

On robust estimation of high-dimensional covariance matrices

Karina Ashurbekova

GIPSA-lab, INRIA

July 18, 2019

No machine learning today! 😊

Sample Covariance Matrix

Warm-up: covariance matrix 😊

- $\mathbf{X} = (X_1, \dots, X_p)$ - p -variate random vector
- Covariance matrix:

$$\Sigma = \begin{bmatrix} \mathbb{E}[(X_1 - \mathbb{E}(X_1))(X_1 - \mathbb{E}(X_1))] & \dots & \mathbb{E}[(X_1 - \mathbb{E}(X_1))(X_p - \mathbb{E}(X_p))] \\ \vdots & \ddots & \vdots \\ \mathbb{E}[(X_p - \mathbb{E}(X_p))(X_1 - \mathbb{E}(X_1))] & \dots & \mathbb{E}[(X_p - \mathbb{E}(X_p))(X_p - \mathbb{E}(X_p))] \end{bmatrix} \quad (1)$$

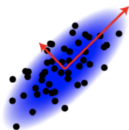
- **Problem: estimate the covariance matrix** from the i.i.d. observations $\mathbf{X}_1 = (X_{11}, \dots, X_{1p}), \dots, \mathbf{X}_n = (X_{n1}, \dots, X_{np})$

Why covariance estimation?

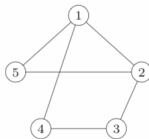
Portfolio selection



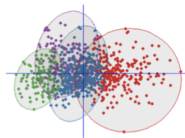
PCA



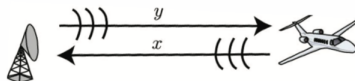
Graphical models



Discriminant Analysis



Radar detection



$$\Sigma^{-1} = \begin{bmatrix} \bullet & \bullet & 0 & \bullet & \bullet \\ \bullet & \bullet & \bullet & 0 & \bullet \\ 0 & \bullet & \bullet & \bullet & 0 \\ \bullet & 0 & \bullet & \bullet & 0 \\ \bullet & \bullet & 0 & 0 & \bullet \end{bmatrix}$$

How to estimate the covariance matrix?

- We have a data matrix $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_n]^T$ (n rows (samples) and p columns (dimension))
- The sample covariance matrix:

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \widehat{\mathbf{X}})(\mathbf{x}_i - \widehat{\mathbf{X}})^T, \quad (2)$$

$$\widehat{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

- **Problem:**
 - The eigenstructure of \mathbf{S} tends to be systematically distorted unless $\frac{p}{n}$ is small \Rightarrow
 - Larger eigenvalues are overestimated; smaller eigenvalues are underestimated

How to estimate the covariance matrix?

- We have a data matrix $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_n]^T$ (n rows (samples) and p columns (dimension))
- The sample covariance matrix:

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \widehat{\mathbf{X}})(\mathbf{x}_i - \widehat{\mathbf{X}})^T, \quad (2)$$

$$\widehat{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

- **Problem:**
 - The eigenstructure of \mathbf{S} tends to be systematically distorted unless $\frac{p}{n}$ is small \Rightarrow
 - Larger eigenvalues are overestimated; smaller eigenvalues are underestimated

How to estimate the covariance matrix?

- We have a data matrix $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_n]^T$ (n rows (samples) and p columns (dimension))
- The sample covariance matrix:

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \widehat{\mathbf{X}})(\mathbf{x}_i - \widehat{\mathbf{X}})^T, \quad (2)$$

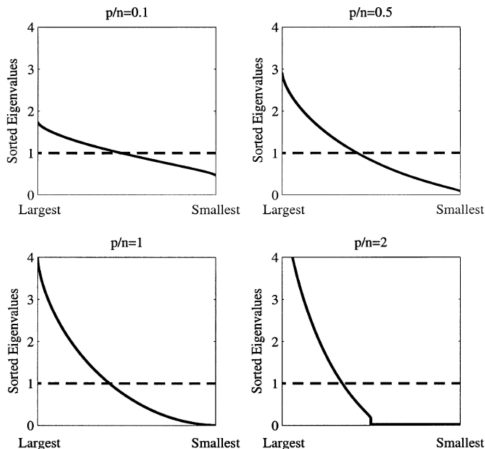
$$\widehat{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

- **Problem:**
 - The eigenstructure of \mathbf{S} tends to be systematically distorted unless $\frac{p}{n}$ is small \Rightarrow
 - Larger eigenvalues are overestimated; smaller eigenvalues are underestimated

Eigenvalues structure for the range $\frac{p}{n}$

Larger eigenvalues are overestimated; smaller eigenvalues are underestimated.

Too bad 😞



$p > n$. What can we do?

Shrink!

- **To ensure non-singularity**, Ledoit and Wolf (2004) proposed a shrinkage estimator of the covariance matrix:

$$\widehat{\Sigma} = \alpha_1 \mathbf{S} + \alpha_2 \mathbf{I}$$

- Use $\widehat{\Sigma}$ that shrinks \mathbf{S} towards to a structure (e.g., a scaled identity matrix) using a tuning (shrinkage) parameter α_2
- Why?
 - Mean Squared Error $\text{MSE}(\widehat{\Sigma}) = \mathbb{E} \left[\left\| \widehat{\Sigma} - \Sigma \right\|_F^2 \right]$ can be reduced by introducing some bias!
 - Positive definiteness of $\widehat{\Sigma}$ can be ensured!

Shrink!

