

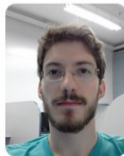
Non-parametric Online Change Point Detection on Riemannian Manifolds

Xiuheng Wang[†], Ricardo Borsoi^{*}, Cédric Richard[†]

[†]Université Côte d'Azur, CNRS, OCA, France

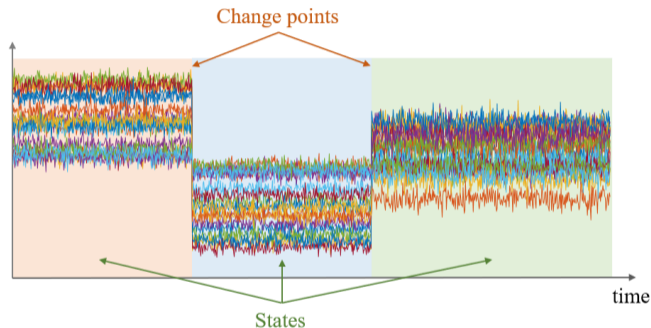
^{*}Université de Lorraine, CNRS, CRAN, France

To appear in ICML 2024.



Change point detection

Change point detection (CPD): detect abrupt changes in the states of time series¹.



- Non-parametric: no prior knowledge of the data distribution.
- Online: process the stream on the fly, ideally without storing raw data.

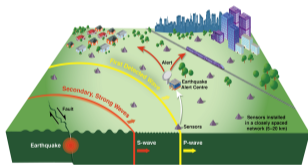
¹Samaneh Aminikhanghahi et al. "A survey of methods for time series change point detection". In: *Knowledge and information systems* 51.2 (2017), pp. 339–367.

CPD on Riemannian manifolds

Many features of signals are lying on different manifolds, e.g.,

- Covariance descriptors
- Subspace representations

Investigate CPD on manifolds can impact many applications, e.g..



earthquake detection²



video change detection³



subspace change detection⁴

²<https://www.earthquakescanada.nrcan.gc.ca/ew-asp/system-en.php>.

³<https://intvo.com/>.

⁴<https://bering-ivis.readthedocs.io/en/stable/>.

CPD on Riemannian manifolds

Developing methods on Riemannian manifolds is challenging:

- the nonlinear geometry;
- lack of vector space structure.

Few works have investigated CPD for manifold-valued data:

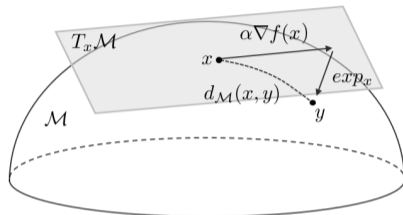
- a parametric algorithm⁵;
- an offline technique⁶.

This work introduces a **general** framework for **non-parametric** and **online** CPD on Riemannian manifolds.

⁵Florent Bouchard et al. "Riemannian geometry for compound Gaussian distributions: Application to recursive change detection". In: *Signal Processing* 176 (2020), p. 107716.

⁶Paromita Dubey et al. "Fréchet change-point detection". In: *The Annals of Statistics* 48.6 (2020), pp. 3312–3335.

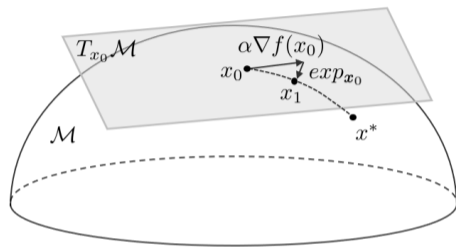
Riemannian optimization: main tools



A few important tools:

- Riemannian gradient: $\nabla f(x) \in T_x\mathcal{M}$
- Exponential mapping: $\exp_x : T_x\mathcal{M} \rightarrow \mathcal{M}$ (maps a vector in the tangent space back to the manifold)
- Riemannian distance: $d_{\mathcal{M}}$ (length of the shortest path between two points on \mathcal{M})

Riemannian optimization: R-SGD, basic structure



Considering a cost $f(\mathbf{x})$, $\mathbf{x} \in \mathcal{M}$ we proceed as⁷:

- compute a stochastic approximation of $\nabla f(\mathbf{x})$ at \mathbf{x}
- “take a step in the negative gradient direction” on \mathcal{M} using the exponential mapping

⁷Silvère Bonnabel. “Stochastic gradient descent on Riemannian manifolds”. In: *IEEE Transactions on Automatic Control* 58.9 (2013), pp. 2217–2229.

