Robust Covariance Matrix Estimators for Sparse Data Using Regularization and RMT

Esa Ollila¹ Frédéric Pascal²

¹ Aalto University, Department of Signal Processing and Acoustics, Finland esa.ollila@aalto.fi http://signal.hut.fi/~esollila/

²CentraleSupelec, Laboratory of Signals and Systems (L2S), France frederic.pascal@centralesupelec.fr http://fredericpascal.blogspot.fr

EUSIPCO'16, Sep 29, 2016



Contents

Part A

Regularized M-estimators of covariance:

- *M*-estimation and geodesic (*g*-)convexity
- Regularization via g-convex penalties
- Application: regularized discriminant analysis

Part B

Regularized M-estimators and RMT

- Robust estimation and RMT
- Regularized *M*-estimators
- Application(s): DoA estimation, target detection

Frederic

Fsa

Covariance estimation problem

- **x** : *p*-variate (centered) random vector
- $\mathbf{x}_1, \ldots, \mathbf{x}_n$ i.i.d. realizations of \mathbf{x}
- Problem: Find an estimate $\hat{\Sigma} = \hat{\Sigma}(\{\mathbf{x}_i\}_{i=1}^n)$ of the positive definite covariance matrix

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

Solution: Maximum likelihood, *M*-estimation.

Conventional estimate: the sample covariance matrix (SCM)

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

Why covariance estimation?



Covariance estimation challenges

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

Problem 1 & 2 = Sparse data \Rightarrow regularization and/or RMT Problem 3 \Rightarrow robust estimation

Why robustness?

- 1 Outliers difficult to glean from high-dimensional data sets
- 2 Impulsive measurement environments in sensing systems (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

Part A : Contents

- I. Ad-hoc shrinkage SCM-s of multiple samples
- II. ML- and $M\mbox{-estimators}$ of scatter matrix
- III. Geodesic convexity
 - Geodesic
 - g-convex functions

IV. Regularized M-estimators

- Shrinkage towards an identity matrix
- Shrinkage towards a target matrix
- Estimation of the regularization parameter
- V. Penalized estimation of multiple covariances
 - Pooling vs joint estimation
 - Regularized discriminant analysis

Acknowledgement

To my co-authors:



David E. Tyler Rutgers University







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

Multiple covariance estimation problem

• We are given K groups of elliptically distributed measurements,

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

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 \u03c6_k of populations is known or that data is *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 $\Sigma_1 = \ldots = \Sigma_K$ 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 S towards \propto 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 then 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:

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

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

Penalized objective function: Add a penalty term and solve

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

where

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

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

Q1: Some formalism to the Ad-Hoc method

The information theoretic Kullback-Leibler (KL) divergence [Cover and Thomas, 2012], distance from N_p(0, A) to N_p(0, B) is

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

As is well known, it verifies $d_{\mathrm{KL}}(\mathbf{A}, \mathbf{B}) \geq 0$ and = 0 for $\mathbf{A} = \mathbf{B}$. Using $d_{\mathrm{KL}}(\boldsymbol{\Sigma}_k, \hat{\boldsymbol{\Sigma}})$ as the penalty, the optimization problem $\mathcal{L}_{\mathrm{G},k}(\boldsymbol{\Sigma}_k) + \lambda d_{\mathrm{KL}}(\boldsymbol{\Sigma}_k, \hat{\boldsymbol{\Sigma}})$ possesses a unique solution given by

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

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

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

Q1: Some formalism to the Ad-Hoc method

The information theoretic Kullback-Leibler (KL) divergence [Cover and Thomas, 2012], distance from N_p(0, A) to N_p(0, B) is

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

As is well known, it verifies $d_{\mathrm{KL}}(\mathbf{A}, \mathbf{B}) \geq 0$ and = 0 for $\mathbf{A} = \mathbf{B}$. Using $d_{\mathrm{KL}}(\boldsymbol{\Sigma}_k, \hat{\boldsymbol{\Sigma}})$ as the penalty, the optimization problem $\mathcal{L}_{\mathrm{G},k}(\boldsymbol{\Sigma}_k) + \lambda d_{\mathrm{KL}}(\boldsymbol{\Sigma}_k, \hat{\boldsymbol{\Sigma}})$ 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, ..., K.

This gives Friedman's Ad-Hoc shrinkage SCM estimators when the shrinkage target matrix $\hat{\Sigma}$ 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 $d_{\mathrm{KL}}(\Sigma_k, \hat{\Sigma})$.

Comments

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

Problems

- The penalized optimization program, L_{G,k}(Σ_k) + λ d_{KL}(Σ_k, Σ̂), does not seem to generalize to using other distance functions or other non-Gaussian cost functions.
- KL-distance $d_{\mathrm{KL}}(\Sigma_k, \Sigma)$ 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

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

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 Part A of the tutorial

- **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

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.

