

# Maximum likelihood estimation over multiple undirected graphs

Yunzhang Zhu, Xiaotong Shen and Wei Pan

University of Minnesota

December 13, 2011

- 1 Introduction
  - Estimation of multiple precision matrices
  - Motivating examples
- 2 Computation
  - DC programming
  - Blockwise minimization
- 3 Theory
- 4 Numerical examples

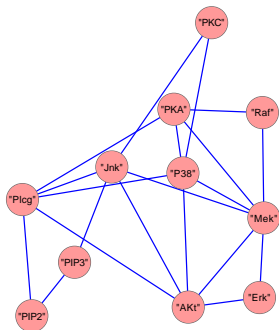
# Estimation of multiple precision matrices

- Multiple Gaussian graphical models:
  - Data:  $\mathbf{Y}_1^{(l)}, \dots, \mathbf{Y}_{n_l}^{(l)} \sim N(\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l); l = 1, \dots, L$ . and  $\boldsymbol{\Sigma}_l$  is  $p \times p$  covariance matrix.
  - Parameter of interest:  $\mathbf{X}_l = \boldsymbol{\Sigma}_l^{-1}; l = 1, \dots, L$ .
- High-dimensional cases:
  - Network size:  $p$  could be large.
  - # networks:  $L$  could be large.
  - # of parameters:  $\frac{p(p-1)}{2}L \gg n = \sum_{l=1}^L n_l$ .

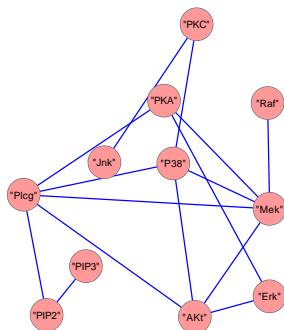
## Why multiple Gaussian graphical models?

- Existing research has focused on single precision matrix estimation through sparseness—Glasso,...
- Networks sharing sparsity pattern: (Guo et al 2011)
  - Entries of matrices tend to be simultaneously zero or not.
- Time-varying networks: (Kolar et al 2009), (Zhou et al 2010), ...
  - Network structures change over time: Adjacent matrices tend to be similar.
- Our goal: Exploit sparse structure within matrices and cluster structure among matrices.

# Motivating examples



(a) condition 1



(b) condition 2

# Estimation of multiple precision matrices

- **Regularization:**

$$S(\mathbf{X}) = \sum_{l=1}^L \left( -\log \det(\mathbf{X}_l) + \text{tr}(\mathbf{S}_l \mathbf{X}_l) \right) + J(\{x_{jkl} : k \neq j\}),$$

Precision matrices:  $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_L)$ ; Sample covariances:  $\mathbf{S}^{(1)}, \dots, \mathbf{S}_L$ ;  $J(\{x_{jkl} : k \neq j\})$  over on off-diagonal elements.

- **Idea:** Sparseness pursuit within matrices and clustering over matrices.

- **Penalty:**  $J(\{x_{jkl} : k \neq j\}) = \lambda_1 \sum_{(l,j,k):k \neq j} J_\tau(|x_{jkl}|) + \lambda_2 \sum_{(j,k,l_1,l_2):|l_1-l_2|=1,k \neq j} J_\tau(|x_{jkl_1} - x_{jkl_2}|)$

- $J_\tau(x) = \min(\frac{x}{\tau}, 1) \rightarrow \mathbb{I}(x \neq 0)$  as  $\tau \rightarrow 0$  (computational surrogate for  $\ell_0$  penalty).

# Multiple precision matrices estimation

- **Parameter estimator:**
  - $p(p+1)L/2$  unknown parameters.
  - Existing methods such as Glasso for single matrix estimation break down.
- **Our strategy:** Difference convex programming + block-wise minimization by sweeping each column of  $\mathbf{X}_l$  across  $l = 1, \dots, L$ .
  - **DC programming:** treat non-convex minimization by solving a sequence of convex relaxations.
  - **Blockwise minimization:** each convex relaxation problem is solved by blockwise coordinate descent method.

