

Regularized M-estimators of the Covariance Matrix

Esa Ollila

Aalto University, Department of Signal Processing and Acoustics, Finland
esa.ollila@aalto.fi <http://signal.hut.fi/~esollila/>

Summer School, Rüdesheim, Sep 22, 2016



Aalto University

Contents

Part A Regularized M -estimators of covariance:

- M -estimation and geodesic (g -)convexity
- Regularization via g -convex penalties
- Penalty (tuning) parameter selection

Part B Applications

- Regularized discriminant analysis
- Radar detection
- Optimal portfolio selection

Covariance estimation problem

- \mathbf{x} : p -variate (centered) random vector
- $\mathbf{x}_1, \dots, \mathbf{x}_n$ i.i.d. realizations of \mathbf{x} , collected in data matrix

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix}$$

- Problem: Find an estimate $\hat{\Sigma} = \hat{\Sigma}(\{\mathbf{x}_i\}_{i=1}^n)$ of the positive definite covariance matrix

$$\Sigma = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] \in \mathcal{S}(p)$$

- Solution: Maximum likelihood, M -estimation.

Conventional estimate: the sample covariance matrix (SCM)

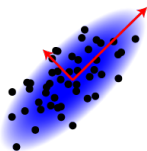
$$\mathbf{S} = \frac{1}{n} \mathbf{X}^\top \mathbf{X} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$$

Why covariance estimation?

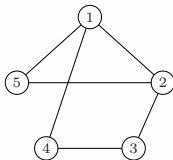
Portfolio selection



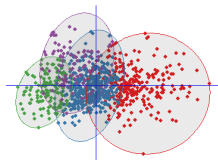
PCA



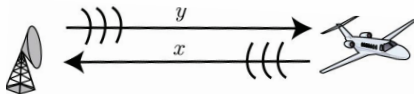
Graphical models



Discriminant Analysis



Radar detection



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

Covariance estimation challenges

- 1 **Insufficient sample support (ISS)** case: $p > n$.
 \Rightarrow Estimate of Σ^{-1} can not be computed!
- 2 **Low sample support (LSS)** (i.e., p of the same magnitude as n)
 $\Rightarrow \hat{\Sigma}$ is estimated with a lot of error.
- 3 **Outliers** or heavy-tailed **non-Gaussian** data
 $\Rightarrow \hat{\Sigma}$ is completely corrupted.

Problem 1 & 2 = **Sparse data**
 \Rightarrow **regularization** (this talk)
 \Rightarrow **RMT** (Frederic's talk)

Problem 3
 \Rightarrow **robust estimation**

Why robustness?

- 1 Outliers difficult to glean from high-dimensional data sets
- 2 Impulsive measurement environments (e.g., fMRI)
- 3 SCM is vulnerable to outliers and inefficient under non-Gaussianity
- 4 Most robust estimators can not be computed in $p > n$ cases

Contents

I. Ad-hoc shrinkage SCM-s of multiple samples

II. ML- and M -estimators of scatter matrix

III. Geodesic convexity

- Geodesic
- g -convex functions

IV. Regularized M -estimators

- Shrinkage towards an identity matrix
- Shrinkage towards a target matrix

V. Penalized estimation of multiple covariances

VI. Estimation of the regularization parameter

- Cross-validation
- Oracle approach

VII. Applications

- Regularized discriminant analysis
- Matched filter detection

Acknowledgement

- To my co-authors:



David E. Tyler

Rutgers University



Ami Wiesel

Hebrew U. Jerusalem



Ilya Soloveychik

Hebrew U. Jerusalem

- and many inspiring people working in this field:

Frederic Pascal, Teng Zhang, Lutz Dümbgen, Romain Couillet,
Matthew R. McKay, Yuri Abramovich, Olivier Besson, Maria
Greco, Fulvio Gini, Daniel Palomar,

Menu

I. Ad-hoc shrinkage SCM-s of multiple samples

II. ML- and M -estimators of scatter matrix

III. Geodesic convexity

IV. Regularized M -estimators

V. Penalized estimation of multiple covariances

VI. Estimation of the regularization parameter

VII. Applications

Multiple covariance estimation problem

- We are given K groups of elliptically distributed measurements,

$$\mathbf{x}_{11}, \dots, \mathbf{x}_{1n_1}, \quad \dots, \quad \mathbf{x}_{K1}, \dots, \mathbf{x}_{Kn_K}$$

- Each group $\mathbf{X}_k = (\mathbf{x}_{k1} \cdots \mathbf{x}_{kn_k})$ containing n_k p -dimensional observation vectors, and

$$N = \sum_{i=1}^K n_k = \text{total sample size}$$

$$\pi_k = \frac{n_k}{N} = \text{relative sample size of the } k\text{-th group}$$

- Sample populations follow elliptical distributions, $\mathcal{E}_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, g_k)$, with different scatter matrices $\boldsymbol{\Sigma}_k$ possessing mutual structure or a **joint center** $\boldsymbol{\Sigma} \Rightarrow$ need to estimate **both** $\{\boldsymbol{\Sigma}_k\}_{k=1}^K$ **and** $\boldsymbol{\Sigma}$.
- We assume that the symmetry center $\boldsymbol{\mu}_k$ of populations is known or that the data sets are *centered*.

Ad-hoc regularization approach

- Gaussian MLE-s of $\Sigma_1, \dots, \Sigma_K$ are the SCM-s $\mathbf{S}_1, \dots, \mathbf{S}_K$
- If n_k small relative to p , common assumption is $\Sigma_1 = \dots = \Sigma_K$ which is estimated by **pooled SCM**

$$\mathbf{S} = \sum_{k=1}^K \pi_k \mathbf{S}_k.$$

- Rather than assume the population covariance matrices are all equal (*hard modeling*), simply shrink them towards equality (*soft modeling*):

$$\mathbf{S}_k(\beta) = \beta \mathbf{S}_k + (1 - \beta) \mathbf{S},$$

e.g., as in [Friedman, 1989], where $\beta \in (0, 1)$ is a regularization parameter, commonly chosen by cross-validation.

- If the the total sample size N is also small relative to dimension p , then Friedman recommends also shrinking the pooled SCM \mathbf{S} towards $\propto \mathbf{I}$.

Regularized covariance matrices

- Q1 Can the Ad-Hoc method be improved or some theory/formalism put behind it?
- Q2 Robustness and resistance, e.g., non-Gaussian models and outliers.
- Q3 Methods other than convex combinations?
- Q4 Shrinkage towards other models?
 - E.g., proportional covariance matrices instead of common covariance matrices?
 - Other types of shrinkage to the structure?

Q1: Some formalism to the Ad-Hoc method

- Gaussian ML cost function ($-2 \times$ neg. log-likelihood) for the k th class:

$$\mathcal{L}_{G,k}(\Sigma_k) = \text{Tr}(\Sigma_k^{-1} \mathbf{S}_k) - \log |\Sigma_k^{-1}|$$

has a unique minimizer at $\hat{\Sigma}_k = \mathbf{S}_k$ (= SCM of the k th sample).

- Penalized objective function: Add a penalty term and solve

$$\min_{\Sigma_k \in \mathcal{S}(p)} \left\{ \mathcal{L}_{G,k}(\Sigma_k) + \lambda d(\Sigma_k, \hat{\Sigma}) \right\}, \quad k = 1, \dots, K,$$

where

- $\lambda > 0$ is a **penalty/regularization** parameter
- $d(\mathbf{A}, \mathbf{B}) : \mathcal{S}(p) \times \mathcal{S}(p) \rightarrow \mathbb{R}_0^+$ is a **penalty/distance function** minimized whenever $\mathbf{A} = \mathbf{B}$

Idea: Penalty shrinks $\hat{\Sigma}_k$ towards a fixed **shrinkage target matrix** $\hat{\Sigma} \in \mathcal{S}(p)$, the amount of shrinkage depends on the magnitude of λ

Q1: Some formalism to the Ad-Hoc method

- The information theoretic **Kullback-Leibler (KL) divergence** [Cover and Thomas, 2012], distance from $\mathcal{N}_p(\mathbf{0}, \mathbf{A})$ to $\mathcal{N}_p(\mathbf{0}, \mathbf{B})$, is

$$d_{\text{KL}}(\mathbf{A}, \mathbf{B}) = \text{Tr}(\mathbf{A}^{-1}\mathbf{B}) - \log |\mathbf{A}^{-1}\mathbf{B}| - p.$$

As is well known, it verifies $d_{\text{KL}}(\mathbf{A}, \mathbf{B}) \geq 0$ and $= 0$ for $\mathbf{A} = \mathbf{B}$.

- Using $d_{\text{KL}}(\Sigma_k, \hat{\Sigma})$ as the penalty, the optimization problem $\mathcal{L}_{G,k}(\Sigma_k) + \lambda d_{\text{KL}}(\Sigma_k, \hat{\Sigma})$ possesses a unique solution given by

$$\mathbf{S}_k(\beta) = \beta \mathbf{S}_k + (1 - \beta) \hat{\Sigma}, \quad k = 1, \dots, K,$$

where $\beta = (1 + \lambda)^{-1} \in (0, 1)$ and $k = 1, \dots, K$.

- This gives Friedman's Ad-Hoc shrinkage SCM estimators when the shrinkage target matrix $\hat{\Sigma}$ is the **pooled SCM \mathbf{S}**

Q1: Some formalism to the Ad-Hoc method

- The information theoretic **Kullback-Leibler (KL) divergence** [Cover and Thomas, 2012], distance from $\mathcal{N}_p(\mathbf{0}, \mathbf{A})$ to $\mathcal{N}_p(\mathbf{0}, \mathbf{B})$, is

$$d_{\text{KL}}(\mathbf{A}, \mathbf{B}) = \text{Tr}(\mathbf{A}^{-1}\mathbf{B}) - \log |\mathbf{A}^{-1}\mathbf{B}| - p.$$

As is well known, it verifies $d_{\text{KL}}(\mathbf{A}, \mathbf{B}) \geq 0$ and $= 0$ for $\mathbf{A} = \mathbf{B}$.

- Using $d_{\text{KL}}(\boldsymbol{\Sigma}_k, \hat{\boldsymbol{\Sigma}})$ as the penalty, the optimization problem $\mathcal{L}_{G,k}(\boldsymbol{\Sigma}_k) + \lambda d_{\text{KL}}(\boldsymbol{\Sigma}_k, \hat{\boldsymbol{\Sigma}})$ possesses a unique solution given by

$$\mathbf{S}_k(\beta) = \beta \mathbf{S}_k + (1 - \beta) \mathbf{S}, \quad k = 1, \dots, K,$$

where $\beta = (1 + \lambda)^{-1} \in (0, 1)$ and $k = 1, \dots, K$.

- This gives Friedman's Ad-Hoc shrinkage SCM estimators when the shrinkage target matrix $\hat{\boldsymbol{\Sigma}}$ is the **pooled SCM S**

Discussion

Note: The Gaussian likelihood $\mathcal{L}_{G,k}(\Sigma_k)$ is convex in Σ_k^{-1} and so is $d_{\text{KL}}(\Sigma_k, \hat{\Sigma})$.

Comments

- Other (non-Gaussian) ML cost functions $\mathcal{L}_k(\Sigma)$ are commonly not convex in Σ^{-1}
- Swapping the order $d_{\text{KL}}(\Sigma_k, \hat{\Sigma})$ to $d_{\text{KL}}(\hat{\Sigma}, \Sigma_k)$ gives a distance function that is non-convex in Σ_k^{-1} .

Problems

- The penalized optimization program, $\mathcal{L}_{G,k}(\Sigma_k) + \lambda d_{\text{KL}}(\Sigma_k, \hat{\Sigma})$, **does not seem to generalize** to using other distance functions or other non-Gaussian cost functions.
- KL-distance $d_{\text{KL}}(\Sigma_k, \Sigma)$ is not so useful when the assumption is $\Sigma_k \propto \Sigma$, i.e., proportional covariance matrices.

How about a robust Ad-hoc method?

- **Plug-In Robust Estimators:** Let $\hat{\Sigma}_k$ and $\hat{\Sigma}$ represent robust estimates of scatter (covariance) matrix for the k th class and the pooled data respectively.
- Then a robust version of Friedman's approach is given by

$$\hat{\Sigma}_k(\beta) = \beta \hat{\Sigma}_k + (1 - \beta) \hat{\Sigma}, \quad k = 1, \dots, K$$

where $\beta \in (0, 1)$.

- **Problems:** This approach fails since many robust estimators of scatter, e.g. M, S, MM, MCD, etc., are not defined or do not vary much from the sample covariance when the data is sparse.

Our approach in this series of lectures

- **Regularization** via jointly g -convex distance functions
- **Robust M -estimation** (robust loss fnc downweights outliers)

Menu

I. Ad-hoc shrinkage SCM-s of multiple samples

II. ML- and M -estimators of scatter matrix

III. Geodesic convexity

IV. Regularized M -estimators

V. Penalized estimation of multiple covariances

VI. Estimation of the regularization parameter

VII. Applications

References



R. A. Maronna (1976).

Robust M-estimators of multivariate location and scatter.
Ann. Stat., 5(1):51–67.



D. E. Tyler (1987).

A distribution-free M-estimator of multivariate scatter.
Ann. Stat., 15(1):234–251.



E. Ollila D. E. Tyler V. Koivunen, and H. V. Poor (2012).

Complex elliptically symmetric distributions: survey, new results and applications.
IEEE Trans. Signal Processing, 60(11):5597 – 5625.

The cone of positive definite matrices $\mathcal{S}(p)$

- A square matrix \mathbf{A} is **positive definite**, denoted $\mathbf{A} \succ 0$, if it is symmetric and satisfies

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0 \quad \forall \mathbf{x} \neq \mathbf{0}$$

Positive semidefinite ($\mathbf{A} \succeq 0$): $\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0 \quad \forall \mathbf{x}$.

- $\mathbf{A} \succ 0$ ($\succeq 0$) if and only if its eigenvalues are positive (non-negative)
- **Eigenvalue decomposition (EVD)** of $\mathbf{A} \succ 0$;

$$\mathbf{A} = \mathbf{E} \mathbf{\Lambda} \mathbf{E}^\top = \sum_{i=1}^p \lambda_i \mathbf{e}_i \mathbf{e}_i^\top$$

$\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_p)$ of positive eigenvalues

$\mathbf{E} = (\mathbf{e}_1 \ \cdots \ \mathbf{e}_p)$ orthonormal eigenvectors ($\mathbf{E}^\top \mathbf{E} = \mathbf{I}$) as columns

- $\mathcal{S}(p) :=$ the set of all $p \times p$ positive definite matrices.

Elliptically symmetric (ES) distribution

$\mathbf{x} \sim \mathcal{E}_p(\mathbf{0}, \Sigma, g)$: p.d.f. is

$$f(\mathbf{x}) \propto |\Sigma|^{-1/2} g(\mathbf{x}^\top \Sigma^{-1} \mathbf{x})$$

- $\Sigma \in \mathcal{S}(p)$, unknown positive definite $p \times p$ **scatter matrix** parameter.
- $g : \mathbb{R}_0^+ \rightarrow \mathbb{R}^+$, fixed **density generator**.

When the covariance matrix exists: $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \Sigma$.

Example: Normal distribution $\mathcal{N}_p(\mathbf{0}, \Sigma)$ has p.d.f.

$$f(\mathbf{x}) = \pi^{-p/2} |\Sigma|^{-1/2} \exp\left(-\frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x}\right).$$

Elliptical distribution with $g(t) = \exp(-t/2)$.

The p -variate t -distribution with ν degrees of freedom

$\mathbf{x} \sim t_{p,\nu}(\mathbf{0}, \Sigma)$: pdf is

$$f(\mathbf{x}|\Sigma) \propto |\Sigma|^{-1} \left(1 + \mathbf{x}^\top \Sigma^{-1} \mathbf{x} / \nu\right)^{-(p+\nu)/2}, \quad \nu > 0$$

so the density generator is $g_\nu(t) = (1 + t/\nu)^{-(p+\nu)/2}$

- $\nu > 0$ is the **degrees of freedom** parameter:
 - $\nu = 1$ is called the **complex Cauchy distribution**
 - $\nu \rightarrow \infty$ yields the p -variate normal distribution
 - finite 2nd-order moments for $\nu > 2$
- Stochastic decomposition:

$$\mathbf{x} =_d \sqrt{\tau} \mathbf{n}, \quad \tau^{-1} \sim \text{Gam}(\nu/2, 2/\nu), \quad \mathbf{n} \sim \mathcal{N}_p(\mathbf{0}, \Sigma), \quad \tau \perp \mathbf{n}.$$

This also provides a straightforward approach to generate a random sample from $t_{p,\nu}(\mathbf{0}, \Sigma)$.

The maximum likelihood estimator (MLE)

- $\{\mathbf{x}_i\} \stackrel{iid}{\sim} \mathcal{E}_p(\mathbf{0}, \Sigma, g)$, where $n > p$.
- The MLE $\hat{\Sigma} \in \mathcal{S}(p)$ minimizes the $(-2/n) \times \log$ -likelihood fnc

$$\mathcal{L}(\Sigma) = \frac{1}{n} \sum_{i=1}^n -2 \ln f(\mathbf{x}_i | \Sigma) = \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}|$$

where $\rho(t) = -2 \ln g(t)$ is the loss function.

- Critical points are solutions to **estimating equations**

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top$$

where $u(t) = \rho'(t)$ is the weight function.

- MLE = "an adaptively weighted sample covariance matrix"

Gaussian MLE

- P.d.f: $f(\mathbf{x}) \propto |\Sigma|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{x}^\top \Sigma^{-1} \mathbf{x}\right)$ with density generator $g(t) = \exp(-t/2)$.
- The Gaussian loss function is $\rho_G(t) = -2 \ln g(t) = t$ and hence the $-(2/n) \times$ log-likelihood function is

$$\begin{aligned}\mathcal{L}_G(\Sigma) &= \frac{1}{n} \sum_{i=1}^n \rho_G(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| \\ &= \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i}_{=\text{Tr}(\Sigma^{-1} \mathbf{x}_i \mathbf{x}_i^\top)} - \ln |\Sigma^{-1}| \\ &= \text{Tr}(\Sigma^{-1} \mathbf{S}) - \ln |\Sigma^{-1}|\end{aligned}$$

where $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$ is the SCM.

Gaussian MLE (cont'd)

- Taking the matrix derivative w.r.t Σ^{-1} yields

$$\begin{aligned}\frac{\partial}{\partial \Sigma^{-1}} \mathcal{L}_G(\Sigma) &= \frac{\partial}{\partial \Sigma^{-1}} \left\{ \sum_{i=1}^n \text{Tr}(\Sigma^{-1} \mathbf{S}) - \ln |\Sigma^{-1}| \right\} \\ &= \mathbf{S} - \Sigma\end{aligned}$$

where we used that $\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}\mathbf{A}) = \mathbf{A}^\top$ and $\frac{\partial}{\partial \mathbf{X}} \ln |\mathbf{X}| = (\mathbf{X}^{-1})^\top$.

- Setting the derivative to $\mathbf{0}$ gives $\hat{\Sigma} = \mathbf{S}$ as the critical point.

Q: Is the found solution, i.e., SCM \mathbf{S} , also the MLE?

- Yes. We have learned this fact *surprisingly recently* from [Watson, 1962] elegant proof.
- Other way to prove this is by showing that $\mathcal{L}_G(\Sigma)$ is convex in Σ^{-1} (Note: $\mathcal{L}_G(\Sigma)$ is non-convex in Σ).
- Or: later we show that $\mathcal{L}_G(\Sigma)$ is *g-convex* both in Σ and Σ^{-1} .

Gaussian MLE (cont'd)

- Taking the matrix derivative w.r.t Σ^{-1} yields

$$\begin{aligned}\frac{\partial}{\partial \Sigma^{-1}} \mathcal{L}_G(\Sigma) &= \frac{\partial}{\partial \Sigma^{-1}} \left\{ \sum_{i=1}^n \text{Tr}(\Sigma^{-1} \mathbf{S}) - \ln |\Sigma^{-1}| \right\} \\ &= \mathbf{S} - \Sigma\end{aligned}$$

where we used that $\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}\mathbf{A}) = \mathbf{A}^\top$ and $\frac{\partial}{\partial \mathbf{X}} \ln |\mathbf{X}| = (\mathbf{X}^{-1})^\top$.

- Setting the derivative to $\mathbf{0}$ gives $\hat{\Sigma} = \mathbf{S}$ as the critical point.

Q: Is the found solution, i.e., SCM \mathbf{S} , also the MLE?

- **Yes.** We have learned this fact *surprisingly recently* from [Watson, 1962] elegant proof.
- Other way to prove this is by showing that $\mathcal{L}_G(\Sigma)$ is convex in Σ^{-1} (Note: $\mathcal{L}_G(\Sigma)$ is non-convex in Σ).
- Or: later we show that $\mathcal{L}_G(\Sigma)$ is *g-convex* both in Σ and Σ^{-1} .

