



**HAL**  
open science

## Consistent change-point detection with kernels

Damien Garreau, Sylvain Arlot

► **To cite this version:**

Damien Garreau, Sylvain Arlot. Consistent change-point detection with kernels. 2016. hal-01416704v1

**HAL Id: hal-01416704**

**<https://hal.science/hal-01416704v1>**

Preprint submitted on 14 Dec 2016 (v1), last revised 28 Jun 2017 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Consistent change-point detection with kernels

Damien Garreau<sup>1</sup> and Sylvain Arlot<sup>2</sup>

<sup>1</sup> *INRIA ; Sierra Project-Team  
Laboratoire d'Informatique de l'Ecole Normale Supérieure  
(CNRS/ENS/INRIA UMR 8548)  
2 rue Simone Iff, 75012 Paris, France  
e-mail: [damien.garreau@ens.fr](mailto:damien.garreau@ens.fr)*

<sup>2</sup> *Laboratoire de Mathématiques d'Orsay  
Univ. Paris-Sud, CNRS, Université Paris-Saclay,  
91405 Orsay, France  
e-mail: [sylvain.arlot@math.u-psud.fr](mailto:sylvain.arlot@math.u-psud.fr)*

**Abstract:** In this paper we study the kernel change-point algorithm (KCP) proposed by Arlot, Celisse and Harchaoui [2], that aims at locating an unknown number of change-points in the distribution of a sequence of independent data taking values in an arbitrary set. The change-points are selected by model selection with a penalized kernel empirical criterion. We provide a non-asymptotic results showing that, with high probability, the KCP procedure retrieves the correct number of change-points, provided that the constant in the penalty is well-chosen; in addition, KCP estimates the change-points location at the minimax rate  $\log(n)/n$ . As a consequence, when using a characteristic kernel, KCP detects all kinds of change in the distribution (not only changes in the mean or the variance), and it is able to do so for complex structured data (not necessarily in  $\mathbb{R}^d$ ). Most of the analysis is conducted assuming that the kernel is bounded; part of the results can be extended when we only assume a finite second-order moment.

**Keywords and phrases:** change-point detection, kernel methods, penalized least-squares.

Received December 2016.

## 1. Introduction

In many situations, some properties of a time series change over time, such as the mean, the variance or higher-order moments. Change-point detection is the long standing question of finding both the number and the localization of such changes. This is an important front-end task in many applications. For instance, detecting changes occurring in comparative genomic hybridization array data (CGH arrays) is crucial to the early diagnosis of cancer [29]. In finance, some intensively examined time series like the volatility process exhibit local homogeneity and it is useful to be able to segment these time series both for modeling and forecasting [34, 40]. Change-point detection can also be used to detect changes in exchange flows — for instance the exchange rate between the US dollar and other currencies [23] — or changes in a sequence of images [25, 16].

Generally speaking, it is of interest to the practitioner to segment a time series in order to calibrate its model on homogeneous sets of datapoints.

Addressing the change-point problem in practice requires to face several important challenges. First, the number of changes cannot be assumed to be known in advance — and cannot be assumed equal to 0 or 1 —, so a practical change-point procedure must be able to learn the number of changes from data. Second, changes do not always occur in the mean or the variance of the data, as assumed by most change-point procedures. We need to be able to detect changes in other aspects of the distribution. Third, parametric assumptions — which are often made for building or for analyzing change-point procedures — are often unrealistic, so that we need a fully non-parametric approach. Fourth, data points in the time series we want to segment can be high-dimensional and/or structured. If the dimensionality is larger than the number of observations, a non-asymptotic analysis is mandatory for theoretical results to be meaningful. When data are structured — for instance, histograms, graphs or strings —, taking their structure into account seems necessary for detecting efficiently the change-points.

Numerous change-point procedures have been proposed since the seminal works of [17] and [6], both parametric and non-parametric. We refer to [42] for an extensive review and some applications. Recent works also consider the case of high-dimensional data when only a few coordinates of the mean change at each change-point [44, and references therein], and the slightly different problem of detecting gradual changes [43]; this paper does not address these slightly different problems.

Nevertheless, no change-point procedure addressed simultaneously all the challenges mentioned above, until the kernel change-point procedure (KCP) was proposed by [2]. In a few words, KCP mixes the penalized least-squares approach to change-point detection [13, 35] with semi-definite positive kernels [3]. On the computational side, the KCP segmentation can be computed efficiently thanks to a dynamic programming algorithm [22, 2], which can be made even faster [12]. An oracle inequality — which is not exactly a result on change-point estimation, but a guarantee on estimation of the “mean” of the time series in the RKHS associated with the kernel chosen — for KCP is proved by [2], who also show its good numerical performance — in terms of change-point estimation — in several experiments.

So, a key theoretical question remains open: does KCP estimate correctly the number of change-points and their locations with a large probability? If yes, at which speed does KCP estimate the change-point locations?

This paper answers these questions, showing that KCP has good theoretical properties for change-point estimation with independent data, under a boundedness assumption (Theorem 1 in Section 3.1). This result is non-asymptotic, hence meaningful for high-dimensional or complex data. In the asymptotic setting — with a fixed true segmentation and more and more data points observed within each segment —, Theorem 1 implies that KCP estimates consistently all changes in the “kernel mean” of the distribution of data, at speed  $\log(n)/n$  which is the minimax rate with respect to the sample size  $n$ . We also provide

a partial result under a weaker moment assumption (Theorem 2 in Section 3.3) and explain in Section 4 how our proofs could be extended to other settings, including the dependent case.

An important case is when KCP is used with a characteristic kernel [18], such as the Gaussian or the Laplace kernel. Then, any change in the distribution of data induces a change in the “kernel mean”. So, Theorem 1 implies that KCP then estimates consistently and at the minimax rate *all changes* in the distribution of the data, without any parametric assumption and without prior knowledge about the number of changes.

Our results also are interesting regarding to the theoretical understanding of least-squares change-point procedures. Indeed, when KCP is used with the linear kernel, it reduces to previously known penalized least-squares change-point procedures [45, 13, 35, for instance]. There are basically two kinds of results on such procedures in the change-point literature: (i) asymptotic statements on change-point estimation [45, 46, 4, 33] and (ii) non-asymptotic oracle inequalities [13, 35, 2], which are based upon concentration inequalities and model selection theory [7] but not directly provide guarantees on the estimated change-point locations. Our results and their proofs show how to reconcile the two approaches when we are interested in change-point locations, which is already new for the case of the linear kernel, and also holds for a general kernel.

### 1.1. Acknowledgments

The authors thank Alain Celisse and Aymeric Dieuleveut for helpful discussions. Damien Garreau PhD scholarship is financed by DGA / Inria. Sylvain Arlot is also member of the Select project-team of Inria Saclay. At the beginning of this work, Sylvain Arlot was financed by CNRS and member of the Sierra team in the Département d’Informatique de l’École normale supérieure (CNRS / ENS / Inria UMR 8548), 45 rue d’Ulm, 75005 Paris, France. This work was also partly done while Sylvain Arlot was supported by Institut des Hautes Études Scientifiques (IHES, Le Bois-Marie, 35, route de Chartres, 91440 Bures-Sur-Yvette, France).

## 2. Kernel change-point detection

This section describes the change-point problem and the kernel change-point procedure of [2].

### 2.1. Change-point problem

Set  $2 \leq n < +\infty$  and consider  $X_1, \dots, X_n$  independent  $\mathcal{X}$ -valued random variables, where  $\mathcal{X}$  is an arbitrary (measurable) space. The goal of change-point detection is to detect abrupt changes in the distribution of the  $X_i$ s. For any  $D \in \{1, \dots, n\}$  and any integers  $0 = \tau_0 < \tau_1 < \dots < \tau_D = 1$ , we define

the *segmentation*  $\tau := [\tau_0, \dots, \tau_D]$  of  $\{1, \dots, n\}$  as the collection of segments  $\lambda_\ell = \{\tau_{\ell-1} + 1, \dots, \tau_\ell\}$ ,  $\ell \in \{1, \dots, D\}$ . We call *change-points* the right-end of the segments, that is the  $\tau_\ell$ ,  $\ell \in \{1, \dots, D\}$ . We denote by  $\mathcal{T}_n^D$  the set of segmentations with  $D$  segments and  $\mathcal{T}_n := \bigcup_{D=1}^n \mathcal{T}_n^D$  the set of all segmentations of  $\{1, \dots, n\}$ . For any  $\tau \in \mathcal{T}_n$ , we write  $D_\tau$  for the number of segments of  $\tau$ . Figure 1 provides a visual example.



FIG 1. We often represent the segmentations as above. The bullet points stand for the elements of  $\{1, \dots, n\}$ . Here,  $n = 10$ ,  $D_\tau = 3$ ,  $\tau_0 = 0$ ,  $\tau_1 = 3$ ,  $\tau_2 = 7$  and  $\tau_3 = 10$ .

An important example we have in mind is the following.

*Example 1* (Asymptotic setting). Let  $K \geq 1$ ,  $0 = b_0 < b_1 < \dots < b_K < b_{K+1} = 1$  and  $P_1, \dots, P_{K+1}$  some probability distributions on  $\mathcal{X}$  be fixed. Then, for any  $n$  and  $i \in \{1, \dots, n\}$ , we set  $t_i := i/n$  and the distribution of  $X_i$  is  $P_{j(i)}$  where  $j(i)$  is such that  $t_i \in [b_j, b_{j+1})$ . In other words, we have a fixed segmentation of  $[0, 1]$ , given by the  $b_j$ , a fixed distribution over each segment, given by the  $P_j$ , and we observe independent realizations from the distributions at discrete times  $t_1, \dots, t_n$ . The corresponding true change-points in  $\{0, \dots, n\}$  are the  $\lfloor nb_j \rfloor$ ,  $j = 1, \dots, K$ . For  $n$  large enough, it has  $K + 1$  segments. Figure 2 shows an example. Let us emphasize that in this setting,  $n$  going to infinity does not mean that new observations are observed over time. We here consider the change-point problem *a posteriori*: a larger  $n$  means that we have been able to observe the phenomenon of interest with a finer time discretization.

## 2.2. Kernel change-point procedure

Let  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a positive semidefinite kernel, that is, a measurable function such that  $(k(x_i, x_j))_{1 \leq i, j \leq m}$  is positive semidefinite for any  $x_1, \dots, x_m \in \mathcal{X}$  [38]. Classical examples of kernels are given by [2, section 3.2], among which:

- the *linear kernel*:  $k^{\text{lin}}(x, y) = \langle x, y \rangle_{\mathbb{R}^p}$  for  $x, y \in \mathcal{X} = \mathbb{R}^p$ .
- the *polynomial kernel* of order  $d \geq 1$ :  $k_d^{\text{poly}}(x, y) = (\langle x, y \rangle_{\mathbb{R}^p} + 1)^d$  for  $x, y \in \mathcal{X} = \mathbb{R}^p$ .
- the *Gaussian kernel* with bandwidth  $h > 0$ :  $k_h^{\text{G}}(x, y) = \exp[-\|x - y\|^2 / (2h^2)]$  for  $x, y \in \mathcal{X} = \mathbb{R}^p$ .
- the *Laplace kernel* with bandwidth  $h > 0$ :  $k_h^{\text{L}}(x, y) = \exp[-\|x - y\| / (2h^2)]$  for  $x, y \in \mathcal{X} = \mathbb{R}^p$ .
- the  $\chi^2$ -kernel:  $k_{\chi^2}(x, y) = \exp\left(-\frac{1}{2} \sum_{i=1}^p \frac{(x_i - y_i)^2}{x_i + y_i}\right)$  for  $x, y \in \mathcal{X}$  the  $p$ -dimensional simplex.

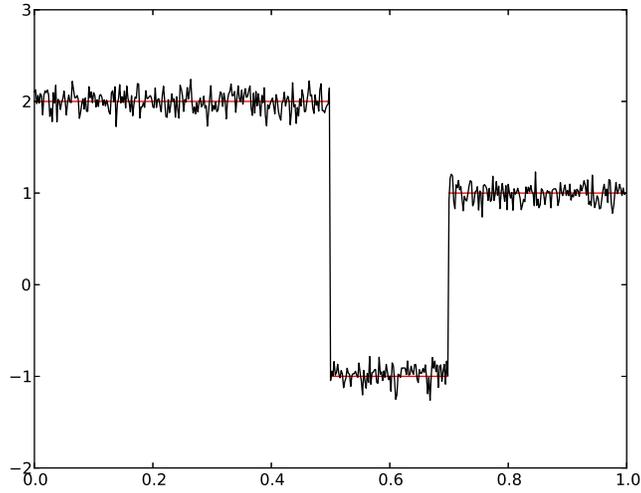


FIG 2. An illustration of the asymptotic setting (Example 1) in the case of changes in the mean of the  $X_i$ . Here,  $\mathcal{X} = \mathbb{R}$ ,  $X_i = f(t_i) + \varepsilon_i$  with  $\varepsilon_1, \dots, \varepsilon_n$  i.i.d. and centered, and  $f : [0, 1] \rightarrow \mathbb{R}$  is a (fixed) piecewise constant function (shown in red). The goal is to recover the number of abrupt changes of  $f$  (here,  $D_{\tau^*} = 2$ ) and their locations ( $b_1 = 0.5$  and  $b_2 = 0.7$ ). Note that other kinds of changes in the distribution of the  $X_i$  can be considered [see 2, for illustrations].

As in [22, 2], for a given segmentation  $\tau \in \mathcal{T}_n^D$ , we assess the adequation of  $\tau$  with the *kernel least-squares criterion*,

$$\widehat{\mathcal{R}}_n(\tau) := \frac{1}{n} \sum_{i=1}^n k(X_i, X_i) - \frac{1}{n} \sum_{\ell=1}^D \left[ \frac{1}{\tau_\ell - \tau_{\ell-1}} \sum_{i=\tau_{\ell-1}+1}^{\tau_\ell} \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} k(X_i, X_j) \right]. \quad (1)$$

Elementary algebra shows that, when  $\mathcal{X} = \mathbb{R}^p$  and  $k = k^{\text{lin}}$ ,  $\widehat{\mathcal{R}}_n$  is the usual least-squares criterion. Minimizing this criterion over the set of all segmentations always outputs the segmentation with  $n$  segments reduced to a point, that is  $[0, \dots, n]$ ; this is a well-known overfitting phenomenon. To counteract this, a classical idea [32, for instance] is to minimize a penalized criterion  $\text{crit}(\tau) := \widehat{\mathcal{R}}_n(\tau) + \text{pen}(\tau)$ , where  $\text{pen} : \mathcal{T}_n \rightarrow \mathbb{R}_+$  is called the penalty. Formally, the kernel change-point procedure (KCP) of [2] selects the segmentation

$$\widehat{\tau} \in \arg \min_{\tau \in \mathcal{T}_n} \{ \text{crit}(\tau) \} \quad \text{where} \quad \text{crit}(\tau) = \widehat{\mathcal{R}}_n(\tau) + \text{pen}(\tau). \quad (2)$$

In the paper, we focus on the classical choice of a penalty proportional to the number of segments:

$$\text{pen}(\tau) = \text{pen}_\ell(\tau) := \frac{CM^2 D_\tau}{n}, \quad (3)$$

similarly to AIC, BIC and  $C_p$  criteria. As mentioned in the introduction, slightly different penalty shapes can be considered, as suggested by [2]. Our results could be extended to the penalty of [2], but we choose to consider the linear penalty (3) only for simplicity.

