Estimation and classification of location and covariance matrix using Riemannian geometry: application to remote sensing

Antoine Collas

Postdoc in the Mind team at Inria Saclay, supervised by Alexandre Gramfort (Meta).

PhD done at SONDRA lab, CentraleSupélec, Univ. Paris-Saclay in 2019-2022.







- 1. Some context around remote sensing
- 2. Riemannian geometry and problematics
- 3. Geodesic triangles for machine learning
- 4. Estimation and classification of non centered and heteroscedastic data
- 5. Probabilistic PCA from heteroscedastic signals

# Some context around remote sensing

## Context

In recent years, many image time series have been taken from the **earth** with different technologies: **SAR**, **multi/hyper spectral imaging**, ...

#### Objective

**To segment semantically** these data using **sensor diversity** (spectral bands, polarization...), and **spatial** and/or **temporal** informations.



Figure 1: Multivariate image time series.

#### Applications

Activity monitoring, land cover mapping, crop type mapping, disaster assessment ...

## Example of multi-spectral time series

## Breizhcrops dataset<sup>1</sup>:

- more than 600 000 crop time series across the whole Brittany,
- 13 spectral bands, 9 classes.



Figure 2: Reflectances  $\rho$  of a time series of meadows.



Figure 3: Reflectances  $\rho$  of a time series of corn.

<sup>&</sup>lt;sup>1</sup>https://breizhcrops.org/

Indian pines dataset:

- $145 \times 145$  pixels, 200 spectral bands,
- 16 classes (corn, grass, wood, ...).



Figure 4: Raw image.



**Figure 5:** Ground truth, one color = one class.

## **Clustering/classification pipeline**



Figure 6: Clustering/classification pipeline.

#### **Assumption:**

 $\boldsymbol{x} \sim f(.; \theta)$ , a probability density function,  $\theta \in \mathcal{M}$ 

#### **Examples of** $\theta$ :

 $\theta = \Sigma$  a covariance matrix,  $\theta = (\mu, \Sigma)$  a vector and a covariance matrix,  $\theta = (\{\tau_i\}, U)$  a scalar and an orthogonal matrix...

 $\mathcal{M}$  can be constrained !

## Step 2: objectives for feature estimation



**Figure 7:** Example of a SAR image (from nasa.gov).



**Figure 8:** Example of a hyperspectral image (from nasa.gov).

#### **Objectives:**

- develop robust estimators, *i.e.* estimators that work well with non Gaussian or heterogeneous data because of the high resolution of images,
- develop **regularized/structured estimators**, *i.e.* estimators that handle the high dimension of hyperspectral images.

## Step 3: objectives for clustering/classification



 $\begin{array}{c} \mathcal{M} \\ \theta_2 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{array}$ 

**Figure 9:** Divergence  $\delta_{\gamma}$ : squared length of the curve  $\gamma$ .



#### **Objectives:**

develop divergences that

- respect the constraints of  $\mathcal{M}$ ,
- are related to the chosen statistical distributions,
- are robust to **mislabelization** and **distribution shifts** between train and test sets.

## Clustering/classification pipeline and Riemannian geometry

Random variable:  $\mathbf{x} \sim f(.; \theta), \ \theta \in \mathcal{M}$ 

Step 2: maximum likelihood estimation

$$\underset{\theta \in \mathcal{M}}{\text{minimize}} \mathcal{L}(\theta | \{ \boldsymbol{x}_i \}_{i=1}^n) = -\log f(\{ \boldsymbol{x}_i \}_{i=1}^n; \theta)$$

Step 3: given  $\delta$ , center of mass of  $\{\theta_i\}_{i=1}^M$ 

$$\underset{\theta \in \mathcal{M}}{\text{minimize}} \sum_{i} \delta(\theta, \theta_i)$$

Use of Riemannian geometry:

- optimization under constraints,
- "Fisher information metric" ⇒ a canonical Riemannian manifold for the parameter space *M* (fast estimators, intrinsic Carmér-Rao bounds...),
- $\delta$ : squared Riemannian distance.

## Riemannian geometry and problematics

## What is a Riemannian manifold ?



Curvature induced by:

- constraints, e.g. the sphere:  $\|\mathbf{x}\| = 1$ ,
- Riemannian metric, *e.g.* on  $S_p^{++}$ :  $\langle \boldsymbol{\xi}, \boldsymbol{\eta} \rangle_{p}^{S_p^{++}} = \operatorname{Tr}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\xi} \boldsymbol{\Sigma}^{-1} \boldsymbol{\eta}).$

#### Some geometric tools:

- tangent space  $T_{\theta}\mathcal{M}$  (vector space): linearization of  $\mathcal{M}$  at  $\theta \in \mathcal{M}$ ,
- **Riemannian metric**  $\langle ., . \rangle_{\theta}^{\mathcal{M}}$ : inner product on  $T_{\theta}\mathcal{M}$ ,
- geodesic  $\gamma$ : curve on  $\mathcal{M}$  with zero acceleration,
- exponential map:  $\exp_{\theta}^{\mathcal{M}}(\xi) = \gamma(1)$  with  $\gamma(0) = \theta$ ,  $\dot{\gamma}(0) = \xi$ ,
- logarithmic map:  $\log_{\theta_1}^{\mathcal{M}}(\theta_2) = \xi$  with  $\exp_{\theta_1}^{\mathcal{M}}(\xi) = \theta_2$ ,

• distance: 
$$d_{\mathcal{M}}(\theta_1, \theta_2) = \left\| \log_{\theta_1}^{\mathcal{M}}(\theta_2) \right\|_{\theta_1}^{\mathcal{M}}$$
.