Gaussian MLE (cont'd)

- Taking the matrix derivative w.r.t Σ^{-1} yields

$$\begin{aligned}\frac{\partial}{\partial \Sigma^{-1}} \mathcal{L}_G(\Sigma) &= \frac{\partial}{\partial \Sigma^{-1}} \left\{ \sum_{i=1}^n \text{Tr}(\Sigma^{-1} \mathbf{S}) - \ln |\Sigma^{-1}| \right\} \\ &= \mathbf{S} - \Sigma\end{aligned}$$

where we used that $\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}\mathbf{A}) = \mathbf{A}^\top$ and $\frac{\partial}{\partial \mathbf{X}} \ln |\mathbf{X}| = (\mathbf{X}^{-1})^\top$.

- Setting the derivative to $\mathbf{0}$ gives $\hat{\Sigma} = \mathbf{S}$ as the critical point.

Q: Is the found solution, i.e., SCM \mathbf{S} , also the MLE?

- **Yes.** We have learned this fact *surprisingly recently* from [Watson, 1962] elegant proof.
- Other way to prove this is by showing that $\mathcal{L}_G(\Sigma)$ is convex in Σ^{-1} (Note: $\mathcal{L}_G(\Sigma)$ is non-convex in Σ).
- Or: later we show that $\mathcal{L}_G(\Sigma)$ is *g-convex* both in Σ and Σ^{-1} .

Gaussian MLE (cont'd)

- Taking the matrix derivative w.r.t Σ^{-1} yields

$$\begin{aligned}\frac{\partial}{\partial \Sigma^{-1}} \mathcal{L}_G(\Sigma) &= \frac{\partial}{\partial \Sigma^{-1}} \left\{ \sum_{i=1}^n \text{Tr}(\Sigma^{-1} \mathbf{S}) - \ln |\Sigma^{-1}| \right\} \\ &= \mathbf{S} - \Sigma\end{aligned}$$

where we used that $\frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X}\mathbf{A}) = \mathbf{A}^\top$ and $\frac{\partial}{\partial \mathbf{X}} \ln |\mathbf{X}| = (\mathbf{X}^{-1})^\top$.

- Setting the derivative to $\mathbf{0}$ gives $\hat{\Sigma} = \mathbf{S}$ as the critical point.

Q: Is the found solution, i.e., SCM \mathbf{S} , also the MLE?

- **Yes.** We have learned this fact *surprisingly recently* from [Watson, 1962] elegant proof.
- Other way to prove this is by showing that $\mathcal{L}_G(\Sigma)$ is convex in Σ^{-1} (Note: $\mathcal{L}_G(\Sigma)$ is non-convex in Σ).
- Or: later we show that $\mathcal{L}_G(\Sigma)$ is *g-convex* both in Σ and Σ^{-1} .

M -estimators of scatter matrix

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top$$

[Maronna, 1976]

- Among the first proposals for robust covariance matrix estimators

Generalizations of ML-estimators:

- $u(t) = \rho'(t)$ non-neg., continuous and non-increasing.
(admits more general ρ fnc's)
- $\psi(t) = tu(t)$ strictly increasing \Rightarrow **unique solution**
- Not too much data lies in some sub-space \Rightarrow **solution exists**

Popular choices of $u(t)$: Huber's, t -likelihood, Tyler's. ...

Huber's M -estimator

- [Maronna, 1976] defined it as an M -estimator with weight fnc

$$u_H(t; c) = \begin{cases} 1/b, & \text{for } t \leq c^2 \\ c^2/(tb), & \text{for } t > c^2 \end{cases}$$

where $c > 0$ is a tuning constant, chosen by the user, and b is a scaling factor used to obtain Fisher consistency at $\mathcal{N}_p(\mathbf{0}, \Sigma)$:

$$b = F_{\chi_{p+2}^2}(c^2) + c^2(1 - F_{\chi_p^2}(c^2))/p.$$

Choose c^2 as q th upper quantile of χ_p^2 : $c^2 = F_{\chi_p^2}^{-1}(q)$.

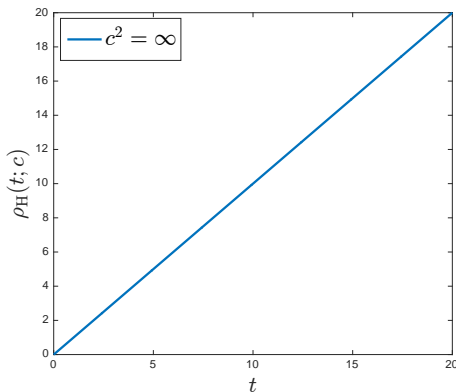
- It is also an MLE with loss function [Ollila et al., 2016]:

$$\rho_H(t; c) = \begin{cases} t/b & \text{for } t \leq c^2, \\ (c^2/b)(\log(t/c^2) + 1) & \text{for } t > c^2. \end{cases}$$

Note: a Gaussian distribution in the middle, but have tails that die down at an inverse polynomial rate. Naturally, $u_H(t; c) = \rho'_H(t; c)$.

Huber's M -estimator

$$\underset{\Sigma \in \mathcal{S}(p)}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \rho_H \left(\underbrace{\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i}_{=t}; c \right) - \ln |\Sigma^{-1}|,$$



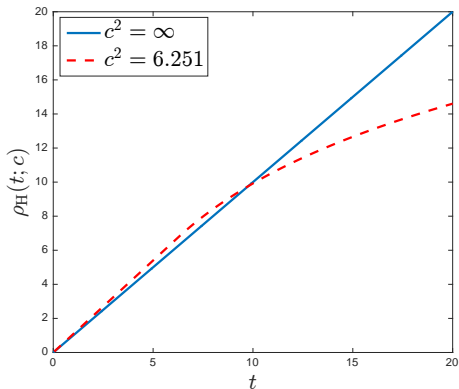
Dimension $p = 3$ and

$$c^2 = F_{\chi_p^2}^{-1}(q)$$

- $q = 1.0 : c^2 = \infty$
- $q = 0.9 : c^2 = 6.251$
- $q = 0.5 : c^2 = 2.366$
- $q = 0.1 : c^2 = 0.584$

Huber's M -estimator

$$\underset{\Sigma \in \mathcal{S}(p)}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \rho_H(\underbrace{\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i}_{=t}; c) - \ln |\Sigma^{-1}|,$$



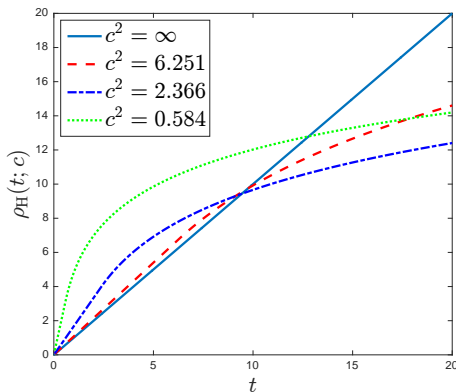
Dimension $p = 3$ and

$$c^2 = F_{\chi_p^2}^{-1}(q)$$

- $q = 1.0 : c^2 = \infty$
- $q = 0.9 : c^2 = 6.251$
- $q = 0.5 : c^2 = 2.366$
- $q = 0.1 : c^2 = 0.584$

Huber's M -estimator

$$\underset{\Sigma \in \mathcal{S}(p)}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \rho_H \left(\underbrace{\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i}_{=t}; c \right) - \ln |\Sigma^{-1}|,$$



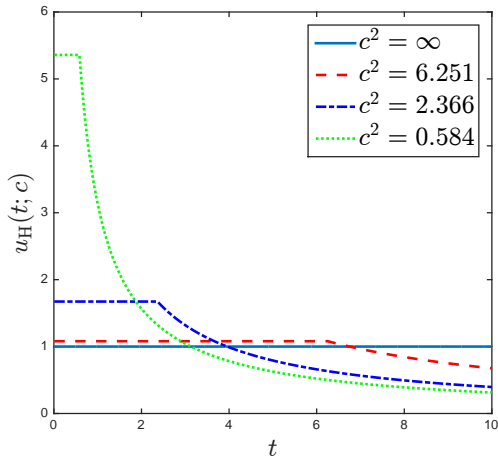
Dimension $p = 3$ and

$$c^2 = F_{\chi_p^2}^{-1}(q)$$

- $q = 1.0 : c^2 = \infty$
- $q = 0.9 : c^2 = 6.251$
- $q = 0.5 : c^2 = 2.366$
- $q = 0.1 : c^2 = 0.584$

Huber's M -estimator

$$\text{solves } \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n u_H\left(\underbrace{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i}_=t; c\right) \mathbf{x}_i \mathbf{x}_i^\top$$



Dimension $p = 3$ and

$$c^2 = F_{\chi_p^2}^{-1}(q)$$

- $q = 1.0$: $c^2 = \infty$
- $q = 0.9$: $c^2 = 6.251$
- $q = 0.5$: $c^2 = 2.366$
- $q = 0.1$: $c^2 = 0.584$

The MLE of $t_{p,\nu}$ -distribution

- Density generator of $t_{p,\nu}(\mathbf{0}, \Sigma)$ is $g_\nu(t) = (1 + t/\nu)^{-\frac{1}{2}(\nu+p)}$.
- M -estimator $\hat{\Sigma}$ based on the respective loss and weight function

$$\begin{aligned}\rho_\nu(t) &= -2 \ln g_\nu(t) = (\nu + p) \ln(1 + t/\nu) \\ u_\nu(t) &= \rho'_\nu(t) = \frac{\nu + p}{\nu + t}\end{aligned}$$

is referred to as $t_\nu M$ -estimator.

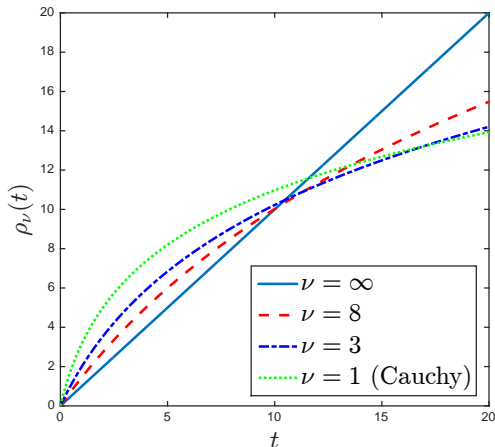
- Naturally an MLE when $\{\mathbf{x}_i\}_{i=1}^n \stackrel{iid}{\sim} t_{p,\nu}(\mathbf{0}, \Sigma)$.
- To obtain a Fisher consistent $t_\nu M$ -estimator $\hat{\Sigma}$ at $\mathcal{N}_p(\mathbf{0}, \Sigma)$, use a scaled loss function:

$$\rho_\nu^*(t) = \frac{1}{b} \rho_\nu(t), \quad b = \{(\nu + p)/p\} \mathbb{E}[\chi_p^2/(\nu + \chi_p^2)].$$

and respective scaled weight function $u_\nu^*(t) = u_\nu(t)/b$,

$t_\nu M$ -estimator

$$\underset{\Sigma \in \mathcal{S}(p)}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \rho_\nu \left(\underbrace{\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i}_{=t} \right) - \ln |\Sigma^{-1}|,$$

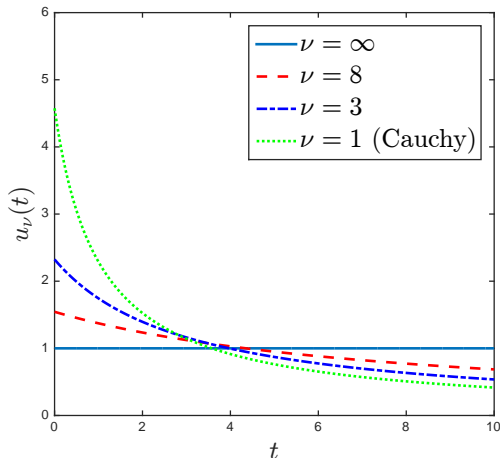


Dimension $p = 3$ and

$$\rho_\nu(t) = (\nu + p) \ln(1 + t/\nu)$$

$t_\nu M$ -estimator

$$\text{solves } \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n u_\nu \left(\underbrace{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i}_{=t} \right) \mathbf{x}_i \mathbf{x}_i^\top$$

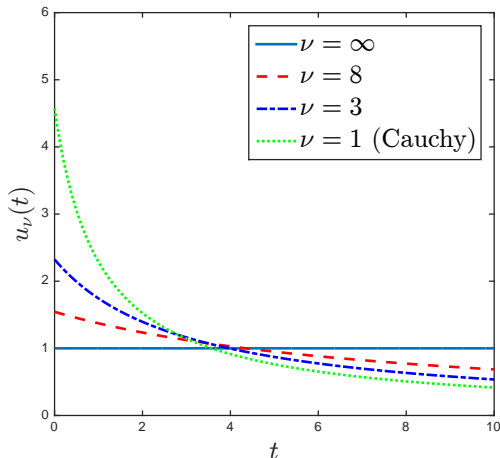


Dimension $p = 3$ and

$$u_\nu(t) = \frac{p + \nu}{\nu + t}$$

$t_\nu M$ -estimator

$$\text{solves } \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \frac{(p + \nu) \mathbf{x}_i \mathbf{x}_i^\top}{\nu + \mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i}$$



Dimension $p = 3$ and

$$u_\nu(t) = \frac{p + \nu}{\nu + t}$$

Tyler's M -estimator

- Distribution-free M -estimator (under elliptical distributions) proposed in [Tyler, 1987].
- Defined as a solution to

$$\hat{\Sigma} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i}$$

\Rightarrow so an M -estimator with Tyler's weight fnc $u(t) = \rho'(t) = p/t$

- Now it is also known that $\hat{\Sigma} \in \mathcal{S}(p)$ minimizes the cost fnc

$$\mathcal{L}_T(\Sigma) = \frac{1}{n} \sum_{i=1}^n \underbrace{p \ln(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i)}_{\rho(t) = p \ln t} - \ln |\Sigma^{-1}|$$

Note: not an MLE for any elliptical density, so $\rho(t) \neq -2 \ln g(t)$!

- Not convex in Σ ! ... or in Σ^{-1}
- Maronna's/Huber's conditions does not apply.

Tyler's M -estimator

- Distribution-free M -estimator (under elliptical distributions) proposed in [Tyler, 1987].
- Defined as a solution to

$$\hat{\Sigma} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i}$$

\Rightarrow so an M -estimator with Tyler's weight fnc $u(t) = \rho'(t) = p/t$

- Now it is also known that $\hat{\Sigma} \in \mathcal{S}(p)$ minimizes the cost fnc

$$\mathcal{L}_T(\Sigma) = \frac{1}{n} \sum_{i=1}^n \underbrace{p \ln(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i)}_{\rho(t) = p \ln t} - \ln |\Sigma^{-1}|$$

Note: not an MLE for any elliptical density, so $\rho(t) \neq -2 \ln g(t)$!

- Not convex in Σ ! ... or in Σ^{-1}

■ Maronna's/Huber's conditions does not apply.

Tyler's M -estimator

- Distribution-free M -estimator (under elliptical distributions) proposed in [Tyler, 1987].
- Defined as a solution to

$$\hat{\Sigma} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i}$$

\Rightarrow so an M -estimator with Tyler's weight fnc $u(t) = \rho'(t) = p/t$

- Now it is also known that $\hat{\Sigma} \in \mathcal{S}(p)$ minimizes the cost fnc

$$\mathcal{L}_T(\Sigma) = \frac{1}{n} \sum_{i=1}^n \underbrace{p \ln(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i)}_{\rho(t) = p \ln t} - \ln |\Sigma^{-1}|$$

Note: not an MLE for any elliptical density, so $\rho(t) \neq -2 \ln g(t)$!

- Not convex in Σ ! ... or in Σ^{-1}
- Maronna's/Huber's conditions does not apply.

Tyler's M -estimator, cont'd

$$\hat{\Sigma} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i}$$

Comments:

- 1 Limiting case of Huber's M -estimator when $c \rightarrow 0$ and of $t_\nu M$ -estimator when $\nu \rightarrow 0$.
- 2 Minimum is a unique *up to a positive scalar*: if $\hat{\Sigma}$ is a minimum then so is $b\hat{\Sigma}$ for any $b > 0$
- $\Rightarrow \hat{\Sigma}$ is a **shape matrix** estimator. We may choose a solution which verifies $|\hat{\Sigma}| = 1$.
- 3 A Fisher consistent estimator at $\mathcal{N}_p(\mathbf{0}, \Sigma)$ can be obtained by scaling any minimum $\hat{\Sigma}$ by

$$b = \text{Median}\{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i; i = 1, \dots, n\} / \text{Median}(\chi_p^2).$$

This scaling is utilized in discriminant analysis later on.

Tyler's M -estimator, cont'd

$$\hat{\Sigma} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i}$$

Comments:

- 1 Limiting case of Huber's M -estimator when $c \rightarrow 0$ and of $t_\nu M$ -estimator when $\nu \rightarrow 0$.
- 2 Minimum is a unique *up to a positive scalar*: if $\hat{\Sigma}$ is a minimum then so is $b\hat{\Sigma}$ for any $b > 0$
- $\Rightarrow \hat{\Sigma}$ is a **shape matrix** estimator. We may choose a solution which verifies $|\hat{\Sigma}| = 1$.
- 3 A Fisher consistent estimator at $\mathcal{N}_p(\mathbf{0}, \Sigma)$ can be obtained by scaling any minimum $\hat{\Sigma}$ by

$$b = \text{Median}\{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i; i = 1, \dots, n\} / \text{Median}(\chi_p^2).$$

This scaling is utilized in discriminant analysis later on.

Tyler's M -estimator, cont'd

$${}_c\hat{\Sigma} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top ({}_c\hat{\Sigma})^{-1} \mathbf{x}_i}$$

Comments:

- 1 Limiting case of Huber's M -estimator when $c \rightarrow 0$ and of $t_\nu M$ -estimator when $\nu \rightarrow 0$.
- 2 Minimum is a unique *up to a positive scalar*: if $\hat{\Sigma}$ is a minimum then so is $b\hat{\Sigma}$ for any $b > 0$
- $\Rightarrow \hat{\Sigma}$ is a **shape matrix** estimator. We may choose a solution which verifies $|\hat{\Sigma}| = 1$.
- 3 A Fisher consistent estimator at $\mathcal{N}_p(\mathbf{0}, \Sigma)$ can be obtained by scaling any minimum $\hat{\Sigma}$ by

$$b = \text{Median}\{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i; i = 1, \dots, n\} / \text{Median}(\chi_p^2).$$

This scaling is utilized in discriminant analysis later on.

Tyler's M -estimator, cont'd

$$\hat{\Sigma} = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i}$$

Comments:

- 1 Limiting case of Huber's M -estimator when $c \rightarrow 0$ and of $t_\nu M$ -estimator when $\nu \rightarrow 0$.
- 2 Minimum is a unique *up to a positive scalar*: if $\hat{\Sigma}$ is a minimum then so is $b\hat{\Sigma}$ for any $b > 0$
- \Rightarrow $\hat{\Sigma}$ is a **shape matrix** estimator. We may choose a solution which verifies $|\hat{\Sigma}| = 1$.
- 3 A Fisher consistent estimator at $\mathcal{N}_p(\mathbf{0}, \Sigma)$ can be obtained by scaling any minimum $\hat{\Sigma}$ by

$$b = \text{Median}\{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i; i = 1, \dots, n\} / \text{Median}(\chi_p^2).$$

This scaling is utilized in discriminant analysis later on.

Tyler's M -estimator as MLE

Tyler's M -estimator $\hat{\Sigma}$ is an MLE of Σ in the following 3 cases

Case I [Kent, 1997]

$\{\mathbf{x}_i\}_{i=1}^n$ i.i.d. from an **angular central Gaussian (ACG) distribution**, $\mathbf{x} \sim \mathcal{AN}_p(\mathbf{0}, \Sigma)$, whose p.d.f. is of the form

$$f(\mathbf{x}) \propto |\Sigma|^{-1/2} (\mathbf{x}^\top \Sigma^{-1} \mathbf{x})^{-p/2}, \quad \|\mathbf{x}\| = 1$$

(Note: Σ is identifiable up to a scale).

■ Stochastic decomposition:

$$\mathbf{x} =_d \frac{\mathbf{n}}{\|\mathbf{n}\|}, \quad \mathbf{n} \sim \mathcal{N}_p(\mathbf{0}, \Sigma)$$

- The term “angular central Gaussian” is a slight misnomer: $\mathcal{N}_p(\mathbf{0}, \Sigma)$ can be replaced by *any* $\mathcal{E}_p(\mathbf{0}, \Sigma, g)$.

