

Riemannian and information geometry in signal processing and machine learning

Part III: Riemannian geometry applied to machine learning

Florent Bouchard, Arnaud Breloy and **Ammar Mian**



Outline

1 Introduction

2 Parameter on a manifold

- General context
- Gaussian mixture models
- Metric Learning
- Deep learning optimization

3 Data on a manifold

- General principles of using Riemannian Geometry
- Tangent-space based approaches
- Distance based approaches
- More complex algorithms

4 Numerical aspects and Toolboxes



Robust Geometric Metric Learning (RGML)

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

Method	Wine $p = 13, n = 178, K = 3$				Vehicle $p = 18, n = 846, K = 4$				Iris $p = 4, n = 150, K = 3$			
	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.

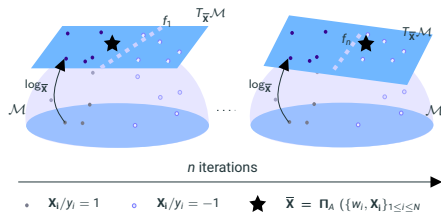
Github: https://github.com/antoinecollas/robust_metric_learning

Another approach: boosting [TPM08a]

Classification based on combining weak-learners (decision trees) $\{f_l : 1 \leq l \leq L\}$ into a classifier with the form $\text{sign}[F(\mathbf{x})] = \text{sign}[\frac{1}{2} \sum_{l=1}^L f_l(\mathbf{x})]$. The probability for feature vector \mathbf{x} of being in class 1 is represented by:

$$p(\mathbf{x}) = \frac{\exp(F(\mathbf{x}))}{\exp(F(\mathbf{x})) + \exp(-F(\mathbf{x}))}, \quad (18)$$

Riemannian equivalent:



Results on DaimerChrysler dataset

		Fold 1	Fold 2	Fold 3	mean
Euclidean	RBF SVM	0.726	0.727	0.727	0.727
	logitboost	0.730	0.734	0.729	0.731
	KNN	0.710	0.708	0.711	0.710
	MDM	0.592	0.590	0.591	0.591
	LogisticRegression	0.700	0.702	0.700	0.701
Riemannian	RBF SVM	0.814	0.814	0.814	0.814
	logitboost	0.741	0.745	0.738	0.741
	KNN	0.727	0.723	0.727	0.726
	MDM	0.638	0.636	0.638	0.638
	LogisticRegression	0.733	0.736	0.735	0.735

Example: Classification of multispectral satellite images

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

Objective

Segment semantically these data using **spatial** information, **temporal** information and **sensor diversity** (spectral bands, polarization...).

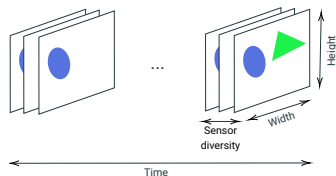


Figure 7: Multivariate image time series.

Applications

Disaster assessment, activity monitoring, land cover mapping, crop type mapping, ...

Example of multi-spectral time series

Breizhcrops dataset¹:

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

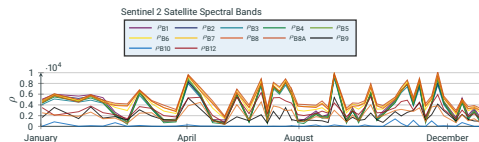


Figure 10: Reflectances ρ of a time series of **meadows**.

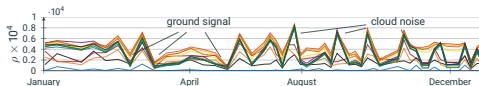


Figure 11: Reflectances ρ of a time series of **corn**.

¹<https://breizhcrops.org/>

Clustering/classification pipeline and Riemannian geometry

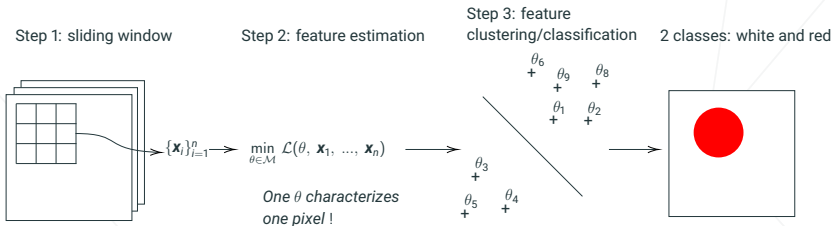


Figure 12: Clustering/classification pipeline.