#### N. Boumal, "An introduction to optimization on smooth manifolds"

## What is a Riemannian manifold ?

#### Examples of ${\mathcal M}$

- linear spaces:  $\mathbb{R}^{p \times k}$ ,  $\mathcal{S}_p = \{ \boldsymbol{X} \in \mathbb{R}^{p \times p} : \boldsymbol{X}^T = \boldsymbol{X} \}$ ,
- norm constraints:  $S^{p^2-1} = \{ \boldsymbol{X} \in \mathbb{R}^{p \times p} : \| \boldsymbol{X} \|_F = 1 \},$
- positivity constraints:  $S_{\rho}^{++} = \{ \boldsymbol{\Sigma} \in S_{\rho} : \forall \boldsymbol{x} \neq \boldsymbol{0} \in \mathbb{R}^{\rho}, \ \boldsymbol{x}^{T} \boldsymbol{\Sigma} \boldsymbol{x} > 0 \},\$
- positivity and scale constraints:  $\mathcal{SS}_p^{++} = \{ \mathbf{\Sigma} \in \mathcal{S}_p^{++} : |\mathbf{\Sigma}| = 1 \},\$
- orthogonality constraints:  $St_{p,k} = \{ \boldsymbol{U} \in \mathbb{R}^{p \times k} : \boldsymbol{U}^T \boldsymbol{U} = \boldsymbol{I}_k \},\$
- quotient spaces:  $Gr_{p,k} = \{ \{ \boldsymbol{U}\boldsymbol{O} : \boldsymbol{O} \in \mathcal{O}_k \} : \boldsymbol{U} \in St_{p,k} \}.$

#### Optimization

 $f:\mathcal{M}
ightarrow\mathbb{R}$ , smooth

 $\underset{\theta \in \mathcal{M}}{\text{minimize } f(\theta)}$ 









## Fisher information metric

#### Random variable, negative log-likelihood

$$oldsymbol{x} \sim f(.; heta), \quad eta \in \mathcal{M}$$
  
 $\mathcal{L}( heta | oldsymbol{x}) = -\log f(oldsymbol{x}; heta)$ 

**Fisher information metric** 

$$\begin{aligned} \langle \boldsymbol{\xi}, \boldsymbol{\eta} \rangle_{\boldsymbol{\theta}}^{\mathsf{FIM}} &= \mathbb{E}_{\mathbf{x} \sim f(.;\boldsymbol{\theta})} \left[ \mathsf{D}^{2} \, \mathcal{L}\left(\boldsymbol{\theta} | \mathbf{x}\right) [\boldsymbol{\xi}, \boldsymbol{\eta}] \right] \\ &= \mathsf{vec}(\boldsymbol{\xi})^{T} I(\boldsymbol{\theta}) \mathsf{vec}(\boldsymbol{\eta}) \end{aligned}$$

where

$$I(\theta) = \mathbb{E}_{\boldsymbol{x} \sim f(.;\theta)} [\text{Hess } \mathcal{L}(\theta|\boldsymbol{x})] \in \mathcal{S}_{p}^{++}$$

is the Fisher information matrix.

(Set of constraints, Fisher information metric) = a Riemannian manifold

## Existing work: centered Gaussian

A well known geometry:

 $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}), \quad \mathbf{\Sigma} \in \mathcal{S}_p^{++}$ 

with the Fisher information metric:

$$\langle \boldsymbol{\xi}, \boldsymbol{\eta} 
angle_{oldsymbol{\Sigma}}^{\mathsf{FIM}} = \mathsf{Tr} \left( oldsymbol{\Sigma}^{-1} oldsymbol{\xi} oldsymbol{\Sigma}^{-1} oldsymbol{\eta} 
ight).$$

#### Induced pipeline

Step 2:

$$\hat{\boldsymbol{\Sigma}}_{\text{SCM}} = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{x}_i \boldsymbol{x}_i^T.$$

Step 3: geodesic distance on  $\mathcal{S}_p^{++}$ 

$$d_{\mathcal{S}_p^{++}}(\boldsymbol{\Sigma}_1,\boldsymbol{\Sigma}_2) = \left\| \log \left( \boldsymbol{\Sigma}_1^{-\frac{1}{2}} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_1^{-\frac{1}{2}} \right) \right\|_2.$$

Riemannian gradient descent to solve:

$$\min_{\mathbf{\Sigma}\in\mathcal{S}_p^{++}}\sum_i d_{\mathcal{S}_p^{++}}^2(\mathbf{\Sigma},\mathbf{\Sigma}_i).$$

Performant and standard pipeline in EEG/MEG signals classification !

A. Barachant et al., "Multi-class Brain Computer Interface Classification by Riemannian Geometry", IEEE Transactions on Biomedical Engineering

## Problematics

Go beyond  $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma})$ 

- $\pmb{x} \sim \mathcal{N}(\pmb{\mu}, \pmb{\Sigma})$  for non-centered data,
- $m{x}_i \sim \mathcal{N}(m{\mu}, au_i m{\Sigma})$  for non-centered data and robustness,
- $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \tau_i \mathbf{U} \mathbf{U}^T + \mathbf{I}_p)$  for high dimensional data and robustness.