Tyler's M -estimator as MLE

Tyler's M -estimator $\hat{\Sigma}$ is an MLE of Σ in the following 3 cases

Case I [Kent, 1997]

$\{\mathbf{x}_i\}_{i=1}^n$ i.i.d. from an **angular central Gaussian (ACG) distribution**, $\mathbf{x} \sim \mathcal{AN}_p(\mathbf{0}, \Sigma)$, whose p.d.f. is of the form

$$f(\mathbf{x}) \propto |\Sigma|^{-1/2} (\mathbf{x}^\top \Sigma^{-1} \mathbf{x})^{-p/2}, \quad \|\mathbf{x}\| = 1$$

(Note: Σ is identifiable up to a scale).

- Stochastic decomposition:

$$\mathbf{x} =_d \frac{\mathbf{n}}{\|\mathbf{n}\|}, \quad \mathbf{n} \sim \mathcal{E}_p(\mathbf{0}, \Sigma, g)$$

- The term “angular central Gaussian” is a slight misnomer: $\mathcal{N}_p(\mathbf{0}, \Sigma)$ can be replaced by *any* $\mathcal{E}_p(\mathbf{0}, \Sigma, g)$.

Case II [Gini and Greco, 2002, Conte et al., 2002]

- $\{\mathbf{x}_i\}_{i=1}^n$ independent random vectors with **Gaussian distributions** of **proportional covariance matrices**, so

$$\mathbf{x}_i \sim \mathcal{N}_p(\mathbf{0}, \eta_i \mathbf{\Sigma}), \quad \eta_i > 0, \quad \text{Tr}(\mathbf{\Sigma}) = p$$

where $\{\eta_i\}_{i=1}^n$ and $\mathbf{\Sigma} \in \mathcal{S}_p$ (with $\text{Tr}(\mathbf{\Sigma}) = p$) are the unknown parameters.

Case III [Ollila and Tyler, 2012] is a generalization of Case II

- $\{\mathbf{x}_i\}_{i=1}^n$ independent from (possibly different) **elliptical distributions** of **proportional covariance matrices**, so

$$\mathbf{x}_i \sim \mathcal{E}_p(\mathbf{0}, \eta_i \mathbf{\Sigma}, g_i), \quad \eta_i > 0, \quad \text{Tr}(\mathbf{\Sigma}) = p$$

- **Note:** density generator (need not be specified or known) and can be different for each \mathbf{x}_i , e.g., $\mathbf{x}_1 \sim \mathcal{N}_p(\mathbf{0}, \eta_1 \mathbf{\Sigma})$, $\mathbf{x}_2 \sim t_{p,\nu}(\mathbf{0}, \eta_2 \mathbf{\Sigma})$, ...

Algorithm

Given an arbitrary initial start $\Sigma_0 \in \mathcal{S}(p)$, the iterations

$$\hat{\Sigma}_{k+1} = \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \hat{\Sigma}_k^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top \quad \text{for } k = 0, 1, \dots$$

converge to the solution

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top$$

under some mild regularity conditions [$\psi(t) = tu(t)$ strictly increasing, not too much data lies in some sub-space]. See details in [\[Ollila et al., 2012\]](#).

Matlab implementation: $t_\nu M$ -estimator

```
function C = tMest(X, v)
[n, p] = size(X);
MAX_ITER = 1000; % Max number of iteration
EPS = 1.0e-5;      % Iteration accuracy
C0 = X'*X/n;       % SCM as initial estimate
invC0 = C0 \ eye(p);
iter=1;
while (iter<MAX_ITER)
    t = sum(X*invC0.*X,2);
    u = (v+p)./(v+t);
    C = X'*((X.*repmat(u,1,p)))/n;
    d = norm(eye(p)-invC0*C,Inf);
    if (d<=EPS),
        break;
    end
    invC0 = C \ eye(p);
    iter = iter+1;
end
```

Effect of an outlier

$\{\mathbf{x}_i\}_{i=1}^n \stackrel{iid}{\sim} \mathcal{N}_2(\mathbf{0}, \Sigma)$, $n = 100$ and $\Sigma = \mathbf{E}\mathbf{\Lambda}\mathbf{E}^\top$, where

$$\mathbf{E} = \begin{pmatrix} \cos(\frac{\pi}{4}) & -\sin(\frac{\pi}{4}) \\ \sin(\frac{\pi}{4}) & \cos(\frac{\pi}{4}) \end{pmatrix}, \quad \mathbf{\Lambda} = \text{diag}(5, 1)$$

Plot the 95% **tolerance ellipses**

$$\{x \in \mathbb{R}^2 : \mathbf{x}^\top \hat{\Sigma}^{-1} \mathbf{x} = F_{\chi_2^2}^{-1}(0.95)\}$$

based on an estimated covariance matrix $\hat{\Sigma}$.

- A single outlier in *south-east*, gradually increasing in magnitude
- t_3M -estimator (**blue ellipse**), computed using `C= tMest(X,3)`
- Sample covariance matrix (**red ellipse**)

Menu

I. Ad-hoc shrinkage SCM-s of multiple samples

II. ML- and M -estimators of scatter matrix

III. Geodesic convexity

- Geodesic
- g -convex functions

IV. Regularized M -estimators

V. Penalized estimation of multiple covariances

VI. Estimation of the regularization parameter

VII. Applications

References



Wiesel, A. (2012a).

Geodesic convexity and covariance estimation.

IEEE Trans. Signal Process., 60(12):6182–6189.



Zhang, T., Wiesel, A., and Greco, M. S. (2013).

Multivariate generalized Gaussian distribution: Convexity and graphical models.

IEEE Trans. Signal Process., 61(16):4141–4148.



Bhatia, R. (2009).

Positive definite matrices.

Princeton University Press.

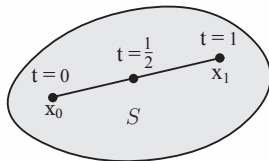
From Euclidean convexity to Riemannian convexity

A set S is **convex** ...

... if $\forall \mathbf{x}_0, \mathbf{x}_1 \in S$ and $t \in [0, 1]$:

$$(1 - t)\mathbf{x}_0 + t\mathbf{x}_1 \in S.$$

convex



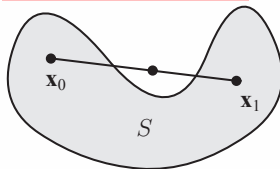
From Euclidean convexity to Riemannian convexity

A set S is **convex** ...

... if $\forall \mathbf{x}_0, \mathbf{x}_1 \in S$ and $t \in [0, 1]$:

$$(1 - t)\mathbf{x}_0 + t\mathbf{x}_1 \in S.$$

nonconvex



From Euclidean convexity to Riemannian convexity

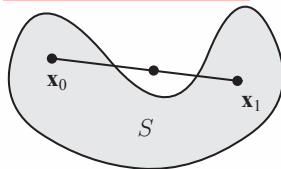
A set S is **convex** ...

... if $\forall \mathbf{x}_0, \mathbf{x}_1 \in S$ and $t \in [0, 1]$:

$$(1 - t)\mathbf{x}_0 + t\mathbf{x}_1 \in S.$$

... if together with \mathbf{x}_0 and \mathbf{x}_1 , it contains the shortest path (geodesic) connecting them

nonconvex

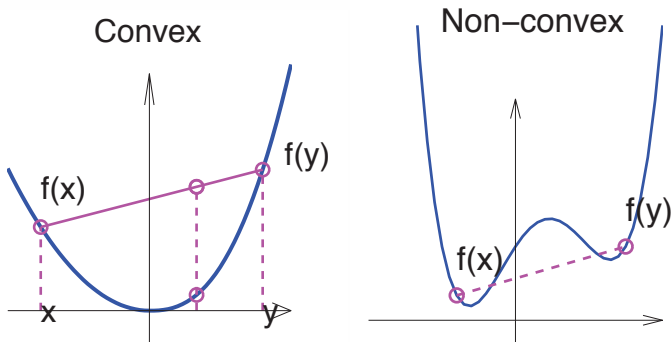


Euclidean convexity

Convex function in $\mathbf{x} \in S$, where S is a convex set

$$f((1-t)\mathbf{x}_0 + t\mathbf{x}_1) \leq (1-t)f(\mathbf{x}_0) + tf(\mathbf{x}_1), \quad \forall t \in [0, 1]$$

strictly convex: strict inequality ($<$) holds $\forall t \in (0, 1)$. **concave**: $-f$ is convex.



Convex function f : the line segment between any two points on the graph of the function f lies above the graph.

Basic examples

Convex:

- x^2 (strictly) , $|x|$, $1/x$ ($x > 0$)
- e^x , $x \log x$
- log-sum-exp function:
 $\log(\sum_{i=1}^n e^{x_i})$.
- $\|\mathbf{x}\|_p$ (any norm, so $p \geq 1$)
- $\mathbf{x}^\top \Sigma \mathbf{x}$ ($\Sigma \in \mathcal{S}(p)$)
- $\mathbf{a}^\top \mathbf{x} + b$
- $\lambda_{\max}(\Sigma)$ ($\Sigma \in \mathcal{S}(p)$)

Concave:

- $x^{1/2}$, $|x|^p$ ($0 \leq p \leq 1$)
- $\log x$
- log determinant:
 $\log |\Sigma|$ ($\Sigma \in \mathcal{S}(p)$)
- $\min(x_1, \dots, x_n)$
- $\lambda_{\min}(\Sigma)$ ($\Sigma \in \mathcal{S}(p)$)
- $\mathbf{a}^\top \mathbf{x} + b$

Geodesic convexity in $p = 1$ variable

convex function in $x \in \mathbb{R}$:

$$f\left(\underbrace{(1-t)x_0 + tx_1}_{\text{line}}\right) \leq (1-t)f(x_0) + tf(x_1)$$

g -convex function in $\sigma^2 \in \mathbb{R}_0^+$:

$$\rho\left(\underbrace{(\sigma_0^2)^{(1-t)}(\sigma_1^2)^t}_{\text{geodesic}}\right) \leq (1-t)\rho(\sigma_0^2) + t\rho(\sigma_1^2)$$

- Convex in $x = \log \sigma^2$ w.r.t. $(1-t)x_0 + tx_1$ is equivalent to g -convex in σ^2 w.r.t. $\sigma_t^2 = (\sigma_0^2)^{(1-t)}(\sigma_1^2)^t$.
- But for $\Sigma \in \mathcal{S}(p)$, $p \neq 1$, the solution is **not** a simple change of variables.

Geodesic convexity in $p = 1$ variable

convex function in $x \in \mathbb{R}$: $f(x) = \rho(e^x)$, $x = \log(\sigma^2)$

$$f\left(\underbrace{(1-t)x_0 + tx_1}_{\text{line}}\right) \leq (1-t)f(x_0) + tf(x_1)$$

g -convex function in $\sigma^2 \in \mathbb{R}_0^+$: $\rho(\sigma^2) = f(\log \sigma^2)$, $\sigma^2 = e^x$

$$\rho\left(\underbrace{(\sigma_0^2)^{(1-t)}(\sigma_1^2)^t}_{\text{geodesic}}\right) \leq (1-t)\rho(\sigma_0^2) + t\rho(\sigma_1^2)$$

- Convex in $x = \log \sigma^2$ w.r.t. $(1-t)x_0 + tx_1$ is equivalent to g -convex in σ^2 w.r.t. $\sigma_t^2 = (\sigma_0^2)^{(1-t)}(\sigma_1^2)^t$.
- But for $\Sigma \in \mathcal{S}(p)$, $p \neq 1$, the solution is **not** a simple change of variables.

Geodesic convexity in $p = 1$ variable

convex function in $x \in \mathbb{R}$:

$$f\left(\underbrace{(1-t)x_0 + tx_1}_{\text{line}}\right) \leq (1-t)f(x_0) + tf(x_1)$$

g -convex function in $\sigma^2 \in \mathbb{R}_0^+$:

$$\rho\left(\underbrace{(\sigma_0^2)^{(1-t)}(\sigma_1^2)^t}_{\text{geodesic}}\right) \leq (1-t)\rho(\sigma_0^2) + t\rho(\sigma_1^2)$$

- Convex in $x = \log \sigma^2$ w.r.t. $(1-t)x_0 + tx_1$ is equivalent to g -convex in σ^2 w.r.t. $\sigma_t^2 = (\sigma_0^2)^{(1-t)}(\sigma_1^2)^t$.
- But for $\Sigma \in \mathcal{S}(p)$, $p \neq 1$, the solution is **not** a simple change of variables.

Geodesic (g -)convexity

On the Riemannian manifold of positive definite matrices, the

geodesic (shortest) path from $\Sigma_0 \in \mathcal{S}(p)$ to $\Sigma_1 \in \mathcal{S}(p)$ is

$$\Sigma_t = \Sigma_0^{1/2} \left(\Sigma_0^{-1/2} \Sigma_1 \Sigma_0^{-1/2} \right)^t \Sigma_0^{1/2} \text{ for } t \in [0, 1].$$

where $\Sigma_t \in \mathcal{S}(p)$ for $0 \leq t \leq 1 \Rightarrow \mathcal{S}(p)$ forms a **g -convex set** (= all geodesic paths Σ_t lie in $\mathcal{S}(p)$).

- Main idea: change the parametric path going from Σ_0 to Σ_1 .
- Midpoint of the path, $\Sigma_{1/2} :=$ **Riemannian (geometric) mean** between Σ_0 and Σ_1 .
- For $p = 1$, the path is $\sigma_t^2 = (\sigma_0^2)^{1-t} (\sigma_1^2)^t$ and the midpoint is the geometric mean

$$\sigma_{1/2}^2 = \sqrt{\sigma_0^2 \sigma_1^2} = \exp \left\{ \frac{1}{2} [\ln(\sigma_0^2) + \ln(\sigma_1^2)] \right\}$$

Riemannian manifold

- Geodesics: informally, shortest paths on a manifold (surface)
- Space of symmetric matrices equipped with inner product

$$\langle \mathbf{A}, \mathbf{B} \rangle = \text{Tr}(\mathbf{AB}) = \text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})$$

and associated Frobenius norm $\| \cdot \|_F = \sqrt{\langle \cdot, \cdot \rangle}$ is a Euclidean space of dimension $p(p+1)/2$.

- Instead, view covariance matrices as elements of a Riemannian manifold
- Endow $\mathcal{S}(p)$ with the **Riemannian metric**
 - local inner product $\langle \mathbf{A}, \mathbf{B} \rangle_\Sigma$ on the tangent space of symmetric matrices

$$\begin{aligned} \langle \mathbf{A}, \mathbf{B} \rangle_\Sigma &= \langle \Sigma^{-1/2} \mathbf{A} \Sigma^{-1/2}, \Sigma^{-1/2} \mathbf{B} \Sigma^{-1/2} \rangle \\ &= \text{Tr}(\mathbf{A} \Sigma^{-1} \mathbf{B} \Sigma^{-1}) = \text{vec}(\mathbf{A})^\top (\Sigma^{-1} \otimes \Sigma^{-1}) \text{vec}(\mathbf{B}) \end{aligned}$$

- Geodesic path Σ_t is the shortest path from Σ_0 to Σ_1 .

Geodesically (g -)convex function

A function $h : \mathcal{S}(p) \rightarrow \mathbb{R}$ is g -convex function if

$$h(\Sigma_t) \leq (1 - t) h(\Sigma_0) + t h(\Sigma_1) \text{ for } t \in (0, 1).$$

If the inequality is strict, then h is strictly g -convex.

Note: Def. of convexity of $h(\Sigma)$ remains the same, i.e., w.r.t. to given path Σ_t . Now geodesic instead of Euclidean path.

g -convexity = convexity w.r.t. geodesic paths

Local is Global

- 1 any local minimum of $h(\Sigma)$ over $\mathcal{S}(p)$ is a global minimum.
- 2 If h is strictly g -convex and a minimum is in $\mathcal{S}(p)$, then it is a unique minimum.
- 3 g -convex + g -convex = g -convex

Useful results on g -convexity: my personal top 3

(not in particular order)



$$\Sigma_t = \Sigma_0^{1/2} \left(\Sigma_0^{-1/2} \Sigma_1 \Sigma_0^{-1/2} \right)^t \Sigma_0^{1/2}$$

1. Joint diagonalization formulation

The geodesic path can be written equivalently as

$$\Sigma_t = \mathbf{E} \mathbf{D}^t \mathbf{E}^\top, \quad t \in [0, 1],$$

where $\Sigma_0 = \mathbf{E} \mathbf{E}^\top$ and $\Sigma_1 = \mathbf{E} \mathbf{D} \mathbf{E}^\top$ by joint diagonalization.

- \mathbf{E} is a *nonsingular* square matrix: row vectors of \mathbf{E}^{-1} are the eigenvectors of $\Sigma_0^{-1} \Sigma_1$
- \mathbf{D} is a diagonal matrix: diagonal elements are the eigenvalues of

$$\Sigma_0^{-1} \Sigma_1 \quad \text{or} \quad \Sigma_0^{-1/2} \Sigma_1 \Sigma_0^{-1/2}.$$

Useful results on g -convexity: my personal top 3

$$\Sigma_t = \Sigma_0^{1/2} \left(\Sigma_0^{-1/2} \Sigma_1 \Sigma_0^{-1/2} \right)^t \Sigma_0^{1/2}$$

2. Convexity w.r.t. t

A continuous function f on a g -convex set \mathcal{M} is g -convex if $f(\Sigma_t)$ is classically convex in $t \in [0, 1]$

3. Midpoint convexity

A continuous function f on a g -convex set \mathcal{M} is g -convex if

$$f(\Sigma_{1/2}) \leq \frac{1}{2} \{f(\Sigma_0) + f(\Sigma_1)\}$$

for any $\Sigma_0, \Sigma_1 \in \mathcal{M}$.

For more results, see [\[Wiesel and Zhang, 2015\]](#)

Some geodesically (g -)convex functions

- 1 if $h(\Sigma)$ is g -convex in Σ , then it is g -convex in Σ^{-1} .

scalar case: if $h(x)$ is convex in $x = \log(\sigma^2) \in \mathbb{R}$, then it is convex in $-x = \log(\sigma^{-2}) = -\log(\sigma^2)$.

- 2 $\pm \log |\Sigma|$ is g -convex. (i.e., **log-determinant is g -linear function**)

scalar case: the scalar g -linear function is the logarithm.

- 3 $\mathbf{a}^\top \Sigma^{\pm 1} \mathbf{a}$ is strictly g -convex ($\mathbf{a} \neq 0$).

- 4 $\log \left| \sum_{i=1}^n \mathbf{H}_i \Sigma^{\pm 1} \mathbf{H}_i \right|$ is g -convex.

scalar case: log-sum-exp function is convex.

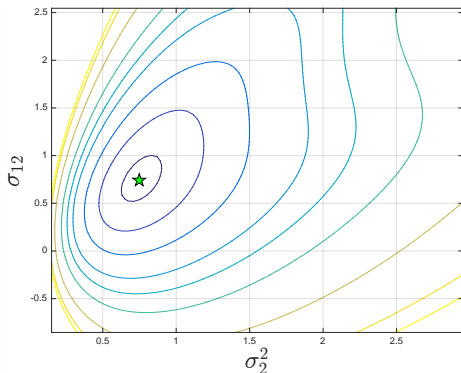
- 5 if $f(\Sigma)$ is g -convex, then $f(\Sigma_1 \otimes \Sigma_2)$ is jointly g -convex.

Example: Tyler's M -estimator of shape

Let's minimize Tyler's cost function $\mathcal{L}_T(\Sigma) = \mathcal{L}_T(\sigma_2^2, \sigma_{12})$ over g -convex set of 2×2 shape matrices:

$$\begin{aligned}\mathcal{M}(2) &= \{\Sigma \in \mathcal{S}(2) : \det(\Sigma) = 1\} \\ &= \left\{ \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} : \sigma_2^2 > 0, \sigma_{12} \in \mathbb{R}, \sigma_1^2 = \frac{1 + \sigma_{12}^2}{\sigma_2^2} \right\}\end{aligned}$$

We generated a Gaussian sample of length $n = 15$ with $\sigma_2^2 = \sigma_{12} = 1$.



$$\min_{\Sigma \in \mathcal{M}(2)} \underbrace{\sum_{i=1}^n \ln(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i)}_{=\mathcal{L}_T(\Sigma)}$$

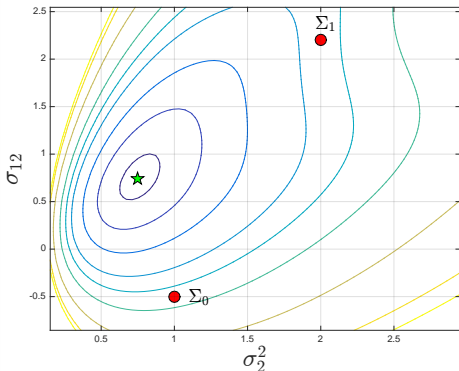
Contours of $\mathcal{L}_T(\Sigma)$
and the solution $\hat{\Sigma}$.