### 2.3. The reproducing kernel Hilbert space

Let  $\mathcal{H}$  be the reproducing kernel Hilbert space (RKHS) associated with  $k$  [3], together with the canonical feature map  $\Phi : \mathcal{X} \rightarrow \mathcal{H}$

$$\begin{aligned} \Phi &: \mathcal{X} &\rightarrow & \mathcal{H} \\ x &\mapsto & \Phi(x) &:= k(\cdot, x). \end{aligned}$$

We write  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  (resp.  $\|\cdot\|_{\mathcal{H}}$ ) for the inner product (resp. the norm) of  $\mathcal{H}$ . For any  $i \in \{1, \dots, n\}$ , define  $Y_i := \Phi(X_i) \in \mathcal{H}$ . In the case where  $k = k^{\text{lin}}$ , then  $Y_i = X_i$  and the empirical risk  $\widehat{\mathcal{R}}_n$  becomes the least-squares criterion

$$\widehat{\mathcal{R}}_n(\tau) = \frac{1}{n} \sum_{\ell=1}^D \sum_{i=\tau_{\ell-1}+1}^{\tau_\ell} (X_i - \bar{X}_\ell)^2,$$

where  $\bar{X}_\ell$  is the empirical mean of the  $X_i$  over the segment  $\{\tau_{\ell-1} + 1, \dots, \tau_\ell\}$ . It is well-known that penalized least-squares procedures detect changes in the

mean of the observations  $X_i$ , see [45]. Hence the kernelized version of this least-squares procedure, KCP, should detect changes in the “mean” of the  $Y_i = \Phi(X_i)$ , which are a nonlinear transformation of the  $X_i$ .

More precisely, assume that  $\mathcal{H}$  is separable and that

$$\forall i \in \{1, \dots, n\}, \quad \mathbb{E} \left[ \sqrt{k(X_i, X_i)} \right] < +\infty.$$

Then  $\mu_i^*$ , the Bochner integral of  $Y_i$  is well-defined [36]. The condition above is satisfied in our setting (when either Assumption 1 or Assumption 2 holds true), and  $\mathcal{H}$  is separable in most cases [15]. The Bochner integral commutes with continuous linear operators, hence the following property holds, which will be of common use:

$$\forall g \in \mathcal{H}, \quad \langle \mu_i^*, g \rangle_{\mathcal{H}} = \mathbb{E}[g(X_i)] = \mathbb{E}[\langle Y_i, g \rangle_{\mathcal{H}}].$$

We now define the “true segmentation”  $\tau^* \in \mathcal{T}_n$  by

$$\begin{aligned} \mu_1^* = \dots = \mu_{\tau_1^*}^*, \quad \mu_{\tau_1^*+1}^* = \dots = \mu_{\tau_2^*}^*, \quad \dots \quad \mu_{\tau_{D_{\tau^*}^*}-1}^*+1 = \dots = \mu_n^* \\ \text{and} \quad \forall i \in \{1, \dots, D_{\tau^*}^* - 1\}, \quad \mu_{\tau_i^*}^* \neq \mu_{\tau_{i+1}^*}^* \end{aligned} \quad (4)$$

with  $1 \leq \tau_1^* < \dots < \tau_{D_{\tau^*}^*}^* \leq n$ . We call the  $\tau_i^*$ s the *true* change-points. It should be clear that it is always possible to define  $\tau^*$ .

A kernel is said to be characteristic if the mapping  $P \mapsto \mathbb{E}_{X \sim P} [\Phi(X)]$  is injective, for  $P$  belonging to the set of Borel probability measures on  $\mathcal{X}$  [41]. In simpler terms, when  $k$  is a characteristic kernel,  $X_i$  and  $X_{i+1}$  have the same distribution if and only if  $\mu_i^* = \mu_{i+1}^*$ , and  $\tau^*$  indeed corresponds to the set of changes in the distribution of the  $X_i$ . For instance, all strictly positive definite kernels are characteristic, including the Gaussian kernel, see [41]. Therefore, in the setting of Example 1,  $D_{\tau^*} = K + 1$  for  $n$  large enough and  $\tau_\ell^* = \lfloor nb_\ell \rfloor$  for  $\ell = 1, \dots, K$ .

For a general kernel, some changes of  $P_{X_i}$ , the distribution of  $X_i$ , might not appear in  $\tau^*$ . For instance, with the linear kernel,  $\tau^*$  only corresponds to changes of the mean of the  $X_i$ . In most cases, a characteristic kernel is known and we can choose to use KCP with a characteristic kernel; then, as we prove in the following, KCP eventually detects any change in the distribution of the observations. But one can also choose a non-characteristic kernel on purpose, hence focusing only in some changes in the distribution of the  $X_i$ . For instance, the polynomial kernel of order  $d$  is not characteristic and leads to the detection of changes in the first  $d$  moments of the distribution; with the linear kernel, KCP detects changes in the mean of the  $X_i$ .

From now on, we focus on the problem of detecting the changes of  $\tau^*$  only, whether the kernel is characteristic or not.

#### 2.4. Rewriting the empirical risk

It is convenient to define  $Y := (Y_1, \dots, Y_n) \in \mathcal{H}^n$ ,  $\mu^* := (\mu_1^*, \dots, \mu_n^*) \in \mathcal{H}^n$ , and  $\varepsilon := Y - \mu^* \in \mathcal{H}^n$ . We identify the elements of  $\mathcal{H}^n$  with the set of the

applications from  $\{1, \dots, n\}$  to  $\mathcal{H}$ , naturally embedded with the inner product and norm given by

$$\forall x, y \in \mathcal{H}^n, \quad \langle x, y \rangle := \sum_{i=j}^n \langle x_j, y_j \rangle_{\mathcal{H}} \quad \text{and} \quad \|x\|^2 := \sum_{j=1}^n \|x_j\|_{\mathcal{H}}^2.$$

We now rewrite the empirical risk (1) as a function of  $\tau$  and  $Y$ . For any segmentation  $\tau \in \mathcal{T}_n$ , define  $F_\tau$  the set of functions from  $\{1, \dots, n\}$  to  $\mathcal{H}$  that are constant over the segments of  $\tau$ . We see  $F_\tau$  as a subspace of  $\mathcal{H}^n$  as a vector space. Take  $f \in \mathcal{H}^n$ , we define  $\Pi_\tau f$  the orthogonal projection of  $f$  onto  $F_\tau$  with respect to  $\|\cdot\|$ :

$$\Pi_\tau f \in \arg \min_{g \in F_\tau} \|f - g\|.$$

It is shown in [2] that for any  $f \in \mathcal{H}^n$ ,

$$\forall 1 \leq \ell \leq D_\tau, \forall i \in \{\tau_{\ell-1}+1, \dots, \tau_\ell\}, \quad (\Pi_\tau f)_i = \frac{1}{|\tau_\ell - \tau_{\ell-1}|} \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} f_j. \quad (5)$$

We are now able to write the empirical risk as

$$\widehat{\mathcal{R}}_n(\tau) = \frac{1}{n} \|Y - \widehat{\mu}_\tau\|^2, \quad (6)$$

where  $\widehat{\mu}_\tau = \Pi_\tau Y$ , following [22, 2].

## 2.5. Assumptions

A key ingredient of our analysis is the concentration of  $\varepsilon$ . Intuitively, the performance of KCP is better when  $\varepsilon$  concentrates strongly around its mean, since without noise we are just given the task to segment a piecewise-constant signal. It is thus natural to make assumptions on  $\varepsilon$  in order to obtain concentration results. We actually formulate assumptions on the kernel  $k$ , which translate automatically onto  $\varepsilon$ .

As in [2], the main hypothesis used in our analysis is the following.

**Assumption 1.** A positive constant  $M$  exists such that

$$\forall i \in \{1, \dots, n\}, \quad k(X_i, X_i) \leq M^2 < +\infty \quad \text{a.s.}$$

If Assumption 1 holds true,

$$\forall i \in \{1, \dots, n\}, \quad \|Y_i\|_{\mathcal{H}} = \sqrt{k(X_i, X_i)} \leq M \quad \text{a.s.}$$

and [2] shows that  $\|\varepsilon_i\|_{\mathcal{H}} \leq 2M$  almost surely.

Assumption 1 is always satisfied for a large class of commonly used kernels, such as the Gaussian, Laplace and  $\chi^2$  kernels.

Note that Assumption 1 is weaker than assuming  $k$  to be bounded — that is,  $k(x, x) \leq M$  for any  $x \in \mathcal{X}$ , which is equivalent to  $k(x, x') \leq M$  for any

$x, x' \in \mathcal{X}$  since  $k$  is positive definite. For instance, if  $\mathcal{X} = \mathbb{R}^p$  and the data are bounded almost surely, Assumption 1 holds true for the linear kernel and all polynomial kernels, which are not bounded on  $\mathbb{R}^p$ .

In the setting of Example 1, Assumption 1 holds true when

$$\forall j \in \{1, \dots, K\}, \quad k(x, x) \leq M^2 \quad \text{for } P_j\text{-a.e. } x \in \mathcal{X} .$$

It is sometimes possible to weaken Assumption 1 into a finite variance assumption

**Assumption 2.** A positive constant  $V$  exists such that

$$\max_{1 \leq i \leq n} \mathbb{E} \left[ \|\varepsilon_i\|_{\mathcal{H}}^2 \right] \leq V .$$

Since  $v_i = \mathbb{E} [k(X_i, X_i)] - \|\mu_i^*\|_{\mathcal{H}}^2$ , Assumption 2 holds true when

$$\forall i \in \{1, \dots, n\}, \quad \mathbb{E} [k(X_i, X_i)] \leq V .$$

As a consequence, Assumption 1 implies Assumption 2 with  $V = M^2$ . Note that Assumption 2 is satisfied for the polynomial kernel of order  $d$  provided that

$$\forall i \in \{1, \dots, n\}, \quad \mathbb{E} \left[ \|X_i\|^{2d} \right] < +\infty .$$

In the setting of Example 1, Assumption 2 holds true with

$$V = \max_{1 \leq k \leq K+1} \mathbb{E}_{X \sim P_k} [k(X, X)] ,$$

provided this maximum is finite.

### 3. Theoretical guarantees for KCP

We are now able to state our main results. The following section is organized as follows. In Section 3.1, we state the main result of the paper, Theorem 1, which provides simple conditions under which KCP recovers the correct number of segments and localizes the true change-points with high probability, under the bounded kernel Assumption 1. Section 3.2 details a few classical losses between segmentations which can be considered in addition to the one used in Theorem 1. Corollary 1 formulates a result on  $\hat{\tau}$  in terms of the Frobenius loss. Finally, Section 3.3 states a partial result on KCP — requiring the number of change-points  $D_{\tau^*}$  to be known — under the weaker Assumption 2.

#### 3.1. Main result

We need first to define some quantities. The size of the smallest jump of  $\mu^*$  in  $\mathcal{H}$  is

$$\underline{\Delta} := \min_{i / \mu_i^* \neq \mu_{i+1}^*} \|\mu_i^* - \mu_{i+1}^*\|_{\mathcal{H}} . \quad (7)$$

Intuitively, the higher  $\underline{\Delta}$  is, the easier it is to detect the smallest jump with our procedure. In the scalar setting (with the linear kernel), the ratio  $\underline{\Delta}/\sigma$  (where  $\sigma^2$  is the variance of the noise) is called the *signal to noise ratio* [5] and is often used as a measure of the magnitude of a change in the signal. In Example 1,

$$\underline{\Delta} = \min_{1 \leq j \leq K} \left\| \mu_{P_j}^* - \mu_{P_{j+1}}^* \right\|_{\mathcal{H}}$$

where  $\mu_{P_j}^*$  denotes the (Bochner) expectation of  $\Phi(X)$  when  $X \sim P_j$ .

For any  $\tau \in \mathcal{T}_n$ , we denote the (normalized) sizes of its smallest and of its largest segment by

$$\underline{\Lambda}_\tau := \frac{1}{n} \min_{1 \leq \ell \leq D_\tau} |\tau_\ell - \tau_{\ell-1}| \quad \text{and} \quad \bar{\Lambda}_\tau := \frac{1}{n} \max_{1 \leq \ell \leq D_\tau} |\tau_\ell - \tau_{\ell-1}|. \quad (8)$$

It should be clear that  $\underline{\Lambda}_{\tau^*}$  should not be too small; otherwise, the corresponding segments of  $\tau^*$  might be undetectable by any change-point procedure [11]. In Example 1,

$$\underline{\Lambda}_\tau \xrightarrow{n \rightarrow +\infty} \min_{0 \leq j \leq K} |b_{j+1} - b_j| \quad \text{and} \quad \bar{\Lambda}_\tau \xrightarrow{n \rightarrow +\infty} \max_{0 \leq j \leq K} |b_{j+1} - b_j|.$$

For any  $\tau^1$  and  $\tau^2 \in \mathcal{T}_n$ , we define

$$d_\infty^{(1)}(\tau^1, \tau^2) := \max_{1 \leq i \leq D_{\tau^1} - 1} \left\{ \min_{1 \leq j \leq D_{\tau^2} - 1} |\tau_i^1 - \tau_j^2| \right\},$$

which is a loss function (a measure of dissimilarity) between the segmentations  $\tau^1$  and  $\tau^2$ . Note that  $d_\infty^{(1)}$  is not a true distance; other possible losses between segmentations and their relationship with  $d_\infty^{(1)}$  are discussed in Section 3.2.

