
Statistical Learning and Data Science

Xinwei Deng

xdeng@vt.edu

Department of Statistics

Gaussian Graphical Model Estimation

- Graphical Model via Covariance Matrix Estimation
 - Modified Cholesky decomposition approach.
 - Structured graphical model estimation.
 - Nonparametric approach for graphical model estimation.
 - Multiple graphical models.

Why Estimating Σ^{-1} : Gaussian Graphical Model

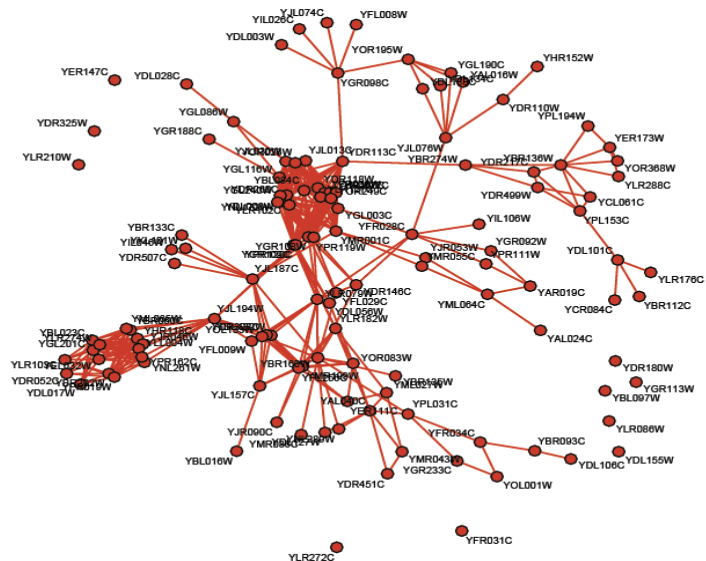
[illegible]

$$\mathbf{\Omega} = \begin{pmatrix} c_{11} & c_{12} & \dots & c_{1q} \\ c_{21} & c_{22} & \dots & c_{2q} \\ \vdots & \ddots & \vdots & \vdots \\ c_{q1} & c_{q2} & \dots & c_{qq} \end{pmatrix}$$

concentration matrix $\Sigma^{-1} \equiv \Omega = (c_{ij})$

Data $\mathbf{Y} = (y_{ij})$ is an $n \times q$ matrix

- describe the **conditional dependency** among variables: if $c_{ij}=0$ zero, then variables i and j are conditionally independent given the other variables.



The Likelihood Viewpoint

- ▶ The multi-response model is simplified as

$$\mathbf{y} = \mu + \epsilon, \epsilon \sim N(0, \Sigma).$$

- ▶ With data $\mathbf{y}_i, i = 1, \dots, n$, it is equivalent to

$$\mathbf{y}_i \sim N(\mu, \Sigma), i = 1, \dots, n.$$

- ▶ To estimate μ and Σ , the log-likelihood function becomes

$$\begin{aligned} L(\mu, \Sigma) &= \log \left\{ \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} |\Sigma^{-1}|^{1/2} \exp\left(-\frac{(\mathbf{y}_i - \mu)' \Sigma^{-1} (\mathbf{y}_i - \mu)}{2}\right) \right\} \\ &\propto \frac{n}{2} \log |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i - \mu)' \Sigma^{-1} (\mathbf{y}_i - \mu) \\ &\propto \log |\Sigma^{-1}| - \text{tr}[\Sigma^{-1} \mathbf{S}], \end{aligned}$$

where

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \mu)(\mathbf{y}_i - \mu)^T.$$

Estimating Σ^{-1}

- Based on

$L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \frac{n}{2} \log |\boldsymbol{\Sigma}^{-1}| - \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \boldsymbol{\mu})$, it is easy to obtain the estimate of $\boldsymbol{\mu}$ as

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i.$$

- To estimate $\boldsymbol{\Sigma}$, take the first derivative with respect to $\boldsymbol{\Omega} = \boldsymbol{\Sigma}^{-1}$

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\Omega}} \{ -\log |\boldsymbol{\Omega}| + \text{tr}[\boldsymbol{\Omega} \mathbf{S}] \} \\ = -\boldsymbol{\Sigma} + \mathbf{S} = 0 \end{aligned}$$

Therefore, the estimate of $\boldsymbol{\Sigma}$ is

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})^T.$$

Graphical Lasso Approach

- ▶ Denote $\Sigma^{-1} = (c_{ij})$. Note that $c_{ij} = 0$ implies that Y_i and Y_j are independent conditional on the other variables.
- ▶ Graphical Lasso (Yuan and Lin, 2007; Friedman et al., 2009) is to estimate $\Omega = \Sigma^{-1}$ by

$$\min_{\Omega} -\log |\Sigma^{-1}| + \text{tr}[\Sigma^{-1} \mathbf{S}] + \lambda \|\Sigma^{-1}\|_1,$$

where λ is a tuning parameter. Here $\|\cdot\|_1$ is defined as $\|\Sigma^{-1}\|_1 = \sum_{i \neq j} |c_{ij}|$.

- ▶ Graphical Lasso can be obtained by recursively solving and updating the lasso regression (R package *glasso*), or coordinate descent method (Witten et al., 2011).

Modified Cholesky Decomposition (MCD)

- ▶ From matrix algebra, various matrix decompositions are used for computing matrix inversion, determinant, etc.
- ▶ The Σ^{-1} can be obtained from the modified Cholesky decomposition.
 - ▶ Advantage: the estimation is transformed into a series of regressions.
- ▶ Let us define

$$\begin{aligned} X_j &= \sum_{t=1}^{j-1} a_{jt} X_t + \epsilon_j \\ &= \mathbf{Z}_j^T \mathbf{a}_j + \epsilon_j, \end{aligned}$$

where $\mathbf{Z}_j = (X_1, \dots, X_{j-1})'$, and $\mathbf{a}_j = (a_{j1}, \dots, a_{j,j-1})'$ are corresponding regression coefficients. Then the vector \mathbf{a}_j can be obtained from

$$\mathbf{a}_j = (\text{Cov}(\mathbf{Z}_j, \mathbf{Z}_j))^{-1} \text{Cov}(X_j, \mathbf{Z}_j).$$

MCD Approach for Σ^{-1}

- ▶ The variance of ϵ_j is defined by

$$d_j^2 = \text{Var}(X_j - \hat{X}_j) = \text{Var}(\epsilon_j).$$

- ▶ Then we can compose a lower triangle matrix \mathbf{A} such that

$$\mathbf{A} = \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ a_{21} & 0 & 0 & \dots & 0 \\ a_{31} & a_{32} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{p1} & a_{p2} & \dots & a_{p,p-1} & 0 \end{pmatrix},$$

which contains all regression coefficients a_{ij} .

- ▶ Let $\mathbf{D} = \text{diag}(d_1^2, \dots, d_p^2)$ as a diagonal matrix. The modified Cholesky decomposition implies that

$$\begin{aligned} \Sigma^{-1} &= (\mathbf{I} - \mathbf{A})' \mathbf{D}^{-1} (\mathbf{I} - \mathbf{A}) \\ &= \mathbf{T}' \mathbf{D}^{-1} \mathbf{T}. \end{aligned} \tag{1}$$

Computational Flexibility

- ▶ The straightforward estimate $\hat{\mathbf{A}}$ and $\hat{\mathbf{D}}$ can be obtained from the least squares estimates of the coefficients a_{ij} and the corresponding residual variances d_{ij} .
- ▶ The above formulation can be viewed as a multiple response regression under a given order of random variables.

$$x_1 = \epsilon_1, \quad \epsilon_1 \sim N(0, d_1^2);$$

$$x_2 = a_{21}x_1 + \epsilon_2, \quad \epsilon_2 \sim N(0, d_2^2);$$

$$x_3 = a_{31}x_1 + a_{32}x_2 + \epsilon_3, \quad \epsilon_3 \sim N(0, d_3^2);$$

$$\vdots$$

$$x_p = a_{p1}x_1 + a_{p2}x_2 + \cdots + a_{p,p-1}x_{p-1} + \epsilon_p, \quad \epsilon_p \sim N(0, d_p^2).$$

- ▶ The penalized regression can be applied to encourage the sparsity on \mathbf{T} , consequently leading to the sparsity of $\mathbf{\Sigma}^{-1}$.

Some Remarks

- ▶ The assumption of modified Cholesky decomposition is that ϵ_j is independent.
- ▶ It can be used to get a banded estimate of Σ^{-1} .
- ▶ Can such an approach be extended for the estimation of sparse Σ ?
- ▶ Recall that the negative log-likelihood function of Σ given the sample, $\mathbf{x}_1, \dots, \mathbf{x}_n$, is proportional to

$$L_n(\Sigma) = -\log |\Sigma^{-1}| + \text{tr}[\Sigma^{-1} \mathbf{S}],$$

where $\mathbf{S} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' / n$ be the sample covariance matrix.