- **To ensure non-singularity**, Ledoit and Wolf (2004) proposed a shrinkage estimator of the covariance matrix:

$$\widehat{\Sigma} = \alpha_1 \mathbf{S} + \alpha_2 \mathbf{I}$$

- Use $\widehat{\Sigma}$ that shrinks \mathbf{S} towards to a structure (e.g., a scaled identity matrix) using a tuning (shrinkage) parameter α_2
- Why?
 - Mean Squared Error $\text{MSE}(\widehat{\Sigma}) = \mathbb{E} \left[\|\widehat{\Sigma} - \Sigma\|_F^2 \right]$ can be reduced by introducing some bias!
 - Positive definiteness of $\widehat{\Sigma}$ can be ensured!

Shrink!

- **To ensure non-singularity**, Ledoit and Wolf (2004) proposed a shrinkage estimator of the covariance matrix:

$$\widehat{\Sigma} = \alpha_1 \mathbf{S} + \alpha_2 \mathbf{I}$$

- Use $\widehat{\Sigma}$ that shrinks \mathbf{S} towards to a structure (e.g., a scaled identity matrix) using a tuning (shrinkage) parameter α_2
- Why?
 - Mean Squared Error $\text{MSE}(\widehat{\Sigma}) = \mathbb{E} \left[\|\widehat{\Sigma} - \Sigma\|_F^2 \right]$ can be reduced by introducing some bias!
 - Positive definiteness of $\widehat{\Sigma}$ can be ensured!

How to estimate α_1 and α_2 ?

Find α_1^* and α_2^* that minimize Mean Squared Error (MSE):

$$MSE(\widehat{\Sigma}) = \mathbb{E} \left[\|\widehat{\Sigma} - \Sigma\|_F^2 \right] = \mathbb{E} \left[\|\alpha_1 \mathbf{S} + \alpha_2 \mathbf{I} - \Sigma\|_F^2 \right]$$

Theorem

The optimal parameters α_1^* and α_2^* are:

$$\alpha_2^* = (1 - \alpha_1^*) \frac{\text{tr}(\Sigma)}{p},$$
$$\alpha_1^* = \frac{p \left(\frac{p \text{tr}(\Sigma^2)}{\text{tr}(\Sigma)^2} - 1 \right) \left(\frac{\text{tr}(\Sigma)}{p} \right)^2}{\mathbb{E}[\text{tr}(\mathbf{S}^2)] - p \left(\frac{\text{tr}(\Sigma)}{p} \right)^2}$$

Thus $\widehat{\Sigma}$ is a convex combination of \mathbf{S} and $\frac{\text{tr}(\Sigma)}{p} \mathbf{I}$

How to estimate α_1 and α_2 ?

Find α_1^* and α_2^* that minimize Mean Squared Error (MSE):

$$MSE(\widehat{\Sigma}) = \mathbb{E} \left[\|\widehat{\Sigma} - \Sigma\|_F^2 \right] = \mathbb{E} \left[\|\alpha_1 \mathbf{S} + \alpha_2 \mathbf{I} - \Sigma\|_F^2 \right]$$

Theorem

The optimal parameters α_1^* and α_2^* are:

$$\alpha_2^* = (1 - \alpha_1^*) \frac{\text{tr}(\Sigma)}{p},$$
$$\alpha_1^* = \frac{p \left(\frac{p \text{tr}(\Sigma^2)}{\text{tr}(\Sigma)^2} - 1 \right) \left(\frac{\text{tr}(\Sigma)}{p} \right)^2}{\mathbb{E} [\text{tr}(\mathbf{S}^2)] - p \left(\frac{\text{tr}(\Sigma)}{p} \right)^2}$$

Thus $\widehat{\Sigma}$ is a convex combination of \mathbf{S} and $\frac{\text{tr}(\Sigma)}{p} \mathbf{I}$

Bias-variance trade-off

LW estimator:

$$\widehat{\Sigma} = \alpha \mathbf{S} + (1 - \alpha) \frac{\text{tr}(\Sigma)}{p} \mathbf{I}$$

- $\frac{\text{tr}(\Sigma)}{p} \mathbf{I}$: all bias no variance.
- \mathbf{S} : all variance no bias (\mathbf{S} is unbiased estimation of Σ i. e. $\mathbb{E}(\mathbf{S}) = \Sigma$)

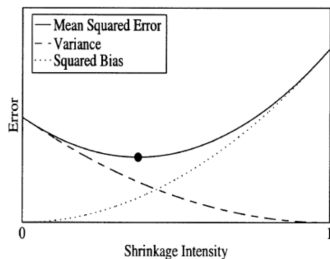


Figure: Shrinkage intensity: $1 - \alpha$

A Bayesian interpretation of shrinkage

Bayes' rule:

$$p(\theta|\mathbf{X}) = \frac{p(\mathbf{X}|\theta) * p(\theta)}{p(\mathbf{X})}$$
$$\text{posterior} = \frac{\text{likelihood} * \text{prior}}{\text{evidence}}$$

Maximum A Posteriori Estimation:

$$\begin{aligned}\theta_{MAP} &= \arg \max_{\theta} p(\mathbf{X}|\theta)p(\theta) = \\ &= \arg \max_{\theta} \log p(\mathbf{X}|\theta) + \log p(\theta)\end{aligned}$$

Thus:

- Laplacian prior $\Rightarrow l_1$ -regularization (lasso)
- Wishart-inverse prior \Rightarrow shrinkage

Difference between lasso and shrinkage:

- lasso impose the sparsity on the elements of covariance matrix!
- shrinkage - no! shrinkage affects only eigenvalues.