**Theorem 1.** *Suppose that Assumption 1 holds true. For any  $y > 0$ , an event  $\Omega$  of probability at least  $1 - e^{-y}$  exists on which the following holds true. For any  $C > 0$ , let  $\hat{\tau}$  be defined as in Eq. (2) with  $\text{pen}_\ell$  defined by Eq. (3). Set*

$$C_{\min} := \frac{74}{3}(D_{\tau^*} + 1)(y + \log n + 1) \quad \text{and} \quad C_{\max} := \frac{\underline{\Delta}^2}{M^2} \frac{\underline{\Lambda}_{\tau^*}}{6D_{\tau^*}} n.$$

Then, if

$$C_{\min} < C < C_{\max}, \quad (9)$$

on  $\Omega$ , we have

$$D_{\hat{\tau}} = D_{\tau^*} \quad \text{and} \quad \frac{1}{n} d_\infty^{(1)}(\tau^*, \hat{\tau}) \leq v_1(y) := \frac{148D_{\tau^*}M^2}{\underline{\Delta}^2} \cdot \frac{y + \log n + 1}{n}.$$

Theorem 1 is proved in Section 5.4. Some remarks follow.

Theorem 1 is a non-asymptotic result: it is valid for any  $n \geq 1$  and there is nothing hidden in  $o(1)$  remainder terms. The latter point is crucial for complex data — for instance,  $\mathcal{X} = \mathbb{R}^p$  with  $p > n$  — since in this case, assuming  $\mathcal{X}$  fixed while  $n \rightarrow +\infty$  is not realistic.

Nevertheless, it is useful to write down what Theorem 1 becomes in the asymptotic setting of Example 1. As previously noticed,  $D_{\tau^*}$ ,  $\underline{\Delta}_{\tau^*}$ ,  $\underline{\Delta}^2$  and  $M^2$  then converge to positive constants as  $n \rightarrow +\infty$ . Therefore,  $C_{\min}$  is of order  $\log(n)$ ,  $C_{\max}$  is of order  $n$  and we always have  $C_{\min} < C_{\max}$  for  $n$  large enough. The upper bound on  $C$  matches classical asymptotic conditions for variable selection [39]. The necessity of taking  $C$  of order at least  $\log(n)$  is shown by [8] in a variable selection setting, which includes change-point detection as a particular example; [8, 1] provide several arguments for the optimality of taking a constant  $C$  of order  $\log(n)$ . When  $C$  satisfies (9), the result of Theorem 1 implies that  $\mathbb{P}(D_{\hat{\tau}} = D_{\tau^*}) \rightarrow 1$ . For the linear kernel in  $\mathbb{R}^d$ , this is a well-known result when the distribution of the  $X_i$  changes only through its mean. The first result dates back to [45, Section 2] for a Gaussian noise, later extended by [37] and [4, Section 3.1] under mixingale hypothesis on the error, and [33] under very mild assumptions satisfied for a large family of zero-mean processes [Section 2.1 33, for the precise statement of the hypothesis, see]. Theorem 1 also shows that the estimated change-points of  $\hat{\tau}$  converge towards the true change-points at speed  $\log(n)/n$ .

This speed matches minimax lower bounds [11, section 4, assuming  $D_{\tau^*} > 2$  and Gaussian noise] and has been obtained previously for various change-point procedures [9, for instance] including least-squares [33]. In the case where  $D_{\tau^*} = 1$  and the change-point is known to be bounded away from 1 and  $n$ , this rate becomes  $n^{-1}$  [28, 27]. This last assumption is in fact unnecessary, as it is shown in [11].

Note finally that KCP also performs well for finite samples, according to the simulation experiments of [2].

Theorem 1 emphasizes the key role of  $\underline{\Delta}^2/M^2$ , which can be seen as a generalization of the signal-to-noise ratio, for the change-point detection performance of KCP. The larger is this ratio, the easier it is to have (9) satisfied and the smaller is  $v_1(y)$ . This suggests to choose  $k$  (theoretically at least) by maximizing  $\underline{\Delta}^2/M^2$ , as we discuss in Section 4. Note that  $\underline{\Delta}^2/M^2$  is invariant by a rescaling of  $k$ , hence the result of Theorem 1 is unchanged when  $k$  is rescaled.

The hypothesis (9) is actually three-fold. First, we use that  $C > C_{\max}$  to get  $D_{\hat{\tau}} \leq D_{\tau^*}$ . We have to assume  $C$  large enough since a too small penalty leads to selecting (with KCP or any other penalized least-squares procedure) the segmentation with  $n$  segments, that is  $D_{\hat{\tau}} = n$ . Second, in the same way,  $C < C_{\max}$  is used to get  $D_{\hat{\tau}} \geq D_{\tau^*}$ . Such an assumption is required since taking a penalty function too large in Eq. (2) would result in selecting the segmentation with only one segment, that is,  $D_{\hat{\tau}} = 1$ . Third,  $C_{\max}$  has to be greater than  $C_{\min}$  for providing a non-empty interval of possible values for  $C$ . This is also used in the proof to obtain the upper bound on  $d_{\infty}^{(1)}(\tau^*, \hat{\tau})$  when we already know that

$D_{\hat{\tau}} = D_{\tau^*}$ . In the asymptotic setting, the  $C_{\min} < C_{\max}$  hypothesis translates into  $\underline{\Delta}_{\tau^*} \succ \log(n)/n$ . That is, the size of the smallest segment has to be of order  $\log n/n$ . This is known to be a necessary condition to obtain the minimax rate in multiple change-point detection [11, section 2].

Theorem 1 helps choosing  $C$ , which is a key parameter of KCP, as in any penalized model selection procedure. However, in practice, we do not recommend to directly use Eq. (2) for choosing  $C$  for two reasons:  $C_{\min}, C_{\max}$  depend on unknown quantities  $D_{\tau^*}, \underline{\Delta}_{\tau^*}, \underline{\Delta}$ , and the exact values of the constants in  $C_{\min}, C_{\max}$  might be pessimistic compared to what we can observe from simulation experiments. We rather suggest to use a data-driven method for choosing  $C$ , see section 4.

If we know  $D_{\tau^*}$ , we can replace  $\hat{\tau}$  by

$$\hat{\tau}(D_{\tau^*}) \in \arg \min_{\tau \in \mathcal{T}_n^{D_{\tau^*}}} \{\widehat{\mathcal{R}}_n(\tau)\}.$$

Then, assuming that  $\underline{\Delta}_{\tau^*} > v_1(y)$  — which is weaker than assuming  $C_{\min} < C_{\max}$  —, the proof of Theorem 1 shows that, on  $\Omega$ , we have

$$\frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}) \leq v_1(y).$$

### 3.2. Loss functions between segmentations

Theorem 1 shows that  $\hat{\tau}$  is close to  $\tau^*$  in terms of  $d_{\infty}^{(1)}$ . Several other loss functions (measures of dissimilarity) can be defined between segmentations [24]. We here consider a few of them, which are often used or natural for the change-point problem.

Let us first consider losses related to the Hausdorff distance. For any  $\tau^1$  and  $\tau^2 \in \mathcal{T}_n$ , we define

$$\begin{aligned} d_{\infty}^{(1)}(\tau^1, \tau^2) &:= \max_{1 \leq i \leq D_{\tau^1} - 1} \left\{ \min_{1 \leq j \leq D_{\tau^2} - 1} |\tau_i^1 - \tau_j^2| \right\} \\ d_{\infty}^{(2)}(\tau^1, \tau^2) &:= \max_{1 \leq i \leq D_{\tau^1} - 1} \left\{ \min_{0 \leq j \leq D_{\tau^2}} |\tau_i^1 - \tau_j^2| \right\} \\ d_{\text{H}}^i(\tau^1, \tau^2) &:= \max\{d_{\infty}^i(\tau^1, \tau^2), d_{\infty}^i(\tau^2, \tau^1)\} \quad \text{for } i \in \{1, 2\}. \end{aligned}$$

Whenever  $D_{\tau^1} = D_{\tau^2}$ , we define

$$d_{\infty}^{(3)}(\tau^1, \tau^2) := \max_{1 \leq i \leq D_{\tau^1} - 1} |\tau_i^1 - \tau_i^2|.$$

Note that  $d_{\infty}^{(3)}$  is symmetric thus there is no need to define  $d_{\text{H}}^3$ . One could also define  $d_{\text{H}}^i$  as the *Hausdorff distance* between the subsets  $\{\tau_1^1, \dots, \tau_{D_{\tau^1} - 1}^1\}$  and

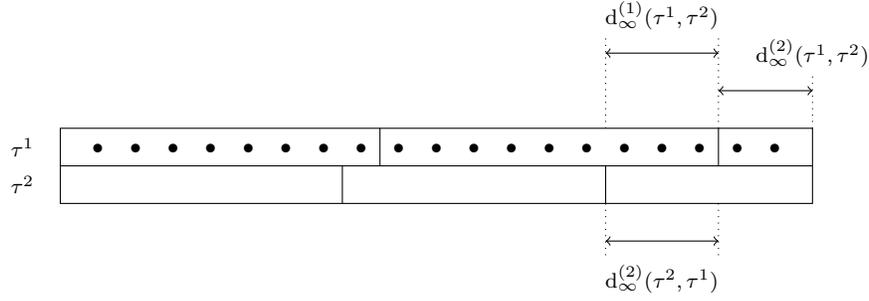


FIG 3. Illustration of the definition of  $d_\infty^i$ , with  $n = 19$ ,  $\tau^1 = [0, 8, 17, 19]$  and  $\tau^2 = [0, 7, 14, 19]$ . In this example,  $D_{\tau^1} = D_{\tau^2} = 3$ . We can compute  $d_\infty^{(1)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^2, \tau^1) = d_\infty^{(2)}(\tau^2, \tau^1) = d_\infty^{(3)}(\tau^1, \tau^2) = 3$  and  $d_\infty^{(2)}(\tau^1, \tau^2) = 2$ .

$\{\tau_1^2, \dots, \tau_{D_{\tau^2}-1}^2\}$  (for  $i = 2$ ) associated to the distance  $\delta(x, y) = |x - y|$  on  $\mathbb{R}$ . These definitions are illustrated by Figure 3.

Interestingly, all these loss functions coincide whenever  $n^{-1}d_\infty^{(1)}(\tau^1, \tau^2)$  is small enough. The following lemma makes this claim rigorous.

**Lemma 1.** *We have the following two properties.*

(i) *For any  $\tau^1, \tau^2 \in \mathcal{T}_n$  such that*

$$\frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) < \frac{1}{2} \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\},$$

*we have  $D_{\tau^1} = D_{\tau^2}$  and*

$$d_\infty^{(1)}(\tau^1, \tau^2) = d_\infty^{(2)}(\tau^1, \tau^2) = d_\infty^{(3)}(\tau^1, \tau^2) = d_H^{(1)}(\tau^1, \tau^2) = d_H^{(2)}(\tau^1, \tau^2).$$

(ii) *For any  $\tau^1, \tau^2 \in \mathcal{T}_n$  such that*

$$D_{\tau^1} = D_{\tau^2} \quad \text{and} \quad \frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) < \frac{\underline{\Lambda}_{\tau^1}}{2},$$

*we have*

$$d_\infty^{(1)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^2, \tau^1) = d_H^{(1)}(\tau^1, \tau^2).$$

Lemma 1 is proved in section A.1. As a direct application of Lemma 1 we see that the statement of Theorem 1 holds true with  $d_\infty^{(1)}$  replaced by *any* of the loss functions that we defined above, at least for  $n$  large enough.

Another loss between segmentations is the *Frobenius* loss [31], which is defined as follows. For any  $\tau^1, \tau^2 \in \mathcal{T}_n$ ,

$$d_F(\tau^1, \tau^2) := \|\Pi_{\tau^1} - \Pi_{\tau^2}\|_F,$$

where  $\Pi_\tau$  is the orthogonal projection onto  $F_\tau$ , as defined in Section 2.4, and  $\|\cdot\|_F$  denotes the Frobenius norm of a matrix:

$$\forall A \in \mathbb{R}^{N \times M}, \quad \|A\|_F^2 := \sum_{i=1}^N \sum_{j=1}^M A_{ij}^2.$$

A closed-form formula for  $d_F$  can be derived from the matrix representation of  $\Pi_\tau$  that is given by Eq. (5):

$$\forall 1 \leq i, j \leq n, \quad (\Pi_\tau)_{i,j} = \begin{cases} \frac{1}{|\lambda|} & \text{if } i \text{ and } j \text{ belong to the same segment } \lambda \text{ of } \tau \\ 0 & \text{otherwise.} \end{cases}$$

An interesting feature of the Frobenius loss is that it is smaller than one only when  $\tau^1$  and  $\tau^2$  have the same number of segments, whereas Hausdorff distances can be small with very different numbers of segments. Indeed, we prove in Section A.2 that

$$|D_{\tau^1} - D_{\tau^2}| \leq d_F(\tau^1, \tau^2)^2 \leq D_{\tau^1} + D_{\tau^2}. \quad (10)$$

The next proposition shows that there is an equivalence (up to constants) between the Hausdorff and Frobenius losses between segmentations, provided that they are close enough.

**Proposition 1.** *Suppose that  $D_{\tau^1} = D_{\tau^2}$  and  $\frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) < \underline{\Lambda}_{\tau^1}/2$ , then*

$$(d_F(\tau^1, \tau^2))^2 \leq \frac{12D_{\tau^1}}{\underline{\Lambda}_{\tau^1}} \frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2).$$

*If in addition  $\frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) < \underline{\Lambda}_{\tau^1}/3$ , then*

$$\frac{2}{3\underline{\Lambda}_{\tau^1}} \frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) \leq (d_F(\tau^1, \tau^2))^2.$$

Proposition 1 was first stated and proved by [30, Theorem B.2]. We prove it in Section A.2 for completeness.

As a corollary of Theorem 1 and Proposition 1, we get the following guarantee on the Frobenius loss between  $\tau^*$  and the segmentation  $\hat{\tau}$  estimated by KCP.