MCD for Estimating Σ

- ▶ The modified Cholesky decomposition (MCD) can be used to estimate Σ from a latent variable regression model.
- ▶ Define $\epsilon = (\epsilon_1, \dots, \epsilon_p)'$. Note that $\epsilon = \mathbf{X} - \hat{\mathbf{X}}$, and hence,

$$\begin{aligned} \text{Var}(\epsilon) &= \text{Var}(\mathbf{X} - \hat{\mathbf{X}}) = \text{Var}(\mathbf{T}\mathbf{X}) \\ &= \mathbf{D} = \mathbf{T}\Sigma\mathbf{T}'. \end{aligned}$$

Therefore, we can write $\mathbf{X} = \mathbf{L}\epsilon$, which leads to

$$\begin{aligned} \text{Var}(\mathbf{X}) &= \text{Var}(\mathbf{L}\epsilon) \\ &= \Sigma = \mathbf{L}\mathbf{D}\mathbf{L}'. \end{aligned} \tag{2}$$

MCD for Estimating Σ (Con't)

- ▶ It implies a modified Cholesky decomposition for Σ itself.
- ▶ we can get an interpretation of latent variable regressions.
 - ▶ The $\epsilon = (\epsilon_1, \dots, \epsilon_p)'$ is unobserved.
 - ▶ The \mathbf{L} is lower triangular, and each variable X_j regresses on the previous latent variables $\epsilon_1, \dots, \epsilon_{j-1}$.
- ▶ It gives a sequence of regression. For $j = 2, \dots, p$, we get

$$\begin{aligned} X_j &= \sum_{t=1}^{j-1} l_{jt} \epsilon_t + \epsilon_j \\ &= \mathbf{U}_j^T \mathbf{l}_j + \epsilon_j, \end{aligned}$$

where $\mathbf{U}_j = (\epsilon_1, \dots, \epsilon_{j-1})'$, and $\mathbf{l}_j = (l_{j1}, \dots, l_{j,j-1})'$ is the vector of regression coefficients.

Sparse Estimate for Σ

- ▶ With data matrix $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)'$, define its j th column to be $\mathbf{x}^{(j)}$. The residual for $\mathbf{x}^{(j)}$ is defined as $\mathbf{e}^{(j)}$, and $\mathbf{Z}^{(j)} = (\mathbf{e}_1, \dots, \mathbf{e}_{j-1})$.
- ▶ Set $\mathbf{e}^{(1)} = \mathbf{x}^{(1)}$. we obtain the estimate of \mathbf{l}_j by least squares

$$\hat{\mathbf{l}}_j = \arg \min_{\mathbf{l}_j} \|\mathbf{x}^{(j)} - \mathbf{Z}^{(j)}\mathbf{l}_j\|_2^2, \quad j = 2, \dots, p,$$

where $\mathbf{e}^{(j)} = \mathbf{x}^{(j)} - \mathbf{Z}^{(j)}\mathbf{l}_j$ constructing the residuals used in $\mathbf{Z}^{(j+1)}$. Get $\hat{\mathbf{D}} = \text{diag}(\hat{d}_1, \dots, \hat{d}_p)$ where $\hat{d}_j = \text{var}(\mathbf{e}^{(j)})$ is sample variance of residuals.

- ▶ Encourage the sparsity on $\hat{\mathbf{L}}$ by using the penalized regression,

$$\hat{\mathbf{l}}_j = \arg \min_{\mathbf{l}_j} \|\mathbf{x}^{(j)} - \mathbf{Z}^{(j)}\mathbf{l}_j\|_2^2 + \lambda \|\mathbf{l}_j\|_1, \quad j = 2, \dots, p.$$

- ▶ Then $\hat{\Sigma} = \hat{\mathbf{L}}\hat{\mathbf{D}}\hat{\mathbf{L}}'$ will be a sparse covariance matrix estimate.

Partial Correlation Approach

- ▶ Assume that $\mathbf{X} \sim N(\mathbf{0}, \Sigma)$ where $\mathbf{X} = (X_1, \dots, X_p)$. Denote $\Sigma^{-1} = (c_{ij})$.
- ▶ **Lemma.** Suppose X_i is expressed as

$$X_i = \sum_{j \neq i} \beta_{ij} X_j + \epsilon_i,$$

Then we can have

$$\beta_{ij} = -c_{ij}/c_{ii}, \text{ and } \text{var}(\epsilon_i) = 1/c_{ii}.$$

Moreover, $\text{cov}(\epsilon_i, \epsilon_j) = c_{ij}/(c_{ii}c_{jj})$.

- ▶ The $\text{cor}(\epsilon_i, \epsilon_j) = \sqrt{c_{ij}/(c_{ii}c_{jj})}$ is the **partial correlation** between X_i and X_j .
- ▶ Now one can consider to estimate $\Sigma^{-1} = (c_{ij})$ through a joint regression effort.
- ▶ Remark: connection to the modified Cholesky decomposition approach?

Partial Correlation Approach: Estimation

- ▶ Let $\mathbf{X}_{(k)}$ be the k th column of the data matrix $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)'$.
- ▶ The joint regression objective function is rewritten as

$$\begin{aligned} \sum_{i=1}^p w_i \|\mathbf{X}_{(i)} - \sum_{j \neq i} \beta_{ij} \mathbf{X}_{(j)}\|^2 \\ = \sum_{i=1}^p \frac{1}{c_{ii}} \|c_{ii} \mathbf{X}_{(i)} + \sum_{j \neq i} c_{ij} \mathbf{X}_{(j)}\|^2, \end{aligned}$$

where $w_i = 1/\text{var}(\epsilon) = c_{ii}$.

- ▶ *Individual penalized regression for estimation VS Joint penalized regressions for estimation*
- ▶ Consider regularization on c_{ij} can lead to the penalized loss function as

$$\min \sum_{i=1}^p \frac{1}{c_{ii}} \|c_{ii} \mathbf{X}_{(i)} + \sum_{j \neq i} c_{ij} \mathbf{X}_{(j)}\|^2 + \lambda \sum_{i \neq j} |c_{ij}|.$$

CLIME Approach

- ▶ The partial correlation approach, Glasso method, MCD approaches are mostly likelihood based methods.
- ▶ CLIME (constrained l_1 -minimization for inverse matrix estimation, Cai et al., 2011) is to estimate $\Sigma^{-1} = (c_{ij})$ by solving

$$\begin{aligned} \min_{\Sigma^{-1}} \sum_{i,j} |c_{ij}| \\ \text{s.t. } |\mathbf{S}\Sigma^{-1} - \mathbf{I}|_{\infty} \leq \lambda, \end{aligned}$$

- ▶ The motivation is from the score function of log-likelihood $-\log |\Sigma^{-1}| + \text{tr}[\Sigma^{-1}\mathbf{S}]$.
- ▶ Such an estimate may not be symmetric, also may not be positive definite.
- ▶ The advantages are relied on few assumptions and computation is convenient.

Computation for CLIME

Approach

- ▶ Note that it is a convex optimization. We can further decompose it into p minimization problems with respect to vectors.
- ▶ Denote by e_j the j th column of \mathbf{I} and β_j the j th column of Σ^{-1} .
- ▶ Then solve a series of minimizations for $j = 1, \dots, p$,

$$\begin{aligned} \min_{\beta_j} & \|\beta_j\|_1 \\ \text{s.t.} & |\mathbf{S}\beta_j - e_j|_\infty \leq \lambda. \end{aligned}$$

- ▶ The above optimization can be solved by linear programming through the linear relaxation technique.

Computation for CLIME Approach (Con't)

- ▶ Because of the convex optimization, it can be decomposed into p minimization problems with respect to vectors.
- ▶ Denote \mathbf{r}'_j to be the j th row of \mathbf{S} . The objective is

$$\min_{\boldsymbol{\beta}_j} \|\boldsymbol{\beta}_j\|_1 \text{ s.t. } |\mathbf{R}\boldsymbol{\beta}_j - \mathbf{e}_j|_\infty \leq \lambda,$$

- ▶ We consider a relaxation of the optimization for $\boldsymbol{\beta}_j$ as

$$\begin{aligned} \min \quad & \sum_{k=j}^p t_k \\ \text{s.t.} \quad & -t_k \leq \beta_j^{(k)} \leq t_k, \quad k = j, \dots, p, \\ & -\lambda \leq \mathbf{r}'_j \boldsymbol{\beta}_j - 1 \leq \lambda, \\ & -\lambda \leq \mathbf{r}'_k \boldsymbol{\beta}_j \leq \lambda, \quad k \neq j. \end{aligned}$$

- ▶ Remark: this linear relaxation technique is also used in Candès and Tao (2007) for the Dantzig selector problem.