Elliptically symmetric (CES) distribution

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

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

Σ ∈ S(p), unknown positive definite p × p scatter matrix parameter.
 g : ℝ⁺₀ → ℝ⁺, fixed density generator.

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

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

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

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

The maximum likelihood estimator (MLE)

•
$$\{\mathbf{x}_i\} \stackrel{iid}{\sim} \mathcal{E}_p(\mathbf{0}, \mathbf{\Sigma}, g)$$
, where $n > p$.

 \blacksquare The MLE $\hat{\mathbf{\Sigma}} \in \mathcal{S}(p)$ minimizes the negative log-likelihood fnc

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

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

Critical points are solutions to estimating equations

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

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

MLE = "an adaptively weighted sample covariance matrix"

M-estimators of scatter matrix

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

[Maronna, 1976]

Among the first proposals for robust covariance matrix estimators

Generalizations of ML-estimators:

- u(t) = ρ'(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 ⇒ solution exists

Huber's *M*-estimator

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

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

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

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

$$\rho_{\rm H}(t;c) = \begin{cases} t/b & \text{for } t \leqslant c^2, \\ (c^2/b) \left(\log(t/c^2) + 1 \right) & \text{for } t > c^2. \end{cases}$$

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

Tyler's (1987) *M*-estimator

- Distribution-free M-estimator (under elliptical distributions)
- Defined as a solution to

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

⇒ so an *M*-estimator with Tyler's weight fnc $u(t) = \rho'(t) = p/t$ ■ Now it is also known that $\hat{\Sigma} \in S(p)$ minimizes the cost fnc

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

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 (1987) *M*-estimator

- Distribution-free *M*-estimator (under elliptical distributions)
- Defined as a solution to

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

⇒ so an *M*-estimator with Tyler's weight fnc $u(t) = \rho'(t) = p/t$ ■ Now it is also known that $\hat{\Sigma} \in S(p)$ minimizes the cost fnc

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

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 (1987) *M*-estimator

- Distribution-free *M*-estimator (under elliptical distributions)
- Defined as a solution to

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

⇒ so an *M*-estimator with Tyler's weight fnc $u(t) = \rho'(t) = p/t$ Now it is also known that $\hat{\Sigma} \in S(p)$ minimizes the cost fnc

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

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.

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

Comments:

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

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

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

Comments:

- **1** Limiting case of Huber's M-estimator when $c \rightarrow 0$
- 2 Minimum is a unique up to a postive scalar: if \$\hlowsymbol{\Sigma}\$ is a minimum then so is b\$\hlowsymbol{\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

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

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

Comments:

- **1** Limiting case of Huber's M-estimator when $c \rightarrow 0$
- 2 Minimum is a unique up to a postive scalar: if îs a minimum then so is bîs 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

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

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

Comments:

- **1** Limiting case of Huber's M-estimator when $c \rightarrow 0$
- 2 Minimum is a unique up to a postive scalar: if îs a minimum then so is bîs 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 $\mathcal{N}_p(0,\Sigma)$ can be obtained by scaling any minimum $\hat{\Sigma}$ by

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

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

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



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

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



From Euclidean convexity to Riemannian convexity



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

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



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 (goedesic) connecting them




Geodesic convexity in p = 1 variable

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

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

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

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

- 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 $\Sigma \in S(p)$, $p \neq 1$, the solution is **not** a simple change of variables.

Geodesic convexity in p = 1 variable

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

$$f\left(\underbrace{(1-t)x_0 + tx_1}_{\text{line}}\right) \le (1-t)f(x_0) + tf(x_1)$$
g-convex function in $\sigma^2 \in \mathbb{R}^+_0$: $\rho(\sigma^2) = f(\log \sigma^2)$, $\sigma^2 = e^x$

$$\rho\left(\underbrace{(\sigma^2_0)^{(1-t)}(\sigma^2_1)^t}_{\text{geodesic}}\right) \le (1-t)\rho(\sigma^2_0) + t\rho(\sigma^2_1)$$

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 $\Sigma \in S(p)$, $p \neq 1$, the solution is **not** a simple change of variables.