**Corollary 1.** *Under the assumptions of Theorem 1, on the event  $\Omega$  defined by Theorem 1, for any  $\hat{\tau}$  satisfying Eq. (2) with  $\text{pen}_\ell$  defined by Eq. (3), we have:*

$$d_F(\tau^*, \hat{\tau}) \leq \frac{43D_{\tau^*}}{\underline{\Lambda}_{\tau^*}} \cdot \frac{M}{\underline{\Delta}} \sqrt{\frac{y + \log n + 1}{n}}.$$

Note that Corollary 1 gives a better result (at least for large  $n$ ) than the obvious bound

$$d_F(\tau^*, \hat{\tau}) \leq D_{\tau^*} + D_{\hat{\tau}} - 2.$$

*Proof.* On the event  $\Omega$ , we have  $\frac{1}{n}d_\infty^{(1)}(\tau^*, \hat{\tau}) < \underline{\Delta}_{\tau^*}/(D_{\tau^*} + 1)$  and  $D_{\tau^*} = D_{\hat{\tau}}$ . Therefore, according to Proposition 1,

$$(d_F(\tau^*, \hat{\tau}))^2 \leq \frac{12D_{\tau^*}}{\underline{\Delta}_{\tau^*}} \frac{1}{n} d_\infty^{(1)}(\tau^*, \hat{\tau}) = \frac{1776D_{\tau^*}^2(y + \log n + 1)}{n\underline{\Delta}_{\tau^*}} \cdot \frac{M^2}{\underline{\Delta}^2}.$$

□

Up to this point, we assessed the quality of the segmentation  $\tau$  by considering the proximity of  $\tau$  with  $\tau^*$ . Another natural idea is to measure the distance between  $\mu^*$  and  $\mu_\tau^*$  in  $\mathcal{H}^n$ . It is closely related to the oracle inequality proved by [2], which implies an upper bound on  $\|\mu^* - \hat{\mu}_{\hat{\tau}}\|^2$ . We can also observe that there is a simple relationship between  $\|\mu^* - \mu_\tau^*\|^2$  and the Frobenius distance between  $\tau$  and  $\tau^*$ . Indeed,

$$\|\mu^* - \mu_\tau^*\|^2 = \|(\Pi_{\tau^*} - \Pi_\tau)\mu^*\|^2 \leq \|\Pi_{\tau^*} - \Pi_\tau\|_2^2 \|\mu^*\|^2 \leq (d_F(\tau^*, \hat{\tau}))^2 \|\mu^*\|^2. \quad (11)$$

Eq. (23) in the proof of Theorem 1 shows that on  $\Omega$ , under the assumptions of Theorem 1,

$$\|\mu^* - \mu_\tau^*\|^2 \leq 74(y + \log(n) + 1)D_{\tau^*}M^2$$

which is slightly better (but similar) to what Corollary 1, Eq. (11) and the bound  $\|\mu^*\|^2 \leq M^2n$  imply.

### 3.3. Extension to the finite variance case

Theorem 1 is valid under a boundedness assumption (Assumption 1). What happens under the weaker Assumption 2? As a first step, we provide a result for

$$\hat{\tau}(D_{\tau^*}, \delta_n) \in \arg \min_{\tau \in \mathcal{T}_n^{D_{\tau^*}} / \underline{\Delta}_\tau \geq \delta_n} \{\hat{\mathcal{R}}_n(\tau)\} \quad (12)$$

for some  $\delta_n > 0$ . In other words, we restrict our search to segmentations  $\tau$  of the correct size — hence  $D_{\tau^*}$  must be known *a priori* — and having no segment with less than  $n\delta_n$  observations. Note that the dynamic programming algorithm of [22, 2] can be used for computing  $\hat{\tau}(D_{\tau^*}, \delta_n)$  efficiently.

We discuss how to relax this restriction right after the statement of Theorem 2. Similarly to  $\underline{\Delta}$ , we define  $\bar{\Delta} := \max_i \|\mu_i^* - \mu_{i+1}^*\|_{\mathcal{H}}$ .

**Theorem 2.** *Suppose that Assumption 2 holds true. For any  $\delta_n, y > 0$ , define:*

$$v_2(y, \delta_n) := 24(D_{\tau^*})^2 \frac{\bar{\Delta}\sqrt{V}}{\underline{\Delta}^2} \frac{y}{\sqrt{n}} + 8D_{\tau^*} \frac{V}{\underline{\Delta}^2} \frac{y^2}{n\delta_n}.$$

*For any  $y > 0$ , an event  $\Omega_2$  exists such that*

$$\mathbb{P}(\Omega_2) \geq 1 - \frac{1}{y^2}$$

and, on  $\Omega_2$ , we have the following: for any  $\delta_n \in (0, \underline{\Delta}_{\tau^*}]$  and any  $\hat{\tau}(D_{\tau^*}, \delta_n)$  satisfying Eq. (12), if  $v_2(y, \delta_n) \leq \underline{\Delta}_{\tau^*}$ ,

$$\frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}(D_{\tau^*}, \delta_n)) \leq v_2(y, \delta_n). \quad (13)$$

Theorem 2 is proved in Section 5.5. Let us make a few remarks.

As for Theorem 1, our result is non-asymptotic. However, it is interesting to write it down in the setting of Example 1. If  $n$  goes to infinity, then the assumption  $\underline{\Delta}_{\tau^*} \geq \delta_n$  is satisfied whenever  $\delta_n \rightarrow 0$ . If we furthermore require that  $n\delta_n \rightarrow \infty$ , then Eq. (13) implies that

$$\frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}(D_{\tau^*}, \delta_n)) \xrightarrow[n \rightarrow +\infty]{(p)} 0,$$

by taking a well-chosen  $y$  of order  $\sqrt{n} + \sqrt{n\delta_n}$ . In the particular case of the linear kernel, this result is known under various hypothesis [33, for instance].

More precisely, if we take  $\delta_n = n^{-1/2}$ , Theorem 2 implies that

$$\frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}(D_{\tau^*}, n^{-1/2}))$$

goes to zero at least as fast as  $\ell_n/\sqrt{n}$ , where  $(\ell_n)_{n \geq 1}$  is any sequence tending to infinity, for instance  $\ell_n = \log(n)$ . This speed seems suboptimal compared to previous results [33, for instance] — which do not consider the case of a general kernel —, but we could not prove tight enough deviation bounds for getting the localization rate  $\log(n)/n$  under Assumption 2.

How does Theorem 2 compares to Theorem 1? First, as noticed by Remark 5 in Section 5.4, the result of Theorem 1 also holds true for  $\hat{\tau}(D_{\tau^*}, \delta_n)$  as long as  $\underline{\Delta}_{\tau^*} \geq \delta_n$ . Second,  $v_1(y)$  is usually smaller than  $v_2(y, \delta_n)$  — its order of magnitude is smaller when  $n \rightarrow +\infty$  —, and the lower bound on the probability of  $\Omega$  is better than the one for  $\Omega_2$ . There is no surprise here: the stronger Assumption 1 helps us proving a stronger result for  $\hat{\tau}(D_{\tau^*}, \delta_n)$ . Nevertheless, these only are upper bounds, so we do not know whether the performance of  $\hat{\tau}(D_{\tau^*}, \delta_n)$  actually changes much depending on the noise assumption. For instance, as already noticed, we do not believe that the localization speed  $\log(n)/n$  requires a boundedness assumption; in particular cases at least, it has been obtained for unbounded data [33, 9].

The dependency in  $k$  of the speed of convergence of  $\hat{\tau}(D_{\tau^*}, \delta_n)$  is slightly less clear than in Theorem 1. The signal-to-noise ratio here appears through  $\underline{\Delta}^2/V$ , as expected, but the size  $\bar{\Delta}$  of the largest true jump also appears in  $v_2$ . At the very least, it is clear that  $\underline{\Delta}^2/V$  should not be too small.

As noted by [33], it may be possible to get rid of the minimal segment length  $\delta_n$ , either by imposing stronger conditions on  $\varepsilon$  — which are not met in our setting — or by constraining the values of  $\hat{\mu}$  to lie in a compact subset  $\Theta \subset \mathcal{H}^{D_{\tau^*}+1}$ .

#### 4. Discussion

Before proving our main results, let us discuss some of their consequences regarding the KCP procedure.

**Fully non-parametric consistent change-point detection** We have proved that for any kernel satisfying some reasonably mild hypotheses, the KCP procedure outputs a segmentation close by the true segmentation with high probability.

An important particular example is the “asymptotic setting” of Example 1, where we have a fixed true segmentation  $\tau^*$  and fixed distributions  $P_1, \dots, P_{K+1}$  from which more and more points are sampled. How fast can KCP recover  $\tau^*$ , without any prior information on the number of segments  $D_{\tau^*}$  or on the distributions  $P_1, \dots, P_{K+1}$ ?

Let us take a bounded characteristic kernel — for instance the Gaussian or the Laplace kernel if  $\mathcal{X} = \mathbb{R}^d$  —, so that Assumption 1 holds true. Then, Theorem 1 shows that KCP detects consistently all changes in the distribution of the  $X_i$ , and localizes them at speed  $\log(n)/n$ . This speed also depends on the adequation between the kernel  $k$  and the differences between the  $P_j$  — through the ratio  $\underline{\Delta}/M$ . In the case of a single change-point, in a non-parametric setting, some consistency result exists for the detection of arbitrary changes in the distribution of the data [10, Th. 3.6.1]. However, obtaining such a fully non-parametric result for multiple change-points with a general set  $\mathcal{X}$  — we only need to know a bounded characteristic kernel on  $\mathcal{X}$  — has never been obtained before, to the best of our knowledge.

**Choice of  $k$**  An important question remains: how to choose the kernel  $k$ ? In Theorem 1,  $k$  only appears through the “signal-to-noise ratio”  $\underline{\Delta}/M$ , leading to better theoretical guarantees when this signal-to-noise ratio is larger: a larger value for  $C_{\max}$  and a smaller bound  $v_1$  on  $d_{\infty}^{(1)}(\tau^*, \hat{\tau})$ . Therefore, a simple strategy for choosing the kernel is to pick  $k$  that maximizes  $\underline{\Delta}/M$ , at least among a family of kernels, for instance Gaussian kernels. This first idea requires to know the distributions of the  $X_i$ , or at least to have prior information on them. Interestingly, when the change-points locations are known,  $\underline{\Delta}^2$  corresponds to the Maximum Mean Discrepancy (MMD) [19] between the distributions of the  $X_i$  over contiguous segments. In this particular setting, it is feasible to estimate and to maximize  $\underline{\Delta}$  with respect to the kernel  $k$ , as it is done in [20]. An interesting future development would be to build an estimator of  $\underline{\Delta}$  without knowing the change-point locations and to maximize this estimator with respect to the kernel  $k$ . We refer to [2, section 7.2] for a complementary discussion about the choice of  $k$  for KCP.

**Choice of  $C$**  Another important parameter of the KCP procedure is the constant  $C$  in front of the penalty. As mentioned below Theorem 1, our theoretical guarantees provide some guidelines for choosing  $C$ , but these are not sufficient

to choose precisely  $C$  in practice. We recommend to follow the advice of [2, section 6.2] on this point, which is to choose  $C$  from data thanks to the “slope heuristic”.

**Modularity of the proofs and possible extensions** Finally, we would like to emphasize what we believe to be an important contribution of this paper. The structure of the proofs of Theorems 1 and 2 — which follow the same strategy — is modular, so that one can easily adapt it to different sets of assumptions.

Our proof strategy is not new, since it is similar to the one of almost all previous papers analyzing the consistency of least-squares change-point detection procedures. In particular, we adapted some ideas of the proofs of [33] to the Hilbert space setting. Nevertheless, these papers formulate their results in asymptotic terms, which can be seen as a limitation — especially when  $n$  is small or  $\mathcal{X}$  is of large dimension. Another approach is the one of [35, 13, 2] where non-asymptotic oracle inequalities — using concentration inequalities and following the model selection results of Birgé and Massart [7] — are provided as theoretical guarantees on some penalized least-squares change-point procedures. Up to now, these two approaches seemed difficult to combine. The proofs of Theorems 1 and 2 show how they can be reconciled, which allows us to mix their strengths.

Indeed, the assumptions on the distributions of the  $X_i$  — Assumption 1 or 2 — are only used for proving bounds on two quantities — a linear term  $L_\tau$  and a quadratic term  $Q_\tau$  —, uniformly over  $\tau \in \mathcal{T}_n$ . Under Assumption 1, this is done thanks to concentration inequalities (Lemmas 8 and 9) which have been proved first by [2] in order to get an oracle inequality. Under Assumption 2, this is done by generalizing the method of [33] to Hilbert-space valued data, through two deterministic bounds (Lemmas 6 and 7) and a deviation inequality for

$$M_n := \max_{1 \leq k \leq n} \left\| \sum_{j=1}^k \varepsilon_j \right\|_{\mathcal{H}}$$

(Lemma 11). The rest of the proofs does not use anything about the distribution of  $X_1, \dots, X_n$ .

As a consequence, if one can generalize these bounds to another setting, a straightforward consequence is that a result similar to Theorem 1 or 2 holds true for the KCP procedure in this new setting. In particular, this could be used for dealing with the case of dependent data  $X_1, \dots, X_n$ . We could also consider an intermediate assumption between Assumption 2 and Assumption 1, of the form:

$$\max_{1 \leq i \leq n} \mathbb{E}[k(X_i, X_i)^\alpha] \leq B_\alpha < +\infty$$

for some  $\alpha \in (1, +\infty)$ .

## 5. Proofs

Let us start by describing our general strategy for proving our main results. Our goal is to build a large probability event on which any  $\hat{\tau} \in \arg \min_{\tau \in \mathcal{T}_n} \text{crit}(\tau)$  belongs to some subset  $\mathcal{E}$  of  $\mathcal{T}_n$ . For proving this, we use the key fact that  $\text{crit}(\tau^*) \geq \text{crit}(\hat{\tau})$ , together with a lower bound on  $\text{crit}(\tau)$  holding simultaneously for all  $\tau \in \mathcal{T}_n$ —hence for  $\tau = \hat{\tau}$ .

In order to get such a lower bound on the empirical penalized criterion, we start by decomposing it in Section 5.1 into terms that are simpler to control individually: two random terms — a linear function of  $\varepsilon$  and a quadratic function of  $\varepsilon$  —, and two deterministic terms — the approximation error and the penalty. Then, we control these terms thanks to deterministic bounds (Section 5.2) and deviation/concentration inequalities (Section 5.3). Finally, we prove Theorem 1 in Section 5.4 and Theorem 2 in Section 5.5.

### 5.1. Decomposition of the empirical risk

The first step in the proofs of Theorems 1 and 2 is to decompose the empirical risk (6).