#### Problems

- Existence of maximum likelihood estimators ?
- Not always closed form estimators: how to get fast iterative algo. ?
- Not always closed form expression of the Riemannian distance: what to do ?
- How to get fast estimators of centers of mass ?



Figure 11: Clustering/classification pipeline.

#### Non-centered multivariate Gaussian distributions

Let  $x_1, \dots, x_n \in \mathbb{R}^p$  distributed as  $x \sim \mathcal{N}(\mu, \Sigma)$  with  $\mu \in \mathbb{R}^p$ ,  $\Sigma \in \mathcal{S}_p^{++}$ . Goal: classify  $\theta = (\mu, \Sigma)$ .

#### Riemannian geometry of Gaussian distributions

Space of  $\theta = (\mu, \Sigma) \in \mathbb{R}^p \times S_p^{++}$  with the Fisher information metric:  $\forall \xi = (\xi_{\mu}, \xi_{\Sigma}), \eta = (\eta_{\mu}, \eta_{\Sigma})$  in the tangent space

$$\langle \xi, \eta \rangle_{\theta}^{\mathsf{FIM}} = \boldsymbol{\xi}_{\mu}^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} \boldsymbol{\eta}_{\mu} + \frac{1}{2} \operatorname{Tr} \left( \boldsymbol{\Sigma}^{-1} \boldsymbol{\xi}_{\boldsymbol{\Sigma}} \boldsymbol{\Sigma}^{-1} \boldsymbol{\eta}_{\boldsymbol{\Sigma}} \right).$$

#### Problem

This Riemannian geometry is not fully known...



M. Calvo and J. M. Oller, "An explicit solution of information geodesic equations for the multivariate normal model", Statistics & Risk Modeling



Riemannian center of mass  $(\mu, \Sigma)$  of  $\{(\mu_i, \Sigma_i)\}$ 

$$(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \argmin_{(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \in \mathbb{R}^{p} \times \mathcal{S}_{p}^{++}} \sum_{i} \delta_{\{c, \bot\}} \left( (\boldsymbol{\mu}, \boldsymbol{\Sigma}), (\boldsymbol{\mu}_{i}, \boldsymbol{\Sigma}_{i}) \right)$$

Computed with a Riemannian gradient descent on  $(\mathbb{R}^{p} \times \mathcal{S}_{p}^{++}, \langle ., . \rangle^{\mathsf{FIM}}).$ 

## Application to the Breizhcrops dataset

Parameter estimation + classification with a Nearest centroid classifier



**Figure 12:** "Overall Accuracy" versus the parameter t of two transformations applied to the test set of the *Breizhcrops* dataset.

## Estimation and classification of non centered and heteroscedastic data

## Non-centered mixtures of scaled Gaussian distributions



Figure 13: Clustering/classification pipeline.

## Non-centered mixtures of scaled Gaussian distributions (NC-MSGs)

Let  $x_1, \dots, x_n \in \mathbb{R}^p$  distributed as  $x_i \sim \mathcal{N}(\mu, \tau_i \Sigma)$  with  $\mu \in \mathbb{R}^p$ ,  $\Sigma \in \mathcal{S}_p^{++}$ , and  $\tau \in (\mathbb{R}^+_*)^n$ . Goal: estimate and classify  $\theta = (\mu, \Sigma, \tau)$ .

Interesting when data are heteroscedastic (e.g. time series) and/or contain outliers.

## Parameter space and cost functions

#### Parameter space: location, scatter matrix, and textures

$$\mathcal{M}_{p,n} = \mathbb{R}^p \times \mathcal{S}_p^{++} \times \mathcal{S}(\mathbb{R}^+_*)^n$$

where

$$\mathcal{S}(\mathbb{R}^+_*)^n = \left\{ oldsymbol{ au} \in (\mathbb{R}^+_*)^n : \prod_{i=1}^n au_i = 1 
ight\}$$

- Positivity constraints:  $\Sigma \succ 0$ ,  $\tau_i > 0$
- Scale constraint:  $\prod_{i=1}^{n} \tau_i = 1$

Need generic optimization algorithms on  $\mathcal{M}_{p,n}$ .

#### Parameter estimation

Minimization of a regularized negative log-likelihood (NLL),  $\beta \geq 0$ 

 $\underset{\theta \in \mathcal{M}_{p,n}}{\text{minimize}} \mathcal{L}\left(\theta | \{\mathbf{x}_i\}_{i=1}^n\right) + \beta \mathcal{R}_{\kappa}(\theta)$ 

#### Center of mass estimation

Averaging parameters  $\{\theta_i\}_{i=1}^M$  with a to be defined divergence  $\delta$ 

$$\underset{\theta \in \mathcal{M}_{p,n}}{\text{minimize}} \ \frac{1}{M} \sum_{i=1}^{M} \delta(\theta, \theta_i)$$

