

# Gaussian Graphical Models and Graphical Lasso

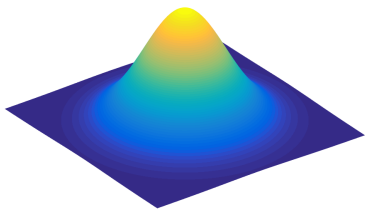


Yuxin Chen

Princeton University, Fall 2020

# Multivariate Gaussians

---

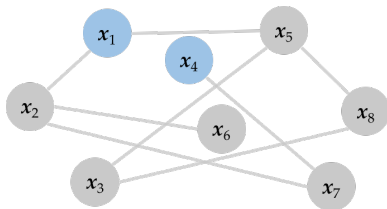
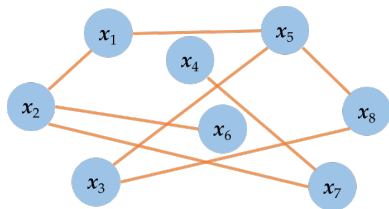


Consider a random vector  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$  with probability density

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

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

# 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

---

## Fact 11.1

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

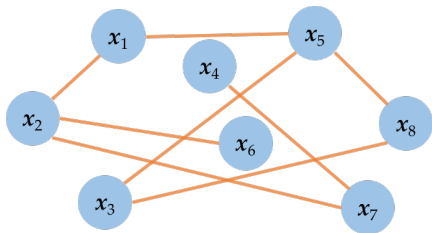
$$x_u \perp\!\!\!\perp x_v \mid \mathbf{x}_{\mathcal{V} \setminus \{u,v\}}$$

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

conditional independence  $\iff$  sparsity

# Gaussian graphical models

---



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

# Likelihoods for Gaussian models

---

Draw  $n$  i.i.d. samples  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ , then the log-likelihood (up to additive constant) is

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

where  $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} \mathbf{x}^{(i)\top}$ : sample covariance;  $\langle \mathbf{S}, \Theta \rangle = \text{tr}(\mathbf{S}\Theta)$

## Maximum likelihood estimation

$$\text{maximize}_{\Theta \succeq \mathbf{0}} \quad \log \det(\Theta) - \langle \mathbf{S}, \Theta \rangle$$

# Challenge in high-dimensional regime

---

Classical theory says MLE converges to the truth as sample size  $n \rightarrow \infty$

Practically, we are often in the regime where the sample size  $n$  is small ( $n < p$ )

- In this regime,  $S$  is rank-deficient, and the MLE does not even exist (why?)

# Graphical lasso (Friedman, Hastie, & Tibshirani '08)

---

In practice, many pairs of variables might be conditionally independent  
 $\iff$  many missing links in the graphical model (sparsity)

**Key idea:** use  $\ell_1$  regularization to promote sparsity

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

- Convex program! (homework)



# Graphical lasso (Friedman, Hastie, & Tibshirani '08)

---

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

- First-order optimality condition

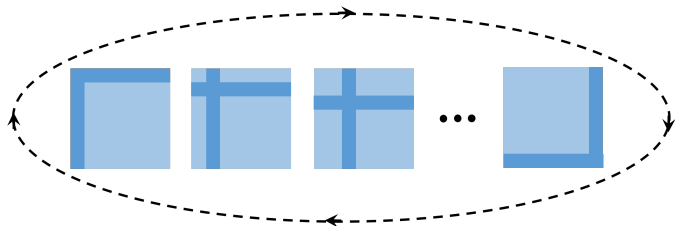
$$\mathbf{0} \in \Theta^{-1} - S - \lambda \underbrace{\partial \|\Theta\|_1}_{\text{subdifferential}} \quad (11.1)$$

- For diagonal entries, one has  $1 \in \partial |\Theta_{i,i}|$  (since  $\Theta_{i,i} > 0$ )

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

## (Optional) 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 a working version of  $\Theta^{-1}$ . Partition all matrices into 1 column / row vs. the rest

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

## (Optional) Blockwise coordinate descent

---

**Blockwise step:** suppose we fix all but the last row / column. It follows from (11.1) that

$$\mathbf{0} \in \mathbf{W}_{11}\boldsymbol{\beta} - \mathbf{s}_{12} - \lambda\partial\|\boldsymbol{\theta}_{12}\|_1 = \mathbf{W}_{11}\boldsymbol{\beta} - \mathbf{s}_{12} + \lambda\partial\|\boldsymbol{\beta}\|_1 \quad (11.2)$$

where  $\boldsymbol{\beta} = -\boldsymbol{\theta}_{12}/\tilde{\theta}_{22}$  (since  $\underbrace{\begin{bmatrix} \Theta_{11} & \boldsymbol{\theta}_{12} \\ \boldsymbol{\theta}_{12}^\top & \theta_{22} \end{bmatrix}^{-1}}_{\text{matrix inverse formula}} = \begin{bmatrix} * & -\frac{1}{\theta_{22}}\Theta_{11}^{-1}\boldsymbol{\theta}_{12} \\ * & * \end{bmatrix}$ ) with

$$\tilde{\theta}_{22} = \theta_{22} - \boldsymbol{\theta}_{12}^\top \Theta_{11}^{-1} \boldsymbol{\theta}_{12} > 0$$

This coincides with the optimality condition for

$$\text{minimize}_{\boldsymbol{\beta}} \quad \frac{1}{2} \|\mathbf{W}_{11}^{1/2}\boldsymbol{\beta} - \mathbf{W}_{11}^{-1/2}\mathbf{s}_{12}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1 \quad (11.3)$$

## (Optional) Blockwise coordinate descent

---

**Algorithm 11.1** 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

$$\text{minimize}_{\boldsymbol{\beta}} \quad \frac{1}{2} \|\mathbf{W}_{11}^{1/2} \boldsymbol{\beta} - \mathbf{W}_{11}^{-1/2} \mathbf{s}_{12}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$$

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

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

---

## (Optional) Blockwise coordinate descent

---

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

### Lemma 11.2 (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

# Reference

---

- "*Sparse inverse covariance estimation with the graphical lasso*," J. Friedman, T. Hastie, and R. Tibshirani, *Biostatistics*, 2008.
- "*The graphical lasso: new insights and alternatives*," R. Mazumder and T. Hastie, *Electronic journal of statistics*, 2012.
- "*Statistical learning with sparsity: the Lasso and generalizations*," T. Hastie, R. Tibshirani, and M. Wainwright, 2015.