^2 e^{-nt^2/2} \right.$$

- Thus, if the performance of an estimator  $\hat{\theta} = \theta(\underline{\Sigma}_s)$  is guaranteed by such large deviation bounds for  $\|\underline{\Sigma}_s - \underline{\Sigma}\|_{\max}$ , then the same guarantee is valid if  $\underline{\Sigma}_s$  is replaced by  $\underline{\Sigma}_{\tau}$  or  $\underline{\Sigma}_{\rho}$  in copula models where  $\mathbf{X}$  is not available.

- Examples of applications

- Thresholding raw  $\underline{\Sigma}$  for sparse  $\Sigma$  (Bickel-Levina, 2008)
- Extension of  $\|\cdot\|_F$  bounds for G-Lasso or SCAD
- Extension of  $\|\cdot\|_S$  for the CLIME

Spectrum error bounds for  $\underline{\Sigma}_T$  (Han-Liu 13; Wegkamp-Zhao 13)

- Write  $\underline{T} = (\tau_{jk})^{p \times p}$ ,  $\underline{\Sigma}_T = (\sin(\pi/2) \tau_{jk})^{p \times p}$  and  $\underline{\Sigma}_T = (\sin(\pi/2) \tau_{jk})^{p \times p}$
- Use the convexity of  $\text{trace}(e^M)$  for  $M = M_T$  to decouple U-statistics,

$$\left( \mathbb{E} \text{trace} \left( \exp(\lambda(\underline{T} - \mathbb{E} \underline{T})) \right) \leq \mathbb{E} \text{trace} \left( \exp \left( \frac{\lambda}{\lfloor n/2 \rfloor} \sum_{l=1}^{\lfloor n/2 \rfloor} \mathbf{h}(\mathbf{x}_{2l}, \mathbf{x}_{2l-1}) \right) \right) \right)$$

where  $\mathbf{h}$  is the kernel of the U-statistics  $\underline{T}$  and  $\mathbf{x}_l$  are the iid rows of  $\mathbf{X}$

- Use a matrix Bernstein inequality (Oliveira 10; Troop 11) to bound the tail

$$\mathbb{P} \left\{ \|\underline{T} - \mathbb{E} \underline{T}\|_S > t \right\} \leq 2p \exp \left( \frac{-\lfloor n/2 \rfloor t^2 / 2}{p \|\underline{\Sigma}\|_S + \|\underline{\Sigma}\|_S^2 + 2pt} \right)$$

- Requirement of  $n \gg p \log p$  for  $\|\cdot\|_S$  convergence, with an extra  $\log p$
- Lack of concentration for controlling  $\|\cdot\|_S$ , e.g. compared with

$$\mathbb{P} \left\{ \mathbf{X}_T^T \mathbf{X} / n - I^S \geq (\sqrt{d/n} + t)(1 + \sqrt{d/n} + t) \right\} \leq 2e^{-n t^2 / 2}$$

for the sample covariance matrix of  $N(0, I)$  data (Davidson-Szarek, 01)

# Spectrum error bounds in Gaussian copula/nonparanormal models (Mitra-Z 14)

- For both nonparametric estimates  $\underline{\Sigma} = \underline{\Sigma}_t$  and  $\underline{\Sigma} = \underline{\Sigma}_t$ 

$$\mathbb{P} \left\{ \underline{\Sigma} \leq \underline{\Sigma} \leq C_0 \left( \sqrt{\frac{d}{n}} + \frac{d}{n} \right) \right\}$$
- For  $\mathcal{A}^{M,d,m} = \{A : \|A\|_S \leq M, |A| \leq d\}$  with  $|\mathcal{A}^{M,d,m}| \leq m$ ,
 
$$\mathbb{P} \left\{ \max_{A \in \mathcal{A}^{M,d,m}} \|(\underline{\Sigma} - A)_{A \times A}\|_S \geq C_M (\Delta_{d,m,t} + \Delta_{d,m,t}^2 + \Delta'_{d,m,t} + \Delta'_{d,m,t}) \right\} \leq e^{-nt}$$
- where  $C_M$  depends on  $M$  only,  $\Delta_{d,m,t} = \sqrt{(d + \log m)/n + t}$  and  $\Delta'_{d,m,t} = \sqrt{(\log m)/n + t + d(\log p)/n + t}$
- Hoeffding and further decompositions; Lipschitz functions of  $N(0, I)$
- Examples of applications
- Consistency of PCA when  $p/n \rightarrow 0$
- $d$ -sparse PCA when  $d(\log p)/n \rightarrow 0$
- Fast convergence for estimating “ $d$ -bandable”  $\Sigma$  (Cai-Z-Zhou, 10)



## Large deviation problems in the elliptical model and in general

- Sub-Gaussian property of

$$(\text{sgn}(X_1), \dots, \text{sgn}(X_p)),$$

the sign sub-Gaussian condition (Han-Liu, 13)

- Concentration in the spectrum norm of

$$\left( n^{-1} \sum_{i=1}^n \left( g(x_{ij}, x_{ik}) - \mathbb{E} g(X_j, X_k) \right) \right)_{p \times p}$$

with a smooth  $g(x, y)$  satisfying  $\|g\|_\infty \leq 1$ . We solved this problem with a strong condition on  $\{g_{jk}\}$  which happens to hold for Gaussian copulas

**Thanks!**