A Bayesian interpretation of shrinkage

Bayes' rule:

$$p(\theta|\mathbf{X}) = \frac{p(\mathbf{X}|\theta) * p(\theta)}{p(\mathbf{X})}$$
$$\text{posterior} = \frac{\text{likelihood} * \text{prior}}{\text{evidence}}$$

Maximum A Posteriori Estimation:

$$\begin{aligned}\theta_{MAP} &= \arg \max_{\theta} p(\mathbf{X}|\theta)p(\theta) = \\ &= \arg \max_{\theta} \log p(\mathbf{X}|\theta) + \log p(\theta)\end{aligned}$$

Thus:

- Laplacian prior $\Rightarrow l_1$ -regularization (lasso)
- Wishart-inverse prior \Rightarrow shrinkage

Difference between lasso and shrinkage:

- lasso impose the sparsity on the elements of covariance matrix!
- shrinkage - no! shrinkage affects only eigenvalues.

A Bayesian interpretation of shrinkage

Bayes' rule:

$$p(\theta|\mathbf{X}) = \frac{p(\mathbf{X}|\theta) * p(\theta)}{p(\mathbf{X})}$$
$$\text{posterior} = \frac{\text{likelihood} * \text{prior}}{\text{evidence}}$$

Maximum A Posteriori Estimation:

$$\begin{aligned}\theta_{MAP} &= \arg \max_{\theta} p(\mathbf{X}|\theta)p(\theta) = \\ &= \arg \max_{\theta} \log p(\mathbf{X}|\theta) + \log p(\theta)\end{aligned}$$

Thus:

- Laplacian prior $\Rightarrow l_1$ -regularization (lasso)
- Wishart-inverse prior \Rightarrow shrinkage

Difference between lasso and shrinkage:

- lasso impose the sparsity on the elements of covariance matrix!
- shrinkage - no! shrinkage affects only eigenvalues.

But...

LW shrinkage approach is based on the sample covariance matrix \mathbf{S} !

- we know how to deal with $p > n$ case only for the **Gaussian data**
- but if the data is corrupted? (outliers, noise)
 - Sample Covariance Matrix is sensitive to outliers. 😞 Why?

$$\begin{aligned}\mathbf{S} &= \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \widehat{\mathbf{X}})(\mathbf{x}_i - \widehat{\mathbf{X}})^T = \\ &= \sum_{i=1}^n w_i (\mathbf{x}_i - \widehat{\mathbf{X}})(\mathbf{x}_i - \widehat{\mathbf{X}})^T, \\ w_i &= \frac{1}{n}\end{aligned}$$

- consider the distance $d_i = \sqrt{(\mathbf{x}_i - \widehat{\mathbf{X}})^T \Sigma^{-1} (\mathbf{x}_i - \widehat{\mathbf{X}})}$
- $w_i = \frac{1}{n}$, normal samples and outliers contribute to \mathbf{S} equally

LW shrinkage approach is based on the sample covariance matrix \mathbf{S} !

- we know how to deal with $p > n$ case only for the **Gaussian data**
- but if the data is corrupted? (outliers, noise)
 - Sample Covariance Matrix is sensitive to outliers. 😞 Why?

$$\begin{aligned}\mathbf{S} &= \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \widehat{\mathbf{X}})(\mathbf{x}_i - \widehat{\mathbf{X}})^T = \\ &= \sum_{i=1}^n w_i (\mathbf{x}_i - \widehat{\mathbf{X}})(\mathbf{x}_i - \widehat{\mathbf{X}})^T, \\ w_i &= \frac{1}{n}\end{aligned}$$

- consider the distance $d_i = \sqrt{(\mathbf{x}_i - \widehat{\mathbf{X}})^T \Sigma^{-1} (\mathbf{x}_i - \widehat{\mathbf{X}})}$
- $w_i = \frac{1}{n}$, normal samples and outliers contribute to \mathbf{S} equally

LW shrinkage approach is based on the sample covariance matrix \mathbf{S} !

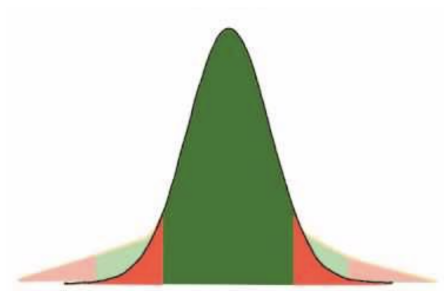
- we know how to deal with $p > n$ case only for the **Gaussian data**
- but if the data is corrupted? (outliers, noise)
 - Sample Covariance Matrix is sensitive to outliers. 😞 Why?

$$\begin{aligned}\mathbf{S} &= \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \widehat{\mathbf{X}})(\mathbf{x}_i - \widehat{\mathbf{X}})^T = \\ &= \sum_{i=1}^n w_i (\mathbf{x}_i - \widehat{\mathbf{X}})(\mathbf{x}_i - \widehat{\mathbf{X}})^T, \\ w_i &= \frac{1}{n}\end{aligned}$$

- consider the distance $d_i = \sqrt{(\mathbf{x}_i - \widehat{\mathbf{X}})^T \Sigma^{-1} (\mathbf{x}_i - \widehat{\mathbf{X}})}$
- $w_i = \frac{1}{n}$, normal samples and outliers contribute to \mathbf{S} equally

What to do in the case of noisy data/data with outliers

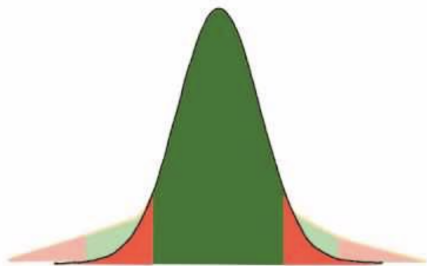
- consider the assumption that data comes from heavy-tailed distribution



- Problem:** estimate the covariance matrix in high-dimensional case for heavy-tailed distributions.