Example: Tyler's M -estimator of shape

Let's minimize Tyler's cost function $\mathcal{L}_T(\Sigma) = \mathcal{L}_T(\sigma_2^2, \sigma_{12})$ over g -convex set of 2×2 shape matrices:

$$\begin{aligned}\mathcal{M}(2) &= \{\Sigma \in \mathcal{S}(2) : \det(\Sigma) = 1\} \\ &= \left\{ \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} : \sigma_2^2 > 0, \sigma_{12} \in \mathbb{R}, \sigma_1^2 = \frac{1 + \sigma_{12}^2}{\sigma_2^2} \right\}\end{aligned}$$

We generated a Gaussian sample of length $n = 15$ with $\sigma_2^2 = \sigma_{12} = 1$.



$$\min_{\Sigma \in \mathcal{M}(2)} \underbrace{\sum_{i=1}^n \ln(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i)}_{=\mathcal{L}_T(\Sigma)}$$

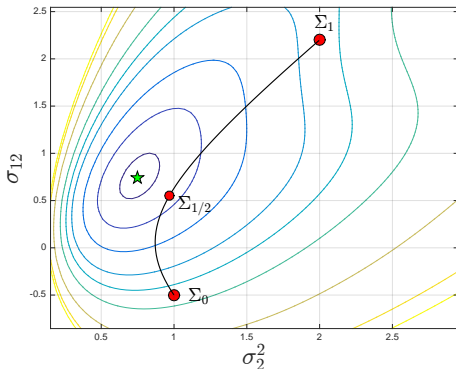
Consider two points Σ_0 and Σ_1 of \mathcal{M} .

Example: Tyler's M -estimator of shape

Let's minimize Tyler's cost function $\mathcal{L}_T(\Sigma) = \mathcal{L}_T(\sigma_2^2, \sigma_{12})$ over g -convex set of 2×2 shape matrices:

$$\begin{aligned}\mathcal{M}(2) &= \{\Sigma \in \mathcal{S}(2) : \det(\Sigma) = 1\} \\ &= \left\{ \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} : \sigma_2^2 > 0, \sigma_{12} \in \mathbb{R}, \sigma_1^2 = \frac{1 + \sigma_{12}^2}{\sigma_2^2} \right\}\end{aligned}$$

We generated a Gaussian sample of length $n = 15$ with $\sigma_2^2 = \sigma_{12} = 1$.



$$\min_{\Sigma \in \mathcal{M}(2)} \underbrace{\sum_{i=1}^n \ln(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i)}_{=\mathcal{L}_T(\Sigma)}$$

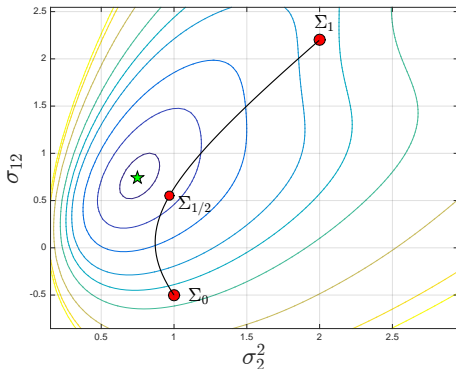
Their geodesic path
 Σ_t and midpoint $\Sigma_{1/2}$

Example: Tyler's M -estimator of shape

Let's minimize Tyler's cost function $\mathcal{L}_T(\Sigma) = \mathcal{L}_T(\sigma_2^2, \sigma_{12})$ over g -convex set of 2×2 shape matrices:

$$\begin{aligned}\mathcal{M}(2) &= \{\Sigma \in \mathcal{S}(2) : \det(\Sigma) = 1\} \\ &= \left\{ \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} : \sigma_2^2 > 0, \sigma_{12} \in \mathbb{R}, \sigma_1^2 = \frac{1 + \sigma_{12}^2}{\sigma_2^2} \right\}\end{aligned}$$

We generated a Gaussian sample of length $n = 15$ with $\sigma_2^2 = \sigma_{12} = 1$.



$$\min_{\Sigma \in \mathcal{M}(2)} \underbrace{\sum_{i=1}^n \ln(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i)}_{=\mathcal{L}_T(\Sigma)}$$

By utilizing the proper (Riemannian) metric, Tyler's cost fnc *is convex*.

Examples of g -convex sets

g -convex set \mathcal{M} = all geodesic paths Σ_t lie in the set, where

$$\Sigma_t = \Sigma_0^{1/2} \left(\Sigma_0^{-1/2} \Sigma_1 \Sigma_0^{-1/2} \right)^t \Sigma_0^{1/2} \text{ for } t \in [0, 1].$$

and Σ_0 and Σ_1 are in \mathcal{M} .

- 1 The set of PDS matrices: $\mathcal{M} = \mathcal{S}_p$
- 2 The set of PDS **shape matrices**: $\mathcal{M} = \{\Sigma \in \mathcal{S}_p : \det(\Sigma) = 1\}$
- 3 The set of PDS **block diagonal matrices**.
- 4 Kronenecker model $\Sigma = \Sigma_1 \otimes \Sigma_2$
- 5 Complex circular symmetric model:

$$\Sigma = \begin{pmatrix} \Sigma_1 & \Sigma_2 \\ -\Sigma_2 & \Sigma_1 \end{pmatrix}$$

- 6 PDS circulant matrices, e.g., $[\Sigma]_{ij} = \rho^{|i-j|}$, $\rho \in (0, 1)$.

Menu

I. Ad-hoc shrinkage SCM-s of multiple samples

II. ML- and M -estimators of scatter matrix

III. Geodesic convexity

IV. Regularized M -estimators

- Shrinkage towards an identity matrix
- Shrinkage towards a target matrix

V. Penalized estimation of multiple covariances

VI. Estimation of the regularization parameter

VII. Applications

References



Ollila, E. and Tyler, D. E. (2014).
Regularized M -estimators of scatter matrix.
IEEE Trans. Signal Process., 62(22):6059–6070.



Ollila, E., Soloveychik, I., Tyler, D. E. and Wiesel, A. (2016).
Simultaneous penalized M-estimation of covariance matrices using
geodesically convex optimization
Journal of Multivariate Analysis (under review), Cite as:
arXiv:1608.08126 [stat.ME]
<http://arxiv.org/abs/1608.08126>

Honey, I shrunk the M -estimator of covariance matrix

Article Abstract

Subscribe Now →

A

Honey, I Shrunk the Sample Covariance Matrix

The Journal of Portfolio Management

Summer 2004, Vol. 30, No. 4: pp. 110-119

DOI: 10.3905/jpm.2004.110

Olivier Ledit and Michael Wolf

View



The central message of this article is that no one should use the sample covariance matrix for portfolio optimization. It is subject to estimation error of the kind most likely to perturb a mean-variance optimizer. Instead, a matrix can be obtained from the sample covariance matrix through a transformation called shrinkage. This tends to pull the most extreme coefficients toward more central values, systematically reducing estimation error when it matters most. Statistically, the challenge is to know the optimal shrinkage intensity. Shrinkage reduces portfolio tracking error relative to a benchmark index, and substantially raises the manager's realized information ratio.

- ⇒ The popular shrinkage SCM estimator (a.k.a Ledit-Wolf estimator) is a member in the large class of regularized M -estimators of [Ollila and Tyler, 2014] that are presented in this section.
- ⇒ **Our message:** for non-Gaussian data or in the presence of outliers, you will do **MUCH BETTER** by shrinking a robust M -estimator of covariance!

Honey, I shrunk the M -estimator of covariance matrix

Article Abstract

Subscribe Now →

A

Honey, I Shrunk the Sample Covariance Matrix

The Journal of Portfolio Management

Summer 2004, Vol. 30, No. 4: pp. 110-119

DOI: 10.3905/jpm.2004.110

Olivier Ledit and Michael Wolf

View



The central message of this article is that no one should use the sample covariance matrix for portfolio optimization. It is subject to estimation error of the kind most likely to perturb a mean-variance optimizer. Instead, a matrix can be obtained from the sample covariance matrix through a transformation called shrinkage. This tends to pull the most extreme coefficients toward more central values, systematically reducing estimation error when it matters most. Statistically, the challenge is to know the optimal shrinkage intensity. Shrinkage reduces portfolio tracking error relative to a benchmark index, and substantially raises the manager's realized information ratio.

⇒ The popular shrinkage SCM estimator (a.k.a Ledit-Wolf estimator) is a member in the large class of regularized M -estimators of [Ollila and Tyler, 2014] that are presented in this section.

⇒ Our message: for non-Gaussian data or in the presence of outliers, you will do MUCH BETTER by shrinking a robust M -estimator of covariance!

Honey, I shrunk the M -estimator of covariance matrix

Article Abstract

Subscribe Now →

A

Honey, I Shrunk the Sample Covariance Matrix

The Journal of Portfolio Management

Summer 2004, Vol. 30, No. 4: pp. 110-119

DOI: 10.3905/jpm.2004.110

Olivier Ledit and Michael Wolf

View



The central message of this article is that no one should use the sample covariance matrix for portfolio optimization. It is subject to estimation error of the kind most likely to perturb a mean-variance optimizer. Instead, a matrix can be obtained from the sample covariance matrix through a transformation called shrinkage. This tends to pull the most extreme coefficients toward more central values, systematically reducing estimation error when it matters most. Statistically, the challenge is to know the optimal shrinkage intensity. Shrinkage reduces portfolio tracking error relative to a benchmark index, and substantially raises the manager's realized information ratio.

- ⇒ The popular shrinkage SCM estimator (a.k.a Ledit-Wolf estimator) is a member in the large class of regularized M -estimators of [Ollila and Tyler, 2014] that are presented in this section.
- ⇒ **Our message:** for non-Gaussian data or in the presence of outliers, you will do **MUCH BETTER** by shrinking a robust M -estimator of covariance!

Regularized M -estimators of scatter matrix

Penalized cost function:

$$\mathcal{L}_\alpha(\Sigma) = \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| + \alpha \mathcal{P}(\Sigma)$$

where $\alpha \geq 0$ is a fixed **regularization parameter**.

Q: Existence, Uniqueness, computation?

Our penalty function

$$\mathcal{P}(\Sigma) = \text{Tr}(\Sigma^{-1})$$

- **Q:** Why is this penalty useful or what is its effect?
- **Q:** Any other penalties ?

Regularized M -estimators of scatter matrix

Penalized cost function:

$$\mathcal{L}_\alpha(\Sigma) = \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| + \alpha \mathcal{P}(\Sigma)$$

where $\alpha \geq 0$ is a fixed **regularization parameter**.

Q: Existence, Uniqueness, computation?

Our penalty function

$$\mathcal{P}(\Sigma) = \text{Tr}(\Sigma^{-1})$$

- **Q:** Why is this penalty useful or what is its effect?
- **Q:** Any other penalties ?

Regularized M -estimators of scatter matrix

Penalized cost function:

$$\mathcal{L}_\alpha(\Sigma) = \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| + \alpha \mathcal{P}(\Sigma)$$

where $\alpha \geq 0$ is a fixed **regularization parameter**.

Q: Existence, Uniqueness, computation?

Our penalty function *pulls* Σ away from singularity

$$\mathcal{P}(\Sigma) = \text{Tr}(\Sigma^{-1})$$

- **Q:** Why is this penalty useful or what is its effect?
- **Q:** Any other penalties ?

Condition 1. [Zhang et al., 2013, Ollila and Tyler, 2014]

- $\rho(t)$ is nondecreasing and continuous for $0 < t < \infty$.
- $\rho(t)$ is g -convex (i.e., $\rho(e^x)$ is convex in $-\infty < x < \infty$)

Are the loss functions presented earlier g -convex?

- Gaussian loss function $\rho_G(t) = t$
as $f(x) = \rho_G(e^x) = e^x$ is strictly convex in x
- Tyler's loss function $\rho_T(t) = p \ln t$
as $f(x) = \rho_T(e^x) = px$ is g -linear.
- t_ν loss function $\rho_\nu(t) = (\nu + p) \ln(\nu + t)$
as $f(x) = \rho_\nu(e^x) = (\nu + p)(\log \nu + x)$
- Huber's loss function $\rho_H(t; c)$
as it is a hybrid of Gaussian and Tyler's loss functions

Condition 1. [Zhang et al., 2013, Ollila and Tyler, 2014]

- $\rho(t)$ is nondecreasing and continuous for $0 < t < \infty$.
- $\rho(t)$ is g -convex (i.e., $\rho(e^x)$ is convex in $-\infty < x < \infty$)

Are the loss functions presented earlier g -convex?

- Gaussian loss function $\rho_G(t) = t$ is strictly g -convex
as $f(x) = \rho_G(e^x) = e^x$ is strictly convex in x
- Tyler's loss function $\rho_T(t) = p \ln t$
as $f(x) = \rho_T(e^x) = px$ is g -linear.
- t_ν loss function $\rho_\nu(t) = (\nu + p) \ln(\nu + t)$
as $f(x) = \rho_\nu(e^x) = (\nu + p)(\log \nu + x)$
- Huber's loss function $\rho_H(t; c)$
as it is a hybrid of Gaussian and Tyler's loss functions

Condition 1. [Zhang et al., 2013, Ollila and Tyler, 2014]

- $\rho(t)$ is nondecreasing and continuous for $0 < t < \infty$.
- $\rho(t)$ is g -convex (i.e., $\rho(e^x)$ is convex in $-\infty < x < \infty$)

Are the loss functions presented earlier g -convex?

- Gaussian loss function $\rho_G(t) = t$ is strictly g -convex
as $f(x) = \rho_G(e^x) = e^x$ is strictly convex in x
- Tyler's loss function $\rho_T(t) = p \ln t$ is g -convex
as $f(x) = \rho_T(e^x) = px$ is g -linear.
- t_ν loss function $\rho_\nu(t) = (\nu + p) \ln(\nu + t)$
as $f(x) = \rho_\nu(e^x) = (\nu + p)(\log \nu + x)$
- Huber's loss function $\rho_H(t; c)$
as it is a hybrid of Gaussian and Tyler's loss functions

Condition 1. [Zhang et al., 2013, Ollila and Tyler, 2014]

- $\rho(t)$ is nondecreasing and continuous for $0 < t < \infty$.
- $\rho(t)$ is g -convex (i.e., $\rho(e^x)$ is convex in $-\infty < x < \infty$)

Are the loss functions presented earlier g -convex?

- Gaussian loss function $\rho_G(t) = t$ is strictly g -convex
as $f(x) = \rho_G(e^x) = e^x$ is strictly convex in x
- Tyler's loss function $\rho_T(t) = p \ln t$ is g -convex
as $f(x) = \rho_T(e^x) = px$ is g -linear.
- t_ν loss function $\rho_\nu(t) = (\nu + p) \ln(\nu + t)$ is g -convex
as $f(x) = \rho_\nu(e^x) = (\nu + p)(\log \nu + x)$
- Huber's loss function $\rho_H(t; c)$ is g -convex
as it is a hybrid of Gaussian and Tyler's loss functions

Main results

$$\mathcal{L}_\alpha(\Sigma) = \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| + \alpha \text{Tr}(\Sigma^{-1}), \quad \alpha > 0$$

Result 1 [Ollila and Tyler, 2014]

Assume $\rho(t)$ satisfies Condition 1.

- (a) Uniqueness: $\mathcal{L}_\alpha(\Sigma)$ is **strictly g -convex** in $\Sigma \in \mathcal{S}(p)$
- (b) Existence: If $\rho(t)$ is bounded below, then the solution to $\mathcal{L}_\alpha(\Sigma)$ **always exists** and is **unique**.
- (c) Furthermore, if $\rho(t)$ is also differentiable, then the minimum corresponds to the unique solution of the regularized M -estimating equation:

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + \alpha \mathbf{I}$$

Main results (cont'd)

Result 1 implies

- $u(t)$ need not be nonincreasing
- Unlike the non-regularized case, no conditions on the data are needed!
→ breakdown point is $= 1$.

Result 1(d) [Ollila and Tyler, 2014, Theorem 2]

Suppose $\rho(t)$ is continuously differentiable, satisfies Condition 1 and that $u(t) = \rho'(t)$ is non-increasing, Then the **Fixed-point (FP) algorithm**

$$\hat{\Sigma}_{k+1} = \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \hat{\Sigma}_k^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + \alpha \mathbf{I}$$

converges to the solution of regularized M -estimating equation given in Result 1(c).

Matlab implementation of regularized $t_\nu M$ -estimator

Regularized $t_\nu M$ -estimator is a solution to

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n u_\nu(\mathbf{x}_i^\top \hat{\Sigma}_k^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + \alpha \mathbf{I}$$

for $\alpha, \beta > 0$.

- Define a new function `C = regtMest(X, v, al, be)`
- Copy-paste the code from Slide [36](#).
- You *only need to change one line*:

replace

```
C = X'*((X.* repmat(u,1,p)))/n;
```

to

```
C = be*X'*((X.* repmat(u,1,p)))/n + al*eye(p);
```


Tuning the $\rho(t)$ function

- Result 1 is general and allows us to tune the $\rho(t)$ function
- For a given ρ -function, a class of tuned ρ -functions are defined as

$$\rho_{\beta}(t) = \beta\rho(t) \quad \text{for } \beta > 0.$$

where β represents *additional tuning constant* which can be used to tune the estimator towards some desirable property.

- Using $\rho_{\beta}(t) = \beta\rho(t)$, our optimization program is

$$\mathcal{L}_{\alpha,\beta}(\Sigma) = \beta \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^{\top} \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| + \alpha \text{Tr}(\Sigma^{-1})$$

- The solution verifies

$$\hat{\Sigma} = \beta \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^{\top} \hat{\Sigma}^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^{\top} + \alpha \mathbf{I}$$

- Special cases: $\alpha = 1 - \beta$ or $\beta = (1 - \alpha)$.

Tuning the $\rho(t)$ function

- Result 1 is general and allows us to tune the $\rho(t)$ function
- For a given ρ -function, a class of tuned ρ -functions are defined as

$$\rho_{\beta}(t) = \beta\rho(t) \quad \text{for } \beta > 0.$$

where β represents *additional tuning constant* which can be used to tune the estimator towards some desirable property.

- Using $\rho_{\beta}(t) = \beta\rho(t)$, our optimization program is

$$\mathcal{L}_{\alpha,\beta}(\Sigma) = \beta \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^{\top} \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| + \alpha \text{Tr}(\Sigma^{-1})$$

- The solution verifies

$$\hat{\Sigma} = \beta \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^{\top} \hat{\Sigma}^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^{\top} + \alpha \mathbf{I}$$

- **Special cases:** $\alpha = 1 - \beta$ or $\beta = (1 - \alpha)$.

A class of regularized SCM's

- Let use *tuned* Gaussian cost fnc $\rho(t) = \beta t$, where $\beta > 0$ is a fixed tuning parameter.
- The penalized cost fnc is then

$$\mathcal{L}_{\alpha,\beta}(\Sigma) = \text{Tr}\{(\beta\mathbf{S} + \alpha\mathbf{I})\Sigma^{-1}\} - \ln|\Sigma^{-1}|$$

where \mathbf{S} denotes the SCM.

- Due to ► Result 1, its unique minimizer $\hat{\Sigma}$ is

$$\mathbf{S}_{\alpha,\beta} = \beta\mathbf{S} + \alpha\mathbf{I}$$

which corresponds to [\[Ledoit and Wolf, 2004\]](#) shrinkage estimator.

- **Note:** Ledoit and Wolf did not show that $\mathbf{S}_{\alpha,\beta}$ solves an penalized Gaussian optimization program.