Problem formulation

There exists a time index $t_r \in \mathbb{N}$ with an abrupt change in the probability measures⁸ of \mathbf{x}_t lying on \mathcal{M} , that is:

$$t < t_r : \mathbf{x}_t \sim P_1(\mathbf{x}), \quad t \geq t_r : \mathbf{x}_t \sim P_2(\mathbf{x}), \quad (1)$$

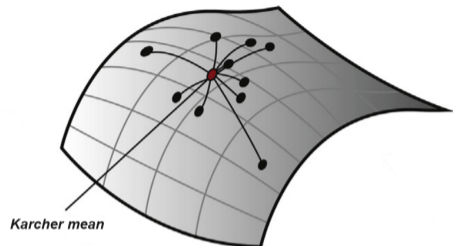
where t_r is the so-called *change point*.

The CPD problem on \mathcal{M} consists of estimating t_r with the following requirements:

- high detection rate;
- low false alarm rate;
- low detection delay.

⁸Xavier Pennec. *Probabilities and statistics on riemannian manifolds: A geometric approach*. Tech. rep. 5093. INRIA, 2004, pp. 1–49.

The algorithm: the Karcher mean



Consider monitoring the **Karcher mean**⁹ on \mathcal{M} , defined as

$$\mathbf{m}^* \in \arg \min_{\mathbf{m}} f(\mathbf{m}). \quad (2)$$

where the Karcher variance

$$f(\mathbf{m}) = \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})} \{ d_{\mathcal{M}}^2(\mathbf{m}, \mathbf{x}) \} = \int d_{\mathcal{M}}^2(\mathbf{m}, \mathbf{x}) dP(\mathbf{x}),$$

⁹Hermann Karcher. "Riemannian center of mass and mollifier smoothing". In: *Communications on pure and applied mathematics* 30.5 (1977), pp. 509–541.

The algorithm: online estimation

To achieve online detection, we consider using the R-SGD algorithm¹⁰ to address problem (2):

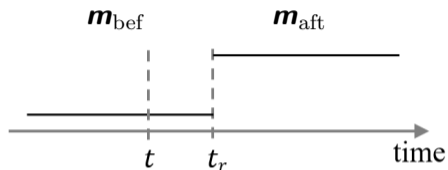
$$\mathbf{m}_{t+1} = \exp_{\mathbf{m}_t} \left(-\alpha H(\mathbf{m}_t, \mathbf{x}_t) \right), \quad (3)$$

where $H(\mathbf{m}, \mathbf{x})$ denotes the unbiased **stochastic gradient** of the loss such that

$$\mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})} \{ H(\mathbf{m}, \mathbf{x}) \} = \int H(\mathbf{m}, \mathbf{x}) dP(\mathbf{x}) = \nabla f(\mathbf{m}).$$

¹⁰Silvère Bonnabel. "Stochastic gradient descent on Riemannian manifolds". In: *IEEE Transactions on Automatic Control* 58.9 (2013), pp. 2217–2229.

The algorithm: an adaptive CPD



To detect change points by monitoring abrupt changes in m ,

- Compute estimates \hat{m}_{bef} and \hat{m}_{aft} ;
- Compare these two quantities using $d_{\mathcal{M}}(\hat{m}_{\text{bef}}, \hat{m}_{\text{aft}})$.

Rationale: The larger the $d_{\mathcal{M}}(\hat{m}_{\text{bef}}, \hat{m}_{\text{aft}})$, the more likely to flag t as a change point.

How to detect change points in an online way?

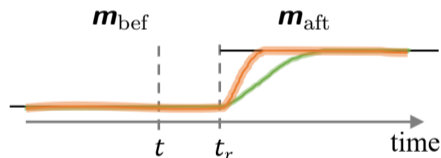
The algorithm: an adaptive CPD

We consider two estimates with two different fixed step sizes $\lambda < \Lambda$ as follows:

$$\mathbf{m}_{\lambda,t+1} = \exp_{\mathbf{m}_{\lambda,t}} \left(-\lambda H(\mathbf{m}_{\lambda,t}, \mathbf{x}_t) \right), \quad (4)$$

$$\mathbf{m}_{\Lambda,t+1} = \exp_{\mathbf{m}_{\Lambda,t}} \left(-\Lambda H(\mathbf{m}_{\Lambda,t}, \mathbf{x}_t) \right). \quad (5)$$

Convergence is directly affected by λ and Λ :

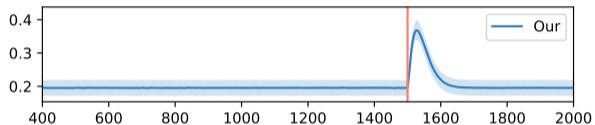


An adaptive CPD statistic is given by:

$$g_t = d_{\mathcal{M}}(\mathbf{m}_{\lambda,t}, \mathbf{m}_{\Lambda,t}). \quad (6)$$

CPD is then performed by comparing g_t to a threshold ξ .