What to do in the case of noisy data/data with outliers

- consider the assumption that data comes from heavy-tailed distribution



- Problem: estimate the covariance matrix in high-dimensional case for heavy-tailed distributions.**

Heavy-tailed distributions

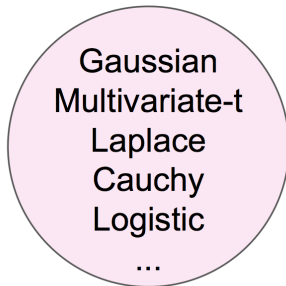
Elliptical distribution

- The probability density function of the elliptical distribution is:

$$f(\mathbf{x}) = C_{p,g} |\Sigma|^{-1/2} g\left((\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right),$$

where $g(t)$ is non-negative generator function, $C_{p,g}$ is the normalisation constant

Elliptical family



Tyler's estimator

- $\mathbf{X}_i \sim \text{elliptical}(0, \Sigma)$ (we don't specify which distribution in elliptical family)
- Normalized sample $\mathbf{Z}_i \triangleq \frac{\mathbf{X}_i - \mu}{\|\mathbf{X}_i - \mu\|_2}$

- pdf, Angular Central Gaussian Distribution:

$$f(\mathbf{z}) = C |\Sigma|^{-1/2} (\mathbf{z}^T \Sigma^{-1} \mathbf{z})^{-p/2}$$

- negative log-likelihood function:

$$\frac{n}{2} \log |\Sigma| + \frac{p}{2} \sum_{i=1}^n \log(\mathbf{Z}_i^T \Sigma^{-1} \mathbf{Z}_i)$$

- Tyler [1987] proposed covariance estimator $\hat{\Sigma}$ as solution to

$$\Sigma = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{Z}_i \mathbf{Z}_i^T}{\mathbf{Z}_i^T \Sigma^{-1} \mathbf{Z}_i}$$

Tyler's estimator

- $\mathbf{X}_i \sim \text{elliptical}(0, \Sigma)$ (we don't specify which distribution in elliptical family)
- Normalized sample $\mathbf{Z}_i \triangleq \frac{\mathbf{X}_i - \mu}{\|\mathbf{X}_i - \mu\|_2}$
 - pdf, Angular Central Gaussian Distribution:

$$f(\mathbf{z}) = C |\Sigma|^{-1/2} (\mathbf{z}^T \Sigma^{-1} \mathbf{z})^{-p/2}$$

- negative log-likelihood function:

$$\frac{n}{2} \log |\Sigma| + \frac{p}{2} \sum_{i=1}^n \log(\mathbf{Z}_i^T \Sigma^{-1} \mathbf{Z}_i)$$

- Tyler [1987] proposed covariance estimator $\hat{\Sigma}$ as solution to

$$\Sigma = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{Z}_i \mathbf{Z}_i^T}{\mathbf{Z}_i^T \Sigma^{-1} \mathbf{Z}_i}$$

Tyler's estimator

- $\mathbf{X}_i \sim \text{elliptical}(0, \Sigma)$ (we don't specify which distribution in elliptical family)
- Normalized sample $\mathbf{Z}_i \triangleq \frac{\mathbf{X}_i - \boldsymbol{\mu}}{\|\mathbf{X}_i - \boldsymbol{\mu}\|_2}$
 - pdf, Angular Central Gaussian Distribution:

$$f(\mathbf{z}) = C |\Sigma|^{-1/2} (\mathbf{z}^T \Sigma^{-1} \mathbf{z})^{-p/2}$$

- negative log-likelihood function:

$$\frac{n}{2} \log |\Sigma| + \frac{p}{2} \sum_{i=1}^n \log(\mathbf{Z}_i^T \Sigma^{-1} \mathbf{Z}_i)$$

- Tyler [1987] proposed covariance estimator $\hat{\Sigma}$ as solution to

$$\Sigma = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{Z}_i \mathbf{Z}_i^T}{\mathbf{Z}_i^T \Sigma^{-1} \mathbf{Z}_i}$$

Tyler's estimator

- $\mathbf{X}_i \sim \text{elliptical}(0, \Sigma)$ (we don't specify which distribution in elliptical family)
- Normalized sample $\mathbf{Z}_i \triangleq \frac{\mathbf{X}_i - \mu}{\|\mathbf{X}_i - \mu\|_2}$
 - pdf, Angular Central Gaussian Distribution:

$$f(\mathbf{z}) = C|\Sigma|^{-1/2}(\mathbf{z}^T \Sigma^{-1} \mathbf{z})^{-p/2}$$

- negative log-likelihood function:

$$\frac{n}{2} \log |\Sigma| + \frac{p}{2} \sum_{i=1}^n \log(\mathbf{Z}_i^T \Sigma^{-1} \mathbf{Z}_i)$$

- Tyler [1987] proposed covariance estimator $\hat{\Sigma}$ as solution to

$$\Sigma = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{Z}_i \mathbf{Z}_i^T}{\mathbf{Z}_i^T \Sigma^{-1} \mathbf{Z}_i}$$