Factor Model Approach

- ▶ Denote $\mathbf{x} = (x_1, \dots, x_p)^T$, and $\mathbf{x} \sim N(\mathbf{0}, \mathbf{\Sigma})$.
- ▶ The factor model is to consider

$$\mathbf{x} = \mathbf{B}\mathbf{z} + \boldsymbol{\epsilon},$$

where $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_p)^T$, and the factors (latent variables) are $\mathbf{z} = (z_1, \dots, z_m)^T$.

- ▶ Assume $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{\Psi})$ where $\mathbf{\Psi} = \text{diag}\{\sigma_1^2, \dots, \sigma_p^2\}$. Then the joint distribution of $(\mathbf{x}_i, \mathbf{z}_i)$ is

$$\begin{pmatrix} \mathbf{z} \\ \mathbf{x} \end{pmatrix} \sim N \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{I}_m & \mathbf{B}' \\ \mathbf{B} & \mathbf{B}\mathbf{B}' + \mathbf{\Psi} \end{bmatrix} \right).$$

- ▶ The underlying assumption is that $\text{cov}(\mathbf{z}, \boldsymbol{\epsilon}) = \mathbf{0}$.

Factor Model Approach for Σ

- ▶ It is easy to see that $\Sigma = \mathbf{B}\mathbf{B}' + \Psi$.
- ▶ The sparse (parsimonious) estimation of \mathbf{B} can imply the sparsity of Σ .
- ▶ With the data matrix $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$, and the latent data matrix $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)^T$, the joint log-likelihood $\ell(\mathbf{B}, \Psi)$ can be written as

$$-\frac{n}{2} \log |\Psi| - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{B}\mathbf{z}_i)' \Psi^{-1} (\mathbf{x}_i - \mathbf{B}\mathbf{z}_i) - \frac{1}{2} \sum_{i=1}^n \mathbf{z}_i^T \mathbf{z}_i,$$

up to some constant.

- ▶ Note that the sparsity is to impose on \mathbf{B} , which is to consider the format of

$$\min -\frac{2}{n} \ell(\mathbf{B}, \Psi) + \lambda J(\mathbf{B}),$$

where $J(\mathbf{B})$ is a penalty function.

Factor Model Approach: Computation

- ▶ However, the latent data points z_i is not observed. The EM-type algorithm can be considered for facilitating the computation.
- ▶ From the joint likelihood of (x, z) , it is known that

$$\begin{aligned}E(z|x) &= B'(BB' + \Psi)^{-1}x \\&= (B'\Psi^{-1}B + I_m)^{-1}B'\Psi^{-1}x \\Var(z|x) &= I_m - B'(BB' + \Psi)^{-1}B \\&= I_m - B'\Psi^{-1}B(B'\Psi^{-1}B + I_m)^{-1}\end{aligned}$$

- ▶ For the M-step, it is to solve

$$\begin{aligned}&\min_{\mathbf{B}} -\frac{2}{n}E[\ell(\mathbf{B}, \Psi)] + \lambda J(\mathbf{B}) \\&\Leftrightarrow \min_{\mathbf{B}} \frac{1}{n} \sum_{i=1}^n E[(x_i - \mathbf{B}z_i)'\Psi^{-1}(x_i - \mathbf{B}z_i)] + \lambda J(\mathbf{B})\end{aligned}$$

Factor Model Approach: Computation (Con't)

- For the estimation of Ψ in M-Step, it is easy to obtain that

$$\hat{\Psi} = \text{diag} \left\{ \frac{1}{n} \sum_{i=1}^n E \left[(\mathbf{x}_i - \mathbf{B} \mathbf{z}_i)(\mathbf{x}_i - \mathbf{B} \mathbf{z}_i)' \right] \right\}$$

- When $\lambda = 0$, the solution of \mathbf{B} becomes the LS estimator that is

$$\hat{\mathbf{B}} = \left\{ \sum_{i=1}^n \mathbf{x}_i E(\mathbf{z}_i') \right\} \left\{ \sum_{i=1}^n E(\mathbf{z}_i \mathbf{z}_i') \right\}^{-1}$$

- Note that the above formulation requires to compute $E(\mathbf{z} \mathbf{z}' | \mathbf{x}) = E(\mathbf{z} | \mathbf{x}) E(\mathbf{z} | \mathbf{x})' + \text{Var}(\mathbf{z} | \mathbf{x})$.
- Remark: Some constraints on the model identification may be needed in the estimation of \mathbf{B} .

Application: Portfolio Optimization

- Consider n stocks whose returns are distributed with mean μ and covariance matrix Σ .
- Markowitz defines the portfolio selection problem as:

$$\begin{aligned} \min_w & w' \Sigma w \\ \text{s.t. } & w' \mathbf{1} = 1, w' \mu = q \end{aligned}$$

where q is the required expected return.

- Using the Lagrange multipliers, the solution is:

$$w = \frac{c - qb}{ac - b^2} \Sigma^{-1} \mathbf{1} + \frac{qa - b}{ac - b^2} \Sigma^{-1} \mu$$

where $a = \mathbf{1}' \Sigma^{-1} \mathbf{1}$, $b = \mathbf{1}' \Sigma^{-1} \mu$, $c = \mu' \Sigma^{-1} \mu$

Log Covariance Matrix Estimation

- Motivation:
 - In univariate case: we can model $\log(\sigma^2)$.
 - In multivariate case, can we still do something as $\log(\Sigma)$

Background

- Covariance matrix estimation is important in multivariate analysis and many statistical applications.
- Suppose x_1, \dots, x_n are i.i.d. p -dimensional random vectors $\sim N(0, \Sigma)$. Let $S = \sum_{i=1}^n x_i x_i' / n$ be the sample covariance matrix. The negative log-likelihood function is proportional to

$$L_n(\Sigma) = -\log |\Sigma^{-1}| + \text{tr}[\Sigma^{-1} S]. \quad (1)$$

- Recent interests on p large or $p \approx n$. S is not a stable estimate.
 - The largest eigenvalue of S overly estimate the true eigenvalue.
 - When $p > n$, S is singular and the smallest eigenvalue is zero. How to estimate Σ^{-1} ?

Motivation

- Estimate of Σ or Σ^{-1} needs to be positive definite.
 - The mathematical restriction makes the covariance matrix estimation problem challenging.
- Matrix transformation/decomposition approaches can relax the restriction: modified Cholesky decomposition, spherical parametrization, Givens parametrization, etc.
- Matrix logarithm: any positive definite Σ can be expressed as a matrix exponential of a real symmetric matrix A .

$$\Sigma = \exp(A) = I + A + \frac{A^2}{2!} + \dots$$

- Expressing the likelihood function in terms of $A \equiv \log(\Sigma)$ releases the mathematical restriction.

Matrix Logarithm Parameterization

- Consider the spectral decomposition of $\Sigma = TDT'$ with $D = \text{diag}(d_1, \dots, d_p)$ where T is an orthonormal matrix consisting of eigenvectors of Σ .
- Then we have

$$A = TMT'$$

with $M = \text{diag}(\log(d_1), \dots, \log(d_p))$.

- In terms of the log-likelihood function,

$$\begin{aligned} -\log |\Sigma^{-1}| &= \text{tr}(A), \\ \text{tr}[\Sigma^{-1}S] &= \text{tr}[\exp(-A)S]. \end{aligned}$$