Geodesic convexity in p = 1 variable

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

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

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

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

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

$$\boldsymbol{\Sigma}_t = \boldsymbol{\Sigma}_0^{1/2} \left(\boldsymbol{\Sigma}_0^{-1/2} \boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_0^{-1/2} \right)^t \boldsymbol{\Sigma}_0^{1/2} \text{ for } t \in [0, 1].$$

where $\Sigma_t \in \mathcal{S}(p)$ for $0 \le t \le 1 \Rightarrow \mathcal{S}(p)$ 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

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

Riemannian manifold

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

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

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**
 - \blacksquare local inner product $\langle A,B\rangle_{\Sigma}$ on the tangent space of symmetric matrices

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

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

Geodesically (g-)convex function

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

$$h(\Sigma_t) \le (1-t) \ h(\Sigma_0) + t \ h(\Sigma_1)$$
 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.

g-convex + g-convex = g-convex

Useful results on g-convexity: my personal top 3

$$\boldsymbol{\Sigma}_t = \boldsymbol{\Sigma}_0^{1/2} \left(\boldsymbol{\Sigma}_0^{-1/2} \boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_0^{-1/2} \right)^t \boldsymbol{\Sigma}_0^{1/2}$$

1. Joint diagonalization formulation

The geodesic path can be written equivalently as

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

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

- E is a nonsingular square matrix: row vectors of E⁻¹ are the eigenvectors of Σ₀⁻¹Σ₁
- **D** is a diagonal matrix: diagonal elements are the eigenvalues of

$$\mathbf{\Sigma}_0^{-1}\mathbf{\Sigma}_1$$
 or $\mathbf{\Sigma}_0^{-1/2}\mathbf{\Sigma}_1\mathbf{\Sigma}_0^{-1/2}.$

Useful results on g-convexity: my personal top 3

$$\boldsymbol{\Sigma}_t = \boldsymbol{\Sigma}_0^{1/2} \left(\boldsymbol{\Sigma}_0^{-1/2} \boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_0^{-1/2} \right)^t \boldsymbol{\Sigma}_0^{1/2}$$

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 on f on a g-convex set $\mathcal M$ is g-convex if

$$f(\boldsymbol{\Sigma}_{1/2}) \leq \frac{1}{2} \{f(\boldsymbol{\Sigma}_0) + f(\boldsymbol{\Sigma}_1)\}$$

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

For more results, see [Wiesel and Zhang, 2015]

Geodesic convexity

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
$$\mathbf{a}^{\top} \mathbf{\Sigma}^{\pm 1} \mathbf{a}$$
 is strictly *g*-convex ($\mathbf{a} \neq 0$).
4 $\log |\sum_{i=1}^{n} \mathbf{H}_{i} \mathbf{\Sigma}^{\pm 1} \mathbf{H}_{i}|$ 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.

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:

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

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



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

 $\begin{array}{l} \mbox{Contours of } \mathcal{L}_T(\boldsymbol{\Sigma}) \\ \mbox{and the solution } \hat{\boldsymbol{\Sigma}}. \end{array}$

Geodesic convexity

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:

$$\mathcal{M}(2) = \{ \mathbf{\Sigma} \in \mathcal{S}(2) : \det(\mathbf{\Sigma}) = 1 \}$$

= $\left\{ \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} : \sigma_2^2 > 0, \sigma_{12} \in \mathbb{R}, \sigma_1^2 = \frac{1 + \sigma_{12}^2}{\sigma_2^2} \right\}$

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



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:

$$\mathcal{M}(2) = \{ \mathbf{\Sigma} \in \mathcal{S}(2) : \det(\mathbf{\Sigma}) = 1 \}$$

= $\left\{ \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} : \sigma_2^2 > 0, \sigma_{12} \in \mathbb{R}, \sigma_1^2 = \frac{1 + \sigma_{12}^2}{\sigma_2^2} \right\}$

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



Geodesic convexity

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:

$$\mathcal{M}(2) = \{ \mathbf{\Sigma} \in \mathcal{S}(2) : \det(\mathbf{\Sigma}) = 1 \}$$

= $\left\{ \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} : \sigma_2^2 > 0, \sigma_{12} \in \mathbb{R}, \sigma_1^2 = \frac{1 + \sigma_{12}^2}{\sigma_2^2} \right\}$

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



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

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

Geodesic convexity

Examples of *g***-convex sets**

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

$$\boldsymbol{\Sigma}_t = \boldsymbol{\Sigma}_0^{1/2} \left(\boldsymbol{\Sigma}_0^{-1/2} \boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_0^{-1/2}
ight)^t \boldsymbol{\Sigma}_0^{1/2} \text{ for } t \in [0,1].$$

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:

$$oldsymbol{\Sigma} = egin{pmatrix} oldsymbol{\Sigma}_1 & oldsymbol{\Sigma}_2 \ -oldsymbol{\Sigma}_2 & oldsymbol{\Sigma}_1 \end{pmatrix}$$

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
- Estimation of the regularization parameter

V. Penalized estimation of multiple covariances

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

Regularized *M*-estimators of scatter matrix: shrinkage towards identity

Penalized cost function:

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

where $\alpha \ge 0$ is a fixed regularization parameter. Q: Existence, Uniqueness, computation?

Our penalty function pulls Σ away from singularity

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

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 $-\infty < x < \infty$)

Note: Tyler's, Huber's, Gaussian loss fnc $\rho(t)$ satisfies Cond. 1.