**Lemma 2.** *Let  $\tau \in \mathcal{T}_n$  be a segmentation. Define  $\mu_\tau^* = \Pi_\tau \mu^*$ . Then we can write*

$$n\widehat{\mathcal{R}}_n(\tau) = \|Y - \widehat{\mu}_\tau\|^2 = \|\mu^* - \mu_\tau^*\|^2 + 2\langle \mu^* - \mu_\tau^*, \varepsilon \rangle - \|\Pi_\tau \varepsilon\|^2 + \|\varepsilon\|^2. \quad (14)$$

*Proof.* First, recall that  $\widehat{\mu}_\tau = \Pi_\tau Y$  and that  $Y = \mu^* + \varepsilon$ , hence

$$\begin{aligned} \|Y - \widehat{\mu}_\tau\|^2 &= \|Y - \Pi_\tau Y\|^2 \\ &= \|\mu^* + \varepsilon - \Pi_\tau(\mu^* + \varepsilon)\|^2 \\ &= \|\mu^* - \Pi_\tau \mu^*\|^2 + \|\varepsilon - \Pi_\tau \varepsilon\|^2 + 2\langle \mu^* - \Pi_\tau \mu^*, \varepsilon - \Pi_\tau \varepsilon \rangle. \end{aligned}$$

Since  $\Pi_\tau$  is an orthogonal projection,

$$\begin{aligned} \|Y - \widehat{\mu}_\tau\|^2 &= \|\mu^* - \mu_\tau^*\|^2 + \|\varepsilon\|^2 - 2\langle \varepsilon, \Pi_\tau \varepsilon \rangle + \|\Pi_\tau \varepsilon\|^2 + 2\langle (\text{Id} - \Pi_\tau)\mu^*, \varepsilon \rangle \\ &= \|\mu^* - \mu_\tau^*\|^2 + \|\varepsilon\|^2 - \|\Pi_\tau \varepsilon\|^2 + 2\langle (\text{Id} - \Pi_\tau)\mu^*, \varepsilon \rangle. \end{aligned}$$

□

Since each term of Eq. (14) behaves differently and is controlled via different techniques depending on the result to be proven, we give a name to each of these terms:

$$L_\tau := \langle \mu^* - \mu_\tau^*, \varepsilon \rangle, \quad Q_\tau := \|\Pi_\tau \varepsilon\|^2 \quad \text{and} \quad A_\tau := \|\mu^* - \mu_\tau^*\|^2. \quad (15)$$

It should be clear that  $L$  stands for “linear”,  $Q$  stands for “quadratic” and  $A$  stands for “approximation error”. We also define

$$\psi_\tau := 2L_\tau - Q_\tau + A_\tau. \quad (16)$$

Thus Lemma 2 states that

$$n\widehat{\mathcal{R}}_n(\tau) = \psi_\tau + \|\varepsilon\|^2.$$

Notice that  $L_{\tau^*} = A_{\tau^*} = 0$  and  $Q_{\tau^*} \geq 0$ , hence  $\psi_{\tau^*} \leq 0$ . Also note that  $\psi$ ,  $L$  and  $Q$  are random quantities depending on  $\varepsilon$ .

### 5.2. Deterministic bounds

In this section, we provide some deterministic bounds that are used in the proofs of Theorems 1 and 2.

#### 5.2.1. Approximation error $A_\tau$

We begin by the following result, which is the reason for the  $\underline{\Delta}_{\tau^*}\underline{\Delta}^2$  term in Theorem 1.

**Lemma 3.** *Let  $\tau \in \mathcal{T}_n$  be a segmentation such that  $D := D_\tau < D_{\tau^*}$ . Then*

$$\frac{1}{n}A_\tau = \frac{1}{n}\|\mu^* - \mu_\tau^*\|^2 \geq \frac{1}{2}\underline{\Delta}_{\tau^*}\underline{\Delta}^2. \quad (17)$$

The proof of Lemma 3 can be found in Section A.3.2.

*Remark 1.* Lemma 3 is tight. Indeed, consider the simple case  $D_\tau = 1$  and  $D_{\tau^*} = 2$ . Assume that  $n = 2m$  is an even number, and let  $\tau_1^* = m$ . It follows from definitions (7) and (8) that, in this case,

$$\underline{\Delta} = \|\mu_1^* - \mu_n^*\|_{\mathcal{H}} \quad \text{and} \quad \underline{\Delta}_{\tau^*} = \frac{1}{2}.$$

According to Eq. (5),  $(\mu_\tau^*)_i = \frac{1}{2}(\mu_1^* + \mu_n^*)$ , which yields

$$\frac{1}{n}A_\tau = \frac{1}{4}\|\mu_1^* - \mu_n^*\|_{\mathcal{H}}^2 = \frac{1}{2}\underline{\Delta}_{\tau^*}\underline{\Delta}^2.$$

Thus, in this particular class of examples, equality holds in Eq. (17).

We next state an analogous result, valid for any  $\tau \in \mathcal{T}_n$ , which plays a key role in the proofs of Theorems 1 and 2.

**Lemma 4.** *For any  $\tau \in \mathcal{T}_n$ ,*

$$\frac{1}{n}A_\tau \geq \frac{1}{2} \min \left\{ \underline{\Delta}_{\tau^*}, \frac{1}{n}d_\infty^{(1)}(\tau^*, \tau) \right\} \underline{\Delta}^2. \quad (18)$$

Lemma 4 is proved in Section A.3.3.

5.2.2. Linear term  $L_\tau$  and quadratic term  $Q_\tau$

The proof of Theorem 2 relies on some deterministic bounds on  $L_\tau$  and  $Q_\tau$ . We start with a preliminary lemma. First define

$$M_n := \max_{1 \leq k \leq n} \left\| \sum_{j=1}^k \varepsilon_j \right\|_{\mathcal{H}}.$$

**Lemma 5.** For any  $\varepsilon_1, \dots, \varepsilon_n \in \mathcal{H}$ ,

$$\max_{1 \leq a < b \leq n} \left\| \sum_{j=a}^b \varepsilon_j \right\|_{\mathcal{H}} \leq 2M_n. \quad (19)$$

*Proof.* For every  $a < b$ , we have:

$$\left\| \sum_{j=a}^b \varepsilon_j \right\|_{\mathcal{H}} = \left\| \sum_{j=1}^b \varepsilon_j - \sum_{j=1}^{a-1} \varepsilon_j \right\|_{\mathcal{H}} \leq \left\| \sum_{j=1}^b \varepsilon_j \right\|_{\mathcal{H}} + \left\| \sum_{j=1}^{a-1} \varepsilon_j \right\|_{\mathcal{H}} \leq 2M_n.$$

□

The following result is a deterministic bound on  $Q_\tau$  in terms of  $M_n$ .

**Lemma 6.** Let  $\tau \in \mathcal{T}_n$  be a segmentation. Then

$$Q_\tau \leq \frac{4D_\tau M_n^2}{n\underline{\Delta}_\tau}.$$

*Proof.* By Eq. (5),

$$\begin{aligned} Q_\tau &= \sum_{\ell=1}^{D_\tau} \frac{1}{|\tau_\ell - \tau_{\ell-1}|} \left\| \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} \varepsilon_j \right\|_{\mathcal{H}}^2 \leq D_\tau \max_{1 \leq \ell \leq D_\tau} \left\{ \frac{1}{|\tau_\ell - \tau_{\ell-1}|} \left\| \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} \varepsilon_j \right\|_{\mathcal{H}}^2 \right\} \\ &\leq \frac{D_\tau}{n\underline{\Delta}_\tau} \max_{1 \leq \ell \leq D_\tau} \left\| \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} \varepsilon_j \right\|_{\mathcal{H}}^2 \leq \frac{4D_\tau}{n\underline{\Delta}_\tau} M_n^2 \end{aligned}$$

where we used Lemma 5 for the last inequality. □

The following result is a deterministic bound on  $L_\tau$ .

**Lemma 7.** For any  $\tau \in \mathcal{T}_n$ ,

$$|L_\tau| \leq 6D_{\tau^*} \max \{D_{\tau^*}, D_\tau\} \overline{\Delta} M_n.$$

Lemma 7 is proved in Section A.4.

### 5.3. Concentration

In this subsection, we present concentration results on  $Q_\tau$ ,  $L_\tau$ , and deviation bounds for  $M_n$  — which will imply deviation bounds on  $Q_\tau$  and  $L_\tau$  by Lemmas 6 and 7). For any  $j \in \{1, \dots, n\}$ ,  $\tau \in \mathcal{T}_n$  and  $\ell \in \{1, \dots, D_\tau\}$ , we define

$$v_j := \mathbb{E} \left[ \|\varepsilon_j\|_{\mathcal{H}}^2 \right] \quad v_{\tau,\ell} := \frac{1}{\tau_\ell - \tau_{\ell-1}} \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} v_j \quad \text{and} \quad v_\tau := \sum_{\ell=1}^D v_{\tau,\ell}.$$

**Concentration under Assumption 1** The first result takes care of the linear term  $L_\tau$  when Assumption 1 is satisfied.

**Lemma 8** (Prop. 3 of [2]). *Suppose that Assumption 1 holds true. Then for any  $x > 0$ , with probability at least  $1 - 2e^{-x}$ , for any  $\theta > 0$ ,*

$$|L_\tau| \leq \theta A_\tau + \left( \frac{4}{3} + \frac{1}{2\theta} \right) M^2 x.$$

The next result deals with the quadratic term  $Q_\tau$  when Assumption 1 is satisfied.

**Lemma 9** (Prop. 1 of [2]). *Suppose that Assumption 1 holds true. Then for any  $x > 0$ , with probability at least  $1 - e^{-x}$ ,*

$$Q_\tau - v_\tau \leq \left( x + 2\sqrt{2xD_\tau} \right) \frac{14M^2}{3}.$$

We merge Lemma 8 and 9 for convenience.

**Lemma 10.** *Suppose that Assumption 1 holds true. Take any  $\lambda > 1$  and  $\tau \in \mathcal{T}_n$  be a segmentation. Then, there exists an event  $\Omega_{\tau,\lambda}^{(0)}$  of probability greater than  $1 - 3e^{-\lambda D_\tau}$  on which:*

$$\psi_\tau \geq \frac{1}{3} A_\tau - \frac{74}{3} \lambda D_\tau M^2.$$

*Proof.* According to Lemma 8 with  $\theta = 1/3$  and  $x = \lambda D_\tau$ , there exists an event  $\Omega_{\tau,\lambda}^{(1)}$  on which  $|L_\tau| \leq \frac{1}{3} A_\tau + \frac{17}{6} \lambda D_\tau M^2$ , with  $\mathbb{P} \left( \Omega_{\tau,\lambda}^{(1)} \right) \geq 1 - 2e^{-\lambda D_\tau}$ . Lemma 9 with  $x = \lambda D_\tau$  gives  $\Omega_{\tau,\lambda}^{(2)}$  on which  $Q_\tau - v_\tau \leq \frac{14}{3} \left( \lambda + 2\sqrt{2\lambda} \right) D_\tau M^2$ , with  $\mathbb{P} \left( \Omega_{\tau,\lambda}^{(2)} \right) \geq 1 - e^{-\lambda D_\tau}$ . Then,  $\Omega_{\tau,\lambda}^{(0)} := \Omega_{\tau,\lambda}^{(1)} \cap \Omega_{\tau,\lambda}^{(2)}$  has a probability larger than  $1 - 3e^{-\lambda D_\tau}$  by the union bound. Since for any  $1 \leq \ell \leq D_\tau$ ,  $v_{\tau,\ell} \leq M^2$ , we have  $v_\tau = \sum_{\ell=1}^{D_\tau} v_{\tau,\ell} \leq D_\tau M^2$ . Hence, by definition (16) of  $\psi_\tau$  and using that  $\lambda \geq 1$ , on the event  $\Omega_{\tau,\lambda}^{(0)}$ , we have:

$$\begin{aligned} \psi_\tau &\geq \frac{1}{3} A_\tau - \left( \frac{31}{3} \lambda + \frac{28}{3} \sqrt{2\sqrt{\lambda}} + 1 \right) D_\tau M^2 \\ &\geq \frac{1}{3} A_\tau - \lambda \left( \frac{31}{3} + \frac{28}{3} \sqrt{2} + 1 \right) D_\tau M^2. \end{aligned}$$

□

*Remark 2.* It is also possible to obtain an upper bound for  $\psi_\tau$ : by Lemma 8, for every  $\lambda \geq 0$ , on the event  $\Omega_{\tau,\lambda}^{(2)} \subset \Omega_{\tau,\lambda}^{(0)}$ ,

$$\psi_\tau \leq \frac{5}{3}A_\tau + \frac{17}{3}\lambda D_\tau M^2.$$

However, we do not need this result thereafter.

**Concentration under Assumption 2** Lemma 6 and 7 directly translate upper bounds on  $M_n$  into controls of  $L_\tau$  and  $Q_\tau$ . Under Assumption 2, this is achieved via the following lemma, a Kolmogorov-like inequality for the noise in the RKHS. This result is a straightforward generalization of the inequality obtained by [26] into the Hilbert setting. A more precise result (for real random variables only) can be found in [21], of which we follow the proof (their proof adapts well in our setting but we do not need the full result).

**Lemma 11.** *If Assumption 2 holds true, then, for any  $x > 0$ ,*

$$\mathbb{P}(M_n \geq x) \leq \frac{1}{x^2} \sum_{j=1}^n v_j. \quad (20)$$

*Remark 3.* We can reformulate Lemma 11 as follows. For any  $y > 0$ , there exists an event of probability at least  $1 - y^{-2}$  on which  $M_n < y\sqrt{\sum_{i=1}^n v_i} \leq y\sqrt{nV}$ . Equivalently, for any  $z \geq 0$ , there exists an event of probability at least  $1 - e^{-z}$  such that  $M_n < e^{z/2}\sqrt{\sum_{i=1}^n v_i} \leq e^{z/2}\sqrt{nV}$ .