### Parameter space with a product metric

#### **Product metric**

Let  $\xi=(\pmb{\xi}_{\mu},\pmb{\xi}_{\pmb{\Sigma}},\pmb{\xi}_{\pmb{ au}}),$   $\eta=(\pmb{\eta}_{\mu},\pmb{\eta}_{\pmb{\Sigma}},\pmb{\eta}_{\pmb{ au}})$  in the tangent space,

$$\langle \xi, \eta \rangle_{\rho, n}^{\mathcal{M}_{\rho, n}^{\text{Prod.}}} = \boldsymbol{\xi}_{\mu}^{T} \boldsymbol{\eta}_{\mu} + \text{Tr}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\xi}_{\boldsymbol{\Sigma}} \boldsymbol{\Sigma}^{-1} \boldsymbol{\eta}_{\boldsymbol{\Sigma}}) + (\boldsymbol{\xi}_{\tau} \odot \boldsymbol{\tau}^{\odot - 1})^{T} (\boldsymbol{\eta}_{\tau} \odot \boldsymbol{\tau}^{\odot - 1})$$

where  $\odot$  is the elementwise operator. Product manifold  $\implies$  easy derivation of the geometric tools  $\implies$  Riemannian gradient descent and conjugate gradient on  $\left(\mathcal{M}_{p,n},\langle.,.\rangle_{\cdot}^{\mathcal{M}_{p,n}^{\text{Prod.}}}\right)$ .

Slow in practice ...



### Parameter space with the Fisher information metric

#### Fisher information metric of NC-MSGs

Let  $\xi=(\pmb{\xi}_{\pmb{\mu}},\pmb{\xi}_{\pmb{\Sigma}},\pmb{\xi}_{\pmb{ au}}),~\eta=(\pmb{\eta}_{\pmb{\mu}},\pmb{\eta}_{\pmb{\Sigma}},\pmb{\eta}_{\pmb{ au}})$  in the tangent space,

$$\langle \xi, \eta \rangle_{\theta}^{\mathcal{M}_{p,n}^{\mathsf{FM}}} = \sum_{i=1}^{n} \frac{1}{\tau_{i}} \xi_{\mu}^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} \boldsymbol{\eta}_{\mu} + \frac{n}{2} \operatorname{Tr}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\xi}_{\boldsymbol{\Sigma}} \boldsymbol{\Sigma}^{-1} \boldsymbol{\eta}_{\boldsymbol{\Sigma}}) + \frac{p}{2} (\boldsymbol{\xi}_{\tau} \odot \boldsymbol{\tau}^{\odot - 1})^{\mathsf{T}} (\boldsymbol{\eta}_{\tau} \odot \boldsymbol{\tau}^{\odot - 1})$$

Most geometric tools remain unknown like geodesics, distance ... But, derivation of the Riemannian gradient and a second order retraction.  $\implies \text{Riemannian gradient descent on } \left(\mathcal{M}_{p,n}, \langle ., \rangle_{-}^{\mathcal{M}_{p,n}^{\text{FIM}}}\right).$ 



## Simulated data



**Figure 16:** MSE over 2000 simulated sets  $\{x_i\}_{i=1}^n \subset \mathbb{R}^{10}$  versus the number samples  $x_i$  for the considered estimators  $\hat{\mu} \in \{\hat{\mu}_{SM}, \hat{\mu}_{Ty}, \hat{\mu}_{IG}\}$  and  $\hat{\Sigma} \in \{\hat{\Sigma}_{SCM}, \hat{\Sigma}_{Ty,\mu}, \hat{\Sigma}_{Ty}, \hat{\Sigma}_{IG}\}$ . The proposed estimators are  $\hat{\mu}_{IG}$  and  $\hat{\Sigma}_{IG}$ .

D. E. Tyler, "A Distribution-Free M-Estimator of Multivariate Scatter", Ann. Statist.

NLL neglecting terms not depending on  $\theta$ :

$$\mathcal{L}(\theta|\{\mathbf{x}_i\}_{i=1}^n) = \frac{1}{2} \sum_{i=1}^n \left[ \log |\tau_i \mathbf{\Sigma}| + \frac{(\mathbf{x}_i - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})}{\tau_i} \right].$$

Existence of the maximum likelihood estimator ?

Observation of sequences  $(\theta^{(\ell)})_{\ell}$  such that

$$\mathcal{L}\left(\theta^{(\ell+1)}|\{\boldsymbol{x}_i\}_{i=1}^n\right) < \mathcal{L}\left(\theta^{(\ell)}|\{\boldsymbol{x}_i\}_{i=1}^n\right) \quad \text{and} \quad \theta^{(\ell)} \xrightarrow[\ell \to +\infty]{} \partial \theta$$

where  $\partial \theta$  is a border of  $\mathcal{M}_{p,n}$  (e.g.  $\tau_i = 0$ ).

The existence of a minimum in  $\mathcal{M}_{p,n}$  depends on  $\{\mathbf{x}_i\}_{i=1}^n$  !!

## Parameter estimation

#### Existence of a regularized maximum likelihood estimator

Under some assumptions on  $\mathcal{R}_{\kappa}$  and  $\beta >$  0, the regularized NLL

 $\theta \mapsto \mathcal{L}\left(\theta | \{\mathbf{x}_i\}_{i=1}^n\right) + \beta \mathcal{R}_{\kappa}(\theta),$ 

admits a minimum in  $\mathcal{M}_{p,n}$ .