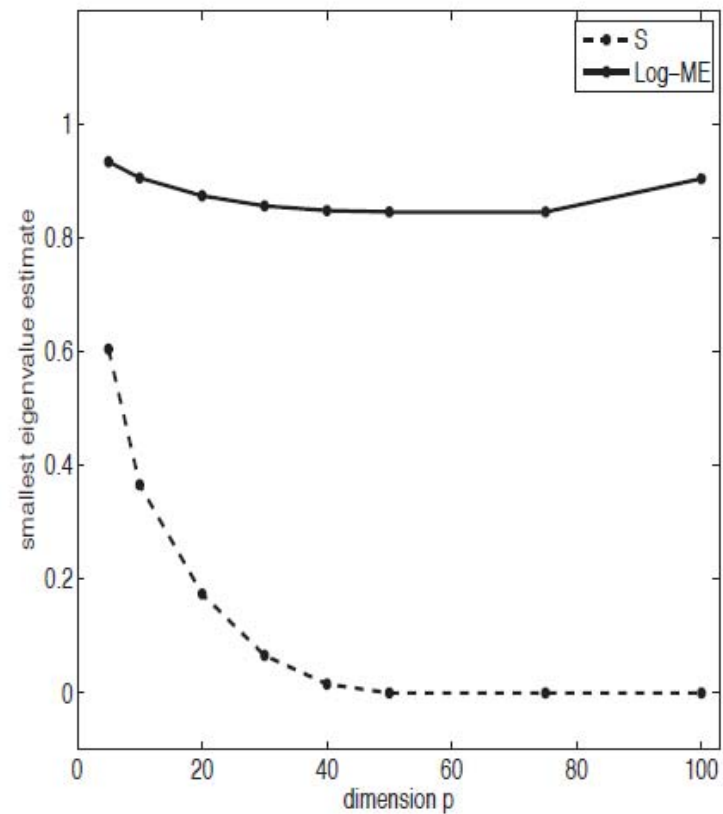
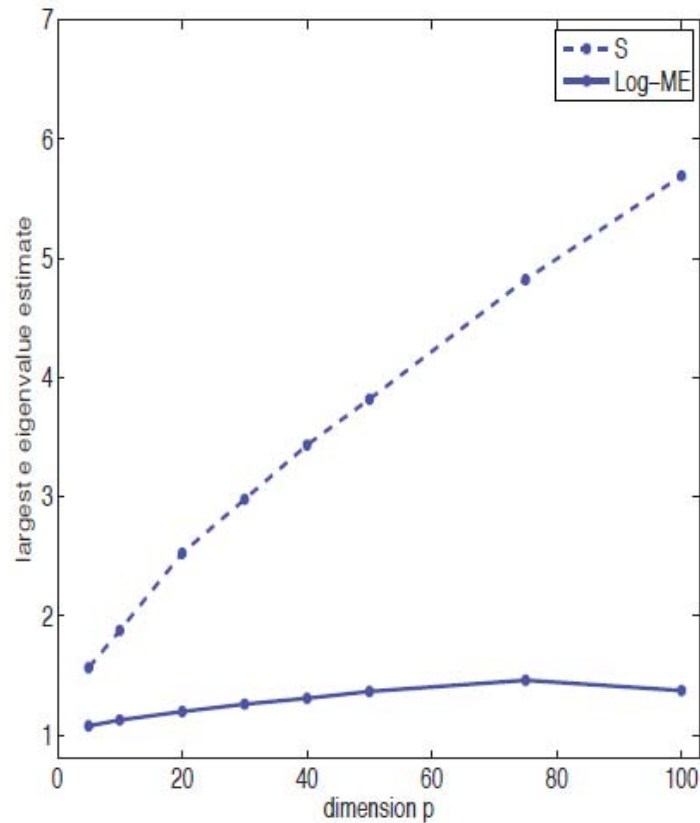
Idea of Proposed Method

- Leonard and Hsu (1992) used this log-transformation method to estimate Σ by approximating the likelihood using Volterra integral equation.
 - Their approximation based on S being nonsingular \Rightarrow not applicable when $p \geq n$.
- We extend the likelihood approximation to the case of singular S .
- Regularize the largest and smallest eigenvalues of Σ *simultaneously*.
- An efficient iterative quadratic programming algorithm to estimate A ($\log \Sigma$).
- The resulting estimate is called “Log-ME”, Logarithm-transformed Matrix Estimate.

A Simple Example

- Experiment: simulate x_i 's from $N(0, I)$, $i = 1, \dots, n$ where $n = 50$.
- For each p varying from 5 to 100, consider the the largest and smallest eigenvalues of the covariance matrix estimate.
- For each p , repeat the experiment 100 times and compute the average of the largest eigenvalues and the average of the smallest eigenvalues for
 - The sample covariance matrix.
 - The Log-ME covariance matrix estimate.

Example Illustration



The averages of the largest and smallest eigenvalues of covariance matrix estimates over the dimension p . The true eigenvalues are all equal to 1.

The Transformed Log-Likelihood

- In terms of the covariance matrix logarithm A , the negative log-likelihood function in (1) becomes

$$L_n(A) = \text{tr}(A) + \text{tr}[\exp(-A)S]. \quad (2)$$

- The problem of estimating a positive definite matrix Σ now becomes a problem of estimating a real symmetric matrix A .
- Because of the matrix exponential term $\exp(-A)S$, estimating A by directly minimizing $L_n(A)$ is nontrivial.
- **Our approach:** Approximate $\exp(-A)S$ using the Volterra integral equation (valid even for S singular case).

The Volterra Integral Equation

- The Volterra integral equation (Bellman, 1970, page 175) is

$$\exp(At) = \exp(A_0t) + \int_0^t \exp(A_0(t-s))(A - A_0) \exp(As) ds. \quad (3)$$

- Repeatedly applying (3) leads to

$$\begin{aligned} \exp(At) = & \exp(A_0t) + \int_0^t \exp(A_0(t-s))(A - A_0) \exp(A_0s) ds \\ & + \int_0^t \int_0^s \exp(A_0(t-s))(A - A_0) \exp(A_0(s-u))(A - A_0) \exp(A_0u) du ds \\ & + \text{cubic and higher order terms,} \end{aligned} \quad (4)$$

where $A_0 = \log(\Sigma_0)$ and Σ_0 is an initial estimate of Σ .

- One can get the expression of $\exp(-A)$ by letting $t = 1$ in (4) and replacing A, A_0 in (4) with $-A, -A_0$.

Approximation to the Log-Likelihood

- The term $\text{tr}[\exp(-A)S]$ can be written as

$$\begin{aligned}\text{tr}[\exp(-A)S] = & \text{tr}(S\Sigma_0^{-1}) - \int_0^1 \text{tr}[(A - A_0)\Sigma_0^{-s}S\Sigma_0^{s-1}]ds \\ & + \int_0^1 \int_0^s \text{tr}[(A - A_0)\Sigma_0^{u-s}(A - A_0)\Sigma_0^{-u}S\Sigma_0^{s-1}]duds \\ & + \text{cubic and higher order terms.}\end{aligned}\tag{5}$$

- By leaving out the higher order terms in (5), we approximate $L_n(A)$ by using $l_n(A)$:

$$\begin{aligned}l_n(A) = & \text{tr}(S\Sigma_0^{-1}) - \left[\int_0^1 \text{tr}[(A - A_0)\Sigma_0^{-s}S\Sigma_0^{s-1}]ds - \text{tr}(A) \right] \\ & + \int_0^1 \int_0^s \text{tr}[(A - A_0)\Sigma_0^{u-s}(A - A_0)\Sigma_0^{-u}S\Sigma_0^{s-1}]duds.\end{aligned}\tag{6}$$

Explicit Form of $I_n(A)$

- The integrations in $I_n(A)$ can be analytically solved through the spectral decomposition of $\Sigma_0 = T_0 D_0 T_0'$.
- **Some Notation:**
 - Here $D_0 = \text{diag}(d_1^{(0)}, \dots, d_p^{(0)})$ with $d_i^{(0)}$'s as the eigenvalues of Σ_0 .
 - $T_0 = (t_1^{(0)}, \dots, t_p^{(0)})$ with $t_i^{(0)}$ as the corresponding eigenvector for $d_i^{(0)}$.
 - Let $B = T_0'(A - A_0)T_0 = (b_{ij})_{p \times p}$, and $\tilde{S} = T_0' S T_0 = (\tilde{s}_{ij})_{p \times p}$.
- In the integration,

$$\begin{aligned} \text{tr}[(A - A_0)\Sigma_0^{-s} S \Sigma_0^{s-1}] &= \text{tr}[(A - A_0)T_0 D_0^{-s} T_0' S T_0 D_0^{s-1} T_0'], \\ \text{tr}[(A - A_0)\Sigma_0^{u-s} (A - A_0)\Sigma_0^{-u} S \Sigma_0^{s-1}] \\ &= \text{tr}[(A - A_0)T_0 D_0^{u-s} T_0' (A - A_0)T_0 D_0^{-u} T_0' S T_0 D_0^{s-1} T_0']. \end{aligned}$$

Explicit Form of $l_n(A)$ (Con't)

- The $l_n(A)$ can be written as a function of b_{ij} :

$$l_n(A) = \sum_{i=1}^p \frac{1}{2} \xi_{ii} b_{ii}^2 + \sum_{i < j} \xi_{ij} b_{ij}^2 + 2 \sum_{i=1}^p \sum_{j \neq i} \tau_{ij} b_{ii} b_{ij} + \sum_{k=1}^p \sum_{i < j, i \neq k, j \neq k} \eta_{kij} b_{ik} b_{kj} - \left[\sum_{i=1}^p \beta_{ii} b_{ii} + 2 \sum_{i < j} \beta_{ij} b_{ij} \right], \quad (7)$$

up to some constant.

- Since $B = T'_0(A - A_0)T_0$, getting $B \leftrightarrow$ getting A .