# Difference convex programming

- **DC decomposition:** a difference of two convex functions:

$$J_\tau(|x|) = |x| - \max(|x| - \tau, 0) \quad (1)$$

- **Convex relaxation:** given current solution, linearize  $\max(|x| - \tau, 0)$  in (1).
- **Convex subproblem:**

$$\sum_{l=1}^L \left( -\log \det(\mathbf{X}_l) + \text{tr}(\mathbf{S}_l \mathbf{X}_l) \right) + J_{\text{con}}^{(m)}(\{x_{jkl}; k \neq j\}), \quad (2)$$

$$J_{\text{con}}^{(m)}(\{x_{jkl}; k \neq j\}) = \sum_{j=1}^p \sum_{k \neq j} \left( \lambda_1 \sum_{l=1}^L \mathbb{I}(|x_{jkl}^{(m)}| \leq \tau) |x_{jkl}| \right. \\ \left. + \lambda_2 \sum_{l=2}^L \mathbb{I}(|x_{jk(l-1)}^{(m)} - x_{jkl}^{(m)}| \leq \tau) |x_{jk(l-1)} - x_{jkl}| \right).$$

- **Blockwise coordinate descent:** For (2), blockwise coordinate descent is used to sweep over one column at a time.



# Blockwise minimization

- **Notation:** decompose the matrix in terms of  $j$ th column

$$\mathbf{X}_l = \begin{pmatrix} \mathbf{X}_{-jl} & \mathbf{x}_{-jl} \\ \mathbf{x}_{-jl}^T & x_{jjl} \end{pmatrix}, \quad \mathbf{S}_l = \begin{pmatrix} \mathbf{S}_{-jl} & \mathbf{s}_{-jl} \\ \mathbf{s}_{-jl}^T & s_{jjl} \end{pmatrix}$$

$$\mathbf{x}_{-j} = (\mathbf{x}_{-j1}, \dots, \mathbf{x}_{-jL}), \quad J_{\text{con}}^{(m)}(\{x_{jkl} : k \neq j\}) = \sum_{j=1}^p J_j^{(m)}(\mathbf{x}_{-j})$$

with

$$J_j^{(m)}(\mathbf{x}_{-j}) = \sum_{k \neq j} \left( \lambda_1 \sum_{l=1}^L \mathbb{I}(|x_{jkl}^{(m)}| \leq \tau) |x_{jkl}| + \lambda_2 \sum_{l=2}^L \mathbb{I}(|x_{jk(l-1)}^{(m)} - x_{jkl}^{(m)}| \leq \tau) |x_{jk(l-1)} - x_{jkl}| \right)$$

# Blockwise minimization

- **One-column problem:** fixing  $\mathbf{X}_{-j1}, \dots, \mathbf{X}_{-jL}$  and using

$$\det(\mathbf{X}_l) = \det\left(\mathbf{X}_{-jl}(\mathbf{x}_{jjl} - (\mathbf{x}_{-jl})^T(\mathbf{X}_{-jl})^{-1}\mathbf{x}_{-jl})\right),$$

- **Negative log-likelihood:**

$$\sum_{l=1}^L \left( -\log\left(\mathbf{x}_{jjl} - (\mathbf{x}_{-jl})^T(\mathbf{X}_{-jl})^{-1}\mathbf{x}_{-jl}\right) + s_{jjl}\mathbf{x}_{jjl} + 2(\mathbf{s}_{-jl})^T\mathbf{x}_{-jl} \right).$$

- **Optimizing out  $x_{jjl}$ :** Plugging in the optimal  $x_{jjl} = 1/s_{jjl} + (\mathbf{x}_{-jl})^T(\mathbf{X}_{-jl})^{-1}\mathbf{x}_{-jl}$ , negative log-likelihood becomes

$$\sum_{l=1}^L \left( s_{jjl}(\mathbf{x}_{-jl})^T(\mathbf{X}_{-jl})^{-1}(\mathbf{x}_{-jl}) + 2(\mathbf{s}_{-jl})^T\mathbf{x}_{-jl} \right),$$

which is quadratic in  $\mathbf{x}_{-j}$  (if fixing  $\mathbf{X}_{-jl}$ )!