A class of regularized Tyler's M -estimators

- Let use *tuned* Tyler's cost fnc $\rho(t) = p\beta \log t$ for fixed $0 < \beta < 1$.
- The penalized Tyler's cost fnc is

$$\mathcal{L}_{\alpha,\beta}(\Sigma) = \beta \frac{p}{n} \sum_{i=1}^n \log(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| + \alpha \text{Tr}(\Sigma^{-1}),$$

- The weight fnc is $u(t) = p\beta/t$, so the regularized M -estimating eq. is

$$\hat{\Sigma} = \beta \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i} + \alpha \mathbf{I}$$

- We commonly use $\alpha = 1 - \beta$, $\beta \in (0, 1]$.
- Target $\mathcal{L}_{\alpha,\beta}(\Sigma)$ is g -convex in Σ , but ρ is not bounded below
 \Rightarrow Result 1(b), for existence does not hold.
- Conditions for existence needs to be considered separately for Tyler's M -estimator;

A class of regularized Tyler's M -estimators

- Let use *tuned* Tyler's cost fnc $\rho(t) = p\beta \log t$ for fixed $0 < \beta < 1$.
- The penalized Tyler's cost fnc is

$$\mathcal{L}_{\alpha,\beta}(\Sigma) = \beta \frac{p}{n} \sum_{i=1}^n \log(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| + \alpha \text{Tr}(\Sigma^{-1}),$$

- The weight fnc is $u(t) = p\beta/t$, so the regularized M -estimating eq. is

$$\hat{\Sigma} = \beta \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i} + (1 - \beta) \mathbf{I}$$

- We commonly use $\alpha = 1 - \beta$, $\beta \in (0, 1]$.
- Target $\mathcal{L}_{\alpha,\beta}(\Sigma)$ is g -convex in Σ , but ρ is not bounded below
 \Rightarrow [Read 10](#), for existence does not hold.
- Conditions for existence needs to be considered separately for Tyler's M -estimator;

A class of regularized Tyler's M -estimators

- Let use *tuned* Tyler's cost fnc $\rho(t) = p\beta \log t$ for fixed $0 < \beta < 1$.
- The penalized Tyler's cost fnc is

$$\mathcal{L}_{\alpha,\beta}(\Sigma) = \beta \frac{p}{n} \sum_{i=1}^n \log(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| + \alpha \text{Tr}(\Sigma^{-1}),$$

- The weight fnc is $u(t) = p\beta/t$, so the regularized M -estimating eq. is

$$\hat{\Sigma} = \beta \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i} + (1 - \beta) \mathbf{I}$$

- We commonly use $\alpha = 1 - \beta$, $\beta \in (0, 1]$.
- Target $\mathcal{L}_{\alpha,\beta}(\Sigma)$ is g -convex in Σ , but ρ is not bounded below
 \Rightarrow ► Result 1(b), for existence does not hold.
- Conditions for existence needs to be considered separately for Tyler's M -estimator;

- **(Sufficient) Condition A.** For any subspace \mathcal{V} of \mathbb{R}^p , $1 \leq \dim(\mathcal{V}) < p$, the inequality

$$\frac{\#\{\mathbf{x}_i \in \mathcal{V}\}}{n} < \frac{\dim(\mathcal{V})}{p\beta}$$

holds. [(**Necessary**) **Condition B:** As earlier but with inequality.]

- Cond A implies $\beta < n/p$ whenever the sample is in “general position” (e.g., when sampling from a continuous distribution)

Result 2 [Ollila and Tyler, 2014]

Consider tuned Tyler's cost $\rho_\beta(t) = p\beta \ln t$ and $\alpha > 0$, $0 \leq \beta < 1$. If Condition A holds, then $\mathcal{L}_{\alpha,\beta}(\Sigma)$ has a unique minimum $\hat{\Sigma}$ in $\mathcal{S}(p)$, the minimum being obtained at the unique solution to

$$\hat{\Sigma} = \beta \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i} + \alpha \mathbf{I},$$

Similar result found independently in [Pascal et al., 2014, Sun et al., 2014].

- For **fixed** $0 < \beta < 1$, consider two different values α_1 and α_2 , and let $\hat{\Sigma}_1$ and $\hat{\Sigma}_2$ represent the respective regularized Tyler's M -estimators.
- It then follows that

$$\hat{\Sigma}_1 = \frac{\alpha_1}{\alpha_2} \cdot \hat{\Sigma}_2$$

\Rightarrow for any fixed $0 < \beta < 1$, the regularized Tyler's M -estimators are *proportional to one another* as α varies.

- Consequently, when the main interest is on estimation of the covariance matrix *up to a scale*, one may set w.l.o.g.

$$\alpha = 1 - \beta \quad [\text{or equivalently } \beta = 1 - \alpha \quad].$$

In these cases, it holds that $\text{Tr}(\hat{\Sigma}^{-1}) = p$.

Related approach for regularizing Tyler's M -estimator

- A related (but *different*) regularized Tyler's M -estimator was proposed by [Abramovich and Spencer, 2007] as the limit of the algorithm

$$\Sigma_{k+1} \leftarrow (1 - \alpha) \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \mathbf{V}_k^{-1} \mathbf{x}_i} + \alpha \mathbf{I}$$
$$\mathbf{V}_{k+1} \leftarrow p \Sigma_{k+1} / \text{Tr}(\Sigma_{k+1}),$$

where $\alpha \in (0, 1]$ is a fixed regularization parameter.

- [Chen et al., 2011] proved that the recursive algorithm above converges to a unique solution regardless of the initialization. [Convergence means convergence in \mathbf{V}_k and not necessarily in Σ_k .]
- **Note 1:** Diagonally loaded version of the fixed-point algorithm for Tyler's M -estimator. Hence we refer to it by **DL-FP**.
- **Note 2:** DL-FP was not shown to be a solution to any penalized form of Tyler's cost function.

Comments

$$\mathcal{L}_\alpha(\Sigma) = \underbrace{\frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}|}_{= \mathcal{L}(\Sigma), \text{ M(L) cost function}} + \alpha \mathcal{P}(\Sigma)$$

- ℓ_1 -penalization $\mathcal{P}(\Sigma) = \|\Sigma^{-1}\|_1$ is convex but not g -convex, so GLASSO [Friedman et al., 2008] does not fit our framework.
- $\mathcal{L}_\alpha(\Sigma)$ is strictly g -convex if either $\mathcal{L}(\Sigma)$ or $\mathcal{P}(\Sigma)$ are.
- Penalties $\mathcal{P}(\Sigma)$, other than $\mathcal{P}(\Sigma) = \text{Tr}(\Sigma^{-1})$, that shrink towards a target matrix are considered in the next section.

Shrinkage towards a target matrix

- Fixed **shrinkage target matrix** $\mathbf{T} \in \mathcal{S}(p)$
- Define penalized M -estimator of scatter matrix as solution to

$$\min_{\Sigma \in \mathcal{S}(p)} \{ \mathcal{L}(\Sigma) + \lambda d(\Sigma, \mathbf{T}) \},$$

or equivalently,

$$\min_{\Sigma \in \mathcal{S}(p)} \{ \beta \mathcal{L}(\Sigma) + (1 - \beta) d(\Sigma, \mathbf{T}) \}, \quad \text{where } \lambda = \frac{1 - \beta}{\beta}$$

where

- $\lambda > 0$ or $\beta \in (0, 1]$ is a **regularization/penalty parameter**
- $d(\mathbf{A}, \mathbf{B}) : \mathcal{S}(p) \times \mathcal{S}(p) \rightarrow \mathbb{R}_0^+$ is **penalty/distance fnc.**

Distance $d(\Sigma, \mathbf{T})$ is used to **enforce similarity** of Σ to target \mathbf{T} and β controls the amount of shrinkage of solution $\hat{\Sigma}$ towards \mathbf{T} .

Properties of the penalty (distance) function

D1 $d(\mathbf{A}, \mathbf{B}) = 0$ if $\mathbf{A} = \mathbf{B}$,

D2 $d(\mathbf{A}, \mathbf{B})$ is **jointly g -convex**

D3 *symmetry*: $d(\mathbf{A}, \mathbf{B}) = d(\mathbf{B}, \mathbf{A})$.

D4 *affine invariance* $d(\mathbf{A}, \mathbf{B}) = d(\mathbf{C}\mathbf{A}\mathbf{C}^\top, \mathbf{C}\mathbf{B}\mathbf{C}^\top)$, \forall nonsingular \mathbf{C}

D5 *scale invariance*: $d(c_1\mathbf{A}, c_2\mathbf{B}) = d(\mathbf{A}, \mathbf{B})$ for $c_1, c_2 > 0$,

Comments:

- D3-D5 are considered optional properties
- Property D5 is needed for shape matrix estimators (e.g. Tyler's). It is also important if Σ_k -s share a common shape matrix only.

Note: Each distance $d(\Sigma_k, \Sigma)$ induce a notion of mean (or center).

\Rightarrow one might expect that a judicious choice of $d(\cdot, \cdot)$ should induce a natural notion of the mean of pos. def. matrices.

Properties of the penalty (distance) function

D1 $d(\mathbf{A}, \mathbf{B}) = 0$ if $\mathbf{A} = \mathbf{B}$,

D2 $d(\mathbf{A}, \mathbf{B})$ is **jointly g -convex**

D3 *symmetry*: $d(\mathbf{A}, \mathbf{B}) = d(\mathbf{B}, \mathbf{A})$.

D4 *affine invariance* $d(\mathbf{A}, \mathbf{B}) = d(\mathbf{C}\mathbf{A}\mathbf{C}^\top, \mathbf{C}\mathbf{B}\mathbf{C}^\top)$, \forall nonsingular \mathbf{C}

D5 *scale invariance*: $d(c_1\mathbf{A}, c_2\mathbf{B}) = d(\mathbf{A}, \mathbf{B})$ for $c_1, c_2 > 0$,

Comments:

- D3-D5 are considered optional properties
- Property D5 is needed for shape matrix estimators (e.g. Tyler's). It is also important if Σ_k -s share a common shape matrix only.

Note: Each distance $d(\Sigma_k, \Sigma)$ induce a notion of mean (or center).

⇒ one might expect that a judicious choice of $d(\cdot, \cdot)$ should induce a natural notion of the mean of pos. def. matrices.

The induced mean or center

- Let $\{\Sigma_k\}_{k=1}^K$ be given matrices in $\mathcal{S}(p)$
- Let weights $\pi = (\pi_1, \dots, \pi_K)$, $\sum_{k=1}^K \pi_k = 1$, be given.

Then

$$\Sigma(\pi) = \arg \min_{\Sigma \in \mathcal{S}(p)} \sum_{k=1}^K \pi_k d(\Sigma_k, \Sigma),$$

is a **weighted mean** associated with distance (penalty) d .

- **Q:** What is a natural mean of positive definite matrices?
- If $p = 1$, so we have $\sigma_1^2, \dots, \sigma_K^2 > 0$, we could consider

$$\sigma = \text{geometric mean } \sigma^2 = (\sigma_1^2, \dots, \sigma_K^2)$$

- **Note:** For a pair σ_0^2, σ_1^2 , the geometric mean is the midpoint of the geodesic $\sigma_t^2 = (\sigma_0^2)^{1-t} (\sigma_1^2)^t$.

The induced mean or center

- Let $\{\Sigma_k\}_{k=1}^K$ be given matrices in $\mathcal{S}(p)$
- Let weights $\pi = (\pi_1, \dots, \pi_K)$, $\sum_{k=1}^K \pi_k = 1$, be given.

Then

$$\Sigma(\pi) = \arg \min_{\Sigma \in \mathcal{S}(p)} \sum_{k=1}^K \pi_k d(\Sigma_k, \Sigma),$$

is a **weighted mean** associated with distance (penalty) d .

- **Q:** What is a natural mean of positive definite matrices?
- If $p = 1$, so we have $\sigma_1^2, \dots, \sigma_K^2 > 0$, we could consider
 - arithmetic mean $\sigma^2 = \frac{1}{K} \sum_{k=1}^K \sigma_k^2$.
 - geometric mean $\sigma^2 = (\sigma_1^2 \cdots \sigma_K^2)^{1/K}$
 - harmonic mean $\sigma^2 = \left(\frac{1}{K} \sum_{k=1}^K (\sigma_k^2)^{-1} \right)^{-1}$
- **Note:** For a pair σ_0^2, σ_1^2 , the geometric mean is the midpoint of the geodesic $\sigma_t^2 = (\sigma_0^2)^{1-t} (\sigma_1^2)^t$.

The induced mean or center

- Let $\{\Sigma_k\}_{k=1}^K$ be given matrices in $\mathcal{S}(p)$
- Let weights $\pi = (\pi_1, \dots, \pi_K)$, $\sum_{k=1}^K \pi_k = 1$, be given.

Then

$$\Sigma(\pi) = \arg \min_{\Sigma \in \mathcal{S}(p)} \sum_{k=1}^K \pi_k d(\Sigma_k, \Sigma),$$

is a **weighted mean** associated with distance (penalty) d .

- **Q:** What is a natural mean of positive definite matrices?
- If $p = 1$, so we have $\sigma_1^2, \dots, \sigma_K^2 > 0$, we could consider
 - arithmetic mean $\sigma^2 = \frac{1}{K} \sum_{k=1}^K \sigma_k^2$.
 - geometric mean $\sigma^2 = (\sigma_1^2 \cdots \sigma_K^2)^{1/K}$
 - harmonic mean $\sigma^2 = \left(\frac{1}{K} \sum_{k=1}^K (\sigma_k^2)^{-1} \right)^{-1}$
- **Note:** For a pair σ_0^2, σ_1^2 , the geometric mean is the midpoint of the geodesic $\sigma_t^2 = (\sigma_0^2)^{1-t} (\sigma_1^2)^t$.

The induced mean or center

- Let $\{\Sigma_k\}_{k=1}^K$ be given matrices in $\mathcal{S}(p)$
- Let weights $\pi = (\pi_1, \dots, \pi_K)$, $\sum_{k=1}^K \pi_k = 1$, be given.

Then

$$\Sigma(\pi) = \arg \min_{\Sigma \in \mathcal{S}(p)} \sum_{k=1}^K \pi_k d(\Sigma_k, \Sigma),$$

is a **weighted mean** associated with distance (penalty) d .

- **Q:** What is a natural mean of positive definite matrices?
- If $p = 1$, so we have $\sigma_1^2, \dots, \sigma_K^2 > 0$, we could consider
 - arithmetic mean $\sigma^2 = \frac{1}{K} \sum_{k=1}^K \sigma_k^2$.
 - geometric mean $\sigma^2 = (\sigma_1^2 \cdots \sigma_K^2)^{1/K}$
 - harmonic mean $\sigma^2 = \left(\frac{1}{K} \sum_{k=1}^K (\sigma_k^2)^{-1} \right)^{-1}$
- **Note:** For a pair σ_0^2, σ_1^2 , the geometric mean is the midpoint of the geodesic $\sigma_t^2 = (\sigma_0^2)^{1-t} (\sigma_1^2)^t$.

The induced mean or center

- Let $\{\Sigma_k\}_{k=1}^K$ be given matrices in $\mathcal{S}(p)$
- Let weights $\pi = (\pi_1, \dots, \pi_K)$, $\sum_{k=1}^K \pi_k = 1$, be given.

Then

$$\Sigma(\pi) = \arg \min_{\Sigma \in \mathcal{S}(p)} \sum_{k=1}^K \pi_k d(\Sigma_k, \Sigma),$$

is a **weighted mean** associated with distance (penalty) d .

- **Q:** What is a natural mean of positive definite matrices?
- If $p = 1$, so we have $\sigma_1^2, \dots, \sigma_K^2 > 0$, we could consider
 - arithmetic mean $\sigma^2 = \frac{1}{K} \sum_{k=1}^K \sigma_k^2$.
 - geometric mean $\sigma^2 = (\sigma_1^2 \cdots \sigma_K^2)^{1/K}$
 - harmonic mean $\sigma^2 = \left(\frac{1}{K} \sum_{k=1}^K (\sigma_k^2)^{-1} \right)^{-1}$
- **Note:** For a pair σ_0^2, σ_1^2 , the geometric mean is the midpoint of the geodesic $\sigma_t^2 = (\sigma_0^2)^{1-t} (\sigma_1^2)^t$.

So for $p > 1$ what penalties could one use?

■ Frobenius distance

$$d_F(\Sigma_k, \Sigma) = \{\text{Tr}[(\Sigma_k - \Sigma)^2]\}^{1/2}$$

gives the standard weighted **arithmetic mean** $\Sigma_F(\pi) = \sum_{k=1}^K \pi_k \Sigma_k$.

... but not g -convex!

- ✓ Riemannian distance $d_R(\mathbf{A}, \mathbf{B})$
- ✓ Kullback-Leibler (KL) divergence $d_{KL}(\mathbf{A}, \mathbf{B})$
- ✓ Ellipticity distance $d_E(\mathbf{A}, \mathbf{B})$

Note: there are also some other distances that are jointly g -convex, and hence fit our framework, e.g., S-divergence of [Sra, 2011].

So for $p > 1$ what penalties could one use?

✗ Frobenius distance

$$d_F(\Sigma_k, \Sigma) = \{\text{Tr}[(\Sigma_k - \Sigma)^2]\}^{1/2}$$

gives the standard weighted **arithmetic mean** $\Sigma_F(\pi) = \sum_{k=1}^K \pi_k \Sigma_k$.
... but not g -convex!

- ✓ Riemannian distance $d_R(\mathbf{A}, \mathbf{B})$
- ✓ Kullback-Leibler (KL) divergence $d_{KL}(\mathbf{A}, \mathbf{B})$
- ✓ Ellipticity distance $d_E(\mathbf{A}, \mathbf{B})$

Note: there are also some other distances that are jointly g -convex, and hence fit our framework, e.g., S-divergence of [Sra, 2011].

So for $p > 1$ what penalties could one use?

✗ Frobenius distance

$$d_F(\Sigma_k, \Sigma) = \{\text{Tr}[(\Sigma_k - \Sigma)^2]\}^{1/2}$$

gives the standard weighted **arithmetic mean** $\Sigma_F(\pi) = \sum_{k=1}^K \pi_k \Sigma_k$.
... but not g -convex!

- ✓ Riemannian distance $d_R(\mathbf{A}, \mathbf{B})$
- ✓ Kullback-Leibler (KL) divergence $d_{KL}(\mathbf{A}, \mathbf{B})$
- ✓ Ellipticity distance $d_E(\mathbf{A}, \mathbf{B})$

Note: there are also some other distances that are jointly g -convex, and hence fit our framework, e.g., S-divergence of [\[Sra, 2011\]](#).

Riemannian distance

■ Riemannian distance

$$d_R(\mathbf{A}, \mathbf{B}) = \|\log(\mathbf{A}^{-1/2}\mathbf{B}\mathbf{A}^{-1/2})\|_F^2,$$

is the length of the geodesic curve between \mathbf{A} and \mathbf{B} .

- The induced mean, called the **Riemannian (or Karcher) mean** is a unique solution to [Bhatia, 2009]

$$\sum_{k=1}^K \pi_k \log(\Sigma_R^{1/2} \Sigma_k^{-1} \Sigma_R^{1/2}) = \mathbf{0}$$

- ☹ No closed-form solution: a number of complex numerical approaches have been proposed in the literature.

Kullback-Leibler (KL) divergence

$$d_{\text{KL}}(\mathbf{A}, \mathbf{B}) = \text{Tr}(\mathbf{A}^{-1}\mathbf{B}) - \log |\mathbf{A}^{-1}\mathbf{B}| - p$$

- KL-distance verifies $d_{\text{KL}}(\mathbf{A}, \mathbf{B}) \geq 0$ and $= 0$ for $\mathbf{A} = \mathbf{B}$.
- utilized as shrinkage penalty in [Sun et al., 2014].

Result 3 [Ollila et al., 2016]

$d_{\text{KL}}(\mathbf{A}, \mathbf{B})$ is **jointly strictly g -convex** and affine invariant and the mean based on it has a unique solution in closed form:

$$\begin{aligned}\Sigma_{\text{I}}(\boldsymbol{\pi}) &= \arg \min_{\boldsymbol{\Sigma} \in \mathcal{S}(p)} \sum_{i=1}^K \pi_k d_{\text{KL}}(\boldsymbol{\Sigma}_k, \boldsymbol{\Sigma}) \\ &= \left(\sum_{k=1}^K \pi_k \boldsymbol{\Sigma}_k^{-1} \right)^{-1},\end{aligned}$$

which is a **weighted harmonic mean** of PDS matrices.

Special case: target matrix $\mathbf{T} = \mathbf{I}$

- If the shrinkage target is $\mathbf{T} = \mathbf{I}$, then the criterion using KL-distance

$$\begin{aligned}\mathcal{L}_{\text{KL},\beta}(\Sigma) &= \beta \mathcal{L}(\Sigma) + (1 - \beta) d_{\text{KL}}(\Sigma, \mathbf{I}) \\ &= \beta \left\{ \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| \right\} + (1 - \beta) \underbrace{\{\text{Tr}(\Sigma^{-1}) - \ln |\Sigma^{-1}|\}}_{d_{\text{KL}}(\Sigma, \mathbf{I})}\end{aligned}$$

looks closely similar to the optimization program which we studied earlier:

$$\mathcal{L}_{\alpha,\beta}(\Sigma) = \beta \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| + \alpha \text{Tr}(\Sigma^{-1}),$$

which utilized the penalty $\mathcal{P}(\Sigma) = \text{Tr}(\Sigma^{-1})$ and a tuned ρ -function $\rho_\beta(t) = \beta \rho(t)$, $\beta > 0$.

Special case: target matrix $\mathbf{T} = \mathbf{I}$ (cont'd)

- Note that

$$\begin{aligned}\mathcal{L}_{\alpha,\beta}(\Sigma) &= \beta \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| + \alpha \text{Tr}(\Sigma^{-1}) \\ &= \beta \underbrace{\left\{ \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) - \ln |\Sigma^{-1}| \right\}}_{=\mathcal{L}(\Sigma)} - (1 - \beta) \ln |\Sigma^{-1}| + \alpha \text{Tr}(\Sigma^{-1})\end{aligned}$$

- This shows that $\mathcal{L}_{\alpha,\beta}(\Sigma) = \mathcal{L}_{\text{KL},\beta}(\Sigma)$ when $\alpha = 1 - \beta$
- Thus results given earlier (e.g. [▶ Result 1\(b\)](#)) transfer directly to penalization using KL-penalty.