The algorithm: preview of the results



1. Can we provide some performance guarantees?
2. How to determine a detection threshold ξ ?

Theoretical analysis: convergence

The performance guarantee of our statistic g_t is based on a non-asymptotic convergence analysis of the R-SGD algorithm:

Theorem

With some assumptions, for any $s \in \mathbb{N}_*$, the stochastic Riemannian gradient descent algorithm with a **constant step size** α satisfies:

$$\mathbb{E}\{f(\mathbf{m}_s) - f(\mathbf{m}^*)\} \leq \frac{(1 - \epsilon)^{(s-1)} D^2}{2\alpha} + \frac{\alpha\sigma^2}{2\epsilon}, \quad (7)$$

with $\epsilon = \min\left\{\frac{1}{\zeta(\kappa, D)}, \alpha\mu\right\}$ and $\zeta(\kappa, D) = \frac{\sqrt{|\kappa|}D}{\tanh(\sqrt{|\kappa|}D)}$.

Theoretical analysis: performance guarantee

Theorem

Under the null hypothesis H_0 , $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{t-1}$ are drawn i.i.d. from $P(\mathbf{x})$ with the Karcher mean \mathbf{m}^* . With some assumptions, at a steady state, *the false alarm rate* can be upper bounded by:

$$\mathbb{P}(g_t \geq \xi | H_0) \leq \frac{2}{\xi} \left(f(\mathbf{m}^*) + \frac{(\lambda + \Lambda)\sigma^2}{4\epsilon} \right)^{\frac{1}{2}}. \quad (8)$$

with $\epsilon = \min \left\{ \frac{1}{\zeta(\kappa, D)}, \lambda\mu \right\}$ and $\xi > 0$ the detection threshold.

This analysis shows that a higher detection threshold ξ and smaller Karcher variance $f(\mathbf{m}^*)$ make this bound tighter.

Theoretical analysis: performance guarantee

Theorem

Under the alternative hypothesis H_1 , $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{t-B-1}$ are drawn i.i.d. from $P_1(\mathbf{x})$ with Karcher mean \mathbf{m}_1^* , and $\mathbf{x}_{t-B}, \mathbf{x}_{t-B+1}, \dots, \mathbf{x}_{t-1}$ are drawn i.i.d. from $P_2(\mathbf{x})$ with Karcher mean \mathbf{m}_2^* . With some assumptions, **the detection rate** can be lower bounded as:

$$\mathbb{P}(g_t > \xi | H_1) \geq \frac{d_{\mathcal{M}}(\mathbf{m}_1^*, \mathbf{m}_2^*) - \psi(\lambda) - \phi(\Lambda) - \xi}{D - \xi}, \quad (9)$$

where $\psi(\lambda) = \left(2f_{\text{bef}}(\mathbf{m}_1^*) + \frac{\lambda\sigma^2}{\epsilon}\right)^{\frac{1}{2}} + \lambda\rho B$ and $\phi(\Lambda) = \left(2f_{\text{aft}}(\mathbf{m}_2^*) + \frac{(1-\epsilon)^B D^2}{\Lambda} + \frac{\Lambda\sigma^2}{\epsilon}\right)^{\frac{1}{2}}$.

This analysis shows that larger values of $d_{\mathcal{M}}(\mathbf{m}_1^*, \mathbf{m}_2^*)$ and smaller values of ξ , Karcher variances $f_{\text{bef}}(\mathbf{m}_1^*)$ and $f_{\text{aft}}(\mathbf{m}_2^*)$ make this bound tighter.

Adaptive threshold selection

Under the null hypothesis, approximate g_t by a Gaussian distribution, set ξ as an estimate of the q -th quantile of g_t by computing only its first two moments¹¹: $\beta_t^g = (1 - \alpha)\beta_{t-1}^g + \alpha g_t$; $\gamma_t^g = (1 - \alpha)\gamma_{t-1}^g + \alpha g_t^2$; $\hat{\xi}_t = \beta_t^g + \sqrt{\gamma_t^g - (\beta_t^g)^2} \sqrt{2} \operatorname{erf}^{-1}(2q - 1)$.