Tyler's estimator

- $\mathbf{X}_i \sim \text{elliptical}(0, \Sigma)$ (we don't specify which distribution in elliptical family)
- Normalized sample $\mathbf{Z}_i \triangleq \frac{\mathbf{X}_i - \boldsymbol{\mu}}{\|\mathbf{X}_i - \boldsymbol{\mu}\|_2}$
 - pdf, Angular Central Gaussian Distribution:

$$f(\mathbf{z}) = C|\Sigma|^{-1/2}(\mathbf{z}^T \Sigma^{-1} \mathbf{z})^{-p/2}$$

- negative log-likelihood function:

$$\frac{n}{2} \log |\Sigma| + \frac{p}{2} \sum_{i=1}^n \log(\mathbf{Z}_i^T \Sigma^{-1} \mathbf{Z}_i)$$

- Tyler [1987] proposed covariance estimator $\hat{\Sigma}$ as solution to

$$\Sigma = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{Z}_i \mathbf{Z}_i^T}{\mathbf{Z}_i^T \Sigma^{-1} \mathbf{Z}_i}$$

Tyler's M-estimator

- Tyler's estimator:

$$\Sigma = \sum_{i=1}^n w_i \mathbf{z}_i \mathbf{z}_i^T$$
$$w_i = \frac{p}{n} \frac{1}{\mathbf{z}_i^T \Sigma^{-1} \mathbf{z}_i}$$

- Why is Tyler's estimator robust to outliers? ☺
 - consider distance $d_i = \sqrt{\mathbf{z}_i^T \Sigma^{-1} \mathbf{z}_i}$
 - $w_i \propto \frac{1}{d_i^2}$, outliers are down-weighted
- Fixed-point equation, iterative algorithm:

$$\tilde{\Sigma}_{t+1} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^T}{\mathbf{z}_i \hat{\Sigma}_t^{-1} \mathbf{z}_i}$$

$$\hat{\Sigma}_{t+1} = \frac{\tilde{\Sigma}_{t+1}}{\text{tr}(\tilde{\Sigma}_{t+1})/p}$$

- existence condition: $n > p$

Tyler's M-estimator

- Tyler's estimator:

$$\Sigma = \sum_{i=1}^n w_i \mathbf{z}_i \mathbf{z}_i^T$$
$$w_i = \frac{p}{n} \frac{1}{\mathbf{z}_i^T \Sigma^{-1} \mathbf{z}_i}$$

- Why is Tyler's estimator robust to outliers? ☺
 - consider distance $d_i = \sqrt{\mathbf{z}_i^T \Sigma^{-1} \mathbf{z}_i}$
 - $w_i \propto \frac{1}{d_i^2}$, outliers are down-weighted
- Fixed-point equation, iterative algorithm:

$$\tilde{\Sigma}_{t+1} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^T}{\mathbf{z}_i \hat{\Sigma}_t^{-1} \mathbf{z}_i}$$

$$\hat{\Sigma}_{t+1} = \frac{\tilde{\Sigma}_{t+1}}{\text{tr}(\tilde{\Sigma}_{t+1})/p}$$

- existence condition: $n > p$

Tyler's M-estimator

- Tyler's estimator:

$$\Sigma = \sum_{i=1}^n w_i \mathbf{z}_i \mathbf{z}_i^T$$
$$w_i = \frac{p}{n} \frac{1}{\mathbf{z}_i^T \Sigma^{-1} \mathbf{z}_i}$$

- Why is Tyler's estimator robust to outliers? ☺
 - consider distance $d_i = \sqrt{\mathbf{z}_i^T \Sigma^{-1} \mathbf{z}_i}$
 - $w_i \propto \frac{1}{d_i^2}$, outliers are down-weighted
- Fixed-point equation, iterative algorithm:

$$\tilde{\Sigma}_{t+1} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^T}{\mathbf{z}_i \hat{\Sigma}_t^{-1} \mathbf{z}_i}$$

$$\hat{\Sigma}_{t+1} = \frac{\tilde{\Sigma}_{t+1}}{\text{tr}(\tilde{\Sigma}_{t+1})/p}$$

- existence condition: $n > p$

Tyler's M-estimator

- Tyler's estimator:

$$\Sigma = \sum_{i=1}^n w_i \mathbf{z}_i \mathbf{z}_i^T$$
$$w_i = \frac{p}{n} \frac{1}{\mathbf{z}_i^T \Sigma^{-1} \mathbf{z}_i}$$

- Why is Tyler's estimator robust to outliers? ☺
 - consider distance $d_i = \sqrt{\mathbf{z}_i^T \Sigma^{-1} \mathbf{z}_i}$
 - $w_i \propto \frac{1}{d_i^2}$, outliers are down-weighted
- Fixed-point equation, iterative algorithm:

$$\tilde{\Sigma}_{t+1} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^T}{\mathbf{z}_i \hat{\Sigma}_t^{-1} \mathbf{z}_i}$$
$$\hat{\Sigma}_{t+1} = \frac{\tilde{\Sigma}_{t+1}}{\text{tr}(\tilde{\Sigma}_{t+1})/p}$$

- existence condition: $n > p$

How to deal with small sample scenario?
Shrink! As for the sample covariance matrix



Shrinkage covariance matrix: modified Tyler's estimator

- Modified Tyler's estimator [Chen *et al.* [2011]]

$$\tilde{\Sigma}_{t+1} = (1 - \rho) \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^T}{\mathbf{z}_i \hat{\Sigma}_t^{-1} \mathbf{z}_i} + \rho \mathbf{I}, \quad \rho \in [0, 1]$$

$$\hat{\Sigma}_{t+1} = \frac{\tilde{\Sigma}_{t+1}}{\text{tr}(\tilde{\Sigma}_{t+1})/p}$$