Ellipticity distance

$$d_E(\mathbf{A}, \mathbf{B}) = p \log \frac{1}{p} \text{Tr}(\mathbf{A}^{-1} \mathbf{B}) - \log |\mathbf{A}^{-1} \mathbf{B}|$$

- d_E is scale invariant. Note: Scale invariance is a useful property for estimators that are scale invariant, e.g., Tyler's M -estimator.
- utilized as shrinkage penalty in [Wiesel, 2012]
- Related to ellipticity factor, $e(\Sigma) = \frac{1}{p} \text{Tr}(\Sigma) / |\Sigma|^{1/p}$, the ratio of the arithmetic and geometric means of the eigenvalues of Σ .

Result 4 [Ollila et al., 2016]

$d_E(\mathbf{A}, \mathbf{B})$ is jointly g -convex and affine and scale invariant. The induced mean is unique (up to a scale) and solves

$$\Sigma_E = \left(\sum_{k=1}^K \pi_k \frac{p \Sigma_k^{-1}}{\text{Tr}(\Sigma_k^{-1} \Sigma_E)} \right)^{-1},$$

which is an (implicitly) weighted harmonic mean of normalized Σ_k -s.

Critical points

$$\min_{\Sigma \in \mathcal{S}(p)} \{ \beta \mathcal{L}(\Sigma) + (1 - \beta) d(\Sigma, \mathbf{T}) \}, \quad \beta \in (0, 1]$$

- Write $\mathcal{P}_0(\Sigma) = d(\Sigma, \mathbf{T})$ and $\mathcal{P}'_0(\Sigma) = \partial \mathcal{P}(\Sigma) / \partial \Sigma^{-1}$.
- The critical points then verify

$$\mathbf{0} = \beta \left\{ \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top - \Sigma \right\} + (1 - \beta) \mathcal{P}'_0(\Sigma)$$

$$\Leftrightarrow \beta \Sigma = \beta \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + (1 - \beta) \mathcal{P}'_0(\Sigma)$$

$$\Leftrightarrow \Sigma = \beta \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + (1 - \beta) \{ \mathcal{P}'_0(\Sigma) + \Sigma \}.$$

- For $\mathcal{P}_0(\Sigma) = d_{\text{KL}}(\Sigma, \mathbf{T}) = \text{Tr}(\Sigma^{-1} \mathbf{T}) - \log |\Sigma^{-1} \mathbf{T}| - p$, this gives

$$\Sigma = \beta \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \Sigma^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + (1 - \beta) \mathbf{T}.$$

Menu

- I. Ad-hoc shrinkage SCM-s of multiple samples
- II. ML- and M -estimators of scatter matrix
- III. Geodesic convexity
- IV. Regularized M -estimators
- V. Penalized estimation of multiple covariances
- VI. Estimation of the regularization parameter
- VII. Applications

Reference



Ollila, E., Soloveychik, I., Tyler, D. E. and Wiesel, A. (2016).
Simultaneous penalized M-estimation of covariance matrices using
geodesically convex optimization
Journal of Multivariate Analysis (under review), Cite as:
arXiv:1608.08126 [stat.ME]
<http://arxiv.org/abs/1608.08126>

Multiple covariance estimation problem

- We are given K groups of elliptically distributed measurements,

$$\mathbf{x}_{11}, \dots, \mathbf{x}_{1n_1}, \quad \dots, \quad \mathbf{x}_{K1}, \dots, \mathbf{x}_{Kn_K}$$

- Each group $\mathbf{X}_k = \{\mathbf{x}_{k1}, \dots, \mathbf{x}_{kn_k}\}$ containing n_k p -dimensional samples, and

$$N = \sum_{i=1}^K n_k = \text{total sample size}$$

$$\pi_k = \frac{n_k}{N} = \text{relative sample size of the } k\text{-th group}$$

- Sample populations follow elliptical distributions, $\mathcal{E}_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, g_k)$, with different scatter matrices $\boldsymbol{\Sigma}_k$ possessing mutual structure or a **joint center** $\boldsymbol{\Sigma} \Rightarrow$ need to estimate **both** $\{\boldsymbol{\Sigma}_k\}_{k=1}^K$ **and** $\boldsymbol{\Sigma}$.
- We assume that symmetry center $\boldsymbol{\mu}_k$ of populations is known or that data is *centered*.

Proposal 1: Regularization towards a pooled center

- A pooled M -estimator of scatter is defined as a minimum of

$$\mathcal{L}(\Sigma) = \sum_{k=1}^K \pi_k \mathcal{L}_k(\Sigma) = \frac{1}{N} \left\{ \sum_{k=1}^K \sum_{i=1}^{n_k} \rho_k(\mathbf{x}_{ki}^\top \Sigma^{-1} \mathbf{x}_{ki}) \right\} - \log |\Sigma^{-1}|$$

over $\Sigma \in \mathcal{S}(p)$.

- Penalized M -estimators of scatter for the individual groups solve

$$\min_{\Sigma_k \in \mathcal{S}(p)} \left\{ \beta \mathcal{L}_k(\Sigma_k) + (1 - \beta) d(\Sigma_k, \hat{\Sigma}) \right\}, \quad k = 1, \dots, K,$$

where

- $\beta \in (0, 1]$ is a regularization/penalty parameter
- $d(\mathbf{A}, \mathbf{B}) : \mathcal{S}(p) \times \mathcal{S}(p) \rightarrow \mathbb{R}_0^+$ is penalty/distance fnc.

Distance $d(\Sigma_k, \hat{\Sigma})$ enforce similarity of Σ_k -s to joint center $\hat{\Sigma}$ and β controls the amount of shrinkage towards $\hat{\Sigma}$.

Proposal 2: Joint regularization enforcing similarity among the group scatter matrices

$$\underset{\{\Sigma_k\}_{k=1}^K, \Sigma \in \mathcal{S}(p)}{\text{minimize}} \quad \sum_{k=1}^K \pi_k \{ \beta \mathcal{L}_k(\Sigma_k) + (1 - \beta) d(\Sigma_k, \Sigma) \}$$

where β is the penalty parameter, $d(\Sigma_k, \Sigma)$ is the distance function as before, and

$$\mathcal{L}_k(\Sigma_k) = \frac{1}{n_k} \sum_{i=1}^{n_k} \rho_k(\mathbf{x}_{ki}^\top \Sigma_k^{-1} \mathbf{x}_{ki}) - \log |\Sigma_k^{-1}|$$

is the M(L)-cost fnc for the k -th class and $\rho_k(\cdot)$ is the loss fnc.

‘Center’ Σ can now be viewed as ‘average’ of Σ_k -s. Namely, for fixed Σ_k -s, the minimum $\hat{\Sigma}$ is found by solving

$$\hat{\Sigma}(\pi) = \arg \min_{\Sigma \in \mathcal{S}(p)} \sum_{k=1}^K \pi_k d(\Sigma_k, \Sigma),$$

which represents the **weighted mean** associated with the distance d .

Modifications to Proposals 1 and 2

- Penalty parameter β can be replaced by **individual tuning constants** β_k , $k = 1, \dots, K$ for each class.

Comment: typically one tends to choose small β_k when sample size is small, but this does not seem to be necessary in our framework

- In Proposal 1, if the total sample size N is small (e.g., $N < p$), then one may add a penalty $\mathcal{P}(\Sigma) = \text{Tr}(\Sigma^{-1})$ and compute pooled center $\hat{\Sigma}$ as a **pooled regularized M -estimator**:

$$\min_{\Sigma} \sum_{k=1}^K \pi_k \mathcal{L}_k(\Sigma) + \gamma \mathcal{P}(\Sigma)$$

where $\gamma > 0$ is the (additional) penalty parameter for the center.

- Such a penalty term can be added to Proposal 2 as well.

- We consider the cases that penalty function $d(\mathbf{A}, \mathbf{B})$ is the KL-distance or ellipticity distance.
- Both distances are *affine invariant*, i.e.

$$d(\mathbf{A}, \mathbf{B}) = d(\mathbf{C}\mathbf{A}\mathbf{C}^\top, \mathbf{C}\mathbf{B}\mathbf{C}^\top), \quad \forall \text{ nonsingular } \mathbf{C}.$$

which is Property D4 in Slide

- If D4 holds, the resulting estimators are **affine equivariant**:

if $\mathbf{x}_{ki} \rightarrow \mathbf{C}\mathbf{x}_{ki}$ for all $k = 1, \dots, K; i = 1, \dots, n_k$
 then $\{\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K, \boldsymbol{\Sigma}\} \rightarrow \{\mathbf{C}\boldsymbol{\Sigma}_1\mathbf{C}^\top, \dots, \mathbf{C}\boldsymbol{\Sigma}_K\mathbf{C}^\top, \mathbf{C}\boldsymbol{\Sigma}\mathbf{C}^\top\}.$

Critical points/algorithm using KL-divergence penalty

Problem:
$$\min_{\{\Sigma_k\}_{k=1}^K, \Sigma} \sum_{k=1}^K \pi_k \{ \beta \mathcal{L}_k(\Sigma_k) + (1 - \beta) d_{\text{KL}}(\Sigma_k, \Sigma) \}$$

Solving

$$\mathbf{0} = \beta \frac{\partial \mathcal{L}_k(\Sigma_k)}{\partial \Sigma_k^{-1}} + (1 - \beta) \frac{\partial d_{\text{KL}}(\Sigma_k, \Sigma)}{\partial \Sigma_k^{-1}}, \quad k = 1, \dots, K$$

$$\mathbf{0} = \sum_{k=1}^K \pi_k \frac{\partial d_{\text{KL}}(\Sigma_k, \Sigma)}{\partial \Sigma}$$

yields **estimating equations**

$$\Sigma_k = \beta \frac{1}{n_k} \sum_{i=1}^{n_k} u_k(\mathbf{x}_{ki}^\top \Sigma_k^{-1} \mathbf{x}_{ki}) \mathbf{x}_{ki} \mathbf{x}_{ki}^\top + (1 - \beta) \Sigma$$

$$\Sigma = \left(\sum_{k=1}^K \pi_k \Sigma_k^{-1} \right)^{-1}$$

where $u_k(t) = \rho'_k(t)$, $k = 1, \dots, K$.

Critical points/algorithm using KL-divergence penalty

Problem:
$$\min_{\{\Sigma_k\}_{k=1}^K, \Sigma} \sum_{k=1}^K \pi_k \{ \beta \mathcal{L}_k(\Sigma_k) + (1 - \beta) d_{\text{KL}}(\Sigma_k, \Sigma) \}$$

Solving

$$\mathbf{0} = \beta \frac{\partial \mathcal{L}_k(\Sigma_k)}{\partial \Sigma_k^{-1}} + (1 - \beta) \frac{\partial d_{\text{KL}}(\Sigma_k, \Sigma)}{\partial \Sigma_k^{-1}}, \quad k = 1, \dots, K$$

$$\mathbf{0} = \sum_{k=1}^K \pi_k \frac{\partial d_{\text{KL}}(\Sigma_k, \Sigma)}{\partial \Sigma}$$

yields **algorithm** that updates covariances cyclically from $\Sigma_1, \dots, \Sigma_K$ to Σ

$$\Sigma_k \leftarrow \beta \frac{1}{n_k} \sum_{i=1}^{n_k} u_k(\mathbf{x}_{ki}^\top \Sigma_k^{-1} \mathbf{x}_{ki}) \mathbf{x}_{ki} \mathbf{x}_{ki}^\top + (1 - \beta) \Sigma$$

$$\Sigma \leftarrow \left(\sum_{k=1}^K \pi_k \Sigma_k^{-1} \right)^{-1}$$

where $u_k(t) = \rho'_k(t)$, $k = 1, \dots, K$.

Critical points/algorithm using ellipticity distance

As for KL-distance, we can easily solve the estimating equations and propose a cyclic algorithm to find the solutions.

Estimating equations

$$\Sigma_k = \beta \frac{1}{n_k} \sum_{i=1}^{n_k} u_k(\mathbf{x}_{ki}^\top \Sigma_k^{-1} \mathbf{x}_{ki}) \mathbf{x}_{ki} \mathbf{x}_{ki}^\top + (1 - \beta) \frac{p \Sigma}{\text{Tr}(\Sigma_k^{-1} \Sigma)},$$
$$\Sigma = \left(\sum_{k=1}^K \pi_k \frac{p \Sigma_k^{-1}}{\text{Tr}(\Sigma_k^{-1} \Sigma)} \right)^{-1}$$

where $u_k(t) = \rho'_k(t)$, $k = 1, \dots, K$.

Critical points/algorithm using ellipticity distance

As for KL-distance, we can easily solve the estimating equations and propose a cyclic algorithm to find the solutions.

Algorithm updates covariances cyclically from $\Sigma_1, \dots, \Sigma_K$ to Σ

$$\Sigma_k \leftarrow \beta \frac{1}{n_k} \sum_{i=1}^{n_k} u_k(\mathbf{x}_{ki}^\top \Sigma_k^{-1} \mathbf{x}_{ki}) \mathbf{x}_{ki} \mathbf{x}_{ki}^\top + (1 - \beta) \frac{p \Sigma}{\text{Tr}(\Sigma_k^{-1} \Sigma)},$$
$$\Sigma \leftarrow \left(\sum_{k=1}^K \pi_k \frac{p \Sigma_k^{-1}}{\text{Tr}(\Sigma_k^{-1} \Sigma)} \right)^{-1}$$

where $u_k(t) = \rho'_k(t)$, $k = 1, \dots, K$.

Menu

I. Ad-hoc shrinkage SCM-s of multiple samples

II. ML- and M -estimators of scatter matrix

III. Geodesic convexity

IV. Regularized M -estimators

V. Penalized estimation of multiple covariances

VI. Estimation of the regularization parameter

- Cross-validation
- Oracle approach

VII. Applications

Estimation of the regularization parameter

- Recall that the regularized M -estimator introduced in Section V solve

$$\hat{\Sigma} = \beta \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + \alpha \mathbf{I}$$

- For simplicity, we often tune only one parameter and set:

$$\beta = (1 - \alpha), \alpha \in (0, 1) \quad \text{or} \quad \alpha = (1 - \beta), \beta \in (0, 1).$$

- We need a disciplined way of choosing β (or α): it determines the amount of shrinkage towards the target matrix (here the identity matrix)
- Also related to a wider topic of **model selection**.
- A judicious choice of β is the one that provides an estimate $\hat{\Sigma}$ that minimizes the mean squared error (or PE).

Approaches:

- 1 Cross-validation
- 2 Oracle/Clairvoyant approach
- 3 Expected likelihood approach
[\[Abramovich and Besson, 2013, Besson and Abramovich, 2013\]](#)
- 4 Random matrix theory (Frederic's talk).

Only the first two topics are addressed in this tutorial.

Cross-validation (CV)

- Ideally, we would split our available data into two portions:
 - **training set** $\{\mathbf{x}_i\}_{i=1}^n$ to estimate $\hat{\Sigma}$ which then gives the **fit** $\hat{F}(\mathbf{x})$.
 - **test set** $\{\mathbf{x}_i^t\}_{i=1}^n$ to validate the model, i.e., to "test" or assess how small is (say) the sum of the fitted values, $\sum_i \hat{F}(\mathbf{x}_i^t)$.
- This may not be plausible (especially if $n \approx p$) and also impractical/inefficient usage of data
- **Note:** The choice of the fit function $\hat{F}(\mathbf{x})$ is not as obvious as it is in the regression setting, where the (squared) prediction error is a natural choice.
- Cross validation solves the problem by splitting the data into K folds, fits the data on $K - 1$ folds, and evaluates the performance (using fit criterion) on the fold that was left out.

Cross-validation (cont'd)

- Partition $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ into Q separate sets of similar size

$$I_1 \cup I_2 \cup \dots \cup I_Q = \{1, \dots, n\} \equiv [n]$$

- Common choices: $Q = 5, 10$ or $Q = n$ (*leave-one-out CV*).
- Taking q th fold out (all \mathbf{x}_i , $i \in I_q$) gives a reduced data set \mathbf{X}_{-q} .
- As the fit function, we use the (non-penalized) criterion function

$$\hat{F}(\mathbf{x}) = \rho(\mathbf{x}^\top \hat{\Sigma}^{-1} \mathbf{x}) - \log |\hat{\Sigma}^{-1}|$$

Test	Train	Train	Train	Train
Train	Test	Train	Train	Train
Train	Train	Test	Train	Train
Train	Train	Train	Test	Train
Train	Train	Train	Train	Test

Figure: $K = 5$ cross validation

CV procedure [Ollila et al., 2016]

For simplicity, tune only one parameter and set $\alpha = 1 - \beta$.

- 1 **for** $\beta \in [\beta]$ (= a grid of β values in $(0, 1)$) and $q \in \{1, \dots, Q\}$ **do**
 - Compute regularized M -estimator based on \mathbf{X}_{-q} , denoted $\hat{\Sigma}(\beta, q)$
 - **CV fit** for β is computed as the sum of fits over the q th folds that were left out:

$$\text{CV}(\beta, q) = \sum_{\tilde{q} \in I_q} \rho(\mathbf{x}_{\tilde{q}}^{\top} [\hat{\Sigma}(\beta, q)]^{-1} \mathbf{x}_{\tilde{q}}) - (\#I_q) \cdot \log |\hat{\Sigma}(\beta, q)^{-1}|$$

end

- 2 Compute the average CV fit: $\text{CV}(\beta) = \frac{1}{Q} \sum_{q=1}^Q \text{CV}(\beta, q)$, $\forall \beta \in [\beta]$.
- 3 Select $\hat{\beta}_{\text{CV}} = \arg \min_{\beta \in [\beta]} \text{CV}(\beta)$.
- 4 Compute $\hat{\Sigma}$ based on the entire data set \mathbf{X} using $\beta = \hat{\beta}_{\text{CV}}$.

Oracle/Clairvoyant approach

- The regularized M -estimator of scatter is a solution of

$$\hat{\Sigma} = \beta \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + \alpha \mathbf{I}$$

for given penalty parameters $\alpha, \beta > 0$.

- Since $\hat{\Sigma}$ does not have a closed form expression, deriving an analytic expression for an expected loss (such as the MSE) is virtually impossible.
- **Solution:** use an approximation for $\hat{\Sigma}$.
- Given the true scatter (covariance) matrix Σ_0 , then

$$\Sigma_{\alpha, \beta} = (1 - \alpha) \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \Sigma_0^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + \alpha \mathbf{I}$$

called the **clairvoyant estimator**, is a closed-form approximation of $\hat{\Sigma}$.

Oracle/Clairvoyant approach

- The regularized M -estimator of scatter is a solution of

$$\hat{\Sigma} = \beta \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + \alpha \mathbf{I}$$

for given penalty parameters $\alpha, \beta > 0$.

- Since $\hat{\Sigma}$ does not have a closed form expression, deriving an analytic expression for an expected loss (such as the MSE) is virtually impossible.
- **Solution:** use an approximation for $\hat{\Sigma}$.
- Given the true scatter (covariance) matrix Σ_0 , then

$$\Sigma_{\alpha,\beta} = (1 - \alpha) \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \Sigma_0^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + \alpha \mathbf{I}$$

called the **clairvoyant estimator**, is a closed-form approximation of $\hat{\Sigma}$.

Oracle/Clairvoyant approach

- The regularized M -estimator of scatter is a solution of

$$\hat{\Sigma} = \beta \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \hat{\Sigma}^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + \alpha \mathbf{I}$$

for given penalty parameters $\alpha, \beta > 0$.

- Since $\hat{\Sigma}$ does not have a closed form expression, deriving an analytic expression for an expected loss (such as the MSE) is virtually impossible.
- **Solution:** use an approximation for $\hat{\Sigma}$.
- Given the true scatter (covariance) matrix Σ_0 , then

$$\Sigma_{\alpha, \beta} = (1 - \alpha) \frac{1}{n} \sum_{i=1}^n u(\mathbf{x}_i^\top \Sigma_0^{-1} \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top + \alpha \mathbf{I}$$

called the **clairvoyant estimator**, is a closed-form approximation of $\hat{\Sigma}$.

Oracle/Clairvoyant approach (cont'd)

Let Σ_0 denote the true unknown scatter (covariance) matrix. Oracle penalty parameters (α_0, β_0) minimize the expected loss,

$$(\alpha_o, \beta_o) = \arg \min_{\alpha, \beta > 0} \mathbb{E}[d(\Sigma_{\alpha, \beta}, \Sigma_0)],$$

for some suitable distance function $d(\mathbf{A}, \mathbf{B})$.