Some Details

- For the linear term,

$$\beta_{ii} = \frac{\tilde{s}_{ii}}{d_i^{(0)}} - 1, \quad \beta_{ij} = \frac{\tilde{s}_{ij}(d_i^{(0)} - d_j^{(0)}) / (d_i^{(0)} d_j^{(0)})}{(\log d_i^{(0)} - \log d_j^{(0)})}.$$

- For the quadratic term,

$$\xi_{ii} = \frac{\tilde{s}_{ii}}{d_i^{(0)}},$$

$$\xi_{ij} = \frac{\tilde{s}_{ii}/d_i^{(0)} - \tilde{s}_{jj}/d_j^{(0)}}{\log d_j^{(0)} - \log d_i^{(0)}} + \frac{(d_i^{(0)}/d_j^{(0)} - 1)\tilde{s}_{ii}/d_i^{(0)} + (d_j^{(0)}/d_i^{(0)} - 1)\tilde{s}_{jj}/d_j^{(0)}}{(\log d_j^{(0)} - \log d_i^{(0)})^2},$$

$$\tau_{ij} = \left[\frac{1/d_j^{(0)} - 1/d_i^{(0)}}{(\log d_j^{(0)} - \log d_i^{(0)})^2} + \frac{1/d_i^{(0)}}{\log d_j^{(0)} - \log d_i^{(0)}} \right] \tilde{s}_{ij},$$

$$\eta_{kij} = \left[\frac{1/d_i^{(0)} - 1/d_j^{(0)}}{\log(d_k^{(0)}/d_j^{(0)}) \log(d_j^{(0)}/d_i^{(0)})} + \frac{1/d_j^{(0)} - 1/d_i^{(0)}}{\log(d_k^{(0)}/d_i^{(0)}) \log(d_i^{(0)}/d_j^{(0)})} + \frac{2/d_k^{(0)} - 1/d_i^{(0)} - 1/d_j^{(0)}}{\log(d_k^{(0)}/d_i^{(0)}) \log(d_k^{(0)}/d_j^{(0)})} \right] \tilde{s}_{ij}$$

The Log-ME Method

- Propose a regularized method to estimate Σ by using the approximate log-likelihood function $l_n(A)$.
- Consider the penalty function $\|A\|_F^2 = \text{tr}(A^2) = \sum_{i=1}^p (\log(d_i))^2$, where d_i is the eigenvalue of the covariance matrix Σ .
 - If d_i goes to zero or diverges to infinity, the value of $\log(d_i)$ goes to infinity in both cases.
 - Such a penalty function can *simultaneously* regularize the largest and smallest eigenvalues of the covariance matrix estimate.
- Estimate Σ , or equivalently A , by minimizing

$$l_{n,\lambda}(B) \equiv l_{n,\lambda}(A) = l_n(A) + \lambda \text{tr}(A^2), \quad (8)$$

where λ is a tuning parameter.

The Log-ME Method (Con't)

- Note that $\text{tr}(A^2) = \text{tr}((T_0 B T_0' + A_0)^2)$ is equivalent to $\text{tr}(B^2) + 2\text{tr}(B\Gamma)$ up to some constant, where $\Gamma = (\gamma_{ij})_{p \times p} = T_0' A_0 T_0$.
- In terms of B , the function $l_{n,\lambda}(A)$ becomes

$$\begin{aligned}
 l_{n,\lambda}(B) = & \sum_{i=1}^p \frac{1}{2} \xi_{ii} b_{ii}^2 + \sum_{i < j} \xi_{ij} b_{ij}^2 + 2 \sum_{i=1}^p \sum_{j \neq i} \tau_{ij} b_{ii} b_{ij} + \sum_{k=1}^p \sum_{i < j, i \neq k, j \neq k} \eta_{kij} b_{ik} b_{kj} \\
 & - \left(\sum_{i=1}^p \beta_{ii} b_{ii} + 2 \sum_{i < j} \beta_{ij} b_{ij} \right) \\
 & + \lambda \left[\frac{1}{2} \sum_{i=1}^p b_{ii}^2 + \sum_{i < j} b_{ij}^2 + \sum_{i=1}^p \gamma_{ii} b_{ii} + 2 \sum_{i < j} \gamma_{ij} b_{ij} \right].
 \end{aligned} \tag{9}$$

- The $l_{n,\lambda}(B)$ is still a quadratic function of $B = (b_{ij})$.

An Iterative Algorithm

- The $l_{n,\lambda}(B)$ depends on an initial estimate Σ_0 , or equivalently, A_0 .
- Propose to iteratively use $l_{n,\lambda}(B)$ to obtain its minimizer \hat{B} :

Algorithm:

Step 1: Set an initial covariance matrix estimate Σ_0 , a positive definite matrix.

Step 2: Use the spectral decomposition $\Sigma_0 = T_0 D_0 T_0'$, and set $A_0 = \log(\Sigma_0)$.

Step 3: Compute \hat{B} by minimizing $l_{n,\lambda}(B)$. Then obtain $\hat{A} = T_0 \hat{B} T_0' + A_0$, and update the estimate of Σ by

$$\hat{\Sigma} = \exp(\hat{A}) = \exp(T_0 \hat{B} T_0' + A_0).$$

Step 4: Check if $\|\hat{\Sigma} - \Sigma_0\|_F^2$ is less than a pre-specified positive tolerance value. Otherwise, set $\Sigma_0 = \hat{\Sigma}$ and go back to **Step 2**.

- Set an initial Σ_0 in **Step 1** to be $S + \varepsilon I$.

Properties of the Log-ME Method

- The iterative algorithm improves \hat{A} step-by-step as the initial estimate Σ_0 is updated in each iteration.

- Recall that $L_n(A) = \text{tr}(A) + \text{tr}[\exp(-A)S]$ and $L_{n,\lambda}(A) \equiv L_n(A) + \lambda \text{tr}(A^2)$.

Proposition. *Suppose \hat{A} is the solution of the proposed iterative algorithm when it converges. Then the value of $l_n(\hat{A})$ in (7) is exactly the same as the value of $L_n(\hat{A})$, i.e., $l_n(\hat{A}) = L_n(\hat{A})$.*

- When the proposed iterative algorithm converges at \hat{A} , $l_n(A)$ can provide a good approximation to the likelihood function $L_n(A)$ around \hat{A} .
- It indicates that $l_{n,\lambda}(A)$ can approximate $L_{n,\lambda}(A)$ well in a neighborhood region of \hat{A} .

Simulation Study

- Consider four different covariance models of $\Sigma = (\sigma_{ij})_{p \times p}$ for comparison,
 - Model 1: Σ is constructed from an MA(2) model with $\sigma_{ii} = 1, \sigma_{i,i-1} = \sigma_{i-1,i} = 0.6, \sigma_{i,i-2} = \sigma_{i-2,i} = 0.3$, and all the other σ_{ij} are zeros.
 - Model 2: $\Sigma = PMP'$, where M is the covariance matrix defined in model 2, and P is a permutation matrix by randomly permuting the rows of the identity matrix.
 - Model 3: $\Sigma^{-1} = (c_{ij})_{p \times p}$ with $c_{ii} = 1$ and $c_{ij} = 0.3$ if $i \neq j$.
 - Model 4: $\Sigma = GG$ where G is the covariance matrix defined in model 2.
- Compare with five covariance matrix estimates: (1) the LW estimate in Ledoit and Wolf (2006), (2) the banding estimate in Rothman et al. (2010), (3) the graphical Lasso estimate (Glasso) in Yuan and Lin (2007), (4) the estimate with a condition number constraint (denoted by CN) in Won et al. (2009), and (5) the CLIME estimate in Cai et al. (2011).

The CN Method

- The CN method is to estimate Σ with a constraint on its condition number (Won et al., 2009).
- They consider $\hat{\Sigma} = T \text{diag}(\hat{u}_1^{-1}, \dots, \hat{u}_p^{-1}) T'$, where T is from the spectral decomposition of $S = T \text{diag}(l_1, \dots, l_p) T'$.
- The $\hat{u}_1, \dots, \hat{u}_p$ are obtained by solving the constraint optimization

$$\begin{aligned} \min_{u, u_1, \dots, u_p} \quad & \sum_i^p (l_i u_i - \log u_i) \\ \text{s.t.} \quad & u \leq u_i \leq \kappa_{\max} u, \quad i = 1, \dots, p, \end{aligned}$$

where κ_{\max} is a tuning parameter.

The CLIME Method

- For the CLIME method, Cai et al. (2011) estimate the inverse covariance matrix $\Sigma^{-1} = (c_{ij})$ by solving

$$\min_{\Sigma^{-1}} \sum_{i,j} |c_{ij}| \text{ s.t. } \|S\Sigma^{-1} - I\|_{\infty} \leq \lambda,$$

where λ is a tuning parameter.