Examples of θ :

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

Clustering/classification pipeline and Riemannian geometry

Clustering/classification and Riemannian geometry

$\theta \in \mathcal{M}$, a *Riemannian manifold* (constraints and non-constant metric):

step 2: minimization of \mathcal{L} over \mathcal{M} ,

step 3: computing distances and centers of mass on \mathcal{M} .

Existing work (e.g. in BCI classification)

$\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ realizations of $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$, $\boldsymbol{\Sigma} \in \mathcal{S}_p^{++}$.

Step 2: maximum likelihood estimator:

$$\theta = \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T. \quad (19)$$

Step 3: Riemannian distance on \mathcal{S}_p^{++} (geodesic distance):

$$d_{\mathcal{S}_p^{++}}(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2) = \left\| \log \left(\boldsymbol{\Sigma}_1^{-\frac{1}{2}} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_1^{-\frac{1}{2}} \right) \right\|_2. \quad (20)$$

Study of a "low rank" statistical model

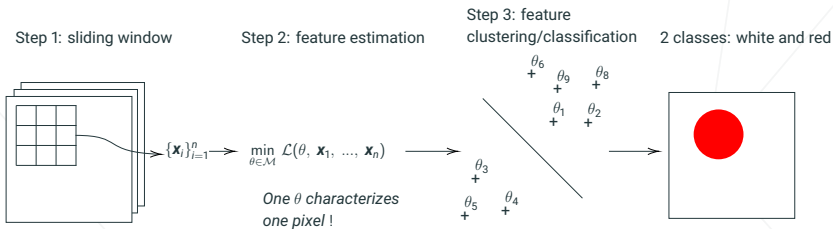


Figure 13: Clustering/classification pipeline.

Statistical model

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

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \tau_i \mathbf{U}\mathbf{U}^T + \mathbf{I}_p) \quad (21)$$

with $\tau_i > 0$ and $\mathbf{U} \in \mathbb{R}^{p \times k}$ is an orthogonal basis ($\mathbf{U}^T \mathbf{U} = \mathbf{I}_k$).

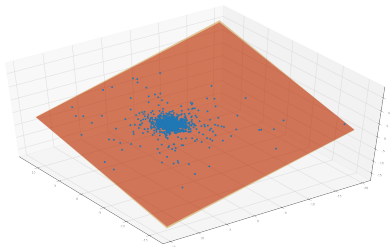
Goal: estimate and classify $\theta = (\mathbf{U}, \tau)$.

Study of a "low rank" statistical model

Statistical model

$$\underbrace{\mathbf{x}_i}_{\in \mathbb{R}^p} \stackrel{d}{=} \underbrace{\sqrt{\tau_i} \mathbf{U} \mathbf{g}_i}_{\text{signal} \in \text{span}(\mathbf{U})} + \underbrace{\mathbf{n}_i}_{\text{noise} \in \mathbb{R}^p} \quad (22)$$

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



Study of a "low rank" statistical model: estimation

Maximum likelihood estimation (MLE)

Minimization of the negative log-likelihood with constraints:

$\mathbf{U} \in \text{Gr}_{p,k}$: orthogonal basis of the subspace (and thus invariant by rotation !)

