

Precision Matrix Estimation

Shihua Zhang

Fall 2019

Contents

- 1 Introduction
- 2 Gaussian Graphical Models
- 3 Main Approaches
- 4 Variants

Outline

- 1 Introduction
- 2 Gaussian Graphical Models
- 3 Main Approaches
- 4 Variants

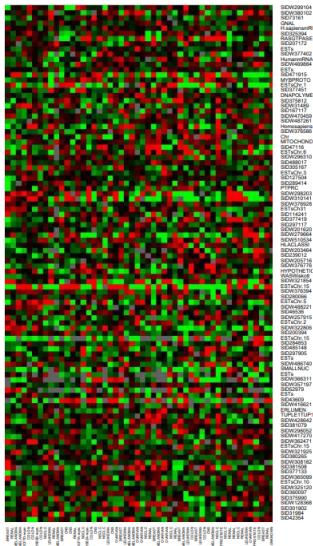
Network of Variables



Gene expression data are typical high-dimensional data.

- Tens of thousands of genes
- Only a few hundreds of samples

Network of Variables



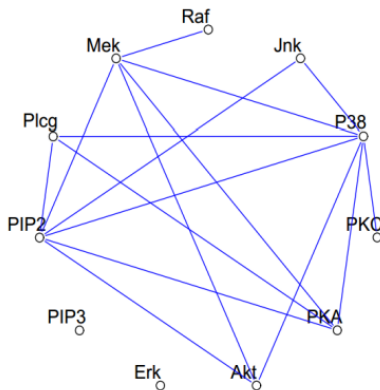
Gene expression data are typical high-dimensional data.

- Tens of thousands of genes
- Only a few hundreds of samples

Network of Variables

- Covariance matrix
- Inverse covariance matrix

Undirected Graphical Model



Consider the simplest undirect graphical model:

- Nodes correspond to random variables
- Edges represent **conditional dependencies** between the variables

Definition of Precision Matrix

In statistics, the **covariance matrix**: $C(i, j)$ is the covariance between the i -th and j -th elements of a random vector.

$$K_{XX} = \begin{bmatrix} E[(X_1 - E[X_1])(X_1 - E[X_1])] & E[(X_1 - E[X_1])(X_2 - E[X_2])] & \cdots & E[(X_1 - E[X_1])(X_n - E[X_n])] \\ E[(X_2 - E[X_2])(X_1 - E[X_1])] & E[(X_2 - E[X_2])(X_2 - E[X_2])] & \cdots & E[(X_2 - E[X_2])(X_n - E[X_n])] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - E[X_n])(X_1 - E[X_1])] & E[(X_n - E[X_n])(X_2 - E[X_2])] & \cdots & E[(X_n - E[X_n])(X_n - E[X_n])] \end{bmatrix}$$

Definition of Precision Matrix

In statistics, the **covariance matrix**: $C(i, j)$ is the covariance between the i -th and j -th elements of a random vector.

$$K_{XX} = \begin{bmatrix} E[(X_1 - E[X_1])(X_1 - E[X_1])] & E[(X_1 - E[X_1])(X_2 - E[X_2])] & \cdots & E[(X_1 - E[X_1])(X_n - E[X_n])] \\ E[(X_2 - E[X_2])(X_1 - E[X_1])] & E[(X_2 - E[X_2])(X_2 - E[X_2])] & \cdots & E[(X_2 - E[X_2])(X_n - E[X_n])] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - E[X_n])(X_1 - E[X_1])] & E[(X_n - E[X_n])(X_2 - E[X_2])] & \cdots & E[(X_n - E[X_n])(X_n - E[X_n])] \end{bmatrix}$$

The **precision matrix** (also known as **concentration matrix**) is the matrix inverse of the covariance matrix.

Challenges of Precision Matrix Estimation

The sample covariance based on the observed data is singular when the dimension is larger than the sample size.

- Many data are of high dimensionality ($p \gg n$)
- Results will be unstable due to the limited number of samples
- The aggregation of a massive amount of estimation errors can lead to considerable adverse impacts on the estimation accuracy

Challenges of Precision Matrix Estimation

The sample covariance based on the observed data is singular when the dimension is larger than the sample size.

- Many data are of high dimensionality ($p \gg n$)
- Results will be unstable due to the limited number of samples
- The aggregation of a massive amount of estimation errors can lead to considerable adverse impacts on the estimation accuracy

Computational complexity is another challenging problem.

- The dimensions of a covariance matrix is $p \times p$.
- The computational complexity of estimating precision matrix directly is $O(p^3)$.

The estimation of large covariance and precision matrices is generally challenging.

Why Precision Matrix is Important?

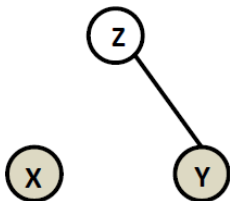
Estimation of a precision matrix is a fundamental problem in many areas of statistical analysis.

- high-dimensional linear discriminant analysis
- complex data visualization
- portfolio allocation
- **graphical model**

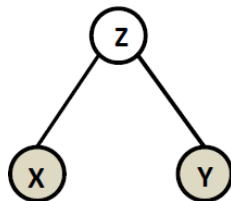
Outline

- 1 Introduction
- 2 Gaussian Graphical Models**
- 3 Main Approaches
- 4 Variants

Independence vs Conditional Independence



$$P(X \cap Y) = P(X)P(Y)$$
$$P(X \cap Y|Z) = P(X|Z)P(Y|Z)$$

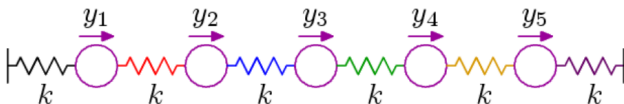


$$P(X \cap Y) \neq P(X)P(Y)$$
$$P(X \cap Y|Z) = P(X|Z)P(Y|Z)$$

Independence: Two events A and B are independent if and only if $P(A \cap B) = P(A)P(B)$.

Conditional independence: Two events A and B are independent if and only if $P(A \cap B|C) = P(A|C)P(B|C)$.

Covariance Matrix vs Precision Matrix



inverse-covariance matrix

or

covariance matrix?

$$\mathbf{K}^{-1} = \frac{k}{T} \begin{bmatrix} 2 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & 2 \end{bmatrix}$$

$$\mathbf{K} = \frac{T}{k} \begin{bmatrix} 0.83 & 0.67 & 0.50 & 0.33 & 0.17 \\ 0.67 & 1.33 & 1.00 & 0.67 & 0.33 \\ 0.50 & 1.00 & 1.50 & 1.00 & 0.50 \\ 0.33 & 0.67 & 1.00 & 1.33 & 0.67 \\ 0.17 & 0.33 & 0.50 & 0.67 & 0.83 \end{bmatrix}$$

For the multivariate Gaussian distribution, **precision matrix** encodes the conditionally independent between variables.

Gaussian Graphical Model

Background:

Multivariate Gaussian Distribution

$$X \sim N_p(\mu, \Sigma)$$

- If Σ is positive definite, distribution has density on \mathbb{R}^p

$$f(x|\mu, \Sigma) = (2\pi)^{-p/2} (\det \Theta)^{1/2} e^{-(x-\mu)^T \Theta (x-\mu)/2}$$

where $\Theta = \Sigma^{-1}$ is the precision matrix of the distribution.

- Conditional distribution:

$$\begin{aligned} X_1|X_2 &\sim \mathcal{N}(\mu_{1|2}, \Sigma_{1|2}) \\ \text{where: } \mu_{1|2} &= \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2) \\ \Sigma_{1|2} &= \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} \end{aligned}$$

Gaussian Graphical Model

Question:

Precision matrix \Leftrightarrow conditional dependencies?

Suppose that observations $x_1, x_2, \dots, x_n \in \mathbb{R}^p$ are i.i.d. $N_p(\mu, \Sigma)$ where $\mu \in \mathbb{R}^p$ and Σ is a $p \times p$ positive definite matrix.

- Partition $X = (Z, Y)$ where $Z = (X_1, \dots, X_{p-1})$ and $Y = X_p$.
- Partitioned Σ and Θ as

$$\Sigma = \begin{pmatrix} \Sigma_{ZZ} & \sigma_{ZY} \\ \sigma_{ZY}^T & \sigma_{YY} \end{pmatrix}, \quad \Theta = \begin{pmatrix} \Theta_{ZZ} & \theta_{ZY} \\ \theta_{ZY}^T & \theta_{YY} \end{pmatrix}$$

Gaussian Graphical Model

Because $\Theta = \Sigma^{-1}$ standards formulas for partitioned inverses given:

- $\theta_{YY} = (\sigma_{YY} - \sigma_{ZY}^T \Sigma_{ZZ}^{-1} \sigma_{ZY})^{-1} \succ 0$
- $\theta_{ZY} = -\theta_{YY} \cdot \Sigma_{ZZ}^{-1} \sigma_{ZY}$

And the conditional distribution of Y given Z :

$$Y(Z = z) \sim N\left(\mu_Y + (z - \mu_Z)^T \Sigma_{ZZ}^{-1} \sigma_{ZY}, \sigma_{YY} - \sigma_{ZY}^T \Sigma_{ZZ}^{-1} \sigma_{ZY}\right)$$

Gaussian Graphical Model

Because $\Theta = \Sigma^{-1}$ standards formulas for partitioned inverses given:

- $\theta_{YY} = (\sigma_{YY} - \sigma_{ZY}^T \Sigma_{ZZ}^{-1} \sigma_{ZY})^{-1} \succ 0$
- $\theta_{ZY} = -\theta_{YY} \cdot \Sigma_{ZZ}^{-1} \sigma_{ZY}$

And the conditional distribution of Y given Z :

$$Y(Z = z) \sim N\left(\mu_Y + (z - \mu_Z)^T \Sigma_{ZZ}^{-1} \sigma_{ZY}, \sigma_{YY} - \sigma_{ZY}^T \Sigma_{ZZ}^{-1} \sigma_{ZY}\right)$$

Comparing these two equations,

- If $(\theta_{ZY})_i = 0$, then $Y(Z = z)$ is nothing to do with the value of z_i

Gaussian Graphical Model

Because $\Theta = \Sigma^{-1}$ standards formulas for partitioned inverses given:

- $\theta_{YY} = (\sigma_{YY} - \sigma_{ZY}^T \Sigma_{ZZ}^{-1} \sigma_{ZY})^{-1} \succ 0$
- $\theta_{ZY} = -\theta_{YY} \cdot \Sigma_{ZZ}^{-1} \sigma_{ZY}$

And the conditional distribution of Y given Z :

$$Y(Z = z) \sim N\left(\mu_Y + (z - \mu_Z)^T \Sigma_{ZZ}^{-1} \sigma_{ZY}, \sigma_{YY} - \sigma_{ZY}^T \Sigma_{ZZ}^{-1} \sigma_{ZY}\right)$$

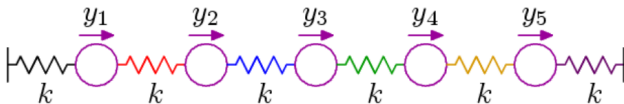
Comparing these two equations,

- If $(\theta_{ZY})_i = 0$, then $Y(Z = z)$ is nothing to do with the value of z_i

$$\theta_{i,p} = 0 \Leftrightarrow i \perp p | V \setminus \{i, p\} \Leftrightarrow \text{Edge } (i, p) \text{ doesn't exist}$$

Gaussian Graphical Model

We can construct a graphical model by estimating the precision matrix.



inverse-covariance matrix

or

covariance matrix?

$$\mathbf{K}^{-1} = \frac{k}{T} \begin{bmatrix} 2 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & 2 \end{bmatrix}$$

$$\mathbf{K} = \frac{T}{k} \begin{bmatrix} 0.83 & 0.67 & 0.50 & 0.33 & 0.17 \\ 0.67 & 1.33 & 1.00 & 0.67 & 0.33 \\ 0.50 & 1.00 & 1.50 & 1.00 & 0.50 \\ 0.33 & 0.67 & 1.00 & 1.33 & 0.67 \\ 0.17 & 0.33 & 0.50 & 0.67 & 0.83 \end{bmatrix}$$

Outline

- 1 Introduction
- 2 Gaussian Graphical Models
- 3 Main Approaches**
 - Sparsity
 - Penalized Likelihood Methods
 - Column-by-column Estimation Methods
 - CLIME
- 4 Variants

1 Introduction

2 Gaussian Graphical Models

3 Main Approaches

- Sparsity
- Penalized Likelihood Methods
- Column-by-column Estimation Methods
- CLIME

4 Variants

The key assumption of the precision matrix estimation is that the target matrix of interest is sparse (i.e., many entries are either zeros or nearly so).

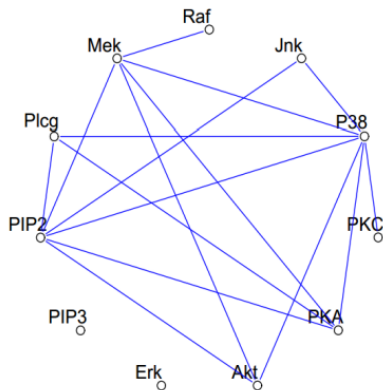
- **Question 1:** why do we need a sparse solution?

The key assumption of the precision matrix estimation is that the target matrix of interest is sparse (i.e., many entries are either zeros or nearly so).

- **Question 1:** why do we need a sparse solution?
 - feature/variable selection
 - better to interpret the data
 - shrink the size of model
 - computational savings
 - discourage overfitting

Sparsity

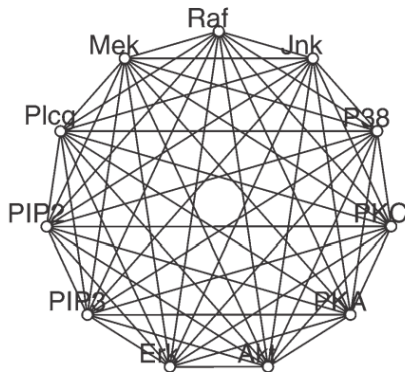
A real network is a set of links with direct dependencies.



- Spare and structured

Sparsity

Estimated network without sparsity constraint.



- Dense and meaningless

Sparsity

The key assumption of precision matrix estimation is that the target matrix of interest is sparse (i.e., many entries are either zeros or nearly so).

- **Question 1:** why do we need a sparse solution?

Sparsity

The key assumption of precision matrix estimation is that the target matrix of interest is sparse (i.e., many entries are either zeros or nearly so).

- **Question 1:** why do we need a sparse solution?

- feature/variable selection
- better to interpret the data
- shrink the size of model
- computational savings
- discourage overfitting

- **Question 2:** how to achieve a sparse solution?

Sparsity

Take the linear regression as an example ($f(x) = w^T * x + b$).

Subset selection: l_0 -norm regularization

$$\min \mathcal{L} = \sum_{i=1}^N |y_i - f(x_i)|^2 + \frac{\lambda}{2} \|w\|_0$$

where

$$\|w\|_0 = \#(i | w_i \neq 0)$$

Sparsity

Take the linear regression as an example ($f(x) = w^T * x + b$).

Subset selection: l_0 -norm regularization

$$\min \mathcal{L} = \sum_{i=1}^N |y_i - f(x_i)|^2 + \frac{\lambda}{2} \|w\|_0$$

where

$$\|w\|_0 = \#(i | w_i \neq 0)$$

- sparse solution
- but non-convex and hard to optimize

Ridge: l_2 -norm regularization

$$\min \mathcal{L} = \sum_{i=1}^N |y_i - f(x_i)|^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Its equivalent form is (constrained optimization):

$$\begin{aligned} \min \mathcal{L} &= \sum_{i=1}^N |y_i - f(x_i)|^2 \\ \text{s.t. } &\|\mathbf{w}\|_2^2 \leq C \end{aligned}$$

Ridge: l_2 -norm regularization

$$\min \mathcal{L} = \sum_{i=1}^N |y_i - f(x_i)|^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

Its equivalent form is (constrained optimization):

$$\begin{aligned} \min \mathcal{L} &= \sum_{i=1}^N |y_i - f(x_i)|^2 \\ \text{s.t. } &\|\mathbf{w}\|_2^2 \leq C \end{aligned}$$

- convex but generate a non-sparse solution (values close to zeros)

Lasso: l_1 -norm regularization

$$\min \mathcal{L} = \sum_{i=1}^N |y_i - f(x_i)|^2 + \frac{\lambda}{2} \|w\|_1$$

Its equivalent form is (constrained optimization):

$$\begin{aligned} \min \mathcal{L} &= \sum_{i=1}^N |y_i - f(x_i)|^2 \\ \text{s.t. } \|w\|_1 &\leq C \end{aligned}$$

Lasso: l_1 -norm regularization

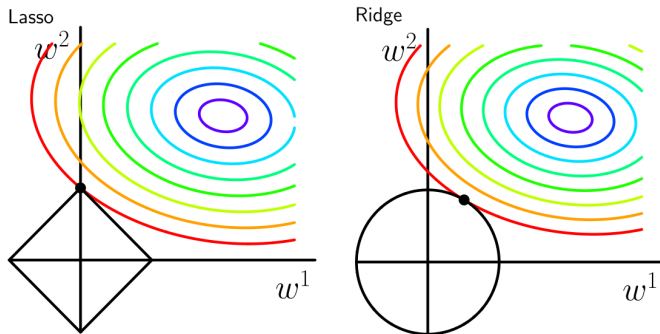
$$\min \mathcal{L} = \sum_{i=1}^N |y_i - f(x_i)|^2 + \frac{\lambda}{2} \|w\|_1$$

Its equivalent form is (constrained optimization):

$$\begin{aligned} \min \mathcal{L} &= \sum_{i=1}^N |y_i - f(x_i)|^2 \\ \text{s.t. } \|w\|_1 &\leq C \end{aligned}$$

- convex
- sparse solution

Why Lasso Leads to Sparsity?



l_1 -norm regularization helps to generate sparse estimation.

1 Introduction

2 Gaussian Graphical Models

3 Main Approaches

- Sparsity
- **Penalized Likelihood Methods**
- Column-by-column Estimation Methods
- CLIME

4 Variants

Penalized Likelihood Methods

One of the most commonly used approaches to estimate sparse precision matrices is the **penalized maximum likelihood**.

- When $x_1, x_2, \dots, x_n \in \mathbb{R}^P$ are i.i.d. $N(\mathbf{0}, \Sigma)$

$$f(x_1, \dots, x_n | \mu, \Sigma) = (2\pi)^{-d/2} (\det \Theta)^{1/2} e^{-\sum_i \text{tr}((x_i - \mu)^T (x_i - \mu) \Theta / 2)}$$

- The negative Gaussian log-likelihood function is given by

$$\ell(\Theta) = \text{tr}(\mathbf{S}\Theta) - \log |\Theta|$$

- Penalized likelihood methods:

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \left\{ \text{tr}(\mathbf{S}\Theta) - \log |\Theta| + \sum_{i \neq j} P(|\theta_{ij}|) \right\}$$

Penalized Likelihood Methods

One of the commonly used convex penalties is the l_1 penalty [Meinshausen et al., 2006].

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \{ \operatorname{tr}(\mathbf{S}\Theta) - \log |\Theta| + \lambda \|\Theta\|_1 \}$$

Solutions

- Interior-point optimization methods [Banerjee et al., 2008]
- **Graphical Lasso (most popular)** [Friedman et al., 2008]
- Alternating direction method of multipliers [Boyd et al., 2011]
- QUIC [Hsieh et al., 2014]

Trick and variants

- Block screen [Mazumder and Hastie, 2012]
- D -trace loss [Zhang and Zou, 2014]

Graphical Lasso

Graphical Lasso [Friedman et al., 2008]

- **Problem:** maximize the l_1 penalized log-likelihood:

$$\log \det \Theta - \text{tr}(\mathbf{S}\Theta) - \lambda \|\Theta\|_1$$

- Optimization by **blockwise coordinate descent**.
- Fast: it solves a 1000-variable problem (about 500,000 parameters) in at most one minute.