- Naturally the found solution (α_o, β_o) will depend on the true scatter matrix Σ_0 which is unknown in practise.
- Replace the unknown true Σ_0 in α_0 and β_0 with some preliminary estimate or guess $\hat{\Sigma}_0$
 $\Rightarrow \hat{\alpha}_o = \alpha_o(\hat{\Sigma}_0)$ and $\hat{\beta} = \beta_o(\hat{\Sigma}_0)$ are the oracle/clairvoyant estimates.

Oracle approach for the regularized SCM

- Let us first consider the regularized M -estimator using (tuned) Gaussian loss function $\rho(t) = \beta t$ ($u(t) = \beta$)
- This gives the regularized SCM (a.k.a. Ledoit-Wolf estimator) :

$$\mathbf{S}_{\alpha,\beta} = \beta \mathbf{S} + \alpha \mathbf{I},$$

where \mathbf{S} is the unbiased SCM $\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, so $\mathbb{E}[\mathbf{S}] = \Sigma$

- **Clairvoyant (oracle) estimator** of $(\alpha, \beta) \in \mathbb{R}_0^+ \times \mathbb{R}^+$ is defined as a minimizer of the mean squared error (MSE),

$$(\alpha_o, \beta_o) = \arg \min_{(\alpha, \beta)} \mathbb{E} \left[\|\Sigma_0 - \mathbf{S}_{\alpha, \beta}\|_F^2 \right].$$

See e.g., [Du et al., 2010].

- The solution will naturally depend on the true unknown covariance matrix Σ_0 .

- The oracle solution is easily found to be

$$\beta_o = \frac{\gamma(\Sigma_0)}{\gamma(\Sigma_0) + \varrho(\Sigma_0)}$$
$$\alpha_o = (1 - \beta_o)\text{Tr}(\Sigma_0)/p.$$

where $\varrho(\Sigma_0) = \mathbb{E}[\|\mathbf{S} - \Sigma_0\|_F^2] = \frac{1}{n}\mathbb{E}[\|\mathbf{x}\|^4] - \frac{1}{n}\text{Tr}(\Sigma_0^2)$

$$\gamma(\Sigma) = \text{Tr}(\Sigma_0^2) - \frac{1}{p}\{\text{Tr}(\Sigma_0)\}^2$$

- **Note:** $\varrho(\Sigma_0)$ is the MSE between \mathbf{S} and Σ_0 and $\beta \in [0, 1]$ and $\alpha \geq 0$.
- **Estimation:** replace the unknown Σ by \mathbf{S} , and $\mathbb{E}[\|\mathbf{x}\|^4]$ by finite sample average.
- For more advanced approach, see [Ledoit and Wolf, 2004].

- The oracle solution is easily found to be

$$\hat{\beta}_o = \frac{\gamma(\mathbf{S})}{\gamma(\mathbf{S}) + \varrho(\mathbf{S})}$$
$$\hat{\alpha}_o = (1 - \hat{\beta}_o) \text{Tr}(\mathbf{S})/p \quad .$$

$$\text{where } \varrho(\mathbf{S}) = \mathbb{E}[\|\mathbf{S} - \Sigma_0\|_F^2] = \frac{1}{n^2} \sum_{i=1}^n \|\mathbf{x}_i\|^4 - \frac{1}{n} \text{Tr}(\mathbf{S}^2)$$

$$\gamma(\mathbf{S}) = \text{Tr}(\mathbf{S}^2) - \frac{1}{p} \{\text{Tr}(\mathbf{S})\}^2$$

- **Note:** $\varrho(\Sigma_0)$ is the MSE between \mathbf{S} and Σ_0 and $\beta \in [0, 1]$ and $\alpha \geq 0$.
- **Estimation:** replace the unknown Σ by \mathbf{S} , and $\mathbb{E}[\|\mathbf{x}\|^4]$ by finite sample average.
- For more advanced approach, see [Ledoit and Wolf, 2004].

- The oracle solution is easily found to be

$$\hat{\beta}_o = \frac{\gamma(\mathbf{S})}{\gamma(\mathbf{S}) + \varrho(\mathbf{S})}$$

$$\hat{\alpha}_o = (1 - \hat{\beta}_o) \text{Tr}(\mathbf{S})/p \quad .$$

$$\text{where } \varrho(\mathbf{S}) = \mathbb{E}[\|\mathbf{S} - \boldsymbol{\Sigma}_0\|_F^2] = \frac{1}{n^2} \sum_{i=1}^n \|\mathbf{x}_i\|^4 - \frac{1}{n} \text{Tr}(\mathbf{S}^2)$$

$$\gamma(\mathbf{S}) = \text{Tr}(\mathbf{S}^2) - \frac{1}{p} \{\text{Tr}(\mathbf{S})\}^2$$

- **Note:** $\varrho(\boldsymbol{\Sigma}_0)$ is the MSE between \mathbf{S} and $\boldsymbol{\Sigma}_0$ and $\beta \in [0, 1]$ and $\alpha \geq 0$.
- **Estimation:** replace the unknown $\boldsymbol{\Sigma}$ by \mathbf{S} , and $\mathbb{E}[\|\mathbf{x}\|^4]$ by finite sample average.
- For more advanced approach, see [[Ledoit and Wolf, 2004](#)].

Oracle approach for regularized Tyler's M -estimator

- Choose $\beta = 1 - \alpha$, $\alpha \in (0, 1)$.
- The clairvoyant estimator of Tyler's M -estimator $\hat{\Sigma}$ is

$$\Sigma_\alpha = (1 - \alpha) \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \Sigma_0^{-1} \mathbf{x}_i} + \alpha \mathbf{I}.$$

- **Q:** What distance d to use in minimizing the expected loss?
- **idea:** Given a shape matrix Σ_0 , verifying $\text{Tr}(\Sigma_0^{-1}) = p$, choose α so that $\Sigma_0^{-1} \Sigma_\alpha$ is as close as possible to $c\mathbf{I}$, for some $c > 0$.
- A natural distance that measures similarity in shape:

$$d(\Sigma_0, \Sigma_\alpha) = \|\Sigma_0^{-1} \Sigma_\alpha - \frac{1}{p} \text{Tr}(\Sigma_0^{-1} \Sigma_\alpha) \mathbf{I}\|^2$$

- The value $\alpha_o \in (0, 1)$ that minimizes the expected loss $\mathbb{E}[d(\Sigma_0, \Sigma_\alpha)]$ is [Ollila and Tyler, 2014, Theorem 5]:

$$\alpha_o = \frac{p - 2 + p \text{Tr}(\Sigma_0)}{p - 2 + p \text{Tr}(\Sigma_0) + n(p + 2) \{p^{-1} \text{Tr}(\Sigma_0^{-2}) - 1\}}$$

Oracle approach for regularized Tyler's M -estimator

- Choose $\beta = 1 - \alpha$, $\alpha \in (0, 1)$.
- The clairvoyant estimator of Tyler's M -estimator $\hat{\Sigma}$ is

$$\Sigma_\alpha = (1 - \alpha) \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \Sigma_0^{-1} \mathbf{x}_i} + \alpha \mathbf{I}.$$

- **Q:** What distance d to use in minimizing the expected loss?
- **idea:** Given a shape matrix Σ_0 , verifying $\text{Tr}(\Sigma_0^{-1}) = p$, choose α so that $\Sigma_0^{-1} \Sigma_\alpha$ is as close as possible to $c\mathbf{I}$, for some $c > 0$.
- A natural distance that measures similarity in shape:

$$d(\Sigma_0, \Sigma_\alpha) = \|\Sigma_0^{-1} \Sigma_\alpha - \frac{1}{p} \text{Tr}(\Sigma_0^{-1} \Sigma_\alpha) \mathbf{I}\|^2$$

- The value $\alpha_o \in (0, 1)$ that minimizes the expected loss $\mathbb{E}[d(\Sigma_0, \Sigma_\alpha)]$ is [Ollila and Tyler, 2014, Theorem 5]:

$$\alpha_o = \frac{p - 2 + p \text{Tr}(\Sigma_0)}{p - 2 + p \text{Tr}(\Sigma_0) + n(p + 2) \{p^{-1} \text{Tr}(\Sigma_0^{-2}) - 1\}}$$

Oracle approach for regularized Tyler's M -estimator

- Choose $\beta = 1 - \alpha$, $\alpha \in (0, 1)$.
- The clairvoyant estimator of Tyler's M -estimator $\hat{\Sigma}$ is

$$\Sigma_\alpha = (1 - \alpha) \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \Sigma_0^{-1} \mathbf{x}_i} + \alpha \mathbf{I}.$$

- **Q:** What distance d to use in minimizing the expected loss?
- **idea:** Given a shape matrix Σ_0 , verifying $\text{Tr}(\Sigma_0^{-1}) = p$, choose α so that $\Sigma_0^{-1} \Sigma_\alpha$ is as close as possible to $c\mathbf{I}$, for some $c > 0$.
- A natural distance that measures similarity in shape:

$$d(\Sigma_0, \Sigma_\alpha) = \left\| \Sigma_0^{-1} \Sigma_\alpha - \frac{1}{p} \text{Tr}(\Sigma_0^{-1} \Sigma_\alpha) \mathbf{I} \right\|^2$$

- The value $\alpha_o \in (0, 1)$ that minimizes the expected loss $\mathbb{E}[d(\Sigma_0, \Sigma_\alpha)]$ is [Ollila and Tyler, 2014, Theorem 5]:

$$\alpha_o = \frac{p - 2 + p \text{Tr}(\Sigma_0)}{p - 2 + p \text{Tr}(\Sigma_0) + n(p + 2) \{p^{-1} \text{Tr}(\Sigma_0^{-2}) - 1\}}$$

Oracle approach for regularized Tyler's M -estimator

- Choose $\beta = 1 - \alpha$, $\alpha \in (0, 1)$.
- The clairvoyant estimator of Tyler's M -estimator $\hat{\Sigma}$ is

$$\Sigma_\alpha = (1 - \alpha) \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\mathbf{x}_i^\top \Sigma_0^{-1} \mathbf{x}_i} + \alpha \mathbf{I}.$$

- **Q:** What distance d to use in minimizing the expected loss?
- **idea:** Given a shape matrix Σ_0 , verifying $\text{Tr}(\Sigma_0^{-1}) = p$, choose α so that $\Sigma_0^{-1} \Sigma_\alpha$ is as close as possible to $c\mathbf{I}$, for some $c > 0$.
- A natural distance that measures similarity in shape:

$$d(\Sigma_0, \Sigma_\alpha) = \left\| \Sigma_0^{-1} \Sigma_\alpha - \frac{1}{p} \text{Tr}(\Sigma_0^{-1} \Sigma_\alpha) \mathbf{I} \right\|^2$$

- The value $\alpha_o \in (0, 1)$ that minimizes the expected loss $\mathbb{E}[d(\Sigma_0, \Sigma_\alpha)]$ is [Ollila and Tyler, 2014, Theorem 5]:

$$\alpha_o = \frac{p - 2 + p \text{Tr}(\Sigma_0)}{p - 2 + p \text{Tr}(\Sigma_0) + n(p + 2) \{p^{-1} \text{Tr}(\Sigma_0^{-2}) - 1\}}$$

Oracle approach for regularized Tyler's M -estimator

- Similar expression is derived for complex-valued case also in [Ollila and Tyler, 2014].
- Estimate $\hat{\alpha}_o = \alpha_o(\hat{\Sigma}_0)$ is obtained by using an estimate $\hat{\Sigma}_0$ which is
 - Conventional (non-regularized) Tyler's M -estimator normalized so that $\text{Tr}(\hat{\Sigma}_0^{-1}) = p$ when $n \geq p$
 - Regularized Tyler's M -estimator using $\beta < n/p$ and $\alpha = 1 - \beta$ when $n < p$.

Oracle approach for DL-FP estimator

$$\Sigma_{\alpha} = (1 - \alpha) \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^{\top}}{\mathbf{x}_i^{\top} \Sigma_0^{-1} \mathbf{x}_i} + \alpha \mathbf{I}.$$

[Chen et al., 2011] proposed an oracle estimator for the tuning parameter of DL-FP estimator ▸ defined in this slide

- Given a shape matrix Σ_0 , verifying $\text{Tr}(\Sigma_0) = p$, find α as

$$\alpha_o = \arg \min_{\alpha} \mathbb{E} \left[\|\Sigma_0 - \Sigma_{\alpha}\|_{\text{F}}^2 \right]$$

- The obtained oracle estimator is (in the real-valued case):

$$\alpha_o = \frac{p^3 + (p - 2)\text{Tr}(\Sigma_0^2)}{\{p^3 + (p - 2)\text{Tr}(\Sigma_0^2)\} + n(p + 2)(\text{Tr}(\Sigma_0^2) - p)}.$$

- Estimate $\hat{\alpha}_o = \alpha_o(\hat{\Sigma}_0)$ is obtained using trace normalized sample sign covariance matrix

$$\hat{\Sigma}_0 = \frac{p}{n} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^{\top}}{\|\mathbf{x}_i\|^2}.$$

Numerical Example

- n realizations from a p -variate (complex) Gaussian distribution with covariance matrix

$$[\Sigma]_{ij} = \rho^{|i-j|}, \quad \rho \in (0, 1).$$

- Distance measure:

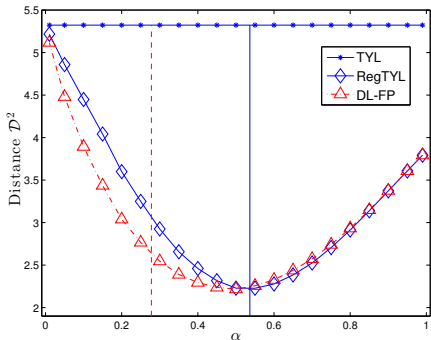
$$\mathcal{D}^2 \equiv \mathcal{D}^2(\Sigma, \hat{\Sigma}) = \|\{p/\text{Tr}(\Sigma^{-1}\hat{\Sigma})\} \Sigma^{-1}\hat{\Sigma} - \mathbf{I}\|_{\text{F}}^2$$

measures the ability $\hat{\Sigma}$ to estimate Σ up to a scale.

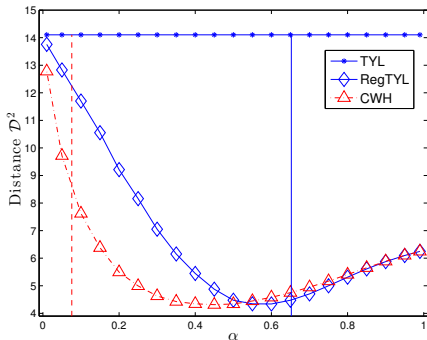
- Simulation report: Average $\mathcal{D}^2(\Sigma, \hat{\Sigma})$ over 1000 MC-trials for both the regularized Tyler M -estimator (using $\beta = 1 - \alpha$) and DL-FP estimator for different values of penalty parameter $\alpha \in (0, 1)$.

Numerical Example (cont'd)

$$n = 48, p = 12$$



$$\rho = 0.5$$



$$\rho = 0.8$$

- Solid vertical line: oracle value α_o for the regularized Tyler's M -estimator given in Slide 92.
- Dotted vertical line: oracle value α_o of DL-FP given in Slide 94.

Menu

I. Ad-hoc shrinkage SCM-s of multiple samples

II. ML- and M -estimators of scatter matrix

III. Geodesic convexity

IV. Regularized M -estimators

V. Penalized estimation of multiple covariances

VI. Estimation of the regularization parameter

VII. Applications

- Regularized discriminant analysis
- Matched filter detection

Quadratic discriminant analysis (QDA)

QDA assigns \mathbf{x} to a group \hat{k} :

$$\hat{k} = \min_{1 \leq k \leq K} \{(\mathbf{x} - \bar{\mathbf{x}}_k)^\top \mathbf{S}_k^{-1} (\mathbf{x} - \bar{\mathbf{x}}_k) + \ln |\mathbf{S}_k|\}.$$

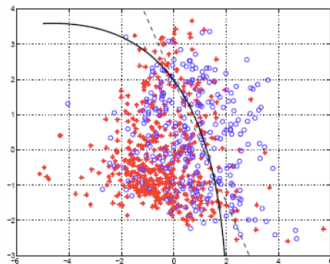
where

$$\mathbf{S}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} (\mathbf{x}_{ki} - \bar{\mathbf{x}}_k)(\mathbf{x}_{ki} - \bar{\mathbf{x}}_k)^\top$$

is the SCM of a training data set \mathbf{X}_k from k th population ($k = 1, \dots, K$).

Assumptions:

- Gaussian populations $\mathcal{N}_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- Covariance matrices can be *different* for each class $\boldsymbol{\Sigma}_i \neq \boldsymbol{\Sigma}_j \quad i \neq j$



Linear discriminant analysis (LDA)

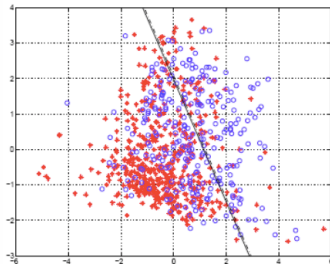
LDA assigns \mathbf{x} to a group \hat{k} :

$$\hat{k} = \min_{1 \leq k \leq K} \{(\mathbf{x} - \bar{\mathbf{x}}_k)^\top \mathbf{S}^{-1}(\mathbf{x} - \hat{\boldsymbol{\mu}}_k)\}.$$

where

$$\mathbf{S} = \sum_{k=1}^K \pi_k \mathbf{S}_k.$$

is the **pooled SCM** estimator.



Assumptions:

- Gaussian populations $\mathcal{N}_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- Covariance matrices are the *same* for each class $\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}_j \quad i \neq j$

Regularized Discriminant Analysis (RDA)

RDA* assigns \mathbf{x} to a group \hat{k} :

$$\hat{k} = \min_{1 \leq k \leq K} \{ (\mathbf{x} - \hat{\boldsymbol{\mu}}_k)^\top [\hat{\boldsymbol{\Sigma}}_k(\beta)]^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}_k) + \ln |\hat{\boldsymbol{\Sigma}}_k(\beta)| \}.$$

where $\hat{\boldsymbol{\Sigma}}_k(\beta)$ are the penalized estimators of scatter matrices obtained either using Proposal 1 or Proposal 2.

Interpretation:

- if $\beta \rightarrow 1$, we do not shrink towards joint center
 \Rightarrow RDA \rightarrow QDA
- if $\beta \rightarrow 0$, we shrink towards joint center
 \Rightarrow RDA \rightarrow LDA
- $0 < \beta < 1 \Rightarrow$ a compromise between LDA and QDA.

For robust loss fnc-s, we use **spatial median** as an estimate $\hat{\boldsymbol{\mu}}_k$ of location

* Inspired by Friedman, "Regularized discriminant Analysis", JASA (1989)

Simulation set-up

- We use the same loss function $\rho = \rho_k$ for each K samples
- **Training data:** \mathbf{X}_k -s generated from $\mathcal{N}_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ or $t_\nu(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, $\nu = 2$. These are used to estimate the discriminant rules.
- **Test data:** generated in exactly the same manner (same size $N = 100$) and classified with the discriminant rules thereby yielding an **estimate of the misclassification risk**.
- RDA rules are computed over a grid of penalty parameters $\beta \in (0, 1)$ and **optimal** (*smallest*) misclassification risk is reported.
- $\text{Prop1}(\rho, d)$ and $\text{Prop2}(\rho, d)$ refer to RDA rules based on Proposal 1 and Proposal 2 estimators, respectively, where
 - ρ refers to the used loss fnc: **Gaussian**, **Huber's**, **Tyler's**)
 - d refers to the used distance fnc: (**KL** or **Ellipticity**).

Unequal spherical covariance matrices ($\Sigma_k = kI$)

- Nr of classes is $K = 3$, total sample size $N = \sum_{k=1}^K n_k = 100$.
- $(n_1, n_2, n_3) \sim \text{Multin}(N; p_1 = p_2 = \frac{1}{4}, p_3 = \frac{1}{2})$.
- $\mu_1 = \mathbf{0}$ and remaining classes μ_k have norm equal to $\delta_k = \|\mu_k\| = 3 + k$ in orthogonal directions

Gaussian case: test misclassification errors %

method	$p = 10$	$p = 20$	$p = 30$
Oracle1	8.8 _(2.6)	6.2 _(2.3)	4.6 _(1.9)
Oracle2	9.8 _(3.1)	7.6 _(2.6)	6.0 _(2.3)
QDA	19.9 _(4.4)	—	—
LDA	17.1 _(3.8)	20.5 _(4.3)	24.0 _(4.9)
Prop1(G,KL)	12.2 _(3.1)	14.6 _(3.5)	17.9 _(4.3)
Prop1(H,KL)	12.4 _(3.2)	14.6 _(3.5)	17.7 _(4.1)
Prop1(T,E)	10.9 _(3.1)	12.1 _(3.3)	16.5 _(3.9)
Prop2(G,E)	10.5 _(3.0)	11.5 _(3.3)	15.9 _(3.8)
Prop2(T,E)	10.9 _(3.1)	12.1 _(3.3)	16.5 _(3.9)
Prop2(H,E)	10.5 _(3.0)	11.6 _(3.3)	15.7 _(3.8)
Prop2(H,KL)	12.3 _(3.2)	14.8 _(3.6)	18.0 _(4.1)

standard deviations inside parantheses in subscript

Oracle1 = QDA rule using true μ_k and Σ_k .