Regularized M-estimators

Main results

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

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 $\mathcal{L}_{\alpha}(\Sigma)$ allways 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:

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

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!
 - \rightarrow 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

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

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

Tuning the $\rho(t)$ function

Result 1 is general and allows us to tune the ρ(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

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

The solution verifies

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

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

Tuning the $\rho(t)$ function

Result 1 is general and allows us to tune the ρ(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

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

The solution verifies

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

• 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

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

where ${\bf S}$ denotes the SCM.

Due to $\mathbf{Pesult 1}$, its unique minimizer $\hat{\Sigma}$ is

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

which corresponds to [Ledoit and Wolf, 2004] shrinkage estimator.

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

A class of regularized Tyler's *M*-estimators

Let use *tuned* Tyler's cost fnc ρ(t) = pβ log t for fixed 0 < β < 1.
The penalized Tyler's cost fnc is

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

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

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

• We commonly use $\alpha = 1 - \beta$.

Target $\mathcal{L}_{\alpha,\beta}(\Sigma)$ is g-convex in Σ , but ρ is not bounded below

 $\Rightarrow \mathbf{P}_{\text{Result 1(b)}}$, 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 ρ(t) = pβ log t for fixed 0 < β < 1.
The penalized Tyler's cost fnc is

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

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

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

• We commonly use $\alpha = 1 - \beta$.

Target $\mathcal{L}_{\alpha,\beta}(\Sigma)$ is g-convex in Σ , but ρ is not bounded below

 $\Rightarrow \bigcirc \text{Result 1(b)}$, for existence does not hold.

 Conditions for existence needs to be considered separately for Tyler's M-estimator;

Regularized *M*-estimators

A class of regularized Tyler's *M*-estimators

Let use *tuned* Tyler's cost fnc ρ(t) = pβ log t for fixed 0 < β < 1.
The penalized Tyler's cost fnc is

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

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

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

- We commonly use $\alpha = 1 \beta$.
- Target $\mathcal{L}_{\alpha,\beta}(\Sigma)$ is *g*-convex in Σ , but ρ is not bounded below ⇒ • 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 β < 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 $\alpha > 0$, $0 \le \beta < 1$. 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

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

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

Regularized M-estimators

- 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

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

 \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.lo.g.

$$\alpha = 1 - \beta$$
 [or equivalently $\beta = 1 - \alpha$].

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

Related approach for regularizing Tyler's *M*-estimator

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

$$\begin{split} \mathbf{\Sigma}_{k+1} &\leftarrow (1-\alpha) \frac{p}{n} \sum_{i=1}^{n} \frac{\mathbf{x}_{i} \mathbf{x}_{i}^{\top}}{\mathbf{x}_{i}^{\top} \mathbf{V}_{k}^{-1} \mathbf{x}_{i}} + \alpha \mathbf{I} \\ \mathbf{V}_{k+1} &\leftarrow p \mathbf{\Sigma}_{k+1} / \text{Tr}(\mathbf{\Sigma}_{k+1}), \end{split}$$

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 V_k and not necessarily in Σ_k.]
- Note 1: essentially a diagonally loaded version of the fixed-point (FP) algorithm for Tyler's *M*-estimator. Hence we call th estimator as DL-FP estimator.
- Note 2: DL-FP was not shown to be a solution to any penalized form of Tyler's cost function.

Regularized M-estimators

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

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

or equivalently,

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

where

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

Regularized M-estimators

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 $d(\mathbf{A}, \mathbf{B}) = d(\mathbf{CAC}^{\top}, \mathbf{CBC}^{\top})$, \forall nonsingular **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 $d(\mathbf{A}, \mathbf{B}) = d(\mathbf{CAC}^{\top}, \mathbf{CBC}^{\top})$, \forall nonsingular **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.

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

$$\boldsymbol{\Sigma}(\boldsymbol{\pi}) = \operatorname*{arg\,min}_{\boldsymbol{\Sigma} \in \mathcal{S}(p)} \sum_{i=1}^{n} \pi_k \, d(\boldsymbol{\Sigma}_k, \boldsymbol{\Sigma}),$$

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 σ²₁,..., σ²_k > 0, we could consider

geometric mean $\sigma^{\perp} = (\sigma_1^{\perp} \cdots \sigma_K^{\perp})^{\perp}$

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

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

$$\boldsymbol{\Sigma}(\boldsymbol{\pi}) = \operatorname*{arg\,min}_{\boldsymbol{\Sigma} \in \mathcal{S}(p)} \sum_{i=1}^{K} \pi_k \, d(\boldsymbol{\Sigma}_k, \boldsymbol{\Sigma}),$$