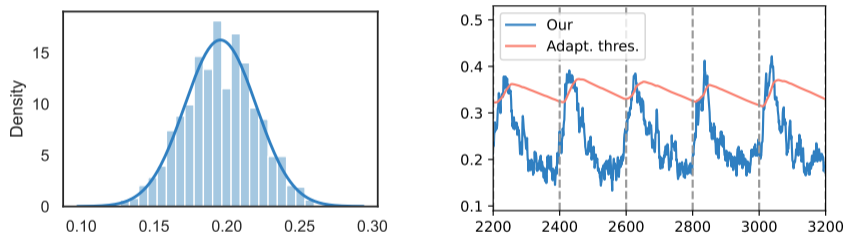


Figure: Distribution of g_t under the null hypothesis (left) and illustration of the adaptive threshold procedure (right).

¹¹Nicolas Keriven et al. "NEWMA: a new method for scalable model-free online change-point detection". In: *IEEE Transactions on Signal Processing* 68 (2020), pp. 3515–3528.

Applications and experiment setups

We apply our strategy to two manifolds as examples:

- The manifold of symmetric positive definite (SPD) matrices: \mathcal{S}_p^{++} ;
- The Grassmann manifold: \mathcal{G}_p^k .

Baselines:

- Scan-B¹², NEWMA¹³ and NODE¹⁴: designed for Euclidean spaces, online;
- F-CPD¹⁵: designed for manifold-valued data, offline.

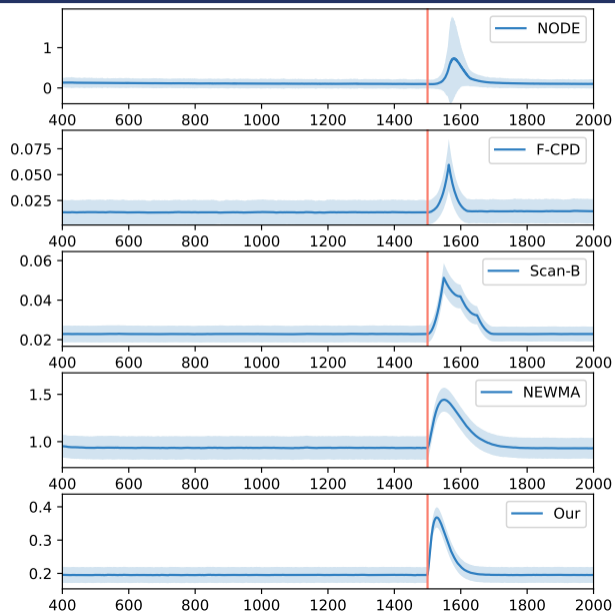
¹²Shuang Li et al. "Scan B-statistic for kernel change-point detection". In: *Sequential Analysis* 38.4 (2019), pp. 503–544.

¹³Nicolas Keriven et al. "NEWMA: a new method for scalable model-free online change-point detection". In: *IEEE Transactions on Signal Processing* 68 (2020), pp. 3515–3528.

¹⁴Xiuheng Wang et al. "Change Point Detection with Neural Online Density-ratio Estimator". In: *IEEE international conference on acoustics, speech and signal processing (ICASSP)*. 2023.

¹⁵Paromita Dubey et al. "Fréchet change-point detection". In: *The Annals of Statistics* 48.6 (2020), pp. 3312–3335.

Experiment with synthetic data on \mathcal{S}_ρ^{++}



Experiment with synthetic data on \mathcal{S}_ρ^{++}

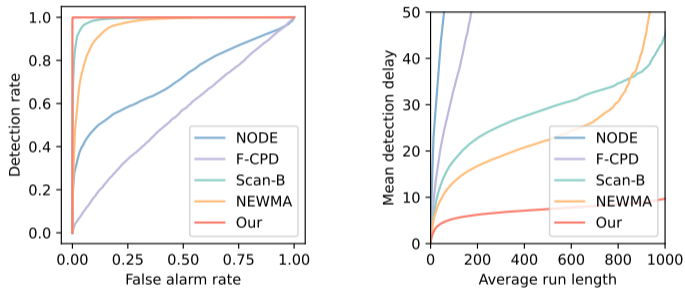
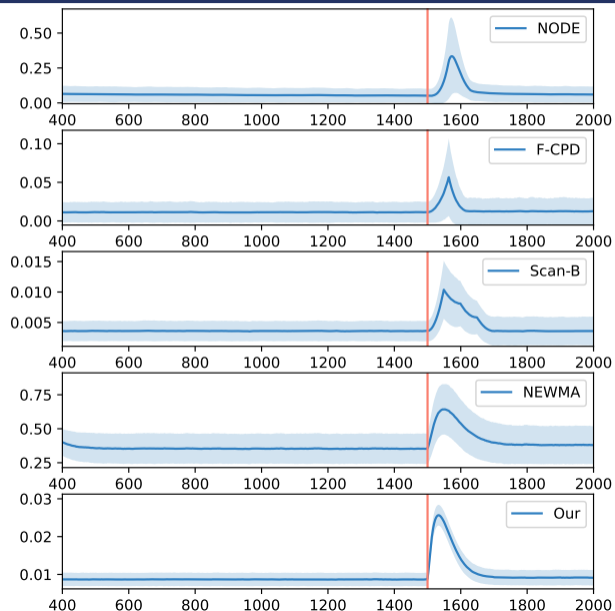


Figure: ROC curves, ARL versus MDD for the compared algorithms.