Name	$\mathcal{R}_{\kappa}$
L1 penalty	$\left\  \left( diag(\boldsymbol{\tau}) \otimes \boldsymbol{\Sigma} \right)^{-1} - \kappa^{-1} \boldsymbol{l}_{n \times p} \right\ _{1} = \sum_{i,j} \left  \left( \tau_{i} \lambda_{j} \right)^{-1} - \kappa^{-1} \right $
L2 penalty	$\left\  \left( diag(\boldsymbol{\tau}) \otimes \boldsymbol{\Sigma} \right)^{-1} - \kappa^{-1} \boldsymbol{I}_{n \times p} \right\ _{2}^{2} = \sum_{i,j} \left( \left( \tau_{i} \lambda_{j} \right)^{-1} - \kappa^{-1} \right)^{2}$
Bures-Wasserstein squared distance	$d_{BW}^{2}\left((diag(\boldsymbol{\tau})\otimes\boldsymbol{\Sigma})^{-1},\kappa^{-1}\boldsymbol{I}_{n\times\rho}\right)=\sum_{i,j}\left(\left(\tau_{i}\lambda_{j}\right)^{-\frac{1}{2}}-\kappa^{-\frac{1}{2}}\right)^{2}$
Gaussian KL divergence	$ \begin{split} \delta_{KL}(\kappa \boldsymbol{I}_{n \times p}, diag(\boldsymbol{\tau}) \otimes \boldsymbol{\Sigma}) = \\ \frac{1}{2} \left[ \sum_{i,j} \left( \kappa \left( \tau_i \lambda_j \right)^{-1} + \log \left( \tau_i \lambda_j \right) \right) - np(1 + \log(\kappa)) \right] \end{split} $

where  $\lambda_j$  are the eigenvalues of  $\Sigma$ .

## Classification

### KL divergence between NC-MSGs

$$\begin{split} \delta_{\mathsf{KL}}(\theta_1, \theta_2) &= \\ & \frac{1}{2} \left[ \sum_{i=1}^n \frac{\tau_{1,i}}{\tau_{2,i}} \operatorname{Tr}\left( \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1 \right) + \sum_{i=1}^n \frac{1}{\tau_{2,i}} \Delta \boldsymbol{\mu}^\top \boldsymbol{\Sigma}_2^{-1} \Delta \boldsymbol{\mu} + n \log\left( \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} \right) - n \boldsymbol{p} \right] \end{split}$$

with  $\Delta \mu = \mu_2 - \mu_1$ .

Symmetrization:

$$\delta_{\mathcal{M}_{p,n}}(\theta_1,\theta_2) = \frac{1}{2} \left( \delta_{\mathsf{KL}}(\theta_1,\theta_2) + \delta_{\mathsf{KL}}(\theta_2,\theta_1) \right)$$

#### **Riemannian center of mass**

$$\begin{array}{l} \underset{\theta \in \mathcal{M}_{p,n}}{\text{minimize}} \; \frac{1}{M} \sum_{i=1}^{M} \delta_{\mathcal{M}_{p,n}}(\theta, \theta_i) \\ \\ \text{Riemannian gradient descent on } \Big( \mathcal{M}_{p,n}, \langle ., . \rangle_{\cdot}^{\mathcal{M}_{p,n}^{\mathsf{FIM}}} \Big). \end{array}$$

## Simulated data

M = 2

M = 100



**Figure 17:** KL variance and its gradient norm versus iterations of the optimizers. The dimensions of the parameter space are p = 10 and n = 150. Two different numbers of points *M* are considered: 2 in the left column and 100 in the right one.

## Application to the Breizhcrops dataset

Parameter estimation + classification with a Nearest centroid classifier



**Figure 18:** "Overall Accuracy" metric versus the parameter *t* associated with transformations applied to the test set. The proposed *Nearest centroid classifier* is " $\theta$  - sym. KL". The regularization is the L2 penalty and  $\beta = 10^{-11}$ .

# Probabilistic PCA from heteroscedastic signals

## Study of a "low rank" statistical model



Figure 19: Clustering/classification pipeline.

**Statistical model** 

 $\boldsymbol{x}_1, \cdots, \boldsymbol{x}_n \in \mathbb{R}^p$ ,  $\forall k < p$ :

$$oldsymbol{x}_i \sim \mathcal{N}(oldsymbol{0}, au_i oldsymbol{U}oldsymbol{U}^T + oldsymbol{I}_p)$$

with  $\tau_i > 0$  and  $\boldsymbol{U} \in \mathbb{R}^{p \times k}$  is an orthogonal basis  $(\boldsymbol{U}^T \boldsymbol{U} = \boldsymbol{I}_k)$ . Goal: estimate and classify  $\theta = (\boldsymbol{U}, \boldsymbol{\tau})$ .

## Study of a "low rank" statistical model

#### **Statistical model**



where  $\boldsymbol{g} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_k)$  and  $\boldsymbol{n} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_p)$  are independent,  $\boldsymbol{\tau} \in (\mathbb{R}^+_*)^n$ , and  $\boldsymbol{U} \in \mathbb{R}^{p \times k}$  is an orthogonal basis  $(\boldsymbol{U}^T \boldsymbol{U} = \boldsymbol{I}_k)$ .



