

# ECE 18-898G: Special Topics in Signal Processing: Sparsity, Structure, and Inference

High-dimensional graphical models

Yuejie Chi

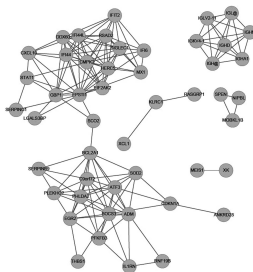
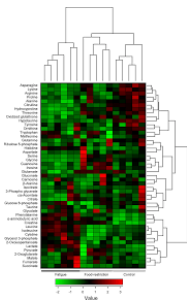
Department of Electrical and Computer Engineering

**Carnegie Mellon University**

Spring 2018

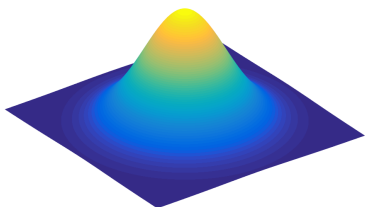
# Identifying Interactions in Data

Given  $n$  data samples,  $x_i \sim \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} \in \mathbb{R}^p$ , how to identify interactions between  $x_i$  and  $x_j$ ?



# Multivariate Gaussians

---



Consider a random vector  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma)$  with pdf

$$\begin{aligned} f(\mathbf{x}) &= \frac{1}{(2\pi)^{p/2} \det(\Sigma)^{1/2}} \exp\left\{-\frac{1}{2}\mathbf{x}^\top \Sigma^{-1} \mathbf{x}\right\} \\ &\propto \det(\Theta)^{1/2} \exp\left\{-\frac{1}{2}\mathbf{x}^\top \Theta \mathbf{x}\right\} \end{aligned}$$

where  $\Sigma = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] \succ \mathbf{0}$  is  $p \times p$  covariance matrix, and  $\Theta = \Sigma^{-1}$  is **inverse covariance matrix / precision matrix**

# Likelihood function for Gaussian models

---

Draw  $n$  i.i.d. samples  $\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathcal{N}(\mathbf{0}, \Sigma)$ , then log-likelihood (up to additive constant) is

$$\begin{aligned}\ell(\Theta) &= \frac{1}{n} \sum_{i=1}^n \log f(\mathbf{x}_i) = \frac{1}{2} \log \det(\Theta) - \frac{1}{2n} \sum_{i=1}^n \mathbf{x}_i^\top \Theta \mathbf{x}_i \\ &= \frac{1}{2} \log \det(\Theta) - \frac{1}{2} \text{tr}(\mathbf{S}\Theta),\end{aligned}$$

where  $\mathbf{S} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$  is sample covariance matrix (SCM).

## Maximum likelihood estimation (MLE)

$$\hat{\Theta} = \operatorname{argmax}_{\Theta \succeq \mathbf{0}} \log \det(\Theta) - \text{tr}(\mathbf{S}\Theta)$$

# The sample-rich regime

---

## Fact 9.1

*If the SCM  $S$  is invertible, the MLE is given as*

$$\hat{\Theta} = S^{-1}.$$

When  $n \gg p$ , the SCM is invertible and classical theory says MLE converges to the truth as sample size  $n \rightarrow \infty$  (consistency).

# High-dimensional / sample-starved regime

---

Practically, we are often in the regime where sample size  $n$  is small, with  $n < p$ . Why?

- Our assumption may only hold for a small window of data collection;
- Our ability may only allow us to collect a few samples;
- The number of features/variables we care is much higher.

In this regime,  $S$  is rank-deficient, and MLE does not even exist.

# High-dimensional / sample-starved regime

---

Practically, we are often in the regime where sample size  $n$  is small, with  $n < p$ . Why?

- Our assumption may only hold for a small window of data collection;
- Our ability may only allow us to collect a few samples;
- The number of features/variables we care is much higher.