*Proof.* Let us put

$$\zeta := \|\varepsilon_1 + \dots + \varepsilon_n\|_{\mathcal{H}}^2.$$

Since for any  $j \neq k$ ,  $\mathbb{E}[\langle \varepsilon_j, \varepsilon_k \rangle_{\mathcal{H}}] = 0$  (see Remark 4), by definition of  $v_j$ ,

$$\mathbb{E}[\zeta] = \mathbb{E}\left[\|\varepsilon_1 + \dots + \varepsilon_n\|_{\mathcal{H}}^2\right] = \sum_{j=1}^n v_j.$$

We recognize the right-hand side of Eq. (20) up to  $1/x^2$ . For any  $r > 1$ , let us denote by  $A_r$  the event

$$\forall 1 \leq s < r, \quad \|\varepsilon_1 + \dots + \varepsilon_s\|_{\mathcal{H}} < x \quad \text{and} \quad \|\varepsilon_1 + \dots + \varepsilon_r\|_{\mathcal{H}} \geq x,$$

and by  $A_1$  the event  $\|\varepsilon_1\|_{\mathcal{H}} \geq x$ . These events are disjoint, thus we can write

$$\mathbb{P}\left(\max_{1 \leq k \leq n} \|\varepsilon_1 + \dots + \varepsilon_k\|_{\mathcal{H}} \geq x\right) = \mathbb{P}\left(\bigcup_{r=1}^n A_r\right) = \sum_{r=1}^n \mathbb{P}(A_r). \quad (21)$$

The law of total expectation and the positiveness of  $\zeta$  yield

$$\mathbb{E}[\zeta] \geq \sum_{r=1}^n \mathbb{E}[\zeta | A_r] \mathbb{P}(A_r).$$

Finally, let  $\ell \leq r < k$  be integers. Since  $\varepsilon_\ell$  is independent from  $\varepsilon_k$  conditionally to  $\sigma(\varepsilon_1, \dots, \varepsilon_r)$ ,  $\varepsilon_\ell$  is independent from  $\varepsilon_k$  conditionally to  $A_r$ . Furthermore,  $\varepsilon_k$  is independent from  $A_r$  and

$$\mathbb{E} [\langle \varepsilon_k, \varepsilon_\ell \rangle_{\mathcal{H}} | A_r] = \langle \mathbb{E} [\varepsilon_k], \mathbb{E} [\varepsilon_\ell | A_r] \rangle_{\mathcal{H}} = 0.$$

Because of this relation and the positivity of the (real) conditional expectation, for any integers  $r \leq k \leq j$ ,

$$\mathbb{E} [\zeta | A_r] = \mathbb{E} [\|\varepsilon_1 + \dots + \varepsilon_n\|_{\mathcal{H}}^2 | A_r] \geq \mathbb{E} [\|\varepsilon_1 + \dots + \varepsilon_r\|_{\mathcal{H}}^2 | A_r] \geq x^2.$$

Therefore,  $\mathbb{E} [\zeta | A_r] \geq x^2$ , which gives  $\mathbb{E} [\zeta] \geq x^2 \sum \mathbb{P}(A_r)$ . This concludes the proof, thanks to Eq. (21).  $\square$

*Remark 4.* The independence between  $\varepsilon_j$  and  $\varepsilon_k$  for  $j \neq k$  yields  $\mathbb{E} [\langle \varepsilon_j, \varepsilon_k \rangle_{\mathcal{H}}] = 0$ . Indeed, we dispose of a conditional expectation on  $\mathcal{H}$  [14, chapter 5], which satisfies the same properties than the conditional expectation with real random variables. Hence we can write

$$\begin{aligned} \mathbb{E} [\langle \varepsilon_j, \varepsilon_k \rangle_{\mathcal{H}}] &= \mathbb{E} [\mathbb{E} [\langle \varepsilon_j, \varepsilon_k \rangle_{\mathcal{H}} | \varepsilon_k]] \\ &= \mathbb{E} [\langle \mathbb{E} [\varepsilon_j | \varepsilon_k], \varepsilon_k \rangle_{\mathcal{H}}] \\ &= \mathbb{E} [\langle \mathbb{E} [\varepsilon_j], \varepsilon_k \rangle_{\mathcal{H}}] = 0. \end{aligned}$$

Note that the  $\varepsilon_j$ s expectation vanishes by hypothesis.

#### 5.4. Proof of Theorem 1

We follow the strategy described at the beginning of Section 5.

**Definition of  $\Omega$ .** Let us define  $\Omega := \bigcap_{\tau \in \mathcal{T}_n} \Omega_{\tau, \lambda}^{(0)}$  with  $\lambda = y + \log n + 1 > 1$ , where we recall that  $\Omega_{\tau, \lambda}^{(0)}$  is defined in Lemma 10. By the union bound, and since the  $\Omega_{\tau, \lambda}^{(0)}$  have probability greater than  $1 - 3e^{-\lambda D_\tau}$ ,

$$\mathbb{P}(\Omega) \geq 1 - 3 \sum_{\tau \in \mathcal{T}_n} e^{-\lambda D_\tau}.$$

The inequality  $\mathbb{P}(\Omega) \geq 1 - e^{-y}$  follows since

$$\begin{aligned} \sum_{\tau \in \mathcal{T}_n} e^{-\lambda D_\tau} &= \sum_{d=1}^n \binom{n-1}{d-1} e^{-\lambda d} = e^{-\lambda} (1 + e^{-\lambda})^{n-1} \\ &\leq e^{-\lambda} \exp((n-1)e^{-\lambda}) \\ &= \frac{e^{-y}}{n e} \exp\left(\frac{n-1}{n} e^{-1-y}\right) \\ &\leq e^{-y} \frac{\exp(e^{-1})}{n e} \leq 0.27 e^{-y}, \end{aligned}$$

where the last inequality uses that  $n \geq 2$ . From now on we work exclusively on  $\Omega$ .

**Key argument.** We now use of the simple (but crucial) remark that  $\text{crit}(\tau^*) \geq \text{crit}(\hat{\tau})$ , hence

$$n \text{pen}(\hat{\tau}) + \psi_{\hat{\tau}} \leq n \text{pen}(\tau^*) + \psi_{\tau^*} \leq n \text{pen}(\tau^*) = CD_{\tau^*}M^2.$$

Since we work on  $\Omega$ , by definition of  $\Omega_{\tau,\lambda}^{(0)}$  in Lemma 10, for any  $\tau \in \mathcal{T}_n$ , we have:

$$\psi_{\tau} \geq \frac{1}{3}A_{\tau} - \frac{74}{3}\lambda D_{\tau}M^2.$$

Therefore, we get:

$$CD_{\tau^*}M^2 \geq \frac{1}{3}A_{\hat{\tau}} + \left(C - \frac{74}{3}\lambda\right) D_{\hat{\tau}}M^2. \quad (22)$$

**Proof that  $D_{\hat{\tau}} \leq D_{\tau^*}$ .** Since  $C > 74\lambda/3$  (by the lower bound in assumption (9)),  $M^2 > 0$  and  $A_{\hat{\tau}} \geq 0$ , Eq. (22) implies that

$$D_{\hat{\tau}} \leq \frac{C}{C - \frac{74}{3}\lambda} D_{\tau^*}.$$

The lower bound in assumption (9) ensures that

$$\frac{C}{C - \frac{74}{3}\lambda} < \frac{D_{\tau^*} + 1}{D_{\tau^*}}$$

hence  $D_{\hat{\tau}} \leq D_{\tau^*}$  on  $\Omega$ .

**Proof that  $D_{\hat{\tau}} \geq D_{\tau^*}$ .** Since  $C > 74\lambda/3$  (by the lower bound in assumption (9)), Eq. (22) implies that

$$A_{\hat{\tau}} \leq 3CD_{\tau^*}M^2.$$

The upper bound in assumption (9) implies that

$$A_{\hat{\tau}} < \frac{n\Lambda_{\tau^*}\Delta^2}{2},$$

hence  $D_{\hat{\tau}} \geq D_{\tau^*}$  by Lemma 3.

**Loss between  $\hat{\tau}$  and  $\tau^*$ .** We have proved that  $D_{\hat{\tau}} = D_{\tau^*}$  on  $\Omega$ , therefore, Eq. (22) can be rewritten

$$A_{\hat{\tau}} \leq 74\lambda D_{\tau^*}M^2. \quad (23)$$

By Lemma 4 and the definition of  $\lambda$ , we get

$$\min\left\{\Lambda_{\tau^*}, \frac{1}{n}d_{\infty}^{(1)}(\tau^*, \hat{\tau})\right\} \leq \frac{148D_{\tau^*}M^2}{\Delta^2} \cdot \frac{y + \log n + 1}{n} = v_1(y). \quad (24)$$

Remark that Assumption (9) implies that

$$\frac{\underline{\Delta}^2}{M^2} \frac{\underline{\Lambda}_{\tau^*}}{6D_{\tau^*}} n > \frac{74}{3} (D_{\tau^*} + 1)(y + \log n + 1)$$

hence

$$\underline{\Lambda}_{\tau^*} > (D_{\tau^*} + 1) \frac{148D_{\tau^*}M^2}{\underline{\Delta}^2} \cdot \frac{y + \log n + 1}{n} > v_1(y).$$

Therefore, Eq. (24) can be simplified into

$$\frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}) \leq v_1(y). \quad \square$$

*Remark 5.* The proof of Theorem 1 generalizes to  $\hat{\tau}$  defined by

$$\hat{\tau} \in \arg \min_{\tau \in \mathcal{T}_n / \underline{\Lambda}_{\tau} \geq \delta_n} \{ \text{crit}(\tau) \}$$

instead of (2), for any  $\delta_n \geq 0$  such that  $\underline{\Lambda}_{\tau^*} \geq \delta_n$ . Indeed, this assumption allows to write  $\text{crit}(\tau^*) \geq \text{crit}(\hat{\tau})$  in the key argument, and the rest of the proof can stay unchanged (with the same event  $\Omega$ ). More generally, any constraint can be added in the argmin defining  $\hat{\tau}$ , provided that  $\tau^*$  satisfies this constraint.

### 5.5. Proof of Theorem 2

We follow the strategy described at the beginning of Section 5. Throughout the proof, we write  $\hat{\tau}_2$  as a shortcut for  $\hat{\tau}(D_{\tau^*}, \delta_n)$ .

**Key argument.** By definition (12) of  $\hat{\tau}_2 = \hat{\tau}(D_{\tau^*}, \delta_n)$ , since we assume  $\underline{\Lambda}_{\tau^*} \geq \delta_n$ ,

$$\hat{\mathcal{R}}_n(\tau^*) \geq \hat{\mathcal{R}}_n(\hat{\tau}_2)$$

hence

$$0 \geq \psi_{\tau^*} \geq \psi_{\hat{\tau}_2} = A_{\hat{\tau}_2} + 2L_{\hat{\tau}_2} - Q_{\hat{\tau}_2}.$$

By Lemma 6, Lemma 7 and the facts that  $D_{\hat{\tau}_2} = D_{\tau^*}$  and  $\underline{\Lambda}_{\hat{\tau}_2} \geq \delta_n$ , we get

$$0 \geq \psi_{\hat{\tau}_2} \geq A_{\hat{\tau}_2} - 12D_{\tau^*}^2 \bar{\Delta} M_n - \frac{4D_{\tau^*}M_n^2}{n\delta_n}$$

hence, using Lemma 4,

$$\min \left\{ \underline{\Lambda}_{\tau^*}, \frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}_2) \right\} \leq \frac{24D_{\tau^*}^2 \bar{\Delta} M_n}{\underline{\Delta}^2} + \frac{8D_{\tau^*} M_n^2}{\underline{\Delta}^2 n^2 \delta_n}. \quad (25)$$

**Definition of  $\Omega_2$ .** We define

$$\Omega_2 := \{M_n \leq y\sqrt{nV}\}.$$

By Lemma 11, under Assumption 2,  $\mathbb{P}(\Omega_2) \geq 1 - y^{-2}$ .

**Conclusion.** By definition of  $\Omega_2$ , Eq. (25) implies that on  $\Omega_2$ :

$$\min\left\{\underline{\Delta}_{\tau^*}, \frac{1}{n}d_{\infty}^{(1)}(\tau^*, \hat{\tau}_2)\right\} \leq 24(D_{\tau^*})^2 \frac{\bar{\Delta}\sqrt{V}}{\underline{\Delta}^2} \frac{y}{\sqrt{n}} + 8D_{\tau^*} \frac{V}{\underline{\Delta}^2} \frac{y^2}{n\delta_n} = v_2(y, \delta_n).$$

Since we assume  $v_2(y, \delta_n) < \underline{\Delta}_{\tau^*}$ , the result follows.  $\square$

## References

- [1] Felix Abramovich, Yoav Benjamini, David L. Donoho, and Iain M. Johnstone. Adapting to unknown sparsity by controlling the false discovery rate. *Ann. Statist.*, 34(2):584–653, 2006.
- [2] S. Arlot, A. Celisse, and Z. Harchaoui. A kernel multiple change-point algorithm via model selection. *arXiv preprint arXiv:1202.3878v2*, 2012.
- [3] N. Aronszajn. Theory of reproducing kernels. *Transactions of the American mathematical society*, pages 337–404, 1950.
- [4] Jushan Bai and Pierre Perron. Estimating and testing linear models with multiple structural changes. *Econometrica*, pages 47–78, 1998.
- [5] M. Basseville and I. Nikiforov. *Detection of abrupt changes: theory and application*. Prentice Hall Englewood Cliffs, 1993.
- [6] Richard Bellman. On the approximation of curves by line segments using dynamic programming. *Communications of the ACM*, 4(6):284, 1961.
- [7] L. Birgé and P. Massart. Gaussian model selection. *Journal of the European Mathematical Society*, 3(3):203–268, 2001.
- [8] Lucien Birgé and Pascal Massart. Minimal penalties for Gaussian model selection. *Probab. Theory Related Fields*, 138(1-2):33–73, 2007.
- [9] Leif Boysen, Angela Kempe, Volkmar Liescher, Axel Munk, and Olaf Witich. Consistencies and rates of convergence of jump-penalized least squares estimators. *Ann. Statist.*, 37(1):157–183, 2009.
- [10] E. Brodsky and Boris S. Darkhovsky. *Nonparametric methods in change point problems*, volume 243. Springer Science & Business Media, 2013.
- [11] Victor-Emmanuel Brunel. Convex set detection. *arXiv preprint arXiv:1404.6224*, 2014.
- [12] Alain Celisse, Guillemette Marot, Morgane Pierre-Jean, and Guillem Rigai. New efficient algorithms for multiple change-point detection with kernels. working paper or preprint, September 2016.
- [13] Fabienne Comte and Yves Rozenholc. A new algorithm for fixed design regression and denoising. *Ann. Inst. Statist. Math.*, 56(3):449–473, 2004.
- [14] Joseph Diestel and John Jerry Uhl. *Vector measures*. Number 15. American Mathematical Soc., 1977.
- [15] Aymeric Dieuleveut and Francis Bach. Non-parametric stochastic approximation with large step sizes. *arXiv preprint arXiv:1408.0361*, 2014.
- [16] Vinicius Ferraris, Nicolas Dobigeon, Qi Wei, and Marie Chabert. Robust fusion of multi-band images with different spatial and spectral resolutions for change detection. *submitted to ICASSP 2017*.