# Algorithm

## Algorithm 1:

- **Initialization:** Set  $X_l^{(0)} = \mathbb{I}_{p \times p}$ ;  $l = 1, \dots, L$ .
- **Step 1:** Iterate over  $m$ , compute  $(X_1^{(m+1)}, \dots, X_L^{(m+1)})$  by solving (2) based on  $(X_1^{(m)}, \dots, X_L^{(m)})$  using **Step 2**.  
 Terminate if  $(X_1^{(m+1)}, \dots, X_L^{(m+1)}) = (X_1^{(m)}, \dots, X_L^{(m)})$ .

- **Step 2:** Iterate over  $j = 1, \dots, p$ , update  $x_{-j}$  by solving

$$\min_{\mathbf{x}_{-j}} \left( \sum_{l=1}^L \left( s_{jpl} (\mathbf{x}_{-jl})^T (\mathbf{X}_{-jl})^{-1} (\mathbf{x}_{-jl}) + 2(s_{-jl})^T \mathbf{x}_{-jl} \right) + J_j^{(m)}(\mathbf{x}_{-j}) \right),$$

and update  $x_{jpl}$  according to

$$x_{jpl} \leftarrow 1/s_{jpl} + (\mathbf{x}_{-jl})^T (\mathbf{X}_{-jl})^{-1} \mathbf{x}_{-jl}.$$

Terminate when certain pre-specified stopping criterion is met.

# Matrix inverse problem

## Remark about Algorithm 1:

- **Matrix inverse:** every iteration requires inverting  $L$   $(p - 1) \times (p - 1)$  sub-matrices  $X_{-jl}$ .
- **Efficient update for inverses:**
  - efficient updates requiring only  $O(p^2)$  operation for each inverse (matrix-vector multiplication and matrix-matrix addition)
  - this is possible because most elements of  $X_l$  are unchanged within each iteration
- **Efficient algorithm for the sub-problem:** fixing  $X_{-jl}$ , the sub-problem is easy to solve.
- **Generalization:** applicable to general form of penalty as long as it is "separable" in column.

# Exponential probability error bound for oracle recovery

- Definition of oracle estimator:
  - MLE given the true sparsity pattern and cluster.
- Simultaneous recovery probability:

$$\mathbb{P}(\text{oracle recovery}) \geq 1 - \exp \left\{ -n \left( c_1 C_{\min}(\boldsymbol{\theta}^0) - 2 \frac{\log d}{n} - \frac{\log(S^*)}{n} \right) \right\}$$

under some regularity conditions.

- $C_{\min}(\boldsymbol{\theta}^0)$ : quantity reflecting the level of difficulty of recovery.
- $d$ : dimension of parameter of interest.
- $S^*$ : quantity related to the graph in the clustering penalty and the true graph.

## The multiple precision matrix case

- Lower bound on entropy:

$$C_{min}(\theta^0) \geq c_2 \min_{1 \leq l \leq L} c_{min}(H_l) \eta_{min}^2,$$

where

$$\eta_{min} = \min \left( \min_{(i,i',l): x_{ii'l}^0 \neq 0} |x_{ii'l}^0|, \frac{1}{\sqrt{2}} \min_{(i,i',l): x_{ii'l}^0 - x_{ii'(l+1)}^0 \neq 0} |x_{ii'l}^0 - x_{ii'(l+1)}^0| \right)$$

- Upper bound on complexity:

$$\log S^* \leq 2g_0 \max(\log(d_0/g_0), 1),$$

where  $g_0 \equiv \sum_{l=1}^{L-1} \sum_{i>i'} \mathbb{I}(x_{ii'l}^0 \neq x_{ii'(l+1)}^0)$  is the number of break points among these clusters.