Experiment with synthetic data on \mathcal{G}_p^k



Experiment with synthetic data on \mathcal{G}_p^k

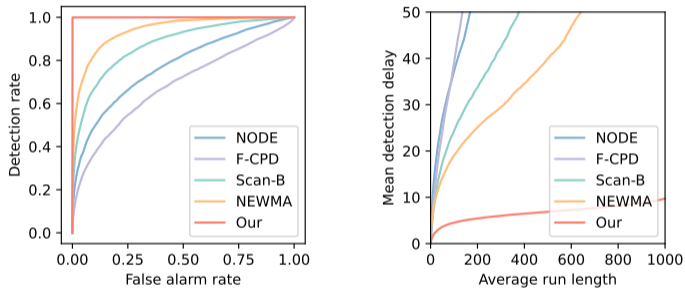


Figure: ROC curves, ARL versus MDD for the compared algorithms.

Voice activity detection

4 seconds of real speech from the TIMIT database¹⁶ was added to 15 seconds of background noises from the QUT-NOISE database¹⁷, with -3 dB Signal-to-Noise Ratio.

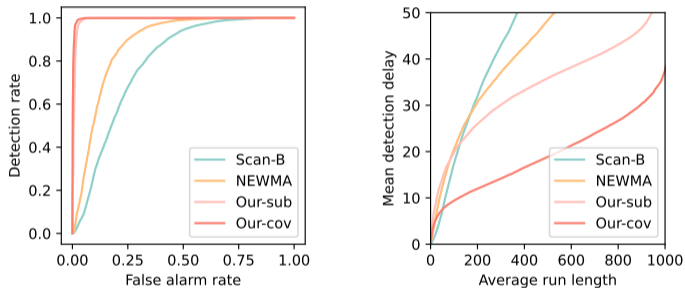


Figure: ROC curves, ARL versus MDD for voice action detection.

¹⁶John S Garofolo. "Timit acoustic phonetic continuous speech corpus". In: *Linguistic Data Consortium, 1993* (1993).

¹⁷David Dean et al. "The QUT-NOISE-TIMIT corpus for evaluation of voice activity detection algorithms". In: *Proceedings of the 11th Annual Conference of the International Speech Communication Association*. International Speech Communication Association. 2010, pp. 3110–3113.

Skeleton-based action recognition

Use the HDM05 motion capture database¹⁸. and generate data points $\mathbf{\Sigma}_t \in \mathcal{S}_p^{++}$ with $p = 93$ by computing the joint covariance descriptor¹⁹ of 3D coordinates of the 31 joints.

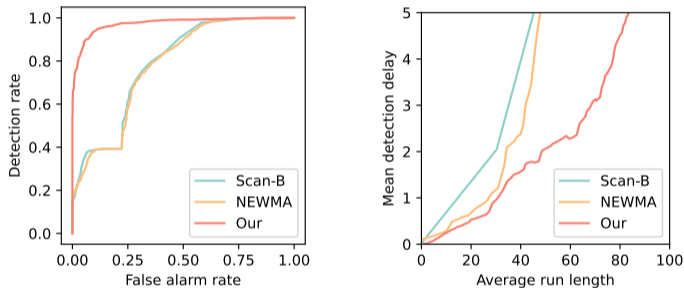


Figure: ROC curves, ARL versus MDD for skeleton-based action recognition.

¹⁸M. Müller et al. *Documentation Mocap Database HDM05*. Tech. rep. CG-2007-2. Universität Bonn, 2007.

¹⁹Mohamed E Hussein et al. "Human action recognition using a temporal hierarchy of covariance descriptors on 3D joint locations". In: *Twenty-third international joint conference on artificial intelligence*. 2013.

Non-parametric Online Change Point Detection on Riemannian Manifolds



Xiuheng Wang, Ricardo Borsoi, Cédric Richard
xiuheng.wang@oca.eu
raborsoi@gmail.com
cedric.richard@unice.fr