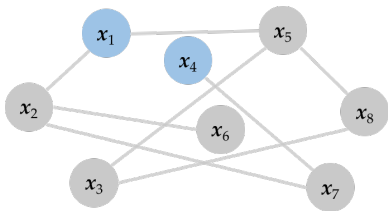
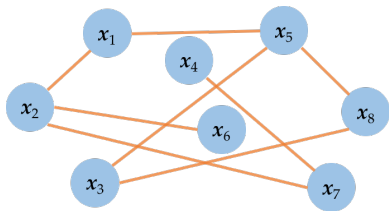
In this regime,  $S$  is rank-deficient, and MLE does not even exist.

Strategy: impose low-dimensional structures.

# Gaussian Graphical Model with Sparsity



# Undirected graphical models



$$x_1 \perp\!\!\!\perp x_4 \mid \{x_2, x_3, x_5, x_6, x_7, x_8\}$$

- Represent a collection of variables  $\mathbf{x} = [x_1, \dots, x_p]^\top$  by a vertex set  $\mathcal{V} = \{1, \dots, p\}$
- Encode conditional independence by a set  $\mathcal{E}$  of edges
  - For any pair of vertices  $u$  and  $v$ ,

$$(u, v) \notin \mathcal{E} \iff x_u \perp\!\!\!\perp x_v \mid \mathbf{x}_{\mathcal{V} \setminus \{u, v\}}$$

# Gaussian graphical models

---

## Lemma 9.2

Consider a Gaussian vector  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ . For any  $u$  and  $v$ ,

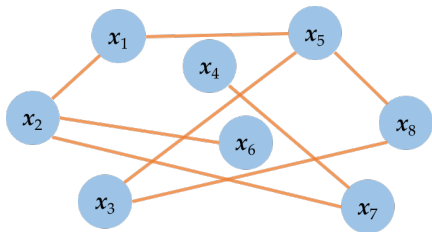
$$x_u \perp\!\!\!\perp x_v \mid \mathbf{x}_{\mathcal{V} \setminus \{u,v\}}$$

iff  $\Theta_{u,v} = 0$ , where  $\Theta = \Sigma^{-1}$ .

Many pairs of variables are conditionally independent  
 $\iff$  many missing links in the graphical model (sparsity)

# Gaussian graphical models

---



$$\underbrace{\begin{bmatrix} * & * & 0 & 0 & * & 0 & 0 & 0 \\ * & * & 0 & 0 & 0 & * & * & 0 \\ 0 & 0 & * & 0 & * & 0 & 0 & * \\ 0 & 0 & 0 & * & 0 & 0 & * & 0 \\ * & 0 & * & 0 & * & 0 & 0 & * \\ 0 & * & 0 & 0 & 0 & * & 0 & 0 \\ 0 & * & 0 & * & 0 & 0 & * & 0 \\ 0 & 0 & * & 0 & * & 0 & 0 & * \end{bmatrix}}_{\Theta}$$

Inverse covariance matrix  $\Theta$  is often (approximately) sparse

# Sparse inverse covariance estimation

---

**Problem definition:** Given  $n$  i.i.d. samples,  $\mathbf{x}_i \sim \mathcal{N}(0, \Sigma)$ , estimate the **sparse** inverse covariance matrix  $\Theta = \Sigma^{-1}$ .

**Two approaches:**

- Graphical Lasso
- CLIME

# Graphical lasso

**Key idea:** regularizing the MLE by imposing  $\ell_1$  regularization (Yuan & Lin'07; Friedman, Hastie, & Tibshirani '08).

## Graphical Lasso (GLasso)

$$\text{maximize}_{\Theta \succeq 0} \quad \log \det(\Theta) - \text{tr}(\mathcal{S}\Theta) - \underbrace{\lambda \|\Theta\|_1}_{\text{lasso penalty}}$$

- It is a convex program! (homework)
- First-order optimality condition

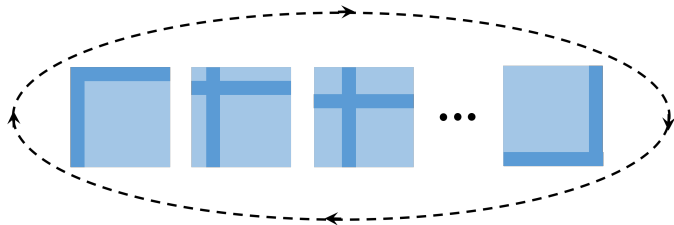
$$\Theta^{-1} - \mathcal{S} - \lambda \underbrace{\partial \|\Theta\|_1}_{\text{subgradient}} = 0 \quad (9.1)$$