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 σ₁²,...,σ_K² > 0, we could consider arithmetic mean σ² = ¹/_K Σ_{k=1}^K σ_k². geometric mean σ² = (σ₁² ··· σ_K²)^{1/K} harmonic mean σ² = ((σ₁² ··· σ_K²))^{1/K}
 Note: For a pair σ₀², σ₁², the geometric mean is the midpole moderic σ² = (σ₁²)¹⁻¹(σ₂²)¹

Regularized M-estimators

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
K

$$\boldsymbol{\Sigma}(\boldsymbol{\pi}) = \operatorname*{arg\,min}_{\boldsymbol{\Sigma} \in \mathcal{S}(p)} \sum_{i=1}^{K} \pi_k \, d(\boldsymbol{\Sigma}_k, \boldsymbol{\Sigma}),$$

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 σ₁²,..., σ_K² > 0, we could consider arithmetic mean σ² = ¹/_K Σ_{k=1}^K σ_k². geometric mean σ² = (σ₁² ··· σ_K²)^{1/K}

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

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
K

$$\boldsymbol{\Sigma}(\boldsymbol{\pi}) = \operatorname*{arg\,min}_{\boldsymbol{\Sigma} \in \mathcal{S}(p)} \sum_{i=1}^{K} \pi_k \, d(\boldsymbol{\Sigma}_k, \boldsymbol{\Sigma}),$$

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 σ₁²,..., σ_K² > 0, we could consider arithmetic mean σ² = 1/K Σ_{k=1}^K σ_k². geometric mean σ² = (σ₁² ··· σ_K²)^{1/K} harmonic mean σ² = (1/K Σ_{k=1}^K (σ_k²)⁻¹)⁻¹

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

Regularized M-estimators

• 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
K

$$\boldsymbol{\Sigma}(\boldsymbol{\pi}) = \operatorname*{arg\,min}_{\boldsymbol{\Sigma} \in \mathcal{S}(p)} \sum_{i=1}^{K} \pi_k \, d(\boldsymbol{\Sigma}_k, \boldsymbol{\Sigma}),$$

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 σ₁²,..., σ_K² > 0, we could consider arithmetic mean σ² = ¹/_K Σ_{k=1}^K σ_k². geometric mean σ² = (σ₁² ··· σ_K²)^{1/K} harmonic mean σ² = (¹/_K Σ_{k=1}^K (σ_k²)⁻¹)⁻¹
 Note: For a pair σ² σ² the geometric mean is the midpoir

• 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

$$d_{\mathrm{F}}(\boldsymbol{\Sigma}_k, \boldsymbol{\Sigma}) = \left\{ \mathrm{Tr}[(\boldsymbol{\Sigma}_k - \boldsymbol{\Sigma})^2] \right\}^{1/2}$$

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

- \checkmark Riemannian distance $d_{
 m R}({f A},{f B})$
- \checkmark Kullback-Leibler (KL) divergence $d_{\mathrm{KL}}(\mathbf{A},\mathbf{B})$
- ✓ Ellipticity distance $d_{\rm E}({\bf A}, {\bf B})$

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

$$d_{\mathrm{F}}(\boldsymbol{\Sigma}_k, \boldsymbol{\Sigma}) = \left\{ \mathrm{Tr}[(\boldsymbol{\Sigma}_k - \boldsymbol{\Sigma})^2] \right\}^{1/2}$$

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

- \checkmark Riemannian distance $d_{
 m R}({f A},{f B})$
- \checkmark Kullback-Leibler (KL) divergence $d_{ ext{KL}}(\mathbf{A},\mathbf{B})$
- ✓ Ellipticity distance $d_{\rm E}({\bf A}, {\bf B})$

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?

X Frobenius distance

$$d_{\mathrm{F}}(\boldsymbol{\Sigma}_k, \boldsymbol{\Sigma}) = \left\{ \mathrm{Tr}[(\boldsymbol{\Sigma}_k - \boldsymbol{\Sigma})^2] \right\}^{1/2}$$

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

- $\checkmark\,$ Riemannian distance $\mathit{d}_{\mathrm{R}}(\mathbf{A},\mathbf{B})$
- ✓ Kullback-Leibler (KL) divergence $d_{\rm KL}({\bf A}, {\bf B})$
- ✓ Ellipticity distance $d_{\rm E}({\bf A}, {\bf B})$

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

$$d_{\rm R}(\mathbf{A}, \mathbf{B}) = \|\log(\mathbf{A}^{-1/2}\mathbf{B}\mathbf{A}^{-1/2})\|_{\rm F}^2,$$

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

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

$$\sum_{k=1}^{K} \pi_k \log(\boldsymbol{\Sigma}_{\mathrm{R}}^{1/2} \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\Sigma}_{\mathrm{R}}^{1/2}) = \boldsymbol{0}$$

Solution: Solution: A number of complex numerical approaches have been proposed in the literature.