#### Maximum likelihood estimation

Minimization of the NLL with constraints,  $heta = (oldsymbol{U}, oldsymbol{ au})$ 

- $U \in Gr_{p,k}$ : orthogonal basis of the subspace (and thus invariant by rotation !)
- $oldsymbol{ au} \in (\mathbb{R}^+_*)^n$  : positivity constraints

 $\underset{\theta \in Gr_{p,k} \times (\mathbb{R}^+_*)^n}{\text{minimize}} \mathcal{L}(\theta | \{\mathbf{x}_i\}_{i=1}^n)$ 

## Study of a "low rank" statistical model: estimation

#### **Fisher information metric**

$$\begin{split} \forall \boldsymbol{\xi} &= \left(\boldsymbol{\xi}_{\boldsymbol{U}}, \boldsymbol{\xi}_{\boldsymbol{\tau}}\right), \boldsymbol{\eta} = \left(\boldsymbol{\eta}_{\boldsymbol{U}}, \boldsymbol{\eta}_{\boldsymbol{\tau}}\right) \text{ in the tangent space} \\ & \left\langle \boldsymbol{\xi}, \boldsymbol{\eta} \right\rangle_{\boldsymbol{\theta}}^{\mathsf{FIM}} = 2nc_{\boldsymbol{\tau}} \operatorname{Tr} \left(\boldsymbol{\xi}_{\boldsymbol{U}}^{\mathsf{T}} \boldsymbol{\eta}_{\boldsymbol{U}}\right) + k \left(\boldsymbol{\xi}_{\boldsymbol{\tau}} \odot (1+\tau)^{\odot -1}\right)^{\mathsf{T}} \left(\boldsymbol{\eta}_{\boldsymbol{\tau}} \odot (1+\tau)^{\odot -1}\right), \end{split}$$
 where  $c_{\boldsymbol{\tau}} = \frac{1}{n} \sum_{i=1}^{n} \frac{\tau_{i}^{2}}{1+\tau_{i}}.$  Derivation of the Riemannian gradient and of a retraction.

To minimize the NLL: Riemannian gradient descent on  $(Gr_{p,k} \times (\mathbb{R}^+_*)^n, \langle ., . \rangle^{FIM})$ .



Figure 20: NLL versus the iterations.

## Study of a "low rank" statistical model: bounds

#### Intrinsic Cramér-Rao bounds

Study of the performance through intrinsic Cramér-Rao bounds:

$$\underbrace{\mathbb{E}[d_{\mathsf{Gr}_{p,k}}^2(\mathsf{span}(\hat{\boldsymbol{U}}),\mathsf{span}(\boldsymbol{U}))]}_{\mathsf{E}[d_{\mathsf{Gr}_{p,k}}^2} \approx \frac{(p-k)k}{nc_{\tau}} \approx \frac{(p-k)k}{n\times\mathsf{SNR}}$$
$$\underbrace{\mathbb{E}[d_{(\mathbb{R}^+_*)^n}^2(\hat{\tau},\tau)]}_{\mathsf{texture estimation error}} \geq \frac{1}{k}\sum_{i=1}^n \frac{(1+\tau_i)^2}{\tau_i^2}$$



Figure 21: Mean squared error versus the number of simulated data.

## Study of a "low rank" statistical model: K-means++





Figure 23: Center of mass  $(U, \tau)$ .





Figure 24:

Euclidean *K-means++*: OA = 31.2%.



Figure 25: Proposed *K*-means++: OA = 47.2%.

Github: https://github.com/antoinecollas/pyCovariance

## Conclusions

## Conclusions



Figure 27: Clustering/classification pipeline.

Interests of using **Riemannian geometry** in this pipeline:

- fast estimators, analysis of constrained estimation problems using intrinsic Cramér-Rao bounds,
- new classifiers robust to distribution shifts.

**Applications** on real datasets: *Indian pines* hyperspectral image, *Breizhcrops* multispectral times series, datasets from the *UCI* repository. Estimation and classification of location and covariance matrix using Riemannian geometry: application to remote sensing

Antoine Collas

Postdoc in the Mind team at Inria Saclay, supervised by Alexandre Gramfort (Meta).

PhD done at SONDRA lab, CentraleSupélec, Univ. Paris-Saclay in 2019-2022.







## Robust Geometric Metric Learning

## Robust Geometric Metric Learning (RGML)

Let be a classification problem with K classes.

#### Metric learning

Find a Mahalanobis distance

$$d_{\boldsymbol{A}}(\boldsymbol{x}_i, \boldsymbol{x}_j) = \sqrt{(\boldsymbol{x}_i - \boldsymbol{x}_j)^T \boldsymbol{A}^{-1} (\boldsymbol{x}_i - \boldsymbol{x}_j)}$$

relevant for classification problems.

#### Metric learning as covariance estimation

Proposed minimization problem,  $\theta = (\mathbf{A}, \{\mathbf{A}_k\})$ 



 $\{\pi_k\}$  are the proportions of the classes and  $\{\mathcal{L}_k\}$  are to be defined.





## Robust Geometric Metric Learning (RGML)

Let  $\mathbf{s}_{ki} = \mathbf{x}_{l} - \mathbf{x}_{m}$  where  $\mathbf{x}_{l}, \mathbf{x}_{m}$  belong to the class k.