Optimization of Graphical Lasso

- Graphical Lasso considers estimation of Σ (rather than Σ^{-1})
- Objective function:

$$\log \det \Sigma^{-1} - \text{tr}(\mathbf{S}\Sigma^{-1}) - \lambda \|\Sigma^{-1}\|_1$$

- Let \mathbf{W} be the estimate of Σ and partitioning \mathbf{W} and \mathbf{S}

$$\mathbf{W} = \begin{pmatrix} W_{11} & w_{12} \\ w_{12}^T & w_{22} \end{pmatrix}, \quad \mathbf{S} = \begin{pmatrix} S_{11} & s_{12} \\ s_{12}^T & s_{22} \end{pmatrix}$$

- Blockwise coordinate descent: Fix W_{11} to optimize w_{12}

Optimization of Graphical Lasso

Equivalence problem [Banerjee et al., 2008]

When fix W_{11} to optimize w_{12} ,

$$\operatorname{argmax}_{w_{12}} \left\{ \log \det \Sigma^{-1} - \operatorname{tr}(\mathbf{S} \Sigma^{-1}) - \lambda \|\Sigma^{-1}\|_1 \right\}$$

equals solving a Lasso problem

$$\min_{\beta} \left\{ \frac{1}{2} \left\| W_{11}^{1/2} \beta - W_{11}^{-1/2} s_{12} \right\|^2 + \lambda \|\beta\|_1 \right\}$$

Optimization of Graphical Lasso

Proof: The subgradient equation of the log-likelihood:

$$w_{12} - s_{12} - \lambda \cdot \gamma_{12} = 0$$

where γ_{12} is the derivative of l_1 -norm.

For the Lasso problem

$$\min_{\beta} \left\{ \frac{1}{2} \left\| W_{11}^{1/2} \beta - W_{11}^{-1/2} s_{12} \right\|^2 + \lambda \|\beta\|_1 \right\}$$

its subgradient equation

$$W_{11} \beta - s_{12} + \lambda \cdot v = 0$$

For (w_{12}, γ_{12}) solves log-likelihood, then $\beta = W_{11}^{-1} w_{12}$ and $v = -\gamma_{12}$ solves the Lasso problem.

Optimization of Graphical Lasso

Graphical Lasso algorithm

1. Start with $W = S + \rho I$. The diagonal of W remains unchanged in what follows.
2. For each $j = 1, 2, \dots, p, 1, 2, \dots, p, \dots$, solve the Lasso problem, which takes as input the inner products W_{11} and s_{12} . This gives a $p - 1$ vector solution $\hat{\beta}$. Fill in the corresponding row and column of W using $w_{12} = W_{11}\hat{\beta}$.
3. Continue until convergence. Obtain estimation of Σ : $\hat{\Sigma} = W$

Optimization of Graphical Lasso

After estimate $\hat{\Sigma} = W$, we can recover $\hat{\Theta} = W^{-1}$ relatively cheaply.

$$\begin{pmatrix} W_{11} & w_{12} \\ w_{12}^T & w_{22} \end{pmatrix} \begin{pmatrix} \Theta_{11} & \theta_{12} \\ \theta_{12}^T & \theta_{22} \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0^T & 1 \end{pmatrix}$$

So

$$\theta_{12} = -W_{11}^{-1} w_{12} \theta_{22}$$

$$\theta_{22} = 1 / (w_{22} - w_{12}^T W_{11}^{-1} w_{12})$$

But since $\hat{\beta} = W_{11}^{-1} w_{12}$

$$\hat{\theta}_{22} = 1 / (w_{22} - w_{12}^T \hat{\beta})$$

$$\hat{\theta}_{12} = -\hat{\beta} \hat{\theta}_{22}$$

Using the stored the coefficients β , we can compute $\hat{\Theta}$ cheaply after convergence.

Graphical Lasso

Graphical Lasso allows each inverse covariance element to be penalized differently,

$$\log \det \Theta - \text{tr}(S\Theta) - \|\Theta * P\|_1$$

where $P = \{\rho_{jk}\}$ with $\rho_{jk} = \rho_{kj}$

Changes the Lasso problem to

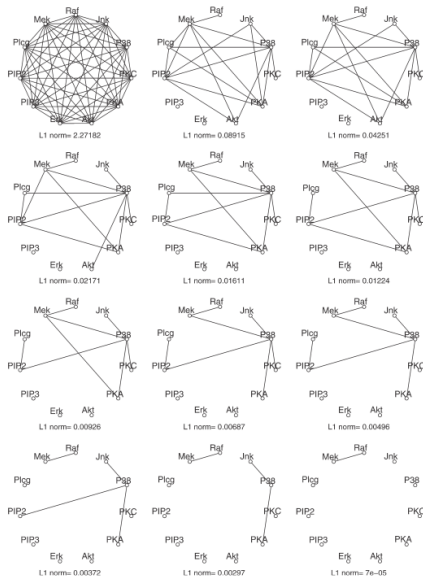
$$\min_{\beta} \left\{ \frac{1}{2} \left\| W_{11}^{1/2} \beta - W_{11}^{-1/2} s_{12} \right\|^2 + P_{12} \|\beta\|_1 \right\}$$

Time Comparison

Table 1. *Timings (seconds) for graphical lasso, Meinhausen–Buhlmann approximation, and COVSEL procedures*

| p | Problem type | (1) Graphical lasso | (2) Approximation | (3) COVSEL | Ratio of (3) to (1) |
|-----|--------------|---------------------|-------------------|------------|---------------------|
| 100 | Sparse | 0.014 | 0.007 | 34.7 | 2476.4 |
| 100 | Dense | 0.053 | 0.018 | 2.2 | 40.9 |
| 200 | Sparse | 0.050 | 0.027 | > 205.35 | > 4107 |
| 200 | Dense | 0.497 | 0.146 | 16.9 | 33.9 |
| 400 | Sparse | 1.23 | 0.193 | > 1616.7 | > 1314.3 |
| 400 | Dense | 6.2 | 0.752 | 313.0 | 50.5 |

Different Penalty Parameters



Penalized Likelihood Methods

One of the commonly used convex penalties is the l_1 penalty [Meinshausen et al., 2006].

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \{ \operatorname{tr}(\mathbf{S}\Theta) - \log |\Theta| + \lambda \|\Theta\|_1 \}$$

Solutions

- Interior-point optimization methods [Banerjee et al., 2008]
- Graphical Lasso (most popular) [Friedman et al., 2008]
- Alternating direction method of multipliers (ADMM) [Boyd et al., 2011]
- QUIC [Hsieh et al., 2014]

Trick and variants

- Block screen [Mazumder and Hastie, 2012]
- D -trace loss [Zhang and Zou, 2014]

ADMM

- ADMM is a method with good robustness of method of multipliers, which can support decomposition.

ADMM problem form (with f, g convex)

$$\begin{array}{ll}\text{minimize} & f(x) + g(z) \\ \text{subject to} & Ax + Bz = c\end{array}$$

The augmented Lagrangian:

$$L_\rho(x, z, y) = f(x) + g(z) + y^T(Ax + Bz - c) + (\rho/2)\|Ax + Bz - c\|_2^2$$

ADMM:

$$x^{k+1} := \operatorname{argmin}_x L_\rho(x, z^k, y^k)$$

$$z^{k+1} := \operatorname{argmin}_z L_\rho(x^{k+1}, z, y^k)$$

$$y^{k+1} := y^k + \rho(Ax^{k+1} + Bz^{k+1} - c)$$

ADMM with scaled dual variables

Combine linear and quadratic terms in the augmented Lagrangian:

$$\begin{aligned} L_{\rho}(x, z, y) &= f(x) + g(z) + y^T (Ax + Bz - c) + (\rho/2) \|Ax + Bz - c\|_2^2 \\ &= f(x) + g(z) + (\rho/2) \|Ax + Bz - c + u\|_2^2 + \text{const.} \end{aligned}$$

with $u^k = (1/\rho)y^k$

ADMM (scaled dual form):

$$\begin{aligned} x^{k+1} &:= \underset{x}{\operatorname{argmin}} \left(f(x) + (\rho/2) \|Ax + Bz^k - c + u^k\|_2^2 \right) \\ z^{k+1} &:= \underset{z}{\operatorname{argmin}} \left(g(z) + (\rho/2) \|Ax^{k+1} + Bz - c + u^k\|_2^2 \right) \\ u^{k+1} &:= u^k + (Ax^{k+1} + Bz^{k+1} - c) \end{aligned}$$

Convergence

Assume (very little!)

- f, g convex, closed, proper
- L_0 has a saddle point

Then ADMM converges:

- Residual convergence: $Ax^k + Bz^k - c \rightarrow 0$
- Objective convergence: $f(x^k) + g(z^k) \rightarrow p^*$
- Dual convergence: $u^k \rightarrow u^*$

Convergence rate: not known in general, theory is currently being developed, e.g., in Hong and Luo (2012), Nishihara *et al.* (2015).
Roughly, it behaves like a first-order method (or a bit faster).