## The multiple precision matrix case

- Sufficient condition for simultaneous pursuit of sparseness and clustering:

$$\text{If } \min_{1 \leq l \leq L} c_{\min}(H_l) \eta_{\min}^2 \geq c_3 \frac{\log(p^2 L) - g_0 \max(\log(d_0/g_0), 1)}{n}$$

for some constant  $c_3 > 0$ , then

$$\mathbb{P}(\hat{\mathbf{X}}^{\ell_0} \neq \hat{\mathbf{X}}^o) \text{ and } \mathbb{P}(\hat{\mathbf{X}}^g \neq \hat{\mathbf{X}}^o) \rightarrow 0 \text{ as } n, d \rightarrow +\infty.$$

## Evaluation metrics

- The following metrics are considered:
  - Average entropy loss (EL):

$$EL = \frac{1}{L} \sum_{l=1}^L \left( \text{tr}((\mathbf{X}_l)^{-1} \hat{\mathbf{X}}_l) - \log \det((\mathbf{X}_l)^{-1} \hat{\mathbf{X}}_l) \right)$$

- Average quadratic loss (QL):

$$QL = \frac{1}{L} \sum_{l=1}^L \text{tr} \left[ ((\mathbf{X}_l)^{-1} \hat{\mathbf{X}}_l - \mathbf{I})^2 \right]$$



# Tuning

- Cross-validation:

$$CV(\boldsymbol{\lambda}) = \frac{1}{L} \sum_{l=1}^L \left( -\log \det (\hat{\mathbf{X}}_l(\boldsymbol{\lambda})) + \text{tr}(\mathbf{S}_l^{\text{test}} \hat{\mathbf{X}}_l(\boldsymbol{\lambda})) \right),$$

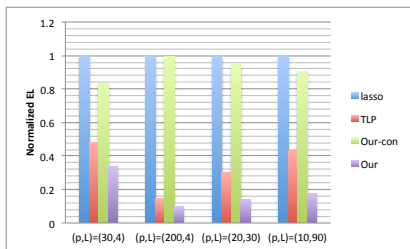
where  $\mathbf{S}_l^{\text{test}}$  is the sample covariance matrix for the testing set;  
 $l = 1, \dots, L$ .

- **Grid search:** get the tuning parameter by minimizing  $CV(\boldsymbol{\lambda})$  through a grid search.

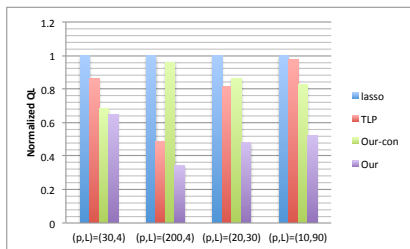
# Example 1

- **Chain Networks:**  $\mathbf{X}_l^{-1} = \Sigma_l$  is AR(1)-structured
  - *ij*-element:  $\sigma_{ijl} = \exp(-|s_{il} - s_{jl}|/2)$  with  $s_{1l} < s_{2l} < \dots < s_{pl}$ ,  $s_{il} - s_{(i-1)l} \sim \text{Unif}(0.5, 1)$ .
- **Scenarios:**
  - $(n, p) = (120, 30)$ ,  $(n, p, L) = (120, 200)$  and  $\mathbf{X}_1 = \mathbf{X}_1, \mathbf{X}_3 = \mathbf{X}_4$ ;
  - $(n, p, L) = (120, 20, 30)$ ,  $(n, p, L) = (120, 10, 90)$  and  $\mathbf{X}_1 = \dots = \mathbf{X}_{L/3}, \mathbf{X}^{(1+L/3)} = \dots = \mathbf{X}_{2L/3}$  and  $\mathbf{X}_{1+2L/3} = \dots = \mathbf{X}_L$ .

# Example 1



(c) example 1: entropy loss



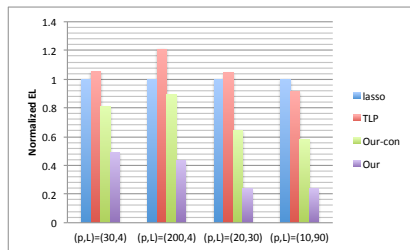
(d) example 1: quadratic loss

- The proposed method outperforms its convex counterpart.
- Simultaneous sparsity and clustering pursuit outperforms that with only sparsity pursuit, and clustering alone.