$$\implies (\Theta^{-1})_{i,i} = S_{i,i} + \lambda, \quad 1 \leq i \leq p$$

# Blockwise coordinate descent

---

**Idea:** repeatedly cycle through all columns/rows and, in each step, optimize only a single column/row



**Notation:** use  $W$  to denote working version of  $\Theta^{-1}$ . Partition all matrices into 1 column/row vs. the rest

$$\Theta = \begin{bmatrix} \Theta_{11} & \theta_{12} \\ \theta_{12}^\top & \theta_{22} \end{bmatrix} \quad S = \begin{bmatrix} S_{11} & s_{12} \\ s_{12}^\top & s_{22} \end{bmatrix} \quad W = \begin{bmatrix} W_{11} & w_{12} \\ w_{12}^\top & w_{22} \end{bmatrix}$$

## Blockwise coordinate descent

---

**Blockwise step:** suppose we fix all but the last row / column. It follows from (9.1) that

$$\mathbf{0} \in \mathbf{W}_{11}\boldsymbol{\beta} - \mathbf{s}_{12} + \lambda\partial\|\boldsymbol{\theta}_{12}\|_1 = \mathbf{W}_{11}\boldsymbol{\beta} - \mathbf{s}_{12} + \lambda\partial\|\boldsymbol{\beta}_{12}\|_1 \quad (9.2)$$

where  $\boldsymbol{\beta} = -\boldsymbol{\theta}_{12} \cdot w_{22}$  (by matrix inverse formula)

This coincides with optimality condition for

$$\text{minimize}_{\boldsymbol{\beta}} \quad \frac{1}{2}\|\mathbf{W}_{11}^{1/2}\boldsymbol{\beta} - \mathbf{W}_{11}^{-1/2}\mathbf{s}_{12}\|^2 + \lambda\|\boldsymbol{\beta}\|_1 \quad (9.3)$$

# Blockwise coordinate descent

---

**Algorithm 9.1** Block coordinate descent for graphical lasso

---

**Initialize**  $\mathbf{W} = \mathbf{S} + \lambda \mathbf{I}$  and fix its diagonals  $\{w_{i,i}\}$ .

**Repeat until convergence:**

**for**  $t = 1, \dots, p$ :

(i) Partition  $\mathbf{W}$  (resp.  $\mathbf{S}$ ) into 4 parts, where the upper-left part consists of all but the  $j$ th row / column

(ii) Solve

$$\text{minimize}_{\boldsymbol{\beta}} \quad \frac{1}{2} \|\mathbf{W}_{11}^{1/2} \boldsymbol{\beta} - \mathbf{W}_{11}^{-1/2} \mathbf{s}_{12}\|^2 + \lambda \|\boldsymbol{\beta}\|_1$$

(iii) Update  $w_{12} = \mathbf{W}_{11} \boldsymbol{\beta}$

**Set**  $\hat{\boldsymbol{\theta}}_{12} = -\hat{\boldsymbol{\theta}}_{22} \boldsymbol{\beta}$  with  $\hat{\boldsymbol{\theta}}_{22} = 1/(w_{22} - \mathbf{w}_{12}^\top \boldsymbol{\beta})$

---



# Blockwise coordinate descent

---

The only remaining thing is to ensure  $\mathbf{W} \succeq \mathbf{0}$ . This is automatically satisfied:

## Lemma 9.3 (Mazumder & Hastie, '12)

*If we start with  $\mathbf{W} \succ \mathbf{0}$  satisfying  $\|\mathbf{W} - \mathbf{S}\|_\infty \leq \lambda$ , then every row/column update maintains positive definiteness of  $\mathbf{W}$ .*

- If we start with  $\mathbf{W}^{(0)} = \mathbf{S} + \lambda \mathbf{I}$ , then  $\mathbf{W}^{(t)}$  will always be positive definite