ADMM for Precision Matrix Estimation

Problem form

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \{ \operatorname{tr}(\mathbf{S}\Theta) - \log |\Theta| + \lambda \|\Theta\|_1 \}$$

ADMM for Precision Matrix Estimation

Problem form

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \{ \operatorname{tr}(\mathbf{S}\Theta) - \log |\Theta| + \lambda \|\Theta\|_1 \}$$

ADMM form

$$\begin{array}{ll} \text{minimize} & \operatorname{tr}(\mathbf{S}X) - \log \det X + \lambda \|Z\|_1 \\ \text{subject to} & X - Z = 0 \end{array}$$

ADMM for Precision Matrix Estimation

Problem form

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \{ \operatorname{tr}(\mathbf{S}\Theta) - \log |\Theta| + \lambda \|\Theta\|_1 \}$$

ADMM form

$$\begin{array}{ll} \text{minimize} & \operatorname{tr}(\mathbf{S}X) - \log \det X + \lambda \|Z\|_1 \\ \text{subject to} & X - Z = 0 \end{array}$$

The augmented Lagrangian:

$$L = \operatorname{tr}(\mathbf{S}X) - \log \det X + \lambda \|Z\|_1 + (\rho/2) \|X - Z + U\|_2^2$$

ADMM for Precision Matrix Estimation

ADMM form

$$L = \text{tr}(\mathbf{S}\mathbf{X}) - \log \det \mathbf{X} + \lambda \|\mathbf{Z}\|_1 + (\rho/2) \|\mathbf{X} - \mathbf{Z} + \mathbf{U}\|_2^2$$

ADMM:

$$\mathbf{X}^{k+1} := \underset{\mathbf{X}}{\text{argmin}} \left(\text{tr}(\mathbf{S}\mathbf{X}) - \log \det \mathbf{X} + (\rho/2) \|\mathbf{X} - \mathbf{Z}^k + \mathbf{U}^k\|_F^2 \right)$$

$$\mathbf{Z}^{k+1} := \underset{\mathbf{Z}}{\text{argmin}} \left((\rho/2) \|\mathbf{X} - \mathbf{Z}^k + \mathbf{U}^k\|_F^2 + \lambda \|\mathbf{Z}\|_1 \right)$$

$$\mathbf{U}^{k+1} := \mathbf{U}^k + (\mathbf{X}^{k+1} - \mathbf{Z}^{k+1})$$

Update for X

Problem

$$X^{k+1} = \underset{X}{\operatorname{argmin}} \left(\operatorname{tr}(SX) - \log \det X + (\rho/2) \|X - Z^k + U^k\|_F^2 \right)$$

Differentiating with respect to X , the minimum solves:

$$S - X^{-1} + \rho(X - Z^k + U^k) = 0$$

that is

$$\rho X - X^{-1} = \rho(Z^k - U^k) - S$$

which is a eigenvalue problem.

Update for X

$$\rho X - X^{-1} = \rho(Z^k - U^k) - S$$

Compute the eigendecomposition:

$$\rho(Z^k - U^k) - S = Q\Lambda Q^T$$

Then the eigendecomposition of X^{k+1} is $Q\tilde{X}Q^T$, where \tilde{X} is a diagonal matrix and $\rho\tilde{X} - \tilde{X}^{-1} = \Lambda$

So $X^{k+1} := Q\tilde{X}Q^T$ with

$$\tilde{X}_{ii} = \frac{\lambda_i + \sqrt{\lambda_i^2 + 4\rho}}{2\rho}$$

Cost of X -update is an eigendecomposition.

Time Cost

For a $1000 \times 1000 \Sigma^{-1}$ with 10^4 nonzeros

- Graphical Lasso (Fortran): 20 seconds – 3 minutes
- ADMM (Matlab): 3 – 10 minutes

It is flexible to extend (such as adding other convex penalty in the log-likelihood function).

Penalized Likelihood Methods

One of the commonly used convex penalties is the l_1 penalty [Meinshausen et al., 2006].

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \{ \operatorname{tr}(\mathbf{S}\Theta) - \log |\Theta| + \lambda \|\Theta\|_1 \}$$

Solutions

- Interior-point optimization methods [Banerjee et al., 2008]
- Graphical Lasso (most popular) [Friedman et al., 2008]
- Alternating direction method of multipliers (ADMM) [Boyd et al., 2011]
- QUIC [Hsieh et al., 2014]

Trick and variants

- Block screen [Mazumder and Hastie, 2012]
- D -trace loss [Zhang and Zou, 2014]

QUIC: QUadratic Approximation of Inverse Covariance Matrices

- Existing methods are first-order iterative methods that mainly use gradient information at each step.
- Disadvantage:** they are at most linear rates of convergence.

Question: Can we achieve superlinearly rate of convergence by considering second-order methods?

QUIC: QUadratic Approximation of Inverse Covariance Matrices

- Existing methods are first-order iterative methods that mainly use gradient information at each step.
- Disadvantage:** they are at most linear rates of convergence.

Question: Can we achieve superlinearly rate of convergence by considering second-order methods?

Difficulties: second-order methods at least in part use the Hessian of the objective function.

- This is too expensive for high-dimensional problem.

QUIC: QUadratic Approximation of Inverse Covariance Matrices

- Existing methods are first-order iterative methods that mainly use gradient information at each step.
- Disadvantage:** they are at most linear rates of convergence.

Question: Can we achieve superlinearly rate of convergence by considering second-order methods?

Difficulties: second-order methods at least in part use the Hessian of the objective function.

- This is too expensive for high-dimensional problem.

QUIC reduces the computational cost of a coordinate descent update from the naive $O(p^2)$ to $O(p)$ complexity.

The Newton Direction

The second-order Taylor expansion of a function f around x^k is

$$f(x) \approx f(x^k) + \nabla f(x^k)^T (x - x^k) + \frac{1}{2} (x - x^k)^T \mathbf{H} (x - x^k)$$

where \mathbf{H} is the Hessian matrix

$$\mathbf{H} = \nabla^2 f(x^k) = \begin{bmatrix} \frac{\partial^2 f(x^k)}{\partial x_1^2} & \cdots & \frac{\partial^2 f(x^k)}{\partial x_1 \partial x_n} \\ \vdots & & \vdots \\ \frac{\partial^2 f(x^k)}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f(x^k)}{\partial x_n^2} \end{bmatrix}$$

The Newton direction is the solution of $\Delta x = x - x^k$ for the second-order expansion

$$D^k = -(\mathbf{H})^{-1} \nabla f(x^k)$$

Let

$$X^* = \arg \min_{X \succ 0} \{-\log \det X + \text{tr}(SX) + \|X\|_1\} = \arg \min_{X \succ 0} f(X)$$

Partition $f(X) = g(X) + h(X)$, where

$$g(X) = -\log \det X + \text{tr}(SX) \quad \text{and} \quad h(X) = \|X\|_1$$

Consider the second-order Taylor expansion of the $g(X)$

$$\begin{aligned} \bar{g}_{X_t}(\Delta) &\equiv g(X_t) + \text{vec}(\nabla g(X_t))^T \text{vec}(\Delta) + \frac{1}{2} \text{vec}(\Delta)^T \nabla^2 g(X_t) \text{vec}(\Delta) \\ &\propto \frac{1}{2} \left\| \left(\nabla^2 g(X_t) \right)^{\frac{1}{2}} \text{vec}(\Delta) + \left(\nabla^2 g(X_t) \right)^{-\frac{1}{2}} \text{vec}(\nabla g(X_t)) \right\|_2^2 \end{aligned}$$

The Newton direction D_t^* for the entire objective $f(X)$

$$D_t^* = \arg \min_{\Delta} \{ \bar{g}_{X_t}(\Delta) + h(X_t + \Delta) \}$$

It can be rewritten as a standard Lasso regression problem

$$\arg \min_{\Delta} \frac{1}{2} \left\| \mathbf{H}^{\frac{1}{2}} \text{vec}(\Delta) + \mathbf{H}^{-\frac{1}{2}} \mathbf{b} \right\|_2^2 + \|X_t + \Delta\|_1$$

where $\mathbf{H} = \nabla^2 g(X_t)$ and $\mathbf{b} = \text{vec}(\nabla g(X_t))$

The Gradient and Hessian of the log-likelihood $g(x)$ are

$$\nabla g(X) = S - X^{-1} \quad \text{and} \quad \nabla^2 g(X) = X^{-1} \otimes X^{-1}$$

The Gradient and Hessian of the log-likelihood $g(x)$ are

$$\nabla g(X) = S - X^{-1} \quad \text{and} \quad \nabla^2 g(X) = X^{-1} \otimes X^{-1}$$

Note $\text{vec}(\Delta)^T (X_t^{-1} \otimes X_t^{-1}) \text{vec}(\Delta) = \text{tr}(X_t^{-1} \Delta X_t^{-1} \Delta)$.

Let $W_t = X_t^{-1}$, the approximation of $g(x)$ can be rewritten as

$$\bar{g}_{X_t}(\Delta) = -\log \det X_t + \text{tr}(S X_t) + \text{tr}((S - W_t)^T \Delta) + \frac{1}{2} \text{tr}(W_t \Delta W_t \Delta)$$

The Gradient and Hessian of the log-likelihood $g(x)$ are

$$\nabla g(X) = S - X^{-1} \quad \text{and} \quad \nabla^2 g(X) = X^{-1} \otimes X^{-1}$$

Note $\text{vec}(\Delta)^T (X_t^{-1} \otimes X_t^{-1}) \text{vec}(\Delta) = \text{tr}(X_t^{-1} \Delta X_t^{-1} \Delta)$.

Let $W_t = X_t^{-1}$, the approximation of $g(x)$ can be rewritten as

$$\bar{g}_{X_t}(\Delta) = -\log \det X_t + \text{tr}(S X_t) + \text{tr}((S - W_t)^T \Delta) + \frac{1}{2} \text{tr}(W_t \Delta W_t \Delta)$$

Thus, the Newton direction can be solved by a Lasso problem, which requires $O(p^4)$.

Speedup Strategy

Use D to denote the current iterate approximating the Newton direction and D' for the updated direction.

- Consider the coordinate descent update for the variable X_{ij} , with $i < j$ that preserves symmetry: $D' = D + \mu (\mathbf{e}_i \mathbf{e}_j^T + \mathbf{e}_j \mathbf{e}_i^T)$
- The solution of the one-variable problem is

$$\arg \min_{\mu} \bar{g}_X (D + \mu (\mathbf{e}_i \mathbf{e}_j^T + \mathbf{e}_j \mathbf{e}_i^T)) + 2\lambda_{ij} |X_{ij} + D_{ij} + \mu|$$

- Omit the terms in $\bar{g}_X (D')$ not dependent on μ

$$\text{tr} \left((S - W)^T D' \right) \propto 2\mu (S_{ij} - W_{ij})$$

$$\text{tr} (W D' W D') = \text{tr} (W D W D) + 4\mu \mathbf{w}_i^T D \mathbf{w}_j + 2\mu^2 (W_{ij}^2 + W_{ji} W_{jj})$$

Speedup Strategy

So the one-variable problem is transformed into minimization of the following function of μ

$$\frac{1}{2} (W_{ij}^2 + W_{ii} W_{jj}) \mu^2 + (S_{ij} - W_{ij} + \mathbf{w}_i^T D \mathbf{w}_j) \mu + \lambda_{ij} |X_{ij} + D_{ij} + \mu|$$

Let $a = W_{ij}^2 + W_{ii} W_{jj}$, $b = S_{ij} - W_{ij} + \mathbf{w}_i^T D \mathbf{w}_j$, and $c = X_{ij} + D_{ij}$, the minimum is achieved for

$$\mu = -c + \mathcal{S}(c - b/a, \lambda_{ij}/a)$$

where

$$\mathcal{S}(z, r) = \text{sign}(z) \max\{|z| - r, 0\}$$

a and c are easy to compute.