Gaussian negative log-likelihood

$$\mathcal{L}_{G,k}(\boldsymbol{A}_k) = \frac{1}{n_k} \sum_{i=1}^{n_k} \boldsymbol{s}_{ki}^{\mathsf{T}} \boldsymbol{A}_k^{-1} \boldsymbol{s}_{ki} + \log |\boldsymbol{A}_k|$$
  
minimized for  $\boldsymbol{A}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} \boldsymbol{s}_{ki} \boldsymbol{s}_{ki}^{\mathsf{T}}$ 

|--|

#### Tyler cost function

$$\mathcal{L}_{T,k}(\boldsymbol{A}_k) = \frac{p}{n_k} \sum_{i=1}^{n_k} \log\left(\boldsymbol{s}_{ki}^T \boldsymbol{A}_k^{-1} \boldsymbol{s}_{ki}\right) + \log|\boldsymbol{A}_k|$$
  
minimized for  $\boldsymbol{A}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} \underbrace{\frac{p}{\boldsymbol{s}_{ki}^T \boldsymbol{A}_k^{-1} \boldsymbol{s}_{ki}}}_{\text{weight of sample } \boldsymbol{s}_{ki}} \boldsymbol{s}_{ki} \boldsymbol{s}_{ki}^T$ 



## Robust Geometric Metric Learning (RGML)

#### **Riemannian metric**

 $orall\xi = \left(oldsymbol{\xi}, \{oldsymbol{\xi}_k\}
ight), \eta = \left(oldsymbol{\eta}, \{oldsymbol{\eta}_k\}
ight)$  in the tangent space

$$\langle \xi, \eta \rangle_{\theta} = \operatorname{Tr} \left( \mathbf{A}^{-1} \boldsymbol{\xi} \mathbf{A}^{-1} \boldsymbol{\eta} \right) + \sum_{k=1}^{K} \operatorname{Tr} \left( \mathbf{A}_{k}^{-1} \boldsymbol{\xi}_{k} \mathbf{A}_{k}^{-1} \boldsymbol{\eta}_{k} \right)$$

 $\implies$  strongly geodesically convexity of the minimization problem  $\implies$  the Riemannian gradient descent is fast



Figure 28: Cost function versus the iterations.

RGML + k-NN on datasets from the UCI Machine Learning Repository

	Wine				Vehicle				Iris			
	p = 13 , $n = 178$ , $K = 3$				p = 18, n = 846, K = 4				p = 4, n = 150, K = 3			
Method	Mislabeling rate				Mislabeling rate				Mislabeling rate			
	0%	5%	10%	15%	0%	5%	10%	15%	0%	5%	10%	15%
Euclidean	30.12	30.40	31.40	32.40	38.27	38.58	39.46	40.35	3.93	4.47	5.31	6.70
SCM	10.03	11.62	13.70	17.57	23.59	24.27	25.24	26.51	12.57	13.38	14.93	16.68
ITML - Identity	3.12	4.15	5.40	7.74	24.21	23.91	24.77	26.03	3.04	4.47	5.31	6.70
ITML - SCM	2.45	4.76	6.71	10.25	23.86	23.82	24.89	26.30	3.05	13.38	14.92	16.67
GMML	2.16	3.58	5.71	9.86	21.43	22.49	23.58	25.11	2.60	5.61	9.30	12.62
LMNN	4.27	6.47	7.83	9.86	20.96	24.23	26.28	28.89	3.53	9.59	11.19	12.22
Proposed - Gaussian	2.07	2.93	5.15	9.20	19.76	21.19	22.52	24.21	2.47	5.10	8.90	12.73
Proposed - Tyler	2.12	2.90	4.51	8.31	19.90	20.96	22.11	23.58	2.48	2.96	4.65	7.83

**Table 1:** Misclassification errors on 3 datasets: Wine, Vehicle and Iris.Mislabeling rate: percentage of labels randomly changed in the training set.

Github: https://github.com/antoinecollas/robust\_metric\_learning

Perspectives

- Beyond KL divergence between non-centered mixtures of scaled Gaussian distributions:
  - Wasserstein distance
  - Entropy-Regularized Optimal Transport

Wasserstein distance between Gaussian distributions (from H. Janati et al., "Entropic Optimal Transport between Unbalanced Gaussian Measures has a Closed Form", Neurips 2020):

 $W_2^2\left(\mathcal{N}(oldsymbol{\mu}_1,oldsymbol{\Sigma}_1),\mathcal{N}(oldsymbol{\mu}_2,oldsymbol{\Sigma}_2)
ight)=$ 

$$\|\mu_1 - \mu_2\|_2^2 + \mathsf{Tr}(\Sigma_1) + \mathsf{Tr}(\Sigma_2) - 2\,\mathsf{Tr}\left(\Sigma_1^{\frac{1}{2}}\Sigma_2\Sigma_1^{\frac{1}{2}}\right)^{\frac{1}{2}}$$

• Apply *RGML* to more recent datasets (larger *p*, larger *n*, larger *K*): SGD, SGA, "low rank + identity" structure, ...

## Perspectives

• Integrate developed estimators in a differentiable pipeline.



**Figure 29:** Tangent Space Mapping Network (from R. J. Kobler et al., "SPD domain-specific batch normalization to crack interpretable unsupervised domain adaptation in EEG", Neurips 2022)

• Domain adaptation: see if developed algorithms are robust to distributions shifts on real datasets; *e.g.* inter-session and -subject in EEG datasets.

## Perspectives

Fisher information metric for optimization:

- Only a good preconditioner ?
- When second order retractions are useful ?