- the second step solves the identifiability issue
- Provable convergence
- Systematic way of choosing parameter ρ

Shrinkage covariance matrix: modified Tyler's estimator

- Modified Tyler's estimator [Chen *et al.* [2011]]

$$\tilde{\Sigma}_{t+1} = (1 - \rho) \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^T}{\mathbf{z}_i \hat{\Sigma}_t^{-1} \mathbf{z}_i} + \rho \mathbf{I}, \quad \rho \in [0, 1]$$

$$\hat{\Sigma}_{t+1} = \frac{\tilde{\Sigma}_{t+1}}{\text{tr}(\tilde{\Sigma}_{t+1})/p}$$

- the second step solves the identifiability issue
- Provable convergence
- Systematic way of choosing parameter ρ

Ok! But what if the mean value is unknown
(cannot use Tyler's estimator anymore).
We still want to deal with heavy-tailed data!

- **Problem: estimate the covariance matrix in high-dimensional case for heavy-tailed distributions with the unknown mean vector**
 - Consider the specific distribution in the class of elliptical

My work 😊

- **Problem: estimate the covariance matrix in high-dimensional case for heavy-tailed distributions with the unknown mean vector**
 - Consider the specific distribution in the class of elliptical

Small Sample Regime & Robust Mean-Covariance Estimators

Elliptical distribution

- fix the distribution in the elliptical family:

$$f(\mathbf{x}) = C_{\rho, g} |\Sigma|^{-1/2} g\left((\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu)\right),$$

- find the MAP estimator:

$$\widehat{\Sigma}_\rho = (1 - \rho) \frac{1}{n} \sum_{i=1}^n u(t) (\mathbf{X}_i - \mu)(\mathbf{X}_i - \mu)^\top + \rho I,$$

$$u(t) = -2 \frac{g'(t)}{g(t)},$$

$$t = (\mathbf{X}_i - \mu)^\top \widehat{\Sigma}_\rho^{-1} (\mathbf{X}_i - \mu)$$

- we find the ρ which minimize $\text{MSE}(\widehat{\Sigma}) = \mathbb{E} \left[\|\widehat{\Sigma} - \Sigma\|_F^2 \right]$
 - no formula here, too long ...

Elliptical distribution

- fix the distribution in the elliptical family:

$$f(\mathbf{x}) = C_{p,g} |\Sigma|^{-1/2} g\left((\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu)\right),$$

- find the MAP estimator:

$$\widehat{\Sigma}_\rho = (1 - \rho) \frac{1}{n} \sum_{i=1}^n u(t) (\mathbf{X}_i - \mu)(\mathbf{X}_i - \mu)^\top + \rho \mathbf{I},$$

$$u(t) = -2 \frac{g'(t)}{g(t)},$$

$$t = (\mathbf{X}_i - \mu)^\top \widehat{\Sigma}_\rho^{-1} (\mathbf{X}_i - \mu)$$

- we find the ρ which minimize $\text{MSE}(\widehat{\Sigma}) = \mathbb{E} \left[\|\widehat{\Sigma} - \Sigma\|_F^2 \right]$
 - no formula here, too long ...

Elliptical distribution

- fix the distribution in the elliptical family:

$$f(\mathbf{x}) = C_{p,g} |\Sigma|^{-1/2} g\left((\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu)\right),$$

- find the MAP estimator:

$$\widehat{\Sigma}_\rho = (1 - \rho) \frac{1}{n} \sum_{i=1}^n u(t) (\mathbf{X}_i - \mu)(\mathbf{X}_i - \mu)^\top + \rho \mathbf{I},$$

$$u(t) = -2 \frac{g'(t)}{g(t)},$$

$$t = (\mathbf{X}_i - \mu)^\top \widehat{\Sigma}_\rho^{-1} (\mathbf{X}_i - \mu)$$

- we find the ρ which minimize $\text{MSE}(\widehat{\Sigma}) = \mathbb{E} \left[\|\widehat{\Sigma} - \Sigma\|_F^2 \right]$
 - no formula here, too long ...

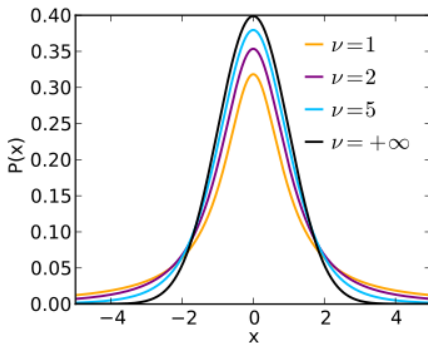
Example: Multivariate t-distribution

- Multivariate t-distribution with degree of freedom ν

$$t_p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu) = \frac{\Gamma(\frac{\nu+p}{2})}{(\pi\nu)^{p/2}\Gamma(\frac{\nu}{2})|\boldsymbol{\Sigma}|^{1/2}} \left[1 + \frac{\delta(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\nu}\right]^{-\frac{\nu+p}{2}},$$

$\delta(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$ is the Mahalanobis distance between \mathbf{x} and $\boldsymbol{\mu}$