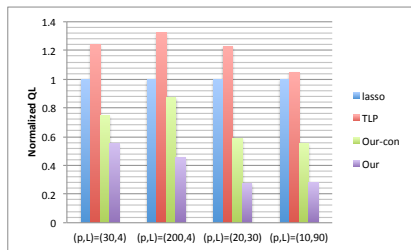
## Example2

- **Nearest Neighbor Networks:** detailed generating scheme in (Li and Gui 2006)
- **Scenarios:**
  - $(n, p, L) = (300, 30, 4)$  and  $(n, p, L) = (300, 200, 4)$ , where  $\mathbf{X}_1 = \mathbf{X}_2, \mathbf{X}_3 = \mathbf{X}_4$  with the second cluster of matrices  $(\mathbf{X}_3, \mathbf{X}_4)$  deleting one edge for each node.
  - $(n, p, L) = (300, 20, 30)$  and  $(n, p, L) = (300, 10, 90)$ , where  $\mathbf{X}_1 = \cdots = \mathbf{X}_{L/3}, \mathbf{X}_{1+L/3} = \cdots = \mathbf{X}_{2L/3}, \mathbf{X}_{1+2L/3} = \cdots = \mathbf{X}_l$  with the second cluster of matrices  $(\mathbf{X}_3, \mathbf{X}_4)$  deleting one edge for each node and the third cluster of matrices  $(\mathbf{X}_{1+2L/3}, \cdots, \mathbf{X}_L)$  adding an edge.

# Example 2



(e) Example 2: entropy loss



(f) Example 2: quadratic loss

## Example 3

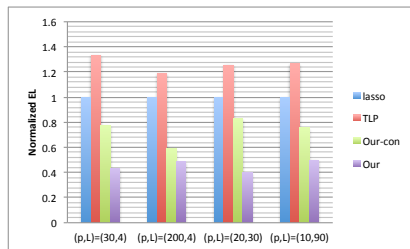
- Exponential Decay Matrix:

- *ij*-element:  $x_{ijl} = \exp(a_l|i - j|)$  with  $a_l$  sampled uniformly over  $[1, 2]$ .

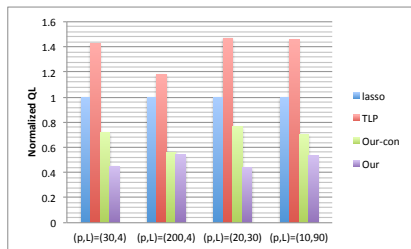
- Scenarios:

- $(p, L) = (30, 4)$ ,  $(p, L) = (200, 4)$  and  $n = 120$  or  $300$  with  $\mathbf{X}_1 = \mathbf{X}_2$ ,  $\mathbf{X}_3 = \mathbf{X}_4$
- $(p, L) = (20, 30)$ ,  $(p, L) = (10, 90)$  and the sample size  $n = 120$  or  $300$  with  $\mathbf{X}_1 = \dots = \mathbf{X}_{L/3}$ ,  $\mathbf{X}_{1+L/3} = \dots = \mathbf{X}_{2L/3}$ ,  $\mathbf{X}_{1+2L/3} = \dots = \mathbf{X}_L$ .

# Example 3: $n = 120$

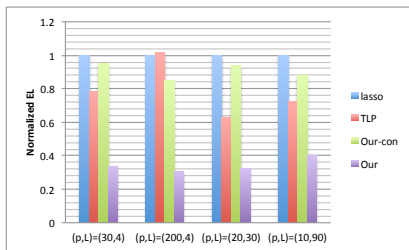


(g) example 3: entropy loss

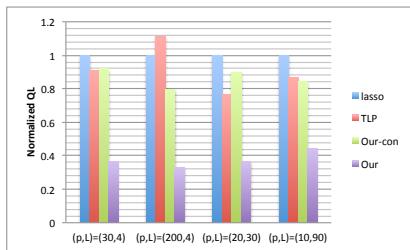


(h) example 3: quadratic loss

# Example 3: $n = 300$



(i) example 3: entropy loss



(j) example 3: quadratic loss



## Proteins data

- Multiple protein network: protein networks under different experimental conditions