Kullback-Leibler (KL) divergence

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

■ KL-distance verifies d_{KL}(A, B) ≥ 0 and = 0 for A = 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:

$$egin{aligned} \mathbf{\Sigma}_{\mathrm{I}}(oldsymbol{\pi}) &= rgmin_{oldsymbol{\Sigma}\in\mathcal{S}(p)}\sum_{i=1}^{K}\pi_k\,d_{\mathrm{KL}}(oldsymbol{\Sigma}_k,oldsymbol{\Sigma}) \ &= \left(\sum_{k=1}^{K}\pi_koldsymbol{\Sigma}_k^{-1}
ight)^{-1}, \end{aligned}$$

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

Special case: target matrix T = I

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

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

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

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

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

Special case: target matrix T = I (cont'd)

Note that

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

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

Ellipticity distance

$$d_{\mathrm{E}}(\mathbf{A}, \mathbf{B}) = p \log \frac{1}{p} \mathrm{Tr}(\mathbf{A}^{-1}\mathbf{B}) - \log |\mathbf{A}^{-1}\mathbf{B}|$$

• $d_{\rm E}$ 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]

 $d_{\rm E}({\bf A},{\bf B})$ is jointly *g*-convex and affine and scale invariant. The induced mean is unique (up to a scale) and solves

$$\mathbf{\Sigma}_{\mathrm{E}} = \left(\sum_{k=1}^{K} \pi_k \frac{p \mathbf{\Sigma}_k^{-1}}{\mathrm{Tr}(\mathbf{\Sigma}_k^{-1} \mathbf{\Sigma}_{\mathrm{E}})}
ight)^{-1},$$

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

53 / 75

Critical points

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

• 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

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

$$\Leftrightarrow \quad \beta \boldsymbol{\Sigma} = \beta \frac{1}{n} \sum_{i=1}^{n} u(\mathbf{x}_{i}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{x}_{i}) \mathbf{x}_{i} \mathbf{x}_{i}^{\top} + (1 - \beta) \mathcal{P}_{0}'(\boldsymbol{\Sigma})$$

$$\Leftrightarrow \quad \boldsymbol{\Sigma} = \beta \frac{1}{n} \sum_{i=1}^{n} u(\mathbf{x}_{i}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{x}_{i}) \mathbf{x}_{i} \mathbf{x}_{i}^{\top} + (1 - \beta) \{\mathcal{P}_{0}'(\boldsymbol{\Sigma}) + \boldsymbol{\Sigma}\}.$$

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

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

Estimation of the regularization parameter

Let us consider the regularized *M*-estimator [Ollila and Tyler, 2014] with shrinkage towards an identity matrix:

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

For simplicity, tune only one parameter and set:

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

Approaches:

- 1 Cross-validation
- 2 Oracle/Clairvoyant approach
- Expected likelihood approach
 [Abramovich and Besson, 2013, Besson and Abramovich, 2013]
- 4 Random matrix theory (more in Frederic's talk after the break).
- Approaches 2 and 3 are especially useful for Tyler's *M*-estimator.

Cross-validation (CV)

- 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 choises: Q = 5, 10 or Q = n (*leave-one-out CV*).
- Taking qth fold out (all \mathbf{x}_i , $i \in I_q$) gives a reduced data set \mathbf{X}_{-q} .

CV procedure (assuming $\alpha = 1 - \beta$) proposed in [Ollila et al., 2016]:

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

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

end

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

Oracle/Clairvoyant approach

 \blacksquare Given the true scatter matrix Σ_0 , define

$$\boldsymbol{\Sigma}_{\alpha} = (1 - \alpha) \frac{1}{n} \sum_{i=1}^{n} u(\mathbf{x}_{i}^{\top} \boldsymbol{\Sigma}_{0}^{-1} \mathbf{x}_{i}) \mathbf{x}_{i} \mathbf{x}_{i}^{\top} + \alpha \mathbf{I}$$

• Choose the oracle α_0 as the value that minimize the expected loss, say

$$\alpha_o = \alpha_o(\mathbf{\Sigma}_0) = \arg\min_{\alpha} \mathbb{E}[d(\mathbf{\Sigma}_{\alpha}, \mathbf{\Sigma}_0)]$$

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

Replace the unknown true Σ_0 in α_0 with some preliminary estimate or guess $\hat{\Sigma}_0$

 $\Rightarrow \hat{\alpha}_o = \alpha_o(\hat{\Sigma}_0)$ is the oracle/clairvoyant estimate

Oracle approach for regularized Tyler's *M*-estimator

$$\boldsymbol{\Sigma}_{\alpha} = (1-\alpha) \frac{p}{n} \sum_{i=1}^{n} \frac{\mathbf{x}_{i} \mathbf{x}_{i}^{\top}}{\mathbf{x}_{i}^{\top} \boldsymbol{\Sigma}_{0}^{-1} \mathbf{x}_{i}} + \alpha \mathbf{I}.$$