- t-distribution has heavier tail than normal



MAP estimation

- $\mathbf{X}_i \sim t_p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu)$
- Algorithm:

$$\hat{\boldsymbol{\mu}}_{t+1} = \frac{\sum_{i=1}^n \tau_i^{t+1} \mathbf{X}_i}{\sum_{i=1}^n \tau_i^{t+1}}, \quad \tau_i^{t+1} = \frac{\nu + \rho}{\nu + (\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t) \hat{\boldsymbol{\Sigma}}_t^{-1} (\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t)}$$
$$\hat{\boldsymbol{\Sigma}}_{t+1} = (1 - \rho) \frac{\rho + \nu}{n} \sum_{i=1}^n \frac{(\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t)(\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t)^T}{(\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t) \hat{\boldsymbol{\Sigma}}_t^{-1} (\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t) + \nu} + \rho \mathbf{I},$$

- ρ can be found by minimizing $\text{MSE}(\hat{\boldsymbol{\Sigma}}_{t+1})$

MAP estimation

- $\mathbf{X}_i \sim t_p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu)$
- Algorithm:

$$\hat{\boldsymbol{\mu}}_{t+1} = \frac{\sum_{i=1}^n \tau_i^{t+1} \mathbf{X}_i}{\sum_{i=1}^n \tau_i^{t+1}}, \quad \tau_i^{t+1} = \frac{\nu + p}{\nu + (\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t) \hat{\boldsymbol{\Sigma}}_t^{-1} (\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t)}$$
$$\hat{\boldsymbol{\Sigma}}_{t+1} = (1 - \rho) \frac{p + \nu}{n} \sum_{i=1}^n \frac{(\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t)(\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t)^T}{(\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t) \hat{\boldsymbol{\Sigma}}_t^{-1} (\mathbf{X}_i - \hat{\boldsymbol{\mu}}_t) + \nu} + \rho \mathbf{I},$$

- ρ can be find by minimizing $\text{MSE}(\hat{\boldsymbol{\Sigma}}_{t+1})$

Optimal parameter

Theorem

The optimal shrinkage parameter ρ for the t -distribution with $\nu > 0$ degrees of freedom:

$$\rho^* = \frac{\text{tr}(\Sigma^2) \left(1 + \frac{\nu}{p} - \frac{2}{p}\right) + p(\nu + p)}{\text{tr}(\Sigma^2) \left((n+1) \left(\frac{\nu}{p} + 1\right) + \frac{2}{p}(n-1)\right) + (p+\nu)(p-n) - 2n}$$

- ρ^* depends on $\text{tr}(\Sigma^2)$, ν which are unknown
 - to estimate $\text{tr}(\Sigma^2)$ we use the normalized sample covariance matrix:

$$\widehat{R} = \frac{p}{n} \sum_{i=1}^n \frac{(\mathbf{X}_i - \widehat{\mathbf{X}})(\mathbf{X}_i - \widehat{\mathbf{X}})^T}{\|\mathbf{X}_i - \widehat{\mathbf{X}}\|^2}$$

- To estimate degrees of freedom parameter ν we use extreme values theory, tail-index

Optimal parameter

Theorem

The optimal shrinkage parameter ρ for the t -distribution with $\nu > 0$ degrees of freedom:

$$\rho^* = \frac{\text{tr}(\Sigma^2) \left(1 + \frac{\nu}{p} - \frac{2}{p}\right) + p(\nu + p)}{\text{tr}(\Sigma^2) \left((n+1) \left(\frac{\nu}{p} + 1\right) + \frac{2}{p}(n-1)\right) + (p+\nu)(p-n) - 2n}$$

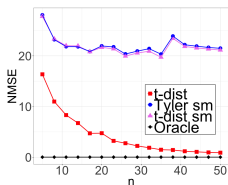
- ρ^* depends on $\text{tr}(\Sigma^2)$, ν which are unknown
 - to estimate $\text{tr}(\Sigma^2)$ we use the normalized sample covariance matrix:

$$\widehat{R} = \frac{p}{n} \sum_{i=1}^n \frac{(\mathbf{X}_i - \widehat{\mathbf{X}})(\mathbf{X}_i - \widehat{\mathbf{X}})^T}{\|\mathbf{X}_i - \widehat{\mathbf{X}}\|^2}$$

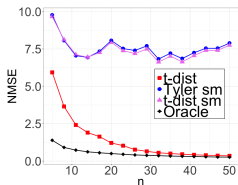
- To estimate degrees of freedom parameter ν we use extreme values theory, tail-index

Experiments

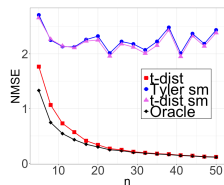
Multivariate t-distribution, $\nu = \{1, 3\}$



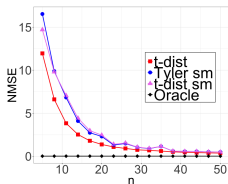
(a) $r = 0.1, \nu = 1$



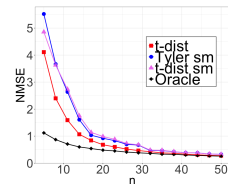
(b) $r = 0.7, \nu = 1$



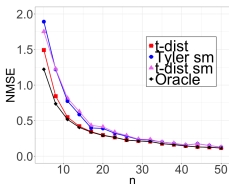
(c) $r = 0.9, \nu = 1$



(d) $r = 0.1, \nu = 3$



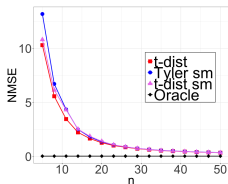
(e) $r = 0.7, \nu = 3$



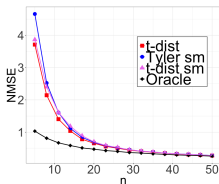
(f) $r = 0.9, \nu = 3$

Figure: AR(r) process: comparison of covariance estimators when $p = 50$ and $r \in \{0.1, 0.5, 0.9\}$ and the samples are from multivariate t-distribution with $\nu \in \{1, 2, 3, 6, 10\}$ degrees of freedom; μ is fixed to be 5 in all simulations