## Proof of Lemma 9.3

---

A key observation for the proof of Lemma 9.3

**Fact 9.4 (Lemma 2, Mazumder & Hastie, '12)**

*Solving (9.3) is equivalent to solving*

$$\text{minimize}_{\gamma} (\mathbf{s}_{12} + \gamma)^\top \mathbf{W}_{11}^{-1} (\mathbf{s}_{12} + \gamma) \quad \text{s.t.} \quad \|\gamma\|_\infty \leq \lambda \quad (9.4)$$

*where solutions to 2 problems are related by  $\hat{\beta} = \mathbf{W}_{11}^{-1} (\mathbf{s}_{12} + \hat{\gamma})$*

- Check that optimality condition of (9.3) and that of (9.4) match

## Proof of Lemma 9.3

---

Suppose in  $t^{\text{th}}$  iteration one has  $\|\mathbf{W}^{(t)} - \mathbf{S}\|_{\infty} \leq \lambda$  and

$$\mathbf{W}^{(t)} \succ \mathbf{0}$$

$$\iff \mathbf{W}_{11}^{(t)} \succ \mathbf{0}; \quad w_{22} - \mathbf{w}_{12}^{(t)\top} \left( \mathbf{W}_{11}^{(t)} \right)^{-1} \mathbf{w}_{12}^{(t)} > 0 \quad (\text{Schur complement})$$

We only update  $\mathbf{w}_{12}$ , so it suffices to show

$$w_{22} - \mathbf{w}_{12}^{(t+1)\top} \left( \mathbf{W}_{11}^{(t)} \right)^{-1} \mathbf{w}_{12}^{(t+1)} > 0 \quad (9.5)$$

Recall that  $\mathbf{w}_{12}^{(t+1)} = \mathbf{W}_{11}^t \boldsymbol{\beta}^{t+1}$ . It follows from Fact 9.4 that and

$$\begin{aligned} \|\mathbf{w}_{12}^{(t+1)} - \mathbf{s}_{12}\|_{\infty} &\leq \lambda; \\ \mathbf{w}_{12}^{(t+1)\top} \left( \mathbf{W}_{11}^{(t)} \right)^{-1} \mathbf{w}_{12}^{(t+1)} &\leq \mathbf{w}_{12}^{(t)\top} \left( \mathbf{W}_{11}^{(t)} \right)^{-1} \mathbf{w}_{12}^{(t)}. \end{aligned}$$

Since  $w_{22} = s_{22} + \lambda$  remains unchanged, we establish (9.5).

# CLIME

---

**Key idea:** Utilize two facts:

- $\Sigma \cdot \Theta = I$ .
- The SCM  $S$  can be used as a surrogate of  $\Sigma$ .

## CLIME (Cai, Liu & Luo, 2011)

$$\text{minimize}_{\Theta} \|\Theta\|_1 \quad \text{s.t.} \quad \|S\Theta - I\|_{\infty} \leq \lambda_n.$$

- Note:  $\|A\|_{\infty} = \max_{i,j} |A_{i,j}|$ .
- Parallelizable for each column of  $\Theta$ , thus very efficient.
- Post-processing step needed to guarantee symmetry and PSD.

# Comparison with GLasso

---

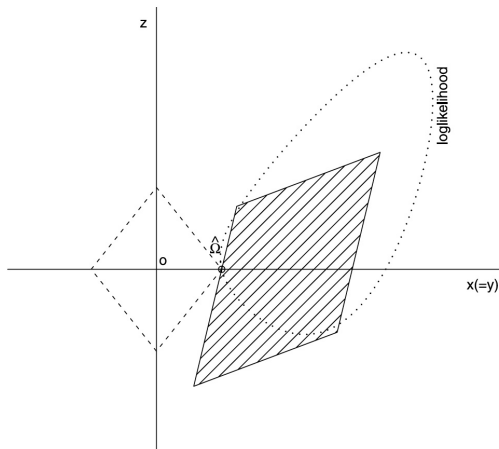


Figure 1. Plot of the elementwise  $\ell_\infty$  constrained feasible set (shaded polygon) and the elementwise  $\ell_1$  norm objective (dashed diamond near the origin) from CLIME. The log-likelihood function as in Glasso is represented by the dotted line.

Figure credit: Cai, Liu & Luo, 2011.