- For all six methods in comparison, the tuning parameter in each method is selected through the independent validation set under the likelihood loss.

Simulation Study (Con't)

- Consider five loss functions to evaluate the performance of each method,
- The first three loss functions are the Kullback-Leibler (KL) loss, the entropy loss (EN), and the Frobenius norm (Fnorm),

$$KL = -\log |\hat{\Sigma}^{-1}| + \text{tr}(\hat{\Sigma}^{-1}\Sigma) - (-\log |\Sigma^{-1}| + p),$$

$$EN = \text{tr}(\Sigma^{-1}\hat{\Sigma}) - \log |\Sigma^{-1}\hat{\Sigma}| - p,$$

and

$$Fnorm = \|\hat{\Sigma} - \Sigma\|_F = \sqrt{\sum_{i,j} (\hat{\sigma}_{ij} - \sigma_{ij})^2}.$$

- The other two loss functions are

$$\Delta_{1/p} = |\hat{d}_1/\hat{d}_p - d_1/d_p| \text{ and } \Delta_1 = |\hat{d}_1 - d_1|.$$

Here d_1 is the largest eigenvalue and d_p is the smallest eigenvalue of the covariance matrix Σ .

Simulation Results

Averages and standard errors from 100 runs for $n = 50, p = 100$.

Method	Model 1					Model 2				
	KL	EN	$Fnorm$	$\Delta_{1/p}$	Δ_1	KL	EN	$Fnorm$	$\Delta_{1/p}$	Δ_1
Log-ME	42.59	46.29	7.80	13.09	0.11	42.89	46.74	7.81	21.95	0.18
	(0.05)	(0.30)	(0.01)	(0.03)	(0.01)	(0.07)	(0.06)	(0.01)	(0.02)	(0.01)
Glasso	28.87	59.59	8.90	11.25	3.41	28.83	59.99	8.89	11.19	3.51
	(0.08)	(0.25)	(0.06)	(0.19)	(0.05)	(0.08)	(0.29)	(0.05)	(0.15)	(0.05)
LW	43.13	86.18	7.92	22.91	0.39	43.11	86.86	7.91	22.94	0.38
	(0.04)	(0.32)	(0.01)	(0.05)	(0.02)	(0.05)	(0.42)	(0.01)	(0.05)	(0.03)
CN	44.70	109.27	8.04	24.69	0.46	45.33	102.76	9.45	22.75	0.73
	(0.04)	(0.40)	(0.03)	(0.03)	(0.01)	(0.05)	(0.39)	(0.02)	(0.01)	(0.01)
Banding	10.12	8.49	3.34	32.35	0.60	57.40	140.69	9.64	24.30	1.19
	(0.10)	(0.07)	(0.02)	(1.15)	(0.03)	(0.05)	(0.50)	(0.00)	(0.16)	(0.02)
Clime	26.13	49.65	6.29	20.72	0.83	25.86	49.34	6.25	20.52	0.81
	(0.12)	(0.46)	(0.02)	(0.08)	(0.01)	(0.11)	(0.39)	(0.02)	(0.08)	(0.01)

Simulation Results (Con't)

Averages and standard errors from 100 runs for $n = 50, p = 100$.

Method	Model 3					Model 4				
	KL	EN	$Fnorm$	$\Delta_{1/p}$	Δ_1	KL	EN	$Fnorm$	$\Delta_{1/p}$	Δ_1
Log-ME	9.34 (0.00)	30.86 (0.02)	4.29 (0.01)	32.63 (0.04)	0.24 (0.04)	80.32 (0.29)	479.83 (5.08)	19.89 (0.04)	750.95 (0.28)	1.26 (0.05)
Glasso	7.05 (0.04)	54.22 (0.21)	6.00 (0.04)	41.12 (0.04)	1.83 (0.03)	77.17 (0.18)	344.15 (2.31)	24.84 (0.14)	679.05 (0.87)	9.97 (0.14)
LW	2.99 (0.01)	36.49 (0.10)	1.64 (0.02)	42.62 (0.01)	0.22 (0.01)	130.51 (0.10)	1377.90 (6.87)	18.98 (0.03)	759.12 (0.13)	2.08 (0.09)
CN	3.72 (0.03)	41.54 (0.57)	2.43 (0.04)	42.85 (0.01)	0.20 (0.01)	158.69 (0.17)	3005.30 (12.68)	29.57 (0.15)	760.32 (0.07)	10.51 (0.15)
Banding	5.14 (0.04)	39.84 (0.09)	3.15 (0.02)	41.07 (0.03)	0.76 (0.02)	175.20 (0.10)	2797.80 (10.56)	26.43 (0.10)	753.26 (0.74)	3.74 (0.07)
Clime	2.90 (0.01)	40.53 (0.02)	1.54 (0.02)	42.68 (0.01)	0.11 (0.03)	85.23 (0.17)	282.74 (2.51)	21.42 (0.03)	747.04 (0.29)	4.16 (0.03)

Classification of Ionosphere Data

- The data contains 351 observations with 34 variables with binary response 1 for “Good” and -1 for “Bad” of radar returns (Sigillito et al., 1989).
- Apply linear discriminant analysis (LDA) for classification under various covariance matrix estimate.
- Examine how different covariance matrix estimates can affect the classification performance of LDA.
- Randomly divide the data into a training set (of size 40), a validation set (of size 40) and a test set.

Misclassification Errors

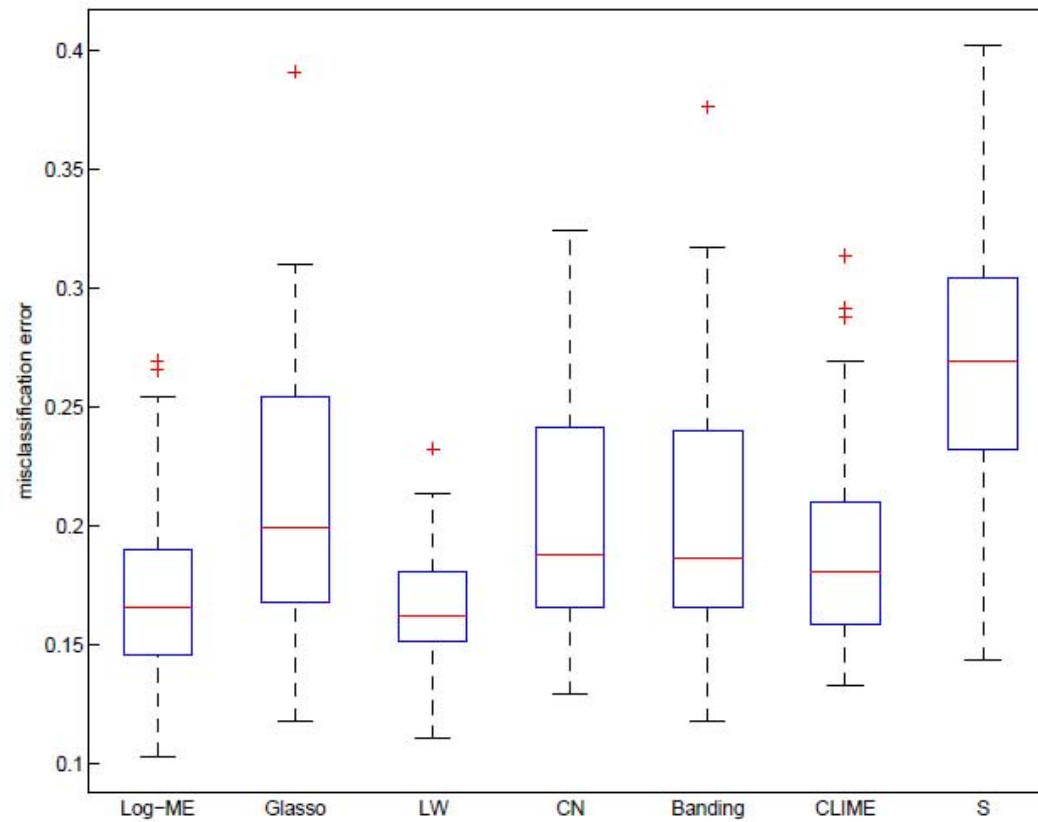


Figure 1. The boxplots of test misclassification errors from 100 replications.

Portfolio Optimization of Stock Data

- Apply the Log-ME method in the portfolio optimization.
- In mean-variance optimization, the risk of a portfolio $w = (w_1, \dots, w_p)$ is measured by the standard deviation $\sqrt{w^T \Sigma w}$, where $\sum_i^p w_i = 1$.
- The estimated minimum variance portfolio optimization problem is

$$\begin{aligned} \min_w w^T \hat{\Sigma} w \\ \text{s.t. } \sum_i^p w_i = 1, \end{aligned} \tag{10}$$

where $\hat{\Sigma}$ is an estimate of the true covariance matrix Σ .