- [17] Walter D Fisher. On grouping for maximum homogeneity. *Journal of the American Statistical Association*, 53(284):789–798, 1958.
- [18] Kenji Fukumizu, Arthur Gretton, Xiaohai Sun, and Bernhard Schölkopf. Kernel measures of conditional dependence. In *Advances in Neural Information Processing Systems 20*, pages 489–496. Curran Associates, Inc., 2008.
- [19] Arthur Gretton, Karsten M Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex J Smola. A kernel method for the two-sample-problem. In *Advances in neural information processing systems*, pages 513–520, 2006.
- [20] Arthur Gretton, Dino Sejdinovic, Heiko Strathmann, Sivaraman Balakrishnan, Massimiliano Pontil, Kenji Fukumizu, and Bharath K Sriperumbudur. Optimal kernel choice for large-scale two-sample tests. In *Advances in Neural Information Processing Systems*, pages 1205–1213, 2012.
- [21] J. Hájek and A. Rényi. Generalization of an inequality of kolmogorov. *Acta Mathematica Hungarica*, 6(3-4):281–283, 1955.
- [22] Zaid Harchaoui and Olivier Cappé. Retrospective multiple change-point estimation with kernels. In *IEEE Workshop on Statistical Signal Processing*, pages 768–772, 2007.
- [23] Lajos Horváth, Marie Hušková, Gregory Rice, and Jia Wang. Estimation of the time of change in panel data. *arXiv preprint arXiv:1503.04455*, 2015.
- [24] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of classification*, 2(1):193–218, 1985.
- [25] Albert Y Kim, Caren Marzban, Donald B Percival, and Werner Stuetzle. Using labeled data to evaluate change detectors in a multivariate streaming environment. *Signal Processing*, 89(12):2529–2536, 2009.
- [26] A.N. Kolmogorov. Über die Summen durch den Zufall bestimmten unabhängigen größen. *Math. Ann*, 99:484–488, 1928.
- [27] Aleksandr Petrovich Korostelev and Alexandre B Tsybakov. *Minimax theory of image reconstruction*, volume 82. Springer Science & Business Media, 2012.
- [28] AP Korostelev. On minimax estimation of a discontinuous signal. *Theory of Probability & Its Applications*, 32(4):727–730, 1988.
- [29] Weil R Lai, Mark D Johnson, Raju Kucherlapati, and Peter J Park. Comparative analysis of algorithms for identifying amplifications and deletions in array cgh data. *Bioinformatics*, 21(19):3763–3770, 2005.
- [30] Rémi Lajugie, Sylvain Arlot, and Francis Bach. Large-margin metric learning for partitioning problems, March 2013. arXiv:1303.1280.
- [31] Rémi Lajugie, Sylvain Arlot, and Francis Bach. Large-margin metric learning for constrained partitioning problems. In *Proceedings of The 31st International Conference on Machine Learning*, pages 297–305, 2014.
- [32] Marc Lavielle. Using penalized contrasts for the change-point problem. *Signal processing*, 85(8):1501–1510, 2005.
- [33] Marc Lavielle and Eric Moulines. Least-squares estimation of an unknown number of shifts in a time series. *Journal of time series analysis*, 21(1):33–59, 2000.
- [34] Marc Lavielle and Gilles Teysiere. Detection of multiple change-points in

- multivariate time series. *Lithuanian Mathematical Journal*, 46(3):287–306, 2006.
- [35] Émilie Lebarbier. Detecting multiple change-points in the mean of a Gaussian process by model selection. *Signal Proces.*, 85:717–736, 2005.
- [36] Michel Ledoux and Michel Talagrand. *Probability in Banach Spaces: isoperimetry and processes*, volume 23. Springer Science & Business Media, 2013.
- [37] Jian Liu, Shiyong Wu, and James V Zidek. On segmented multivariate regression. *Statistica Sinica*, 7(2):497–525, 1997.
- [38] Bernhard Schölkopf and Alexander J Smola. *Learning with kernels: Support vector machines, regularization, optimization, and beyond*. MIT press, 2002.
- [39] Jun Shao. An asymptotic theory for linear model selection. *Statistica Sinica*, 7(2):221–242, 1997.
- [40] Vladimir Spokoiny. Multiscale local change point detection with applications to value-at-risk. *The Annals of Statistics*, pages 1405–1436, 2009.
- [41] Bharath K. Sriperumbudur, Kenji Fukumizu, Arthur Gretton, Gert R. G. Lanckriet, and Bernhard Schölkopf. Kernel choice and classifiability for rkhs embeddings of probability distributions. In *Advances in Neural Information Processing Systems*, volume 21. NIPS Foundation (<http://books.nips.cc>), 2009.
- [42] Alexander Tartakovsky, Igor Nikiforov, and Basseville Michèle. *Sequential Analysis: Hypothesis Testing and Changepoint Detection*, volume 136 of *Monographs on Statistics and Applied Probability*. Chapman and Hall/CRC, Boca Raton, FL, 2014.
- [43] Michael Vogt and Holger Dette. Detecting gradual changes in locally stationary processes. *arXiv preprint arXiv:1403.3808*, 2014.
- [44] Tengyao Wang and Richard J. Samworth. High-dimensional change-point estimation via sparse projection. Technical report, arXiv, 2016. arXiv:1606.06246.
- [45] Yi-Ching Yao. Estimating the number of change-points via schwarz’criterion. *Statistics & Probability Letters*, 6(3):181–189, 1988.
- [46] Yi-Ching Yao and S. T. Au. Least-squares estimation of a step function. *Sankhyā: The Indian Journal of Statistics, Series A*, pages 370–381, 1989.

## Appendix A: Proofs

The following additional notation are used in the appendix.

We denote by  $\lambda_1^*, \dots, \lambda_{D_{\tau^*}}^*$  the segments of  $\tau^*$ , that is,  $\lambda_i^* = \{\tau_{i-1}^* + 1, \dots, \tau_i^*\}$ .

For any segment  $\lambda$  of  $\tau \in \mathcal{T}_n$ , we denote by  $\mu_\lambda^*$  the value of  $\mu_\tau^*$  on  $\lambda$ , which does not depend on  $\tau$  and is given by Eq. (5):

$$\mu_\lambda^* = \frac{1}{|\lambda|} \sum_{j \in \lambda} \mu_j^*. \quad (26)$$

**A.1. Proof of Lemma 1**

**Proof of (i)** We set  $D^i := D_{\tau^i}$  for  $i \in \{1, 2\}$ . Let us show first that  $d_\infty^{(2)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^1, \tau^2)$ . Take any  $i \in \{1, \dots, D^1 - 1\}$ , by the definition of  $\underline{\Lambda}_{\tau^1}$ ,

$$|\tau_i^1 - \tau_{D^2}^2| = |\tau_i^1 - n| \geq n\underline{\Lambda}_{\tau^1} > n\underline{\Lambda}_{\tau^1}/2 \geq n \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\}/2 > d_\infty^{(1)}(\tau^1, \tau^2).$$

The same holds for  $|\tau_i^1 - \tau_0^2|$ . Hence, for any  $i \in \{1, \dots, D^1 - 1\}$ ,

$$\min_{0 \leq j \leq D^2} |\tau_i^1 - \tau_j^2| = \min_{1 \leq j \leq D^2 - 1} |\tau_i^1 - \tau_j^2|,$$

which proves that  $d_\infty^{(2)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^1, \tau^2)$ .

Next, we prove that  $D^1 = D^2$  and  $d_\infty^{(3)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^1, \tau^2)$ . Define  $\phi : \{1, \dots, D^1 - 1\} \rightarrow \{1, \dots, D^2 - 1\}$  such that  $\{\phi(i)\} = \arg \min_{1 \leq j \leq D^2 - 1} |\tau_i^1 - \tau_j^2|$  for all  $i \in \{1, \dots, D^1 - 1\}$ . This mapping is well-defined: indeed, suppose that  $j, k \in \{1, \dots, D^2 - 1\}$  both realize the minimum for some  $i \in \{1, \dots, D^1 - 1\}$ . Since we assumed  $\frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) < \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\}/2$ ,

$$|\tau_i^1 - \tau_j^2| = |\tau_i^1 - \tau_k^2| \leq d_\infty^{(1)}(\tau^1, \tau^2) < n \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\}/2.$$

By the triangle inequality,

$$|\tau_j^2 - \tau_k^2| < n \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\} \leq n\underline{\Lambda}_{\tau^2},$$

hence  $j = k$ . Next, we show that  $\phi$  is increasing. Take  $i, j \in \{1, \dots, D^1 - 1\}$  such that  $i < j$ . Recall that  $\tau^k$  is increasing ( $k = 1, 2$ ). Then

$$\begin{aligned} \tau_{\phi(i)}^2 - \tau_{\phi(j)}^2 &= \tau_{\phi(i)}^2 - \tau_i^1 + \tau_i^1 - \tau_j^1 + \tau_j^1 - \tau_{\phi(j)}^2 \\ &= \tau_{\phi(i)}^2 - \tau_i^1 - |\tau_i^1 - \tau_j^1| + \tau_j^1 - \tau_{\phi(j)}^2 \\ &\leq \left| \tau_{\phi(i)}^2 - \tau_i^1 \right| - |\tau_i^1 - \tau_j^1| + \left| \tau_j^1 - \tau_{\phi(j)}^2 \right| \\ &\leq 2d_\infty^{(1)}(\tau^1, \tau^2) - |\tau_i^1 - \tau_j^1| \\ &< n \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\} - n\underline{\Lambda}_{\tau^1} \leq 0. \end{aligned}$$

Hence  $\phi(i) < \phi(j)$ , so  $\phi$  is increasing. As a consequence,  $\phi$  is injective and we get  $D^1 \leq D^2$ . The same argument, exchanging  $\tau^1$  and  $\tau^2$ , shows that  $D^2 \leq D^1$ . Therefore,  $D^1 = D^2$  and  $\phi$  is an increasing permutation of  $\{1, \dots, D^1 - 1\}$ , hence it is the identity. As a consequence,  $d_\infty^{(3)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^1, \tau^2)$ .

Finally, since  $d_\infty^{(3)}$  is symmetric,  $d_\infty^i(\tau^1, \tau^2) = d_{\text{H}}^i(\tau^1, \tau^2)$  for any  $i \in \{1, 2, 3\}$ .

**Proof of (ii)** Since  $D_{\tau^1} = D_{\tau^2}$ , we can set  $D = D_{\tau^1} = D_{\tau^2}$ . Next, define  $\phi(i) := \arg \min_{1 \leq j \leq D-1} |\tau_i^1 - \tau_j^2|$  and  $C_\phi(i) := |\phi(i)|$  for all  $i \in \{1, \dots, D-1\}$ . Clearly,  $C_\phi(i) \geq 1$  for any  $i$ . Let us show that we actually have  $C_\phi(i) = 1$ .

Take  $i$  and  $j$  distincts elements of  $\{1, \dots, D-1\}$ , and suppose that  $\phi(i) \cap \phi(j)$  is non-empty. Let  $k$  be any element of  $\phi(i) \cap \phi(j)$ . By the triangle inequality and the definition of  $d_\infty^{(1)}$ ,

$$n\underline{\Lambda}_{\tau^1} \leq |\tau_i^1 - \tau_j^1| \leq |\tau_i^1 - \tau_k^2| + |\tau_k^2 - \tau_j^1| \leq 2d_\infty^{(1)}(\tau^1, \tau^2) < n\underline{\Lambda}_{\tau^1}.$$

Hence, the  $\phi(i)$  are disjoint and we can write  $\sum_{i=1}^{D-1} C_\phi(i) = D-1$ , which clearly implies that  $C_\phi(i) = 1$ .

From now on, we identify  $\phi(i)$  with its unique element. Let us show that  $\phi$  is increasing similarly to what we have done for proving (i). Take  $i, j \in \{1, \dots, D-1\}$  such that  $i < j$ . We showed that

$$\tau_{\phi(i)}^2 - \tau_{\phi(j)}^2 \leq 2d_\infty^{(1)}(\tau^1, \tau^2) - |\tau_i^1 - \tau_j^1|,$$

thus according to the definition of  $\underline{\Lambda}_{\tau^1}$ , and our assumption,

$$\tau_{\phi(i)}^2 - \tau_{\phi(j)}^2 < n\underline{\Lambda}_{\tau^1} - n\underline{\Lambda}_{\tau^1} \leq 0.$$

Hence  $\phi(i) < \phi(j)$ :  $\phi$  is increasing. As a consequence,  $d_\infty^{(1)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^2, \tau^1) = d_{\mathbb{H}}^{(1)}(\tau^1, \tau^2)$ .  $\square$

## A.2. Proofs about the Frobenius loss

### A.2.1. A formula for $d_{\mathbb{F}}^2$

We start by proving a general formula for  $d_{\mathbb{F}}$ , which is stated by [31]. We prove it here for completeness.

$$\forall \tau^1, \tau^2 \in \mathcal{T}_n, \quad d_{\mathbb{F}}(\tau^1, \tau^2)^2 = D_{\tau^1} + D_{\tau^2} - 2 \sum_{k=1}^{D_{\tau^1}} \sum_{\ell=1}^{D_{\tau^2}} \frac{|\lambda_k^1 \cap \lambda_\ell^2|^2}{|\lambda_k^1| \times |\lambda_\ell^2|}. \quad (27)$$

Indeed, by definition, we have

$$d_{\mathbb{F}}(\tau^1, \tau^2)^2 = \text{Tr}((\Pi_{\tau^1} - \Pi_{\tau^2})^2) = \underbrace{\text{Tr}(\Pi_{\tau^1})}_{=D_{\tau^1}} + \underbrace{\text{Tr}(\Pi_{\tau^2})}_{=D_{\tau^2}} - 2 \text{Tr}(\Pi_{\tau^1} \Pi_{\tau^2})$$

$$\text{and} \quad \text{Tr}(\Pi_{\tau^1} \Pi_{\tau^2}) = \sum_{i=1}^n \sum_{j=1}^n \frac{\mathbf{1}_{\{\lambda_1(i)=\lambda_1(j) \text{ and } \lambda_2(i)=\lambda_2(j)\}}}{|\lambda_1(i)| |\lambda_2(i)|} = \sum_{k=1}^{D_{\tau^1}} \sum_{\ell=1}^{D_{\tau^2}} \frac{|\lambda_k^1 \cap \lambda_\ell^2|^2}{|\lambda_k^1| \times |\lambda_\ell^2|},$$

where we denoted by  $\lambda_k(i)$  the segment of  $\tau^k$  to which  $i \in \{1, \dots, n\}$  belongs.