Figure 1: A typical situation encountered when performing large discrete updates in the original parameter space. The red arrow is the natural gradient direction given by the vector \$\overline{\nu}\$ in parameter space) and the black arrow is the path generated by taking θ - a\overline{\nu}\$ h for a \overline{\nu}\$.

Figure 30: Euclidean retraction (from James Martens, "New Insights and Perspectives on the Natural Gradient Method", JMLR 2020)

• Iterates defined up to reparametrizations of the sample and parameter spaces:

$$\begin{array}{l} \underset{\boldsymbol{\Sigma} \in \mathcal{S}_{\rho}^{++}}{\text{minimize }} f(\boldsymbol{\Sigma} \mid \{\boldsymbol{x}_i\}) = f(\boldsymbol{A}\boldsymbol{\Sigma}\boldsymbol{A}^T \mid \{\boldsymbol{A}\boldsymbol{x}_i + \mu\}) \quad (\boldsymbol{A} \text{ invertible}) \\ \boldsymbol{\Sigma}^{(\ell+1)} = \exp_{\boldsymbol{\Sigma}^{(\ell)}}^{\mathcal{S}_{\rho}^{++}} (-\alpha \nabla f(\boldsymbol{\Sigma}^{(\ell)})) \\ \exp_{\boldsymbol{A}\boldsymbol{\Sigma}^{(\ell)}\boldsymbol{A}^T}^{\mathcal{S}_{\rho}^{++}} (-\alpha \nabla f(\boldsymbol{A}\boldsymbol{\Sigma}^{(\ell)}\boldsymbol{A}^T)) = \boldsymbol{A} \exp_{\boldsymbol{\Sigma}}^{\mathcal{S}_{\rho}^{++}} (-\alpha \nabla f(\boldsymbol{\Sigma}^{(\ell)})) \boldsymbol{A}^T = \boldsymbol{A}\boldsymbol{\Sigma}^{(\ell+1)}\boldsymbol{A}^T \end{aligned}$$

# References, publications and additional figures

#### Optimization on Riemannian manifolds:

- P.-A. Absil, R. Mahony, and R. Sepulchre. Optimization Algorithms on Matrix Manifolds. Princeton, NJ: Princeton University Press, 2008, pp. xvi+224. ISBN: 978-0-691-13298-3.
- [9] Nicolas Boumal. An introduction to optimization on smooth manifolds. 2022. URL: http://www.nicolasboumal.net/book.

#### Robust statistics:

- [10] Ricardo Antonio Maronna. Robust M-Estimators of Multivariate Location and Scatter. 1976. DOI: 10.1214/aos/1176343347.
- [11] David E. Tyler. A Distribution-Free M-Estimator of Multivariate Scatter. 1987. DOI: 10.1214/aos/1176350263.
- [12] Esa Ollila et al. Complex Elliptically Symmetric Distributions: Survey, New Results and Applications. 2012. DOI: 10.1109/TSP.2012.2212433.

#### Software:

- [13] N. Boumal et al. Manopt, a Matlab Toolbox for Optimization on Manifolds. 2014. URL: https://www.manopt.org.
- [14] J. Townsend, N. Koep, and S. Weichwald. Pymanopt: A Python Toolbox for Optimization on Manifolds Using Automatic Differentiation. 2016.

## Publications

#### Conferences:

- A. Collas et al. "A Tyler-Type Estimator of Location and Scatter Leveraging Riemannian Optimization". In: ICASSP 2021 -2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Toronto, Canada (Virtual). 2021.
- [2] A. Collas et al. "On The Use of Geodesic Triangles Between Gaussian Distributions for Classification Problems". In: ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Singapore. 2022.
- [3] A. Collas et al. "Robust Geometric Metric Learning". In: 2022 30th European Signal Processing Conference (EUSIPCO), Belgrade, Serbia. 2022.
- [4] A. Collas et al. "Apprentissage robuste de distance par géométrie riemannienne". In: GRETSI 2022 XXVIIIème colloque, Nancy, France. 2022.

#### Journals:

- [5] A. Mian, A. Collas et al. "Robust Low-Rank Change Detection for Multivariate SAR Image Time Series". In: IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing (2020).
- [6] A. Collas et al. "Probabilistic PCA From Heteroscedastic Signals: Geometric Framework and Application to Clustering". In: IEEE Transactions on Signal Processing (2021).
- [7] A. Collas et al. "Riemannian optimization for non-centered mixture of scaled Gaussian distributions". In: submitted to IEEE Transactions on Signal Processing (2022).

#### Award:

2022: Best Student Paper Award at the EUSIPCO 2022 conference.

## Study of a "low rank" statistical model: estimation



Figure 31: Negative log-likelihood versus the iterations.

## NC-MSGs: application to the Breizhcrops dataset



**Figure 32:** "Overall Accuracy" metric achieved by the proposed *Nearest centroïd classifier* on the *Breizhcrops* dataset versus  $\beta$ . The regularization is the L2 penalty.

### NC-MSGs: numerical experiments for the estimation



1000

10

100

Iterations

1000

1000

10

Iterations

100

Figure 33: Regularized NLL and its gradient norm versus iterations of the optimizers. The chosen regularization is the L2 penalty. Each estimation is performed on n = 150 samples in  $\mathbb{R}^{10}$  sampled from a NC-MSG.