[Ollila and Tyler, 2014]

- idea: Given a shape matrix Σ₀, verifying Tr(Σ₀⁻¹) = p, choose α so that Σ₀⁻¹Σ_α is as close as possible to cI, for some c > 0.
- A natural distance that measures similarity in shape:

$$d(\boldsymbol{\Sigma}_0, \boldsymbol{\Sigma}_\alpha) = \|\boldsymbol{\Sigma}_0^{-1}\boldsymbol{\Sigma}_\alpha - \frac{1}{p}\mathrm{Tr}(\boldsymbol{\Sigma}_0^{-1}\boldsymbol{\Sigma}_\alpha)\mathbf{I}\|^2$$

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

$$\alpha_o = \frac{p - 2 + p \operatorname{Tr}(\mathbf{\Sigma}_0)}{p - 2 + p \operatorname{Tr}(\mathbf{\Sigma}_0) + n(p + 2) \{ p^{-1} \operatorname{Tr}(\mathbf{\Sigma}_0^{-2}) - 1 \}}$$

Estimate $\hat{\alpha}_o = \alpha_o(\hat{\mathbf{\Sigma}}_0)$ is obtained by using $\hat{\mathbf{\Sigma}}_0$ that is

- Tyler's *M*-estimator normalized s.t. $\operatorname{Tr}(\hat{\Sigma}_0^{-1}) = p$ when $n \ge p$
- \bullet regularized Tyler's estimator using $\beta < n/p$ & $\alpha = 1-\beta$ when n < p

Oracle approach for DL-FP estimator

$$\boldsymbol{\Sigma}_{\alpha} = (1-\alpha) \frac{p}{n} \sum_{i=1}^{n} \frac{\mathbf{x}_{i} \mathbf{x}_{i}^{\top}}{\mathbf{x}_{i}^{\top} \boldsymbol{\Sigma}_{0}^{-1} \mathbf{x}_{i}} + \alpha \mathbf{I}.$$

[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 $\operatorname{Tr}(\Sigma_0) = p$, find α as

$$\alpha_o = \arg\min_{\alpha} \mathbb{E}[\|\boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_{\alpha}\|^2]$$

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

$$\alpha_o = \frac{p^3 + (p-2)\operatorname{Tr}(\boldsymbol{\Sigma}_0^2)}{\{p^3 + (p-2)\operatorname{Tr}(\boldsymbol{\Sigma}_0^2)\} + n(p+2)(\operatorname{Tr}(\boldsymbol{\Sigma}_0^2) - p)}$$

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

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

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
 - Pooling vs joint estimation
 - Regularized discriminant analysis

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,

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

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 \u03c6_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

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

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

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

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

where

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

Penalized estimation of multiple covariances

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

$$\min_{\{\boldsymbol{\Sigma}_k\}_{k=1}^K, \boldsymbol{\Sigma} \in \mathcal{S}(p)} \quad \sum_{k=1}^K \pi_k \left\{ \beta \mathcal{L}_k(\boldsymbol{\Sigma}_k) + (1-\beta) \, d(\boldsymbol{\Sigma}_k, \boldsymbol{\Sigma}) \right\}$$

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

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

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 minumum $\hat{\Sigma}$ is found by solving

$$\hat{\boldsymbol{\Sigma}}(\boldsymbol{\pi}) = \operatorname*{arg\,min}_{\boldsymbol{\Sigma}\in\mathcal{S}(p)} \sum_{i=1}^{K} \pi_k \, d(\boldsymbol{\Sigma}_k, \boldsymbol{\Sigma}),$$

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 $\beta_k, k = 1, \dots, K$ 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 P(Σ) = Tr(Σ⁻¹) and compute pooled center Σ̂ as a pooled regularized M-estimator:

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

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(A, B) is the KL-distance or ellipticity distance.
- Both distances are *affine invariant*, i.e.

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

which is Property D4 in Slide

If D4 holds, the resulting estimators are affine equivariant:

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

Solving

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

yields estimating equations

$$\boldsymbol{\Sigma}_{k} = \beta \frac{1}{n_{k}} \sum_{i=1}^{n_{k}} u_{k} (\mathbf{x}_{ki}^{\top} \boldsymbol{\Sigma}_{k}^{-1} \mathbf{x}_{ki}) \mathbf{x}_{ki} \mathbf{x}_{ki}^{\top} + (1 - \beta) \boldsymbol{\Sigma}$$
$$\boldsymbol{\Sigma} = \left(\sum_{k=1}^{K} \pi_{k} \boldsymbol{\Sigma}_{k}^{-1} \right)^{-1}$$