The main computational cost is $\mathbf{w}_i^T D \mathbf{w}_j$.

Speedup Strategy

Calculating $\mathbf{w}_i^T D \mathbf{w}_j$ requires $O(p^2)$ times.

- Instead, we maintain matrix $U = DW$, and then compute $\mathbf{w}_i^T D \mathbf{w}_j$ by $\mathbf{w}_i^T \mathbf{u}_j$ using $O(p)$ flops.
- The maintain of $U = DW$ needs to update $2p$ elements

$$\mathbf{u}_{i.} \leftarrow \mathbf{u}_{i.} + \mu \mathbf{w}_j.$$

$$\mathbf{u}_{j.} \leftarrow \mathbf{u}_{j.} + \mu \mathbf{w}_i.$$

where $\mathbf{u}_{i.}$ refers to the i -th row vector of U .

Guarantee of Convergence

Theorem 1

Algorithm converges to the unique global optimum Y^* .

Theorem 2

The sequence $\{X_t\}$ generated by the QUIC algorithm converges quadratically to X^* , that is for some constant $\kappa > 0$,

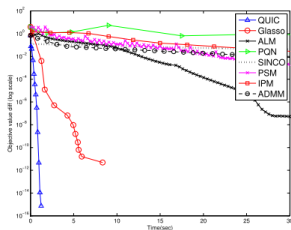
$$\lim_{t \rightarrow \infty} \frac{\|X_{t+1} - X^*\|_F}{\|X_t - X^*\|_F^2} = \kappa$$

Time Comparison using Synthetic Data Sets

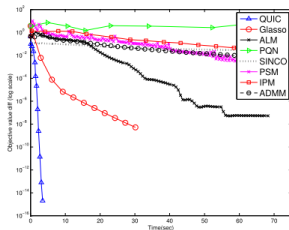
| Parameters | | | | Time (in seconds) | | | | | | | |
|------------|-------|-----------|------------|-------------------|-------|--------|-------|-------|-------|-------|------|
| pattern | p | λ | ϵ | QUIC | ALM | Glasso | PSM | IPM | SINCO | PQN | ADMM |
| chain | 1000 | 0.4 | 10^{-2} | < 1 | 19 | 9 | 16 | 86 | 120 | 110 | 62 |
| | | | 10^{-6} | 2 | 42 | 20 | 35 | 151 | 521 | 210 | 281 |
| chain | 4000 | 0.4 | 10^{-2} | 11 | 922 | 460 | 568 | 3458 | 5246 | 672 | 1028 |
| | | | 10^{-6} | 54 | 1734 | 1371 | 1258 | 5754 | * | 10525 | 2584 |
| chain | 10000 | 0.4 | 10^{-2} | 217 | 13820 | 10250 | 8450 | * | * | * | * |
| | | | 10^{-6} | 987 | 28190 | * | 19251 | * | * | * | * |
| random | 1000 | 0.12 | 10^{-2} | < 1 | 42 | 7 | 20 | 72 | 61 | 33 | 35 |
| | | | 10^{-6} | 1 | 28250 | 15 | 60 | 117 | 683 | 158 | 252 |
| | | 0.075 | 10^{-2} | 1 | 66 | 14 | 24 | 78 | 576 | 15 | 56 |
| | | | 10^{-6} | 7 | * | 43 | 92 | 146 | 4449 | 83 | * |
| random | 4000 | 0.08 | 10^{-2} | 23 | 1429 | 864 | 1479 | 4928 | 7375 | 2052 | 1025 |
| | | | 10^{-6} | 160 | * | 1743 | 4232 | 8097 | * | 4387 | * |
| | | 0.05 | 10^{-2} | 66 | * | 2514 | 2963 | 5621 | * | 2746 | * |
| | | | 10^{-6} | 479 | * | 5712 | 9541 | 13650 | * | 8718 | * |
| random | 10000 | 0.08 | 10^{-2} | 338 | 26270 | 14296 | * | * | * | * | * |
| | | | 10^{-6} | 1125 | * | * | * | * | * | * | * |
| | | 0.04 | 10^{-2} | 804 | * | * | * | * | * | * | * |
| | | | 10^{-6} | 2951 | * | * | * | * | * | * | * |

- Graphical Lasso (Glasso) is without block screen.

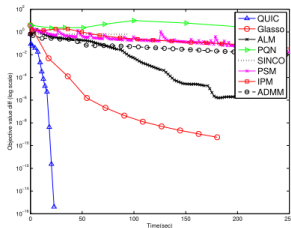
Time Comparison using Real Data Sets



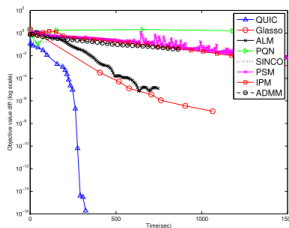
(a) Time taken on ER data set, $p = 692$, $\frac{\|X^*\|_0}{p^2} = 0.0222$



(b) Time taken on Arabidopsis data set, $p = 834$, $\frac{\|X^*\|_0}{p^2} = 0.0296$



(c) Time taken on Leukemia data set, $p = 1,255$, $\frac{\|X^*\|_0}{p^2} = 0.0221$



(d) Time taken on hereditarybc data set, $p = 1,869$, $\frac{\|X^*\|_0}{p^2} = 0.0198$

Penalized Likelihood Methods

One of the commonly used convex penalties is the l_1 penalty [Meinshausen et al., 2006].

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \{ \operatorname{tr}(\mathbf{S}\Theta) - \log |\Theta| + \lambda \|\Theta\|_1 \}$$

Solutions

- Interior-point optimization methods [Banerjee et al., 2008]
- Graphical Lasso (most popular) [Friedman et al., 2008]
- Alternating direction method of multipliers (ADMM) [Boyd et al., 2011]
- QUIC [Hsieh et al., 2014]

Trick and variants

- Block screen [Mazumder and Hastie, 2012]
- D -trace loss [Zhang and Zou, 2014]

Block Screen (Important)

Motivation

Suppose $\hat{\Theta}$ has the following sparse pattern

$$\hat{\Theta} = \begin{pmatrix} \hat{\Theta}_1 & 0 & \cdots & 0 \\ 0 & \hat{\Theta}_2 & 0 & \cdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \hat{\Theta}_{k(\lambda)} \end{pmatrix}$$

The log-likelihood problem can be decomposed to subproblems:

$$\hat{\Theta}_\ell = \arg \min_{\Theta_\ell} \left\{ -\log \det (\Theta_\ell) + \text{tr} (\mathbf{S}_\ell \Theta_\ell) + \lambda \sum_{ij} |(\Theta_\ell)_{ij}| \right\}$$

Can we learn such a sparse pattern?

Block Screen (Important)

Define the sparsity pattern of the solution $\hat{\Theta}^{(\lambda)}$ as follows

$$\mathcal{E}_{ij}^{(\lambda)} = \begin{cases} 1 & \text{if } \hat{\Theta}_{ij}^{(\lambda)} \neq 0, i \neq j \\ 0 & \text{otherwise} \end{cases}$$

Block Screen [Mazumder and Hastie, 2012]

Given λ , block screen performs a thresholding on the entries of the sample covariance matrix \mathbf{S} and obtain a graph edge skeleton

$$E_{ij}^{(\lambda)} = \begin{cases} 1 & \text{if } |\mathbf{S}_{ij}| > \lambda, i \neq j \\ 0 & \text{otherwise} \end{cases}$$

Block Screen (Important)

Proof: The log-likelihood problem

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \{ \operatorname{tr}(\mathbf{S}\Theta) - \log |\Theta| + \lambda \|\Theta\|_1 \}$$

The KKT conditions of optimality of log-likelihood problem is:

$$\left| \mathbf{s}_{ij} - \hat{\mathbf{W}}_{ij} \right| \leq \lambda, \quad \forall \hat{\Theta}_{ij} = 0$$

$$\hat{\mathbf{W}}_{ij} = \mathbf{s}_{ij} + \lambda, \quad \forall \hat{\Theta}_{ij} > 0$$

$$\hat{\mathbf{W}}_{ij} = \mathbf{s}_{ij} - \lambda, \quad \forall \hat{\Theta}_{ij} < 0$$

Where $\hat{\mathbf{W}} = (\hat{\Theta})^{-1}$

For the $E_{ij}^{(\lambda)} = 0$, set $\hat{\Theta}_{ij} = 0$, so $\hat{\mathbf{W}}_{ij} = 0$, the KKT condition

$\left| \mathbf{s}_{ij} - \hat{\mathbf{W}}_{ij} \right| = |\mathbf{s}_{ij}| \leq \lambda$ is satisfied.

Block Screen (Important)

- Block screen is not a specific algorithm for the penalized likelihood problem.
- It can be used as a wrapper around existing algorithms leads to enormous performance boosts.
- The optimization problem is completely separated into $k(\lambda)$ separated optimization sub-problems of the form. Help to solve high-dimensional problem.
- Easy to compute in a distributed manner.