- An accurate covariance matrix estimate $\hat{\Sigma}$ can lead to a better portfolio strategy.

The Setting-up

- Consider the weekly returns of $p = 30$ components of the Dow Jones Industrial Index from January 8th, 2007 to June 28th, 2010.
- Use the first $n = 50$ observations as a training set, the next 50 observations as a validation set, and *the remaining* 83 observations as a test set.
- Let X_{ts} be the test set and S_{ts} be the sample covariance matrix of X_{ts} . The performance of a portfolio w is measured by the *realized return*

$$R(w) = \sum_{x \in X_{ts}} w^T x,$$

and the *realized risk*

$$\sigma(w) = \sqrt{w^T S_{ts} w}.$$

- The optimal portfolio \tilde{w} is computed with $\hat{\Sigma}$ from the Log-ME method and the other five methods, separately.

The Comparison Results

Table 1. The comparison of the realized return and the realized risk.

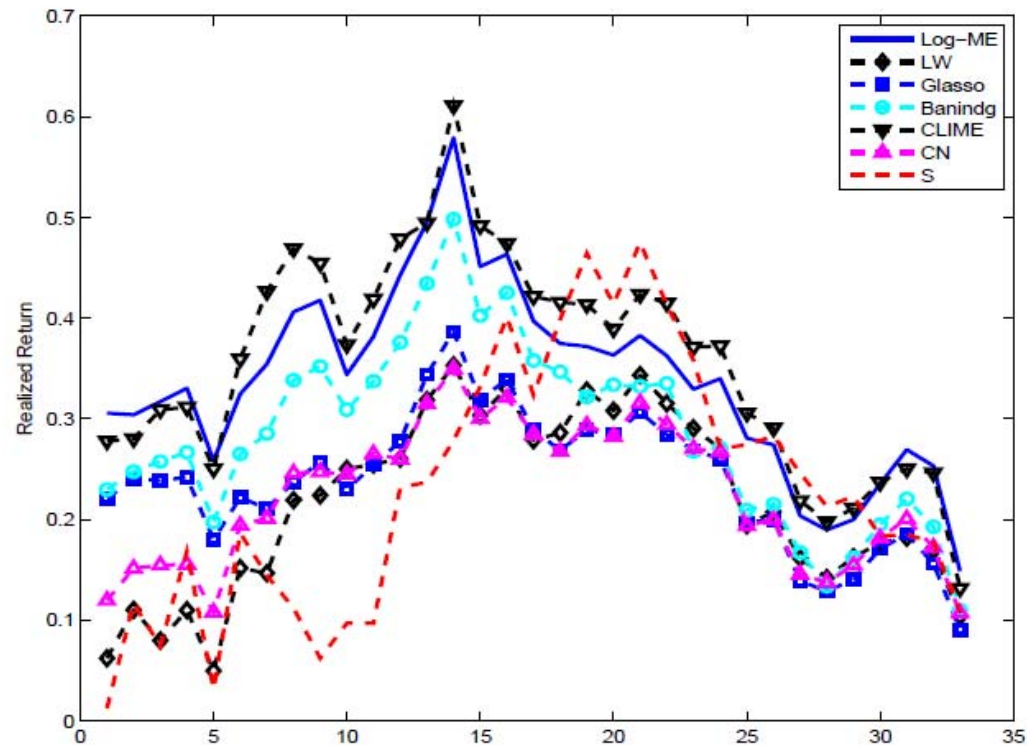
	Log-ME	CN	S	LW	Glasso	Banding	CLIME
Realized return $R(\tilde{w})$	0.218	0.128	0.059	0.062	0.193	0.192	0.211
Realized risk $\sigma(\tilde{w})$	0.029	0.024	0.035	0.025	0.028	0.029	0.029

- The Log-ME method produced a portfolio with a larger realized return but smaller realized risk.

Comparison in Different Periods

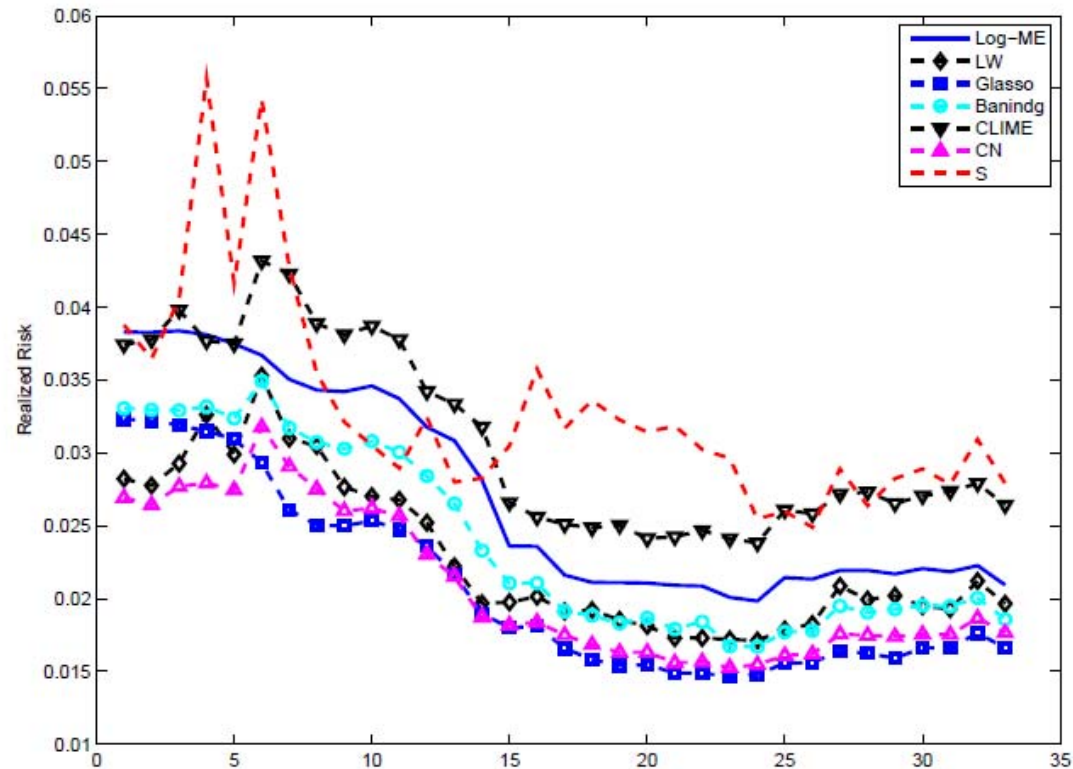
- Given a starting week, use the first 50 observations as the training set, the next 50 observations as a validation set, and *the third 50 observations* as a test set.
- Shift the starting week one ahead every time, and evaluate the portfolio strategy of 33 different consecutive test periods.
- The optimal portfolio \tilde{w} is computed with $\hat{\Sigma}$ obtained from the Log-ME method and other methods in comparison.

The Realized Returns



The proposed Log-ME covariance matrix estimate can lead to higher returns.

The Realized Risks



The log-ME method has relatively higher risks than other methods, but it provides larger realized returns.

Extension: Other Types of Penalties

- One can consider other types of penalties for regularization.
- Consider L_1 penalty on $A = (a_{ij})$ such that

$$L_{n,\lambda}(A) = \text{tr}(A) + \text{tr}[\exp(-A)S] + \lambda \sum_{i,j} |a_{ij}|.$$

- Need to be more careful on the interpretation of L_1 penalty.
- Alternatively, one can consider a penalty term with respect to a target,

$$L_{n,\lambda}(A) = \text{tr}(A) + \text{tr}[\exp(-A)S] + \lambda \|A - A^*\|_F^2.$$

Extension: Bayesian Perspective

- The approximated log-likelihood function $l_n(A)$ is a quadratic function of A .
- Denote $a = \text{vec}(A) = (a_{11}, \dots, a_{pp}, a_{12}, \dots, a_{p-1,p}, \dots, a_{1,p})'$, Then we have

$$l_n(a) = \text{Const} + (a - a^{(0)} - Q^{-1}\beta)^T Q (a - a^{(0)} - Q^{-1}\beta). \quad (11)$$

where

$$Q = \frac{1}{2} \sum_{i=1}^p \xi_{ii} f_{ii} f_{ii}^T + \sum_{i < j} \xi_{ij} f_{ij} f_{ij}^T + 2 \sum_{i=1}^p \sum_{j \neq i} \tau_{ij} f_{ii} f_{ij}^T + \sum_{k=1}^p \sum_{i < j, i \neq k, j \neq k} \eta_{kij} f_{ik} f_{kj}^T,$$

and

$$\beta = \sum_{i=1}^p \frac{1}{2} \beta_{ii} f_{ii} + \sum_{i < j} \beta_{ij} f_{ij}.$$

Here $f_{ij} = t_i^{(0)} * t_j^{(0)}$ the product of the eigenvalues $t_i^{(0)}$ and $t_j^{(0)}$ such that $(a - a^{(0)})^T (t_i^{(0)} * t_j^{(0)}) = (t_i^{(0)})^T (A - A_0) t_j^{(0)} = b_{ij}$.