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

Penalized estimation of multiple covariances

Solving

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

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

$$\boldsymbol{\Sigma}_{k} \leftarrow \beta \frac{1}{n_{k}} \sum_{i=1}^{n_{k}} u_{k} (\mathbf{x}_{ki}^{\top} \boldsymbol{\Sigma}_{k}^{-1} \mathbf{x}_{ki}) \mathbf{x}_{ki} \mathbf{x}_{ki}^{\top} + (1 - \beta) \boldsymbol{\Sigma}$$
$$\boldsymbol{\Sigma} \leftarrow \left(\sum_{k=1}^{K} \pi_{k} \boldsymbol{\Sigma}_{k}^{-1} \right)^{-1}$$

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

Penalized estimation of multiple covariances

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

$$\boldsymbol{\Sigma}_{k} = \beta \frac{1}{n_{k}} \sum_{i=1}^{n_{k}} u_{k} (\mathbf{x}_{ki}^{\top} \boldsymbol{\Sigma}_{k}^{-1} \mathbf{x}_{ki}) \mathbf{x}_{ki} \mathbf{x}_{ki}^{\top} + (1 - \beta) \frac{p \boldsymbol{\Sigma}}{\operatorname{Tr}(\boldsymbol{\Sigma}_{k}^{-1} \boldsymbol{\Sigma})},$$
$$\boldsymbol{\Sigma} = \left(\sum_{k=1}^{K} \pi_{k} \frac{p \boldsymbol{\Sigma}_{k}^{-1}}{\operatorname{Tr}(\boldsymbol{\Sigma}_{k}^{-1} \boldsymbol{\Sigma})} \right)^{-1}$$

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

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 $\mathbf{\Sigma}_1, \dots \mathbf{\Sigma}_K$ to $\mathbf{\Sigma}$

$$\boldsymbol{\Sigma}_{k} \leftarrow \beta \frac{1}{n_{k}} \sum_{i=1}^{n_{k}} u_{k} (\mathbf{x}_{ki}^{\top} \boldsymbol{\Sigma}_{k}^{-1} \mathbf{x}_{ki}) \mathbf{x}_{ki} \mathbf{x}_{ki}^{\top} + (1-\beta) \frac{p \boldsymbol{\Sigma}}{\operatorname{Tr}(\boldsymbol{\Sigma}_{k}^{-1} \boldsymbol{\Sigma})},$$
$$\boldsymbol{\Sigma} \leftarrow \left(\sum_{k=1}^{K} \pi_{k} \frac{p \boldsymbol{\Sigma}_{k}^{-1}}{\operatorname{Tr}(\boldsymbol{\Sigma}_{k}^{-1} \boldsymbol{\Sigma})} \right)^{-1}$$

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

Quadratic discriminant analysis (QDA)

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

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

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 kth population (k = 1, ..., K).

Assumptions:

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

Linear discriminant analysis (LDA)

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

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

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 $\Sigma_i = \Sigma_j \ i \neq j$

Regularized Discriminant Analysis (RDA)

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

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

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

Interpretation:

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

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

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

Simulation set-up

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

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

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

Gaussian cas	e: 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

 $\begin{array}{l} \textbf{Oracle1} = \mathsf{QDA} \ \mathsf{rule} \\ \mathsf{using true} \ \boldsymbol{\mu}_k \ \mathsf{and} \ \boldsymbol{\Sigma}_k. \end{array}$

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

- sample means in Gaussian case
- spatial median
 in t₂ case

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

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

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)

Gaussian case: test misclassification errors %

standard deviations inside parantheses in subscript

 $\begin{array}{l} \textbf{Oracle1} = \mathsf{QDA} \ \mathsf{rule} \\ \mathsf{using true} \ \boldsymbol{\mu}_k \ \mathsf{and} \ \boldsymbol{\Sigma}_k. \end{array}$

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

- sample means in Gaussian case
- spatial median
 in t₂ 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) S_k(β) with shrinkage towards pooled SCM 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 $d({f A},{f B})$: Kullback-Leibler , Ellipticity, Riemannian distance, . . .
- $\checkmark\,$ many different g-convex loss fnc's $\rho(t):$ Gaussian, Tyler's, Huber's, \ldots
- ✓ 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) S_k(β) with shrinkage towards pooled SCM 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:
- $\checkmark\,$ many different g-convex penalty fnc's $d({\bf A},{\bf B}):\,$ Kullback-Leibler , Ellipticity, Riemannian distance, \ldots
- $\checkmark\,$ many different g-convex loss fnc's $\rho(t)$: Gaussian, Tyler's, Huber's, \ldots
- $\checkmark\,$ robust: good performance under non-Gaussianity or outliers

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.



Cover, T. M. and Thomas, J. A. (2012). *Elements of information theory.* John Wiley & Sons.



Friedman, J. H. (1989). Regularized discriminant analysis. J. Amer. Stat. Assoc., 84(405):165–175.

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. (2014).

Regularized $M\mbox{-estimators}$ 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.



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.