Synthetic Examples

| K | p_1 / p | λ | Algorithm | Algorithm Timings (sec) | | Ratio Speedup factor | Time (sec) graph partition |
|---|------------|----------------|-----------|-------------------------|-------------------|----------------------------|----------------------------------|
| | | | | with screen | without screen | | |
| 2 | 200 / 400 | λ_I | GLASSO | 11.1 | 25.97 | 2.33 | 0.04 |
| | | | SMACS | 12.31 | 137.45 | 11.16 | |
| | | λ_{II} | GLASSO | 1.687 | 4.783 | 2.83 | 0.066 |
| | | | SMACS | 10.01 | 42.08 | 4.20 | |
| 2 | 500 / 1000 | λ_I | GLASSO | 305.24 | 735.39 | 2.40 | 0.247 |
| | | | SMACS | 175 | 2138* | 12.21 | |
| | | λ_{II} | GLASSO | 29.8 | 121.8 | 4.08 | 0.35 |
| | | | SMACS | 272.6 | 1247.1 | 4.57 | |
| 5 | 300 / 1500 | λ_I | GLASSO | 210.86 | 1439 | 6.82 | 0.18 |
| | | | SMACS | 63.22 | 6062* | 95.88 | |
| | | λ_{II} | GLASSO | 10.47 | 293.63 | 28.04 | 0.123 |
| | | | SMACS | 219.72 | 6061.6 | 27.58 | |
| 5 | 500 / 2500 | λ_I | GLASSO | 1386.9 | - | - | 0.71 |
| | | | SMACS | 493 | - | - | |
| | | λ_{II} | GLASSO | 17.79 | 963.92 | 54.18 | 0.018 |
| | | | SMACS | 354.81 | - | - | |
| 8 | 300 / 2400 | λ_I | GLASSO | 692.25 | - | - | 0.713 |
| | | | SMACS | 185.75 | - | - | |
| | | λ_{II} | GLASSO | 9.07 | 842.7 | 92.91 | 0.023 |
| | | | SMACS | 153.55 | - | - | |

Penalized Likelihood Methods

One of the commonly used convex penalties is the l_1 penalty [Meinshausen et al., 2006].

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \{ \operatorname{tr}(\mathbf{S}\Theta) - \log |\Theta| + \lambda \|\Theta\|_1 \}$$

Solutions

- Interior-point optimization methods [Banerjee et al., 2008]
- Graphical Lasso (most popular) [Friedman et al., 2008]
- Alternating direction method of multipliers (ADMM) [Boyd et al., 2011]
- QUIC [Hsieh et al., 2014]

Trick and variants

- Block screen [Mazumder and Hastie, 2012]
- *D*-trace loss [Zhang and Zou, 2014]

Variants for Penalized Likelihood Methods

Motivation

- Existing methods for precision matrix estimation do not always guarantee that the final estimator is positive definite.
- The log-determinant term is hard to be analyzed theoretically.
- The assumption of multivariate Gaussian is not appropriate for non-Gaussian data.

Variants for Penalized Likelihood Methods

Motivation

- Existing methods for precision matrix estimation do not always guarantee that the final estimator is positive definite.
- The log-determinant term is hard to be analyzed theoretically.
- The assumption of multivariate Gaussian is not appropriate for non-Gaussian data.

Zhang *et al.* (2014) proposed the D -trace loss to estimate precision matrix.

D-trace Loss Estimator

D-trace loss estimator:

$$\hat{\Theta} = \arg \min_{\Theta \succeq \epsilon I} \frac{1}{2} \left\langle \Theta^2, \hat{\Sigma} \right\rangle - \text{tr}(\Theta) + \lambda \|\Theta\|_{1,\text{off}}$$

$A \succeq B$ represents $A - B$ is positive semidefinite.

D-trace Loss Estimator

D-trace loss estimator:

$$\hat{\Theta} = \arg \min_{\Theta \succeq \epsilon I} \frac{1}{2} \left\langle \Theta^2, \hat{\Sigma} \right\rangle - \text{tr}(\Theta) + \lambda \|\Theta\|_{1,\text{off}}$$

$A \succeq B$ represents $A - B$ is positive semidefinite.

- **Condition 1:** It is a smooth convex function of Θ .
- **Condition 2:** The unique minimizer occurs at $\hat{\Sigma}^{-1}$.

ADMM for D -trace Loss Estimator

ADMM form

$$\arg \min_{\Theta_1 \succeq \epsilon I} \frac{1}{2} \left\langle \Theta^2, \hat{\Sigma} \right\rangle - \text{tr}(\Theta) + \lambda \|\Theta_0\|_{1,\text{off}} \quad \text{s.t. } [\Theta, \Theta] = [\Theta_0, \Theta_1]$$

The augmented Lagrangian

$$\begin{aligned} L(\Theta, \Theta_0, \Theta_1, \Lambda_0, \Lambda_1) = & \frac{1}{2} \left\langle \Theta^2, \hat{\Sigma} \right\rangle - \text{tr}(\Theta) + \lambda \|\Theta_0\|_{1,\text{off}} + h(\Theta_1 \succeq \epsilon I) \\ & + \langle \Lambda_0, \Theta - \Theta_0 \rangle + \langle \Lambda_1, \Theta - \Theta_1 \rangle \\ & + (\rho/2) \|\Theta - \Theta_0\|_F^2 + (\rho/2) \|\Theta - \Theta_1\|_F^2 \end{aligned}$$

where

$$h(\Theta_1 \succeq \epsilon I) = \begin{cases} 0, & \Theta_1 \succeq \epsilon I \\ \infty, & \text{otherwise} \end{cases}$$

ADMM for D -trace Loss Estimator

1. Update Θ

$$\Theta^{k+1} = \arg \min_{\Theta = \Theta^T} \frac{1}{2} \left\langle \Theta^2, \hat{\Sigma} + 2\rho I \right\rangle - \left\langle \Theta, I + \rho\Theta_0^k + \rho\Theta_1^k - \Lambda_0^k - \Lambda_1^k \right\rangle$$

2. Update Θ_0

$$\Theta_0^{k+1} = \arg \min_{\Theta_0 = \Theta_0^T} \frac{\rho}{2} \left\langle \Theta_0^2, I \right\rangle - \left\langle \Theta_0, \rho\Theta^{k+1} + \Lambda_0^k \right\rangle + \lambda \|\Theta_0\|_{1,\text{off}}$$

3. Update Θ_1

$$\Theta_1^{k+1} = \arg \min_{\Theta_1 \succeq \epsilon I} \frac{\rho}{2} \left\langle \Theta_1^2, I \right\rangle - \left\langle \Theta_1, \rho\Theta^{k+1} + \Lambda_1^k \right\rangle$$

$$4. [\Lambda_0^{k+1}, \Lambda_1^{k+1}] = [\Lambda_0^k, \Lambda_1^k] + \rho [\Theta^{k+1} - \Theta_0^{k+1}, \Theta^{k+1} - \Theta_1^{k+1}]$$

1 Introduction

2 Gaussian Graphical Models

3 Main Approaches

- Sparsity
- Penalized Likelihood Methods
- Column-by-column Estimation Methods
- CLIME

4 Variants

Column-by-column Estimation Methods

Column-by-column regression is another approach to estimate the precision matrix.

- **Main idea:** exploit the relationship between the conditional distribution of multivariate Gaussian and linear regressions.

Column-by-column Estimation Methods

Column-by-column regression is another approach to estimate the precision matrix.

- **Main idea:** exploit the relationship between the conditional distribution of multivariate Gaussian and linear regressions.

Problem

1. Precision matrix \Leftrightarrow conditional dependencies ✓
2. Precision matrix \Leftrightarrow conditional dependencies \Leftrightarrow linear regressions?

Gaussian Graphical Models

Suppose that observations $x_1, x_2, \dots, x_n \in \mathbb{R}^p$ are i.i.d. $N_p(\mu, \Sigma)$ where $\mu \in \mathbb{R}^p$ and Σ is a $p \times p$ positive definite matrix.

- Partition $X = (Z, Y)$ where $Z = (X_1, \dots, X_{p-1})$ and $Y = X_p$.
- Partitioned Σ and Θ as

$$\Sigma = \begin{pmatrix} \Sigma_{ZZ} & \sigma_{ZY} \\ \sigma_{ZY}^T & \sigma_{YY} \end{pmatrix}, \Theta = \begin{pmatrix} \Theta_{ZZ} & \theta_{ZY} \\ \theta_{ZY}^T & \theta_{YY} \end{pmatrix}$$

Gaussian Graphical Models

Suppose that observations $x_1, x_2, \dots, x_n \in \mathbb{R}^p$ are i.i.d. $N_p(\mu, \Sigma)$ where $\mu \in \mathbb{R}^p$ and Σ is a $p \times p$ positive definite matrix.

- Partition $X = (Z, Y)$ where $Z = (X_1, \dots, X_{p-1})$ and $Y = X_p$.
- Partitioned Σ and Θ as

$$\Sigma = \begin{pmatrix} \Sigma_{ZZ} & \sigma_{ZY} \\ \sigma_{ZY}^T & \sigma_{YY} \end{pmatrix}, \Theta = \begin{pmatrix} \Theta_{ZZ} & \theta_{ZY} \\ \theta_{ZY}^T & \theta_{YY} \end{pmatrix}$$

Because $\Theta = \Sigma^{-1}$ standards formulas for partitioned inverses given:

- $\theta_{YY} = (\sigma_{YY} - \sigma_{ZY}^T \Sigma_{ZZ}^{-1} \sigma_{ZY})^{-1} > 0$
- $\theta_{ZY} = -\theta_{YY} \cdot \Sigma_{ZZ}^{-1} \sigma_{ZY}$

Gaussian Graphical Models

If we perform multiple linear regression of Y on Z

$$\beta_Y = \arg \min_{\beta} \|Y - Z\beta_Y\|_2^2$$

Then

$$\beta_Y = (Z^T Z)^{-1} Z^T Y = \Sigma_{ZZ}^{-1} \sigma_{ZY} = -\theta_{ZY} / \theta_{YY}$$

Gaussian Graphical Models

If we perform multiple linear regression of Y on Z

$$\beta_Y = \arg \min_{\beta} \|Y - Z\beta_Y\|_2^2$$

Then

$$\beta_Y = (Z^T Z)^{-1} Z^T Y = \Sigma_{ZZ}^{-1} \sigma_{ZY} = -\theta_{ZY} / \theta_{YY}$$

We can learn about this dependence structure through multiple linear regression $((\beta_Y)_i = 0 \Leftrightarrow (\theta_{ZY})_i = 0)$.