$\boldsymbol{\tau} \in (\mathbb{R}_*^+)^n$: positivity constraints

$$\underset{(\mathbf{U}, \boldsymbol{\tau}) \in \text{Gr}_{p,k} \times (\mathbb{R}_*^+)^n}{\text{minimize}} \quad \mathcal{L}(\mathbf{U}, \boldsymbol{\tau}) \quad (23)$$

Study of a "low rank" statistical model: estimation

Fisher information metric

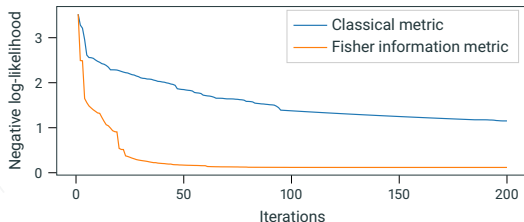
$\forall \xi = (\xi_U, \xi_\tau), \eta = (\eta_U, \eta_\tau)$ in the tangent space

$$\langle \xi, \eta \rangle_{(U, \tau)}^{\text{FIM}} = \mathbb{E} [D \mathcal{L}(\theta)[\xi] D \mathcal{L}(\theta)[\eta]] \quad (24)$$

$$= 2nc_\tau \text{Tr} \left(\xi_U^T \eta_U \right) + k \left(\xi_\tau \odot (\mathbf{1} + \tau)^{\odot -1} \right)^T \left(\eta_\tau \odot (\mathbf{1} + \tau)^{\odot -1} \right), \quad (25)$$

where $c_\tau = \frac{1}{n} \sum_{i=1}^n \frac{\tau_i^2}{1+\tau_i}$.

To solve (23) : Riemannian gradient descent on $(\text{Gr}_{p,k} \times (\mathbb{R}_*^+)^n, \langle \cdot, \cdot \rangle^{\text{FIM}})$.



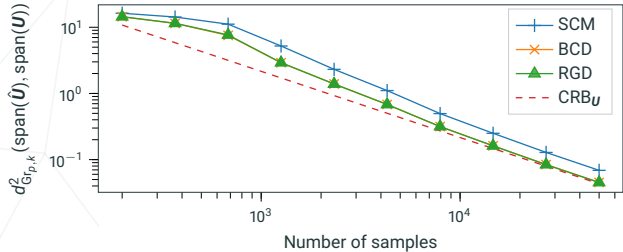
Study of a "low rank" statistical model: bounds

Intrinsic Cramér-Rao bounds

Study of the performance through intrinsic Cramér-Rao bounds:

$$\overbrace{\mathbb{E}[d_{Gr_{p,k}}^2(\text{span}(\hat{\mathbf{U}}), \text{span}(\mathbf{U}))]}^{\text{subspace estimation error}} \geq \frac{(p-k)k}{nc_\tau} \approx \frac{(p-k)k}{n \times \text{SNR}} \quad (26)$$

$$\underbrace{\mathbb{E}[d_{(\mathbb{R}^+)^n}^2(\hat{\tau}, \tau)]}_{\text{texture estimation error}} \geq \frac{1}{k} \sum_{i=1}^n \frac{(1 + \tau_i)^2}{\tau_i^2} \quad (27)$$



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

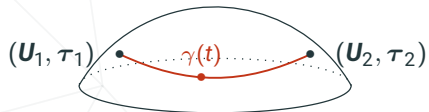


Figure 16: Distance.

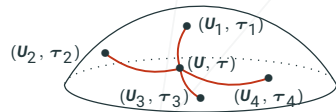


Figure 17: Center of mass (U, τ) .

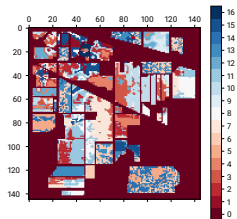


Figure 18: Euclidean *K-means++*:
OA = 31.2%.

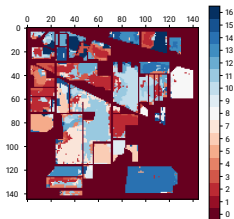


Figure 19: Proposed *K-means++*:
OA = 47.2%.

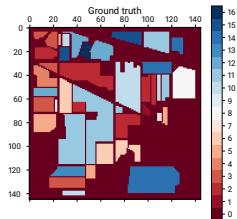


Figure 20: Ground truth.

Approaches when data is on a manifold



Solution 1: Tangent space mapping

Anything Goes

SVM: p.30

Logistic Regression: p.32

MLP: p.30

Logitboost: p.35

Solution 2: Account for the true distances



distance-based algorithm: p. 38

KNN

MDM

K-means

more complex algorithms

GMM: p.58

Kernels: p.54

PCA: p.66

Neural Networks p.59