Bayesian Covariance Matrix Estimation (Con't)

- The $l_n(a)$ can be viewed as the log-density function of a multivariate normal distribution with mean vector $\alpha^{(0)} + Q^{-1}\beta$ and covariance matrix Q^{-1} .
- The proposed framework enables to apply flexible prior structures.
- Conjugate priors: $a = \text{vec}(A) \sim N(\gamma, R)$.
 - The posterior distribution of a is a multivariate normal with mean a^* and covariance R^* as

$$a^* = (Q + R^{-1})^{-1}(Q\alpha^{(0)} + \beta + R^{-1}\gamma),$$

$$R^* = (Q + R^{-1})^{-1}.$$

Hierarchical Prior Structures

- Consider a hierarchical structure on the different bands of A . We set the prior of $a = \text{vec}(A)$ as

$$\begin{aligned}a_{11}, a_{22}, \dots, a_{pp} &\sim N(\mu_1, \sigma_1^2), \\a_{12}, a_{23}, \dots, a_{p-1,p} &\sim N(\mu_2, \sigma_2^2), \\&\vdots \\a_{1,p-1}, a_{2p} &\sim N(\mu_{p-1}, \sigma_{p-1}^2), \\a_{1,p} &\sim N(\mu_p, \sigma_p^2).\end{aligned}$$

- The hierarchical structure lies in that $\mu = (\mu_1, \dots, \mu_p)$ and $\sigma^2 = (\sigma_1^2, \dots, \sigma_p^2)$ with prior distribution $p(\mu, \sigma^2)$ as

$$\begin{aligned}\mu_1, \mu_2, \dots, \mu_p &\sim N(\mu_0, \tau_0^2), \\ \sigma_1^2, \sigma_2^2, \dots, \sigma_p^2 &\sim \text{Inv-}\chi^2(\nu_0, \sigma_0^2),\end{aligned}$$

where μ_i 's and σ_i^2 's are independent each other.

Hierarchical Prior Structures (Con't)

- The posterior distribution of $a = \text{vec}(A)$ is $f(a|x_1, \dots, x_n) \propto$

$$\exp\left\{-\frac{1}{2}l_n(a)\right\} \times \prod_{k=1}^p \prod_{i-j+1=k} N(a_{ij}|\mu_k, \sigma_k^2) \times \prod_{k=1}^p N(\mu_k|\mu_0, \tau_0^2) \prod_{k=1}^p \text{Inv-}\chi^2(\sigma_k^2|v_0, \sigma_0^2).$$

- The conditional posteriors have close forms:
 - The conditional posterior distribution $p(a|u, \sigma^2, X) \propto N(a^*, R^*)$ with

$$a^* = (Q + R^{-1})^{-1}(Qa^{(0)} + \beta + R^{-1}\gamma), \quad R^* = (Q + R^{-1})^{-1}.$$

where $\gamma = (\mu_1 1_p, \dots, \mu_{p-1} 1_1)^T$ and $R = \text{diag}(\sigma_1 1_p, \dots, \sigma_{p-1} 1_1)$.

- The conditional posterior distribution $(\mu_k|a, \sigma^2, X) \propto N(\mu_k^*, (\tau_k^2)^*)$, where

$$\mu_k^* = \frac{\frac{\sum_{i-j+1=k} a_{ij}}{\sigma_k^2} + \frac{\mu_0}{\tau_0^2}}{\frac{p-k+1}{\sigma_k^2} + \frac{1}{\tau_0^2}}, \quad (\tau_k^2)^* = \frac{1}{\frac{\sum_{i-j+1=k} a_{ij}}{\sigma_k^2} + \frac{1}{\tau_0^2}}.$$

- The conditional posterior of $(\sigma^2|a, \mu, X)$ follows scaled inverse χ^2 distributions.

Remarks

- Estimate the covariance matrix through its matrix logarithm based on a penalized likelihood function approximation.
- The Log-ME method regularizes the largest and smallest eigenvalues simultaneously by imposing a convex penalty.
- Can be extended to Bayesian covariance matrix estimation with flexible prior structures.

Multiple Graphical Models

- ▶ Suppose there is a heterogeneous data set with p variables and K categories. The k -th category 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 -dimensional row vector.
- ▶ Assume $\mathbf{x}^{(k)} \sim N(\mathbf{0}, \Sigma^{(k)})$. It means there are multiple graphical models for estimation
- ▶ Denote $\Omega^{(k)} = (\Sigma^{(k)})^{-1}$. The negative log-likelihood can be written as

$$\min_{\Omega^{(1)}, \dots, \Omega^{(K)}} \sum_{k=1}^K n_k \left[-\log |\Omega^{(k)}| + \text{tr}(\Omega^{(k)} \mathbf{S}^{(k)}) \right],$$

up to some constant, where

$$\mathbf{S}^{(k)} = \frac{1}{n_k} \sum_{i=1}^{n_k} \mathbf{x}_i^{(k)} (\mathbf{x}_i^{(k)})^T.$$

Multiple Graphical Model: Joint Estimation

- ▶ For a joint estimation, we can impose consistent sparse structure on the $\mathbf{\Omega}^{(k)} = (\omega_{ij}^{(k)})_{p \times p}$.
- ▶ Remark: The idea is similar as the group lasso. while now we focus group structure for graphical lasso.
- ▶ The penalized log-likelihood estimation can be considered by minimizing

$$\sum_{k=1}^K n_k \left[-\log |\mathbf{\Omega}^{(k)}| + \text{tr}(\mathbf{\Omega}^{(k)} \mathbf{S}^{(k)}) \right] + \lambda \sum_{i \neq j} \sqrt{\sum_{k=1}^K |\omega_{ij}^{(k)}|},$$

where λ is a tuning parameter.

- ▶ Other meaningful penalties can also be considers. The reparameterization technique is also an alternative approach to address the structured penalty.

Joint Estimation Tools

- ▶ To solve the above optimization efficiently, we can apply local linear approximation to the penalty term.
- ▶ Denote $(\omega_{ij}^{(k)})_{(m)}$ the estimate of $\omega_{ij}^{(k)}$ in $f_k^{-1} = (\omega_{ij}^{(k)})_{p \times p}$ at m th iteration. Then using local linear approximation, we write

$$\sqrt{\sum_{k=0}^n |\omega_{ij}^{(k)}|} = \frac{\sum_{k=0}^n |\omega_{ij}^{(k)}|}{\sqrt{\sum_{k=0}^n |(\omega_{ij}^{(k)})_{(m)}|}} \equiv \alpha_{ij}^{(m)} \sum_{k=0}^n |\omega_{ij}^{(k)}|,$$

where $\alpha_{ij}^{(m)} = 1/\sqrt{\sum_{k=0}^n |(\omega_{ij}^{(k)})_{(m)}|}$.

- ▶ Now the glasso (Friedman et al., 2008) can be used for efficient computation. At the $(m+1)$ th iteration, we minimize

$$\sum_{k=1}^K n_k \left[-\log |\mathbf{\Omega}^{(k)}| + \text{tr}(\mathbf{\Omega}^{(k)} \mathbf{S}^{(k)}) \right] + \lambda \sum_{i \neq j} \alpha_{ij}^{(m)} |\omega_{ij}^{(k)}|.$$

- ▶ The above is to estimate each $\mathbf{\Omega}^{(k)}$ with positive-definiteness.

Some References

- ▶ Huang, J. Z., Liu, N., Pourahmadi, M., and Liu, L. (2006). Covariance matrix selection and estimation via penalised normal likelihood. *Biometrika*, **93**, 85–98.
- ▶ Bickel, P. J. and Levina, E. (2008). Regularized estimation of large covariance matrices. *Annals of Statistics*, **36**, 199–227.
- ▶ Rothman, A. J., Levina, E., and Zhu, J. (2010). A new approach to Cholesky-based estimation of high-dimensional covariance matrices. *Biometrika*, **97**, 539–550.
- ▶ Cai, T., Liu, W., and Luo, X. (2011). A constrained l_1 minimization approach to sparse precision matrix estimation. *Journal of American Statistical Association*, **106**, 594–607.
- ▶ Candès, E., and Tao, T. (2007), The Dantzig Selector, Statistical Estimation When p is Much Larger Than n . *The Annals of Statistics*, **35**, 2313–2351.

Remarks

Thank you!

- **Questions and Comments?**