Column-by-column Estimation Methods

Inspired by the linear regression model in and the fact that regression coefficient is sparse [Meinshausen et al., 2006]

$$\hat{\beta}_j = \operatorname{argmin}_{\alpha_j \in \mathbb{R}^{p-1}} \frac{1}{2n} \|Y_{*j} - Y_{*/j} \beta_j\|_2^2 + \lambda_j \|\beta_j\|_1$$

- Once $\hat{\beta}_j$ is obtained, we obtain the neighbourhood edges of node j by reading out the non-zero coefficients of $\hat{\beta}_j$.
- To estimate Θ

$$\hat{\theta}_{jj}^2 = \frac{1}{n} \|Y_{*j} - Y_{*/j} \theta_{jj}\|_2^2$$

and plug it into $\theta_{ij} = -\theta_{jj} \beta_{ij}$

Time Costs

Table 1. *Timings (seconds) for graphical lasso, Meinhausen–Buhlmann approximation, and COVSEL procedures*

| p | Problem type | (1) Graphical lasso | (2) Approximation | (3) COVSEL | Ratio of (3) to (1) |
|-----|--------------|---------------------|-------------------|------------|---------------------|
| 100 | Sparse | 0.014 | 0.007 | 34.7 | 2476.4 |
| 100 | Dense | 0.053 | 0.018 | 2.2 | 40.9 |
| 200 | Sparse | 0.050 | 0.027 | > 205.35 | > 4107 |
| 200 | Dense | 0.497 | 0.146 | 16.9 | 33.9 |
| 400 | Sparse | 1.23 | 0.193 | > 1616.7 | > 1314.3 |
| 400 | Dense | 6.2 | 0.752 | 313.0 | 50.5 |

- Fast and easy to implement and parallelize
- It is a solution to the quadratic approximation of the log-likelihood function [Banerjee et al., 2008].

Problem:

How to choose tuning parameters?

Tuning-Insensitive Graph Estimation and Regression (TIGER)
[Liu et al., 2017]

- based on the square-root-Lasso (SQRT-Lasso):

$$\hat{\beta}_j = \operatorname{argmin}_{\alpha_j \in \mathbb{R}^{p-1}} \frac{1}{2\sqrt{n}} \|Y_{*j} - Y_{* \setminus j} \beta_j\|_2^2 + \lambda_j \|\beta_j\|_1$$

which is asymptotically tuning-free.

- **tuning-insensitive**

1 Introduction

2 Gaussian Graphical Models

3 Main Approaches

- Sparsity
- Penalized Likelihood Methods
- Column-by-column Estimation Methods
- **CLIME**

4 Variants

Problem:

Most methods are designed for the Gaussian distribution.

For some classes of non-Gaussian distributions, the problem of estimating the graph can also be reduced to estimating the precision matrix (e.g., the nonparanormal distribution [Liu et al., 2009]).

Aim:

Estimating the precision matrix for both sparse and nonsparse matrices, without restricting to a specific sparsity pattern.

CLIME: a constrained l_1 minimization approach to sparse precision matrix estimation [Cai et al., 2011].

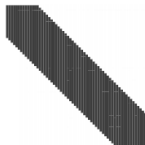
$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \|\Theta\|_1 \text{ s.t. } \left\| \hat{\Sigma}\Theta - I \right\|_{\infty} \leq \delta_j, \Theta \in R^{p \times p}$$

CLIME is equivalent to solving the p optimization sub-problems:

$$\hat{\Theta}_{*j} = \underset{\Theta_{*j}}{\operatorname{argmin}} \|\Theta_{*j}\|_1 \text{ s.t. } \left\| \hat{\Sigma}\Theta_{*j} - \mathbf{e}_j \right\|_{\infty} \leq \delta_j, \text{ for } j = 1, \dots, p$$

Simulated Experiments

Model 1



(a) Truth



(b) CLIME

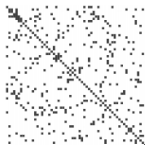


(c) Glasso

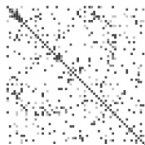


(d) SCAD

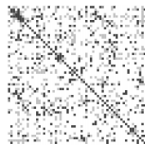
Model 2



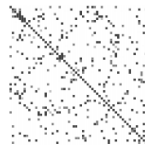
(e) Truth



(f) CLIME



(g) Glasso



(h) SCAD

Real Data Experiments

Table 4. Comparison of average (SE) pCR classification errors over 100 replications. Glasso, Adaptive lasso, and SCAD results are taken from Fan, Feng, and Wu (2009, table 2)

| 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) |

Outline

- 1 Introduction
- 2 Gaussian Graphical Models
- 3 Main Approaches
- 4 Variants**
 - Graphical Model with Hubs
 - Joint Graphical Lasso
 - Differential Network Estimation

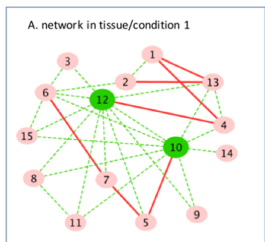
- 1 Introduction
- 2 Gaussian Graphical Models
- 3 Main Approaches
- 4 Variants**
 - **Graphical Model with Hubs**
 - Joint Graphical Lasso
 - Differential Network Estimation

Graphical Model with Hubs

Motivation

In many applications, there are a few **hub nodes** that are densely-connected to many other nodes.

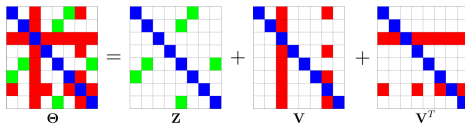
- such as critical genes for organisms to complete their life cycle.



Hub Graphical Lasso was proposed to estimate networks with hub nodes [Tan et al., 2014].

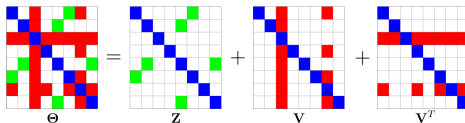
Hub Graphical Lasso

Decomposition of a symmetric matrix Θ into $\mathbf{Z} + \mathbf{V} + \mathbf{V}^T$, where \mathbf{Z} is sparse, and most columns of \mathbf{V} are entirely zeros.



Hub Graphical Lasso

Decomposition of a symmetric matrix Θ into $\mathbf{Z} + \mathbf{V} + \mathbf{V}^T$, where \mathbf{Z} is sparse, and most columns of \mathbf{V} are entirely zeros.



Hub Graphical Lasso [Tan et al., 2014]

$$\begin{aligned} & \text{minimize}_{\Theta \in \mathcal{S}, \mathbf{V}, \mathbf{Z}} \left\{ \ell(\mathbf{X}, \Theta) + \lambda_1 \|\mathbf{Z} - \text{diag}(\mathbf{Z})\|_1 + \lambda_2 \|\mathbf{V} - \text{diag}(\mathbf{V})\|_1 \right. \\ & \quad \left. + \lambda_3 \sum_{j=1}^p \|(\mathbf{V} - \text{diag}(\mathbf{V}))_{*j}\|_2 \right\} \\ & \text{subject to} \quad \Theta = \mathbf{V} + \mathbf{V}^T + \mathbf{Z} \end{aligned}$$

where $\ell(\mathbf{X}, \Theta)$ is the negative log-likelihood.

ADMM for Hub Graphical Lasso

ADMM form

Let $\mathbf{B} = (\Theta, \mathbf{V}, \mathbf{Z})$, $\tilde{\mathbf{B}} = (\tilde{\Theta}, \tilde{\mathbf{V}}, \tilde{\mathbf{Z}})$

$$f(\mathbf{B}) = \ell(\mathbf{X}, \Theta) + \lambda_1 \|\mathbf{Z} - \text{diag}(\mathbf{Z})\|_1 + \lambda_2 \|\mathbf{V} - \text{diag}(\mathbf{V})\|_1 \\ + \lambda_3 \sum_{j=1}^p \|(\mathbf{V} - \text{diag}(\mathbf{V}))_{*j}\|_2$$

$$g(\tilde{\mathbf{B}}) = \begin{cases} 0, & \text{if } \tilde{\Theta} = \tilde{\mathbf{V}} + \tilde{\mathbf{V}}^T + \tilde{\mathbf{Z}} \\ \infty, & \text{otherwise} \end{cases}$$

The ADMM form:

$$\underset{\mathbf{B}, \tilde{\mathbf{B}}}{\text{minimize}} \{f(\mathbf{B}) + g(\tilde{\mathbf{B}})\} \quad \text{s.t. } \mathbf{B} = \tilde{\mathbf{B}}$$

ADMM for Hub Graphical Lasso

The scaled augmented Lagrangian

$$\begin{aligned} L(\mathbf{B}, \tilde{\mathbf{B}}, \mathbf{W}) = & \ell(\mathbf{X}, \Theta) + \lambda_1 \|\mathbf{Z} - \text{diag}(\mathbf{Z})\|_1 + \lambda_2 \|\mathbf{V} - \text{diag}(\mathbf{V})\|_1 \\ & + \lambda_3 \sum_{j=1}^p \|(\mathbf{V} - \text{diag}(\mathbf{V}))_{*j}\|_2 + g(\tilde{\mathbf{B}}) + \frac{\rho}{2} \|\mathbf{B} - \tilde{\mathbf{B}} + \mathbf{W}\|_F^2 \end{aligned}$$

ADMM:

1. Update \mathbf{B} (include $\Theta, \mathbf{V}, \mathbf{Z}$)
2. Update $\tilde{\mathbf{B}}$ (include $\tilde{\Theta}, \tilde{\mathbf{V}}, \tilde{\mathbf{Z}}$)
3. Update \mathbf{W}

- 1 Introduction
- 2 Gaussian Graphical Models
- 3 Main Approaches
- 4 Variants**
 - Graphical Model with Hubs
 - Joint Graphical Lasso**
 - Differential Network Estimation

Joint Graphical Lasso (JGL)

Motivation

The standard formulation for estimating a precision matrix assumes that each observation is drawn from the same distribution.

However, in many datasets the observations may correspond to several distinct classes.

Joint Graphical Lasso (JGL)

Motivation

The standard formulation for estimating a precision matrix assumes that each observation is drawn from the same distribution.

However, in many datasets the observations may correspond to several distinct classes.

Consider the gene expression data of a set of cancer samples and normal samples, respectively.

- **Solution 1:** estimating graphical model using all sample
 - Ignore the heterogeneity, inappropriate
- **Solution 2:** estimating separate graphical models for the cancer and normal samples
 - This strategy does not exploit the similarity between the true graphical models.

Joint Graphical Lasso

Assume a heterogeneous data with p variables and K classes.