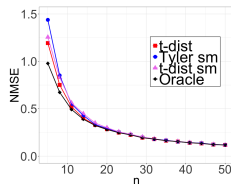
Multivariate t-distribution, $\nu = \{6, 10\}$



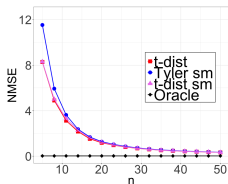
(a) $r = 0.1$, $\nu = 6$



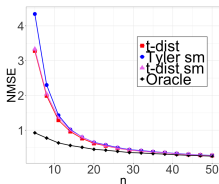
(b) $r = 0.7$, $\nu = 6$



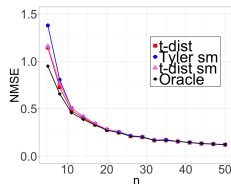
(c) $r = 0.9$, $\nu = 6$



(d) $r = 0.1$, $\nu = 10$



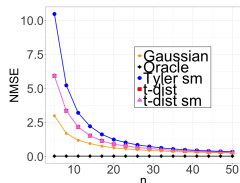
(e) $r = 0.7$, $\nu = 10$



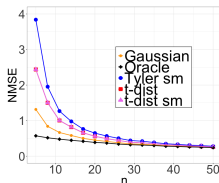
(f) $r = 0.9$, $\nu = 10$

Figure: AR(r) process: comparison of covariance estimators when $p = 50$ and $r \in \{0.1, 0.5, 0.9\}$ and the samples are from multivariate t-distribution with $\nu \in \{1, 2, 3, 6, 10\}$ degrees of freedom; μ is fixed to be 5 in all simulations

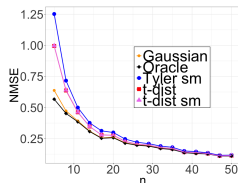
Gaussian distribution



(a) $r = 0.1$



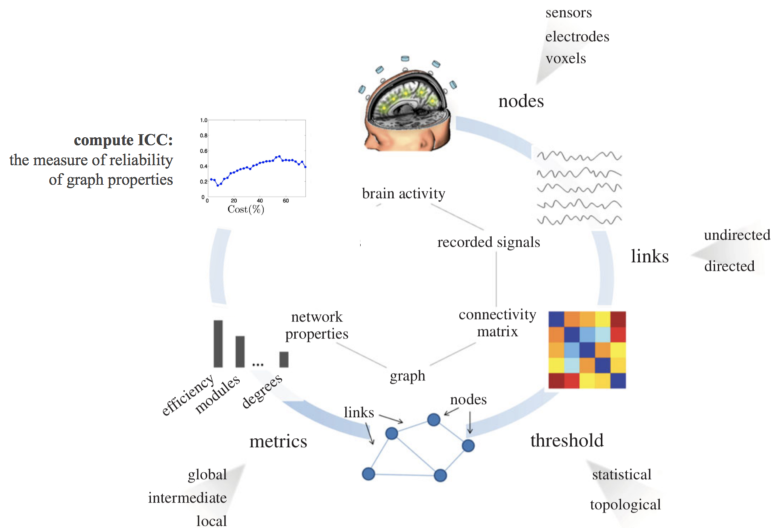
(b) $r = 0.7$



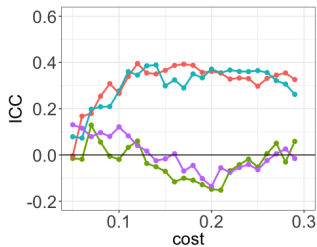
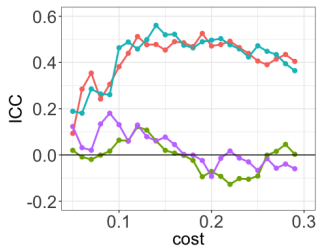
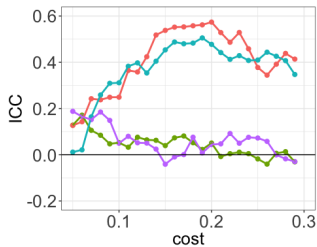
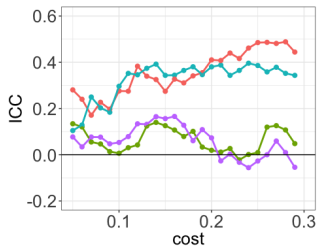
(c) $r = 0.9$

Figure: AR(r) process: $[\Sigma]_{ij} = \varrho^{|i-j|}$, $\varrho \in (0, 1)$. Comparison of covariance estimators when $p = 50$ and $r \in \{0.1, 0.5, 0.9\}$ and the samples are from multivariate normal distribution; μ is fixed to be 5 in all simulations

Application



ICC values



lw sample pc shrink t-dist t-dist

- CHEN, YILUN, WIESEL, AMI, & HERO, ALFRED O. 2011. Robust shrinkage estimation of high-dimensional covariance matrices. *IEEE Transactions on Signal Processing*, **59**(9), 4097–4107.
- TYLER, DAVID E. 1987. A distribution-free M-estimator of multivariate scatter. *The Annals of Statistics*, 234–251.