# Gaussian Graphical Model with Latent Variables

# Latent variables in graphical models

---

**Motivation:** some of the variables are not directly observable.



medical/biological



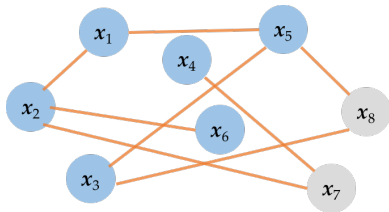
economy

We call the unobserved/missing variables the **latent variables**.

# Graphical models with latent variables

What if one only observes a subset of variables?

$$\begin{bmatrix} \mathbf{x}_o \\ \mathbf{x}_h \end{bmatrix} \quad \begin{array}{l} \text{(observed variables)} \\ \text{(hidden variables)} \end{array}$$



$$\mathbf{x}_o = [x_1, \dots, x_6]^\top, \mathbf{x}_h = [x_7, x_8]^\top$$

Covariance and precision matrices can be partitioned as

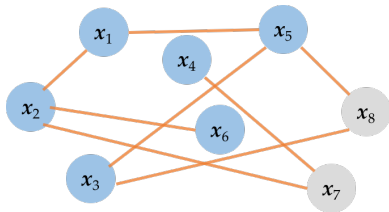
$$\Sigma = \begin{bmatrix} \overbrace{\Sigma_o}^{\text{observed part}} & \Sigma_{o,h} \\ \Sigma_{o,h}^\top & \Sigma_h \end{bmatrix} = \begin{bmatrix} \Theta_o & \Theta_{o,h} \\ \Theta_{o,h}^\top & \Theta_h \end{bmatrix}^{-1}$$



# Graphical models with latent variables

What if one only observes a subset of variables?

$$\begin{bmatrix} \mathbf{x}_o \\ \mathbf{x}_h \end{bmatrix} \quad \begin{array}{l} \text{(observed variables)} \\ \text{(hidden variables)} \end{array}$$



$$\mathbf{x}_o = [x_1, \dots, x_6]^\top, \mathbf{x}_h = [x_7, x_8]^\top$$

$$\Theta_o = \underbrace{\Sigma_o^{-1}}_{\text{observed}} = \underbrace{\Theta_o}_{\text{sparse}} - \underbrace{\Theta_{o,h} \Theta_h^{-1} \Theta_{h,o}}_{\text{low-rank if \# latent vars is small}}$$

sparse + low-rank decomposition

# Inverse covariance estimation for LVGGM

**Problem definition:** Given  $n$  i.i.d. samples,  $\mathbf{x}_i \sim \mathcal{N}(0, \Sigma)$ , estimate the **sparse - low-rank** inverse covariance matrix  $\Theta = \Sigma^{-1}$ .

First write

$$\Theta = \Psi - L$$

where  $\Psi \succeq 0$ ,  $L \succeq 0$ .

**LVGGM (Chandrasekaran, Parrilo, Willsky, 2012)**

$$\begin{aligned} & \text{maximize}_{\Phi, L} \underbrace{\log \det(\Theta) - \text{tr}(\mathcal{S}(\Phi - L))}_{\text{log-likelihood}} - \lambda_n (\|\Psi\|_1 + \eta \text{tr}(L)) \\ & \text{s.t.} \quad \Phi - L \succeq 0, \quad L \succeq 0. \end{aligned}$$

# Reference

---

- [1] "Sparse inverse covariance estimation with the graphical lasso," J. Friedman, T. Hastie, and R. Tibshirani, *Biostatistics*, 2008.
- [2] "The graphical lasso: new insights and alternatives," R. Mazumder and T. Hastie, *Electronic journal of statistics*, 2012.
- [3] "Statistical learning with sparsity: the Lasso and generalizations," T. Hastie, R. Tibshirani, and M. Wainwright, 2015.
- [4] "A constrained  $\ell_1$  minimization approach to sparse precision matrix estimation," T. T. Cai, W. Liu, and X. Luo, *JASA*, 2011.
- [5] "Latent variable graphical model selection via convex optimization," V. Chandrasekaran, P. A. Parrilo, and A. S. Willsky, *The Annals of Statistics*, 2012.