- The k -th class contains n_k observations $(\mathbf{x}_1^{(k)}, \dots, \mathbf{x}_{n_k}^{(k)})$, where each $\mathbf{x}_i^{(k)} = (x_{i,1}^{(k)}, \dots, x_{i,p}^{(k)})$ is a p -dim row vector.
- $\mathbf{S}^{(k)}$ is the sample covariance matrix of the k -th class.
- $\Theta^{(k)}$ is the inverse covariance matrix of the k -th class.

Joint Graphical Lasso

Assume a heterogeneous data with p variables and K classes.

- The k -th class contains n_k observations $(\mathbf{x}_1^{(k)}, \dots, \mathbf{x}_{n_k}^{(k)})$, where each $\mathbf{x}_i^{(k)} = (x_{i,1}^{(k)}, \dots, x_{i,p}^{(k)})$ is a p -dim row vector.
- $\mathbf{S}^{(k)}$ is the sample covariance matrix of the k -th class.
- $\Theta^{(k)}$ is the inverse covariance matrix of the k -th class.

General formulation for JGL [Danaher et al., 2014]

$$\text{maximize}_{\{\Theta\}} \left(\sum_{k=1}^K n_k [\log \{\det (\Theta^{(k)})\} - \text{tr} (\mathbf{S}^{(k)} \Theta^{(k)})] - P(\{\Theta\}) \right)$$

$$P(\{\Theta\}) = \lambda_1 \sum_k \sum_{i \neq j} |\theta_{ij}^{(k)}| + \hat{P}(\{\Theta\})$$

where \hat{P} is a convex function.

Joint Graphical Lasso

Fused graphical Lasso [Danaher et al., 2014]

The fused graphical lasso (FGL) is the solution to joint graphical model with the penalty:

$$P(\{\Theta\}) = \lambda_1 \sum_{k=1}^K \sum_{i \neq j} \left| \theta_{ij}^{(k)} \right| + \lambda_2 \sum_{k < k'} \sum_{i,j} \left| \theta_{ij}^{(k)} - \theta_{ij}^{(k')} \right|$$

- encourages similar edge values.

Joint Graphical Lasso

Group graphical Lasso (GGL) [Danaher et al., 2014]

GGL is the solution to joint graphical model with the penalty:

$$P(\{\Theta\}) = \lambda_1 \sum_{k=1}^K \sum_{i \neq j} \left| \theta_{ij}^{(k)} \right| + \lambda_2 \sum_{i \neq j} \left(\sum_{k=1}^K \left(\theta_{ij}^{(k)} \right)^2 \right)^{1/2}$$

- encourages a similar pattern of sparsity (i.e. there will be a tendency for the 0s in the K estimated precision matrices to occur in the same places)

ADMM for JGL

ADMM form

$$\begin{aligned} \underset{\{\Theta\}, \{\mathbf{Z}\}}{\text{minimize}} \quad & \left(- \sum_{k=1}^K n_k [\log \{\det (\Theta^{(k)})\} - \text{tr} (\mathbf{S}^{(k)}(k))] + P(\{\mathbf{Z}\}) \right) \\ \text{s.t. } & \mathbf{Z}^{(k)} = \Theta^{(k)} \text{ for } k = 1, \dots, K \end{aligned}$$

The scaled augmented Lagrangian

$$\begin{aligned} L_p(\{\Theta\}, \{\mathbf{Z}\}, \{\mathbf{U}\}) = & - \sum_{k=1}^K n_k [\log \det (\Theta^{(k)}) - \text{tr} (\mathbf{S}^{(k)}(k))] + P(\{\mathbf{Z}\}) \\ & + \frac{\rho}{2} \sum_{k=1}^K \|\Theta^{(k)} - \mathbf{Z}^{(k)} + \mathbf{U}^{(k)}\|_F^2 - \frac{\rho}{2} \sum_{k=1}^K \|\mathbf{U}^{(k)}\|_F^2 \end{aligned}$$

ADMM for JGL

ADMM for JGL

1. Update $\Theta^{(k)}$, $k = 1, \dots, K$

$$\begin{aligned}\hat{\Theta}^{(k)} = \operatorname{argmin}(& -n_k [\log \{\det(\Theta^{(k)})\} - \operatorname{tr}(\mathcal{S}^{(k)}\Theta^{(k)})] \\ & + (\rho/2) \left\| \Theta^{(k)} - Z_{i-1}^{(k)} + U_{i-1}^{(K)} \right\|_F^2)\end{aligned}$$

- Eigenvalue problem

2. Update $\mathbf{Z}^{(k)}, k = 1, \dots, K$

$$\hat{\mathbf{Z}}^{(k)} = \operatorname{argmin} \sum_{k=1}^K \left\| \mathbf{Z}^{(k)} - \left(\Theta_i^{(k)} + \mathbf{U}_{i-1}^{(k)} \right) \right\|_F^2 + P(\mathbf{Z})$$

- Fused Graphical Lasso

- The subproblem is a group lasso problem \rightarrow explicit solution [Friedman et al., 2010].

- Fused Graphical Lasso

- The subproblem is a fused lasso problem \rightarrow for each (i, j) , costs is $O\{K \log(K)\}$ [Hoeffling, 2010].

3. Update $\mathbf{U} \left\{ \mathbf{U}_{(i)} \right\} \leftarrow \left\{ \mathbf{U}_{(i-1)} \right\} + \left(\left\{ \Theta_{(i)} \right\} - \left\{ \mathbf{Z}_{(i)} \right\} \right)$

- 1 Introduction
- 2 Gaussian Graphical Models
- 3 Main Approaches
- 4 Variants**
 - Graphical Model with Hubs
 - Joint Graphical Lasso
 - Differential Network Estimation**

Differential Network Estimation

Motivation

Still consider the gene expression data of a set of cancer tissue samples and a set of normal tissue samples.

- A complete understanding of the molecular basis of cancer will require characterization of the differential network.
- Direct estimation of differential networks will help us to advance the understanding of cancer development.

Differential Network Estimation

Motivation

Still consider the gene expression data of a set of cancer tissue samples and a set of normal tissue samples.

- A complete understanding of the molecular basis of cancer will require characterization of the differential network.
- Direct estimation of differential networks will help us to advance the understanding of cancer development.

Zhao *et al.* applied precision matrix estimation to differential network analysis [Zhao et al., 2014].

Differential Network Estimation

Suppose we have two classes of data X, Y

- S^X, S^Y are the sample covariance matrix of \mathbf{X} and \mathbf{Y} .
- Θ^X, Θ^Y are the inverse covariance matrix of \mathbf{X} and \mathbf{Y} .

The differential network can be represented by $\Delta = \Theta^Y - \Theta^X$.

Differential Network Estimation

Suppose we have two classes of data X, Y

- S^X, S^Y are the sample covariance matrix of \mathbf{X} and \mathbf{Y} .
- Θ^X, Θ^Y are the inverse covariance matrix of \mathbf{X} and \mathbf{Y} .

The differential network can be represented by $\Delta = \Theta^Y - \Theta^X$.

Zhao *et al.* (2014) applied CLIME to estimate the Δ .

- $\Delta = \Theta^Y - \Theta^X$
- $S^X \Delta S^Y - (S^X - S^Y) = 0$

Thus, Δ can be estimated by

$$\min \|\Delta\|_1 \text{ s.t. } |S^X \Delta S^Y - (S^X - S^Y)|_\infty \leq \lambda_n$$

References I



Banerjee, O., Ghaoui, L. E., and d'Aspremont, A. (2008).

Model selection through sparse maximum likelihood estimation for multivariate gaussian or binary data.
Journal of Machine learning research, 9(Mar):485–516.



Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J., et al. (2011).

Distributed optimization and statistical learning via the alternating direction method of multipliers.
Foundations and Trends® in Machine learning, 3(1):1–122.



Cai, T., Liu, W., and Luo, X. (2011).

A constrained l_1 minimization approach to sparse precision matrix estimation.
Journal of the American Statistical Association, 106(494):594–607.



Danaher, P., Wang, P., and Witten, D. M. (2014).

The joint graphical lasso for inverse covariance estimation across multiple classes.
Journal of the Royal Statistical Society: Series B (Statistical Methodology), 76(2):373–397.



Friedman, J., Hastie, T., and Tibshirani, R. (2008).

Sparse inverse covariance estimation with the graphical lasso.
Biostatistics, 9(3):432–441.



Friedman, J., Hastie, T., and Tibshirani, R. (2010).

A note on the group lasso and a sparse group lasso.
arXiv preprint arXiv:1001.0736.



Hoefling, H. (2010).

A path algorithm for the fused lasso signal approximator.
Journal of Computational and Graphical Statistics, 19(4):984–1006.

References II



Hsieh, C.-J., Sustik, M. A., Dhillon, I. S., and Ravikumar, P. (2014).
Quic: quadratic approximation for sparse inverse covariance estimation.
The Journal of Machine Learning Research, 15(1):2911–2947.



Liu, H., Lafferty, J., and Wasserman, L. (2009).
The nonparanormal: Semiparametric estimation of high dimensional undirected graphs.
Journal of Machine Learning Research, 10(Oct):2295–2328.



Liu, H., Wang, L., et al. (2017).
Tiger: A tuning-insensitive approach for optimally estimating gaussian graphical models.
Electronic Journal of Statistics, 11(1):241–294.



Mazumder, R. and Hastie, T. (2012).
Exact covariance thresholding into connected components for large-scale graphical lasso.
Journal of Machine Learning Research, 13(Mar):781–794.



Meinshausen, N., Bühlmann, P., et al. (2006).
High-dimensional graphs and variable selection with the lasso.
The annals of statistics, 34(3):1436–1462.



Tan, K. M., London, P., Mohan, K., Lee, S.-I., Fazel, M., and Witten, D. (2014).
Learning graphical models with hubs.
The Journal of Machine Learning Research, 15(1):3297–3331.



Zhang, T. and Zou, H. (2014).
Sparse precision matrix estimation via lasso penalized d-trace loss.
Biometrika, 101(1):103–120.

References III



Zhao, S. D., Cai, T. T., and Li, H. (2014).
Direct estimation of differential networks.
Biometrika, 101(2):253–268.