Oracle2 = QDA rule using true Σ_k , but estimated $\hat{\mu}_k$.

- = sample means in Gaussian case
- = spatial median in t_2 case

Unequal spherical covariance matrices ($\Sigma_k = k\mathbf{I}$)

- Nr of classes is $K = 3$, total sample size $N = \sum_{k=1}^K n_k = 100$.
- $(n_1, n_2, n_3) \sim \text{Multin}(N; p_1 = p_2 = \frac{1}{4}, p_3 = \frac{1}{2})$.
- $\mu_1 = \mathbf{0}$ and remaining classes μ_k have norm equal to $\delta_k = \|\mu_k\| = 4 + k$ in orthogonal directions

t_2 case: test misclassification errors %

Oracle1	15.7 _(3.8)	18.2 _(3.9)	21.1 _(4.0)
Oracle2	16.2 _(3.5)	19.1 _(4.2)	21.9 _(4.1)
QDA	26.9 _(5.2)	—	—
LDA	21.8 _(4.9)	25.3 _(5.3)	27.7 _(5.3)
Prop1(G,KL)	19.7 _(4.8)	22.7 _(5.2)	24.7 _(5.1)
Prop1(H,KL)	15.5 _(3.7)	17.9 _(4.0)	20.3 _(4.1)
Prop1(T,E)	16.8 _(4.0)	20.4 _(4.3)	23.4 _(4.7)
Prop2(G,E)	22.3 _(5.9)	24.3 _(5.1)	25.9 _(4.8)
Prop2(T,E)	16.8 _(4.0)	20.4 _(4.4)	23.5 _(4.8)
Prop2(H,E)	16.6 _(3.9)	20.2 _(4.4)	23.6 _(4.6)
Prop2(H,KL)	15.5 _(3.7)	17.9 _(4.0)	20.5 _(4.1)

standard deviations inside parantheses in subscript

Oracle1 = QDA rule using true μ_k and Σ_k .

Oracle2 = QDA rule using true Σ_k , but estimated $\hat{\mu}_k$:

- = sample means in Gaussian case
- = spatial median in t_2 case

Comments

- I do not want to bug you with more simulations...
- I just mention that, we can perform *much better* than estimators regularized sample covariance matrices (SCM-s) $\mathbf{S}_k(\beta)$ with shrinkage towards pooled SCM \mathbf{S} (as in Friedman's RDA) even when the clusters follow Gaussian distributions.

Why?

- We use more natural Riemannian geometry and our class of joint regularized estimators is **huge**:
 - ✓ many different g -convex penalty fnc's $d(\mathbf{A}, \mathbf{B})$: Kullback-Leibler , Ellipticity, Riemannian distance, ...
 - ✓ many different g -convex loss fnc's $\rho(t)$: Gaussian, Tyler's, Huber's,
 - ✓ robust: good performance under non-Gaussianity or outliers

Comments

- I do not want to bug you with more simulations...
- I just mention that, we can perform *much better* than estimators regularized sample covariance matrices (SCM-s) $\mathbf{S}_k(\beta)$ with shrinkage towards pooled SCM \mathbf{S} (as in Friedman's RDA) even when the clusters follow Gaussian distributions.

Why?

- We use more natural Riemannian geometry and our class of joint regularized estimators is **huge**:
 - ✓ many different g -convex penalty fnc's $d(\mathbf{A}, \mathbf{B})$: Kullback-Leibler , Ellipticity, Riemannian distance, ...
 - ✓ many different g -convex loss fnc's $\rho(t)$: Gaussian, Tyler's, Huber's,
 - ✓ robust: good performance under non-Gaussianity or outliers

A short recap on complex numbers/calculus

We consider *complex-valued* measurements. Our results are also valid for the real-valued case.

- A complex vector $\mathbf{a} = \mathbf{a}_R + j\mathbf{a}_I$ in complex Euclidean p -space \mathbb{C}^p , where $j = \sqrt{-1}$ is the imaginary unit.
- A complex conjugate: $\mathbf{a}^* = \mathbf{a}_R - j\mathbf{a}_I$
- (Hermitian) *norm*: $\|\mathbf{a}\| = \sqrt{\mathbf{a}^\top \mathbf{a}}$, where $(\cdot)^\top = [(\cdot)^*]^\top$ denotes the *Hermitian (conjugate) transpose*.
- Modulus of $a = a_R + ja_I \in \mathbb{C}$ is $|a|^2 = \sqrt{aa^*} = \sqrt{a_R^2 + a_I^2}$.
- A complex matrix $\mathbf{A} = \mathbf{A}_R + j\mathbf{A}_I$ is *Hermitian* if $\mathbf{A}^\top = \mathbf{A}$ which implies that \mathbf{A}_R is symmetric ($\mathbf{A}_R^\top = \mathbf{A}_R$) and \mathbf{A}_I is skew-symmetric ($\mathbf{A}_I^\top = -\mathbf{A}_I$).

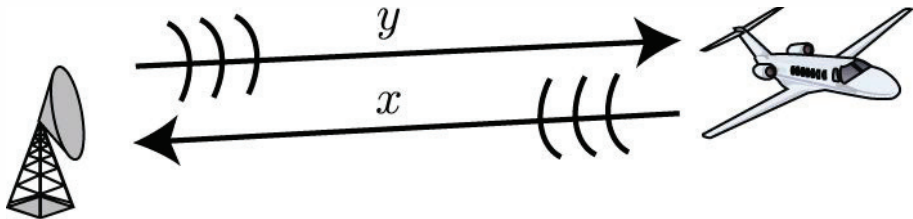
Matched filter detection

- **Matched filter (MF)** is a linear filter that correlates a known "signal" p with a received (measured) data x with the purpose of testing

\mathcal{H}_0 : signal-absent (noise only)

vs \mathcal{H}_1 : signal-present (signal + noise) .

- Optimal = maximizes the signal to noise ratio (SNR).
- Decide (absence/presence) by comparing the magnitude of the (normalized) matched filter output to a *threshold*.
- Threshold is set to attain the desired **probability of false alarm**).



- The output of a linear filter with **filter weights** c_1, \dots, c_p is

$$y = \mathbf{c}^\top \mathbf{x} = \sum_{i=1}^p c_i^* z_i$$

where \mathbf{x} denotes the measured p -variate data \mathbf{x} .

- The MF filter weights are

$$\mathbf{c} = \frac{\Sigma^{-1} \mathbf{p}}{(\mathbf{p}^H \Sigma^{-1} \mathbf{p})^{1/2}}$$

where \mathbf{p} denotes the (known) transmitted p -variate signal.

- In practise, the **noise covariance matrix** Σ is unknown and is estimated from a set of *secondary data* $\mathbf{x}_1, \dots, \mathbf{x}_n$ acquired under noise-only (\mathcal{H}_0) environment.
- **Problem:** $p > n$ or $n \approx p$.

Radar detection using NMF

$$\boxed{\mathcal{H}_0 : \mathbf{x} = \mathbf{c}} \quad \text{vs.} \quad \boxed{\mathcal{H}_1 : \mathbf{x} = \gamma \mathbf{p} + \mathbf{c}}$$

- \mathbf{x} : p -variate received data.
- \mathbf{p} : known complex **signal vector**
- \mathbf{c} : **noise** r.v., $\mathbf{c} \sim \mathcal{E}_p(\mathbf{0}, \Sigma, g)$. [complex elliptically symmetric random vector]
- $\gamma \in \mathbb{C}$: signal parameter

EX In radar, γ is accounting for both channel propagation effect and target backscattering; \mathbf{p} is the transmitted known radar pulse vector.

The normalized matched filter (NMF) detector

$$\Lambda = \frac{|\mathbf{c}^H \mathbf{x}|^2}{\|\Sigma^{-1/2} \mathbf{x}\|^2} = \frac{|\mathbf{p}^H \Sigma^{-1} \mathbf{x}|^2}{(\mathbf{x}^H \Sigma^{-1} \mathbf{x})(\mathbf{p}^H \Sigma^{-1} \mathbf{p})} \underset{\mathcal{H}_0}{\overset{\mathcal{H}_1}{\gtrless}} \lambda$$

- **Null distribution:** $\Lambda \sim_{\mathcal{H}_0} \text{Beta}(1, p - 1)$ i.e., distribution-free under the class of CES distributed noise.
- **Threshold:** To obtain a desired level P_{FA} , threshold λ can be set as $(1 - P_{\text{FA}})$ th quantile of $\text{Beta}(1, p - 1)$:

$$P_{\text{FA}} = \Pr(\Lambda > \lambda | \mathcal{H}_0) = (1 - \lambda)^{p-1}$$

$$\Rightarrow \lambda = 1 - P_{\text{FA}}^{1/(p-1)}.$$

- **Adaptive NMF detector** $\hat{\Lambda} =$ replace Σ by its estimate $\hat{\Sigma}$.
 - **Problems:** $p > n$ or $n \approx p$. Furthermore, noise (clutter) in radar applications often heavy-tailed non-Gaussian and hence accurate estimation of Σ is crucial and SCM will do a poor job.
- \Rightarrow Adaptive NMF detector often does not retain the constant false alarm rate (CFAR) of NMF nor its probability of detection
- **Solution:** Adaptive NMF based on regularized M -estimator

- **Null distribution:** $\Lambda \sim_{\mathcal{H}_0} \text{Beta}(1, p - 1)$ i.e., distribution-free under the class of CES distributed noise.
- **Threshold:** To obtain a desired level P_{FA} , threshold λ can be set as $(1 - P_{\text{FA}})$ th quantile of $\text{Beta}(1, p - 1)$:

$$P_{\text{FA}} = \Pr(\Lambda > \lambda | \mathcal{H}_0) = (1 - \lambda)^{p-1}$$

$$\Rightarrow \lambda = 1 - P_{\text{FA}}^{1/(p-1)}.$$

- **Adaptive NMF detector** $\hat{\Lambda} =$ replace Σ by its estimate $\hat{\Sigma}$.
 - **Problems:** $p > n$ or $n \approx p$. Furthermore, noise (clutter) in radar applications often heavy-tailed non-Gaussian and hence accurate estimation of Σ is crucial and SCM will do a poor job.
- \Rightarrow Adaptive NMF detector often does not retain the constant false alarm rate (CFAR) of NMF nor its probability of detection
- **Solution:** Adaptive NMF based on regularized M -estimator

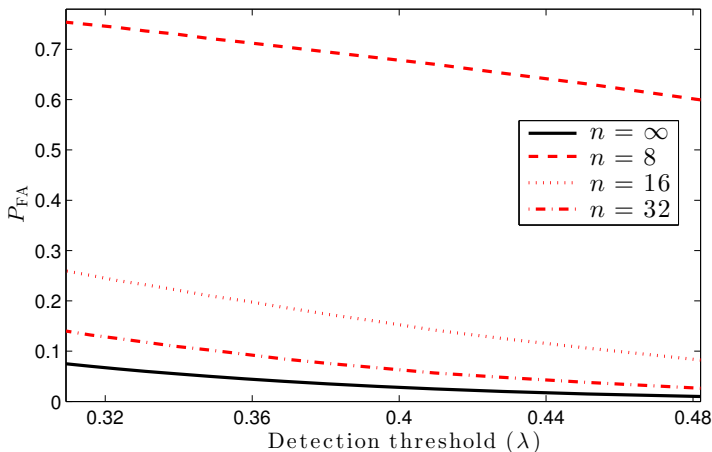
Simulation set-up

Q: How well does the adaptive NMF detector $\hat{\Lambda}$ based on estimated Σ maintain the preset PFA?

Simulation set-up: \mathcal{H}_0 holds, and $\mathbf{x} = \mathbf{c} \sim \mathbb{C}K_{p,\nu}(\mathbf{0}, \Sigma)$

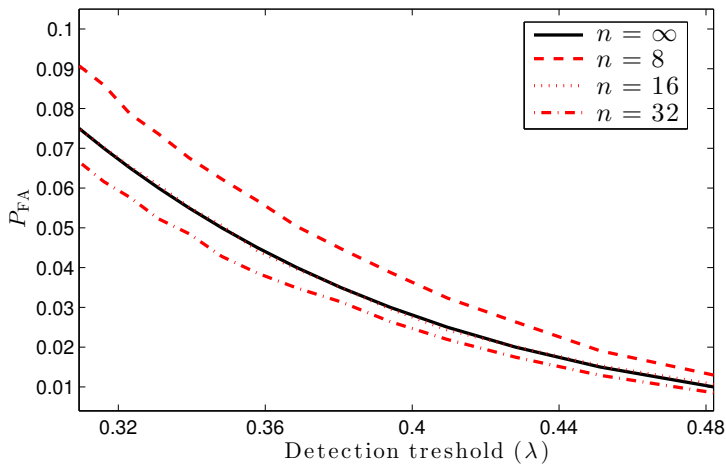
- received data \mathbf{x} (input for the NMF detector)
- secondary data $\mathbf{x}_1, \dots, \mathbf{x}_n$ (input to estimate $\hat{\Sigma}$).
- $p = 8$ and the shape parameter of the K -distribution is $\nu = 4.5$.

For 10000 trials, we calculated the empirical P_{FA} for fixed λ where the true covariance matrix Σ was generated *randomly for each trial*.



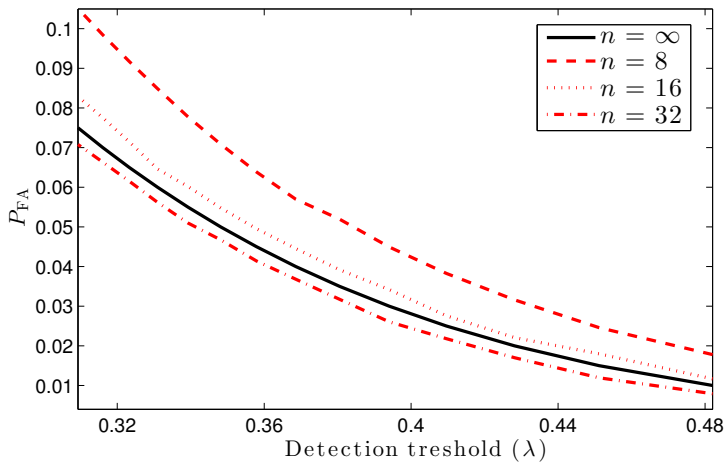
Adaptive NMF detector $\hat{\Lambda}$ using non-regularized Tyler's M -estimator as an estimate for Σ .

$n := \#$ of secondary samples $\{\mathbf{x}_i\}_{i=1}^n$ available for estimation



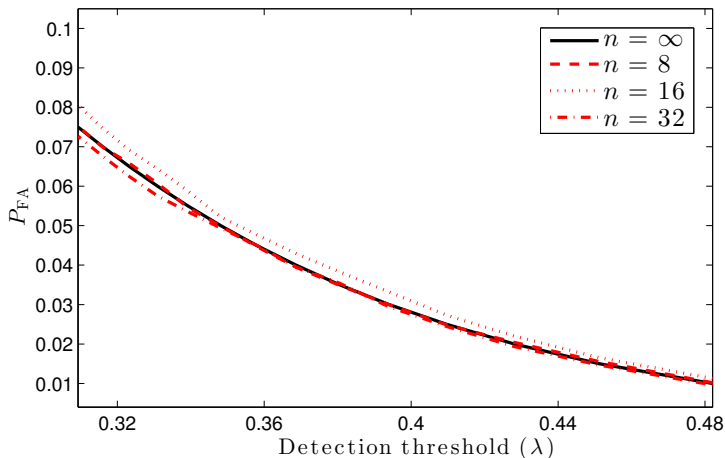
Adaptive NMF detector $\hat{\Lambda}$ using regularized SCM $S_{\alpha,\beta}$ and plug-in oracle estimator $\hat{\alpha}_o$ and $\hat{\beta}_0$ given in Slide 91.

$n := \#$ of secondary samples $\{\mathbf{x}_i\}_{i=1}^n$ available for estimation



Adaptive NMF detector $\hat{\Lambda}$ using DL-FP estimator and plug-in oracle estimator $\hat{\alpha}_0$ given in Slide 94

$n := \#$ of secondary samples $\{\mathbf{x}_i\}_{i=1}^n$ available for estimation



Adaptive NMF detector $\hat{\Lambda}$ using regularized Tyler's M -estimator and plug-in oracle estimator $\hat{\alpha}_o$ (and $\hat{\beta}_o = 1 - \hat{\alpha}_o$) given in Slide 92.

$n := \#$ of secondary samples $\{\mathbf{x}_i\}_{i=1}^n$ available for estimation

Thank you !!!

References: see the next slides



Abramovich, Y. and Besson, O. (2013).

Regularized covariance matrix estimation in complex elliptically symmetric distributions using the expected likelihood approach-part 1: The over-sampled case.

IEEE Trans. Signal Process., 61(23):5807–5818.



Abramovich, Y. I. and Spencer, N. K. (2007).

Diagonally loaded normalised sample matrix inversion (LNSMI) for outlier-resistant adaptive filtering.

In *Proc. Int'l Conf. on Acoustics, Speech and Signal Processing (ICASSP'07)*, pages 1105–1108.



Besson, O. and Abramovich, Y. (2013).

Regularized covariance matrix estimation in complex elliptically symmetric distributions using the expected likelihood approach-part 2: The under-sampled case.

IEEE Trans. Signal Process., 61(23):5819–5829.



Bhatia, R. (2009).

Positive definite matrices.

Princeton University Press.



Chen, Y., Wiesel, A., and Hero, A. O. (2011).

Robust shrinkage estimation of high-dimensional covariance matrices.

IEEE Trans. Signal Process., 59(9):4097 – 4107.



Conte, E., De Maio, A., and Ricci, G. (2002).

Recursive estimation of the covariance matrix of a compound-Gaussian process and its application to adaptive CFAR detection.

IEEE Trans. Signal Process., 50(8):1908 – 1915.



Cover, T. M. and Thomas, J. A. (2012).

Elements of information theory.

John Wiley & Sons.



Du, L., Li, J., and Stoica, P. (2010).

Fully automatic computation of diagonal loading levels for robust adaptive beamforming.

IEEE Trans. Aerosp. Electron. Syst., 46(1):449–458.



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

Sparse inverse covariance estimation with the graphical lasso.

Biostatistics, 9(3):432–441.



Friedman, J. H. (1989).

Regularized discriminant analysis.

J. Amer. Stat. Assoc., 84(405):165–175.



Gini, F. and Greco, M. (2002).

Covariance matrix estimation for CFAR detection in correlated heavy-tailed clutter.
Signal Processing, 82(12):1847–1859.



Kent, J. T. (1997).

Data analysis for shapes and images.
J. Statist. Plann. Inference, 57(2):181–193.



Ledoit, O. and Wolf, M. (2004).

A well-conditioned estimator for large-dimensional covariance matrices.
J. Mult. Anal., 88:365–411.



Maronna, R. A. (1976).

Robust M-estimators of multivariate location and scatter.
Ann. Stat., 5(1):51–67.



Ollila, E., Soloveychik, I., Tyler, D. E., and Wiesel, A. (2016).

Simultaneous penalized M-estimation of covariance matrices using geodesically convex optimization.
Journal of Multivariate Analysis, submitted.



Ollila, E. and Tyler, D. E. (2012).

Distribution-free detection under complex elliptically symmetric clutter distribution.

In *Proc. IEEE Sensor Array and Multichannel Signal Processing Workshop (SAM'12)*, Hoboken, NJ, USA.



Ollila, E. and Tyler, D. E. (2014).

Regularized M -estimators of scatter matrix.

IEEE Trans. Signal Process., 62(22):6059–6070.



Pascal, F., Chitour, Y., and Quek, Y. (2014).

Generalized robust shrinkage estimator and its application to stap detection problem.

IEEE Trans. Signal Process., 62(21):5640–5651.



Sra, S. (2011).

Positive definite matrices and the S-divergence.

arXiv preprint arXiv:1110.1773.



Sun, Y., Babu, P., and Palomar, D. P. (2014).

Regularized tyler's scatter estimator: Existence, uniqueness, and algorithms.

IEEE Trans. Signal Process., 62(19):5143–5156.



Wiesel, A. (2012).

Unified framework to regularized covariance estimation in scaled gaussian models.

IEEE Trans. Signal Process., 60(1):29–38.



Wiesel, A. and Zhang, T. (2015).

Structured robust covariance estimation.

Foundations and Trends in Signal Processing, 8(3):127–216.



Zhang, T., Wiesel, A., and Greco, M. S. (2013).

Multivariate generalized Gaussian distribution: Convexity and graphical models.

IEEE Trans. Signal Process., 61(16):4141–4148.