### A.2.2. Proof of Eq. (10)

Eq. (10) is stated by [31]. The upper bound is a straightforward consequence of Eq. (27). We prove the lower bound here for completeness. We remark that

$$\sum_{k=1}^{D_{\tau^1}} \sum_{\ell=1}^{D_{\tau^2}} \frac{|\lambda_k^1 \cap \lambda_\ell^2|^2}{|\lambda_k^1| \times |\lambda_\ell^2|} \leq \sum_{k=1}^{D_{\tau^1}} \sum_{\ell=1}^{D_{\tau^2}} \frac{|\lambda_k^1 \cap \lambda_\ell^2|}{|\lambda_k^1|} = D_{\tau^1},$$

hence Eq. (27) shows that

$$d_{\mathbb{F}}(\tau^1, \tau^2)^2 \geq D_{\tau^2} - D_{\tau^1}.$$

The lower bound follows since  $\tau^1$  and  $\tau^2$  play symmetric roles.  $\square$

### A.2.3. Proof of Proposition 1

Throughout the proof, we write  $D = D_{\tau^1} = D_{\tau^2}$ ,  $\epsilon = n^{-1}d_{\infty}^{(1)}(\tau^1, \tau^2)$  and we denote by  $(\lambda_k^1)_{1 \leq k \leq D}$  and  $(\lambda_k^2)_{1 \leq k \leq D}$  the segments of  $\tau^1$  and  $\tau^2$ , respectively.

**Preliminary remark.** Since we assume that  $D_{\tau^1} = D_{\tau^2}$  and  $\frac{1}{n}d_{\infty}^{(1)}(\tau^1, \tau^2) = \epsilon < \frac{\underline{\Lambda}_{\tau^1}}{2}$ , point (ii) in Lemma 1 shows that  $d_{\infty}^{(1)}(\tau^1, \tau^2) = d_{\mathbb{H}}^{(1)}(\tau^1, \tau^2) = d_{\infty}^3(\tau^1, \tau^2)$ . In other words, for every  $k \in \{1, \dots, D-1\}$ , we have  $|\tau_k^1 - \tau_k^2| \leq n\epsilon$ , and some  $k_0 \in \{1, \dots, D-1\}$  exists such that  $|\tau_{k_0}^1 - \tau_{k_0}^2| = n\epsilon$ . As a consequence, for every  $k \in \{1, \dots, D-1\}$ ,

$$\left| |\lambda_k^1| - |\lambda_k^2| \right| \leq 2n\epsilon \quad \text{and} \quad |\lambda_k^1 \cap \lambda_k^2| \geq |\lambda_k^1| - 2n\epsilon. \quad (28)$$

**Upper bound for  $d_{\mathbb{F}}(\tau^1, \tau^2)^2$ .** We focus on the sum appearing in the right-hand side of Eq. (27). Using Eq. (28), we get:

$$\begin{aligned} \sum_{k=1}^D \sum_{\ell=1}^D \frac{|\lambda_k^1 \cap \lambda_{\ell}^2|^2}{|\lambda_k^1| \times |\lambda_{\ell}^2|} &\geq \sum_{k=1}^D \frac{|\lambda_k^1 \cap \lambda_k^2|^2}{|\lambda_k^1| \times |\lambda_k^2|} \\ &\geq \sum_{k=1}^D \left[ \frac{(|\lambda_k^1| - 2n\epsilon)^2}{|\lambda_k^1| \times (|\lambda_k^1| + 2n\epsilon)} \right] = \sum_{k=1}^D \frac{\left(1 - \frac{2n\epsilon}{|\lambda_k^1|}\right)^2}{1 + \frac{2n\epsilon}{|\lambda_k^1|}} \\ &\geq \sum_{k=1}^D \left(1 - \frac{6n\epsilon}{|\lambda_k^1|}\right) \geq D - \frac{6\epsilon D}{\underline{\Lambda}_{\tau^1}}, \end{aligned}$$

since for any  $x \geq 0$ ,  $\frac{(1-x)^2}{1+x} \geq 1 - 3x$ . The upper bound follows, using Eq. (27).

**Lower bound for  $d_{\mathbb{F}}(\tau^1, \tau^2)^2$ .** As shown in the preliminary remark, there exists some  $k_0 \in \{1, \dots, D-1\}$  such that  $|\tau_{k_0}^1 - \tau_{k_0}^2| = n\epsilon$ . First consider the case where  $\tau_{k_0}^1 < \tau_{k_0}^2$ . Then, by definition of  $d_{\mathbb{F}}$  and  $\Pi_{\tau}$ , we have:

$$\begin{aligned} d_{\mathbb{F}}(\tau^1, \tau^2)^2 &:= \sum_{1 \leq i, j \leq n} (\Pi_{\tau^1} - \Pi_{\tau^2})_{i,j}^2 \\ &\geq \sum_{i \in \lambda_{k_0+1}^1 \cap \lambda_{k_0}^2} \sum_{j \in \lambda_{k_0+1}^1 \cap \lambda_{k_0+1}^2} \frac{1}{|\lambda_{k_0+1}^1|^2} + \sum_{i \in \lambda_{k_0+1}^1 \cap \lambda_{k_0+1}^2} \sum_{j \in \lambda_{k_0+1}^1 \cap \lambda_{k_0}^2} \frac{1}{|\lambda_{k_0+1}^1|^2} \\ &= \frac{2|\lambda_{k_0+1}^1 \cap \lambda_{k_0}^2| \cdot |\lambda_{k_0+1}^1 \cap \lambda_{k_0+1}^2|}{|\lambda_{k_0+1}^1|^2}. \end{aligned}$$

Now, remark that  $|\lambda_{k_0+1}^1 \cap \lambda_{k_0}^2| = n\epsilon$ , by the preliminary remark and our assumption  $\tau_{k_0}^2 > \tau_{k_0}^1$ . Using also Eq. (28), we get:

$$d_{\mathbb{F}}(\tau^1, \tau^2)^2 \geq \frac{2n\epsilon(|\lambda_{k_0+1}^1| - 2n\epsilon)}{|\lambda_{k_0+1}^1|^2} \geq \frac{2n\epsilon}{3\bar{\Lambda}_{\tau^1}},$$

since  $|\lambda_{k_0+1}^1| - 2n\epsilon \geq |\lambda_{k_0+1}^1|/3$  and  $|\lambda_{k_0+1}^1| \leq \bar{\Lambda}_{\tau^1}$ . When  $\tau_{k_0}^1 > \tau_{k_0}^2$ , we apply the same reasoning, restricting the sum over  $i, j$  in the definition of  $d_{\mathbb{F}}$  to  $i \in \lambda_{k_0}^1 \cap \lambda_{k_0}^2$  and  $j \in \lambda_{k_0}^1 \cap \lambda_{k_0+1}^2$  (plus its symmetric). We obtain the same lower bound, which concludes the proof.  $\square$

### A.3. Lower bounds on the approximation error

This section provides the proofs of Lemmas 3 and 4.

#### A.3.1. Preliminary lemma

We start by a lemma useful in the two proofs.

**Lemma 12.** *If a segment  $\lambda \subset \{1, \dots, n\}$  intersects only two segments of  $\tau^*$ ,  $\lambda_i^*$  and  $\lambda_{i+1}^*$ , then we have:*

$$\sum_{j \in \lambda} \|\mu_j^* - \mu_\lambda^*\|_{\mathcal{H}}^2 = \frac{|\lambda \cap \lambda_i^*| \cdot |\lambda \cap \lambda_{i+1}^*|}{|\lambda \cap \lambda_i^*| + |\lambda \cap \lambda_{i+1}^*|} \|\mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^*\|_{\mathcal{H}}^2 \quad (29)$$

$$\geq \left( \frac{|\lambda \cap \lambda_i^*|}{|\lambda_i^*|} \wedge \frac{|\lambda \cap \lambda_{i+1}^*|}{|\lambda_{i+1}^*|} \right) \cdot \frac{|\lambda_i^*| \cdot |\lambda_{i+1}^*|}{|\lambda_i^*| + |\lambda_{i+1}^*|} \cdot \|\mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^*\|_{\mathcal{H}}^2. \quad (30)$$

*Proof.* We first prove Eq. (29). Since  $\lambda$  only intersects  $\lambda_i^*$  and  $\lambda_{i+1}^*$ , we have:

$$\begin{aligned} \sum_{j \in \lambda} \|\mu_j^* - \mu_\lambda^*\|_{\mathcal{H}}^2 &= \sum_{j \in \lambda \cap \lambda_i^*} \|\mu_j^* - \mu_\lambda^*\|_{\mathcal{H}}^2 + \sum_{j \in \lambda \cap \lambda_{i+1}^*} \|\mu_j^* - \mu_\lambda^*\|_{\mathcal{H}}^2 \\ &= |\lambda \cap \lambda_i^*| \cdot \|\mu_{\lambda_i^*}^* - \mu_\lambda^*\|_{\mathcal{H}}^2 + |\lambda \cap \lambda_{i+1}^*| \cdot \|\mu_{\lambda_{i+1}^*}^* - \mu_\lambda^*\|_{\mathcal{H}}^2. \end{aligned} \quad (31)$$

Since  $\mu_\lambda^*$  is given by Eq. (26), we obtain

$$\begin{aligned} \|\mu_{\lambda_i^*}^* - \mu_\lambda^*\|_{\mathcal{H}}^2 &= \left\| \frac{1}{|\lambda|} \sum_{j \in \lambda} (\mu_{\lambda_i^*}^* - \mu_j^*) \right\|_{\mathcal{H}}^2 = \left\| \frac{1}{|\lambda|} \sum_{j \in \lambda \cap \lambda_{i+1}^*} (\mu_{\lambda_i^*}^* - \mu_{\lambda_{i+1}^*}^*) \right\|_{\mathcal{H}}^2 \\ &= \frac{|\lambda \cap \lambda_{i+1}^*|^2}{|\lambda|^2} \|\mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^*\|_{\mathcal{H}}^2. \end{aligned}$$

The same computation on  $\lambda \cap \lambda_{i+1}^*$  yields

$$\left\| \mu_{\lambda_{i+1}^*}^* - \mu_{\lambda}^* \right\|_{\mathcal{H}}^2 = \frac{|\lambda \cap \lambda_i^*|^2}{|\lambda|^2} \left\| \mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^* \right\|_{\mathcal{H}}^2.$$

Therefore, Eq. (31) and the fact that  $|\lambda| = |\lambda \cap \lambda_i^*| + |\lambda \cap \lambda_{i+1}^*|$  yield Eq. (29).

Now, we remark that for any  $a, b, c, d > 0$ ,

$$\frac{abcd}{ab + cd} = \frac{1}{\frac{ab}{\max(a,c)} + \frac{cd}{\max(a,c)}} \times \min(a, c) \times bd \geq \min(a, c) \frac{bd}{b + d}.$$

Taking  $a = |\lambda \cap \lambda_i^*| / |\lambda_i^*|$ ,  $b = |\lambda_i^*|$ ,  $c = |\lambda \cap \lambda_{i+1}^*| / |\lambda_{i+1}^*|$  and  $d = |\lambda_{i+1}^*|$ , we get Eq. (30).  $\square$

### A.3.2. Proof of Lemma 3

In fact, we prove a slightly stronger statement. We show that, for any  $n \geq 2$ , for any  $D_{\tau^*} \in \{2, \dots, n\}$ , for any  $D \in \{1, \dots, D_{\tau^*} - 1\}$  and any  $\tau \in \mathcal{T}_n^D$ ,

$$\left\| \mu^* - \mu_{\tau}^* \right\|^2 \geq \min_{1 \leq i \leq D_{\tau^*} - 1} \left\{ \frac{|\lambda_i^*| \cdot |\lambda_{i+1}^*|}{|\lambda_i^*| + |\lambda_{i+1}^*|} \cdot \left\| \mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^* \right\|_{\mathcal{H}}^2 \right\}. \quad (32)$$

Then,

$$\left\| \mu^* - \mu_{\tau}^* \right\|^2 \geq \underline{\Gamma} \cdot \underline{\Delta}^2 \quad \text{where} \quad \underline{\Gamma} = \left( n \max_{1 \leq i \leq D_{\tau^*} - 1} \left\{ \frac{1}{|\lambda_i^*|} + \frac{1}{|\lambda_{i+1}^*|} \right\} \right)^{-1}.$$

Since we always have

$$\underline{\Lambda}_{\tau^*} \geq \underline{\Gamma} \geq \frac{1}{2} \underline{\Lambda}_{\tau^*},$$

Eq. (17) follows.

**Proof of Eq. (32) by induction.** We show by strong induction on  $D_{\tau^*}$  that, for any  $D_{\tau^*} \geq 2$ , for any  $D \in \{1, \dots, D_{\tau^*} - 1\}$ , any  $n \geq D_{\tau^*}$  and any  $\tau \in \mathcal{T}_n^D$ , Eq. (32) holds true.

First, if  $D_{\tau^*} = 2$ , the result follows by Eq. (30) in Lemma 12 since we then have  $i = 1$  and

$$\frac{|\lambda \cap \lambda_1^*|}{|\lambda_1^*|} = \frac{|\lambda \cap \lambda_2^*|}{|\lambda_2^*|} = 1.$$

Suppose now that the result is proved for all  $D_{\tau^*} \in \{2, \dots, p\}$  and consider a change-point problem  $(\tau^*, \mu^*)$  with  $D_{\tau^*} = D_{\tau^*} = p + 1$  and  $n \geq p + 1$ . Let  $D < p + 1$  and some segmentation  $\tau \in \mathcal{T}_n^D$  be fixed. Then one of these two scenarios occurs: (i) there exists  $\lambda_i^*$  with  $2 \leq i \leq D_{\tau^*} - 1$  that does not contain any change-point of  $\tau$ , or (ii)  $\lambda_2^*, \dots, \lambda_{D_{\tau^*} - 1}^*$  all contain a change-point of  $\tau$ .

**Case (i)** Suppose that there exists an inner segment  $\lambda_i^*$  of  $\tau^*$ ,  $2 \leq i \leq D_{\tau^*} - 1$ , that does not contain any change-point of  $\tau$  (see Figure 4). Therefore, there exists  $k \in \{1, \dots, D\}$  such that  $\lambda_i^* \subsetneq \lambda_k$ . By definition, there are  $i - 1$  change-points of  $\tau^*$  to the left of  $\lambda_i^*$  and  $k - 1$  change-points of  $\tau$  to the left of  $\lambda_i^*$ . Suppose that  $k < i$ . We define  $\tau^\circ$  as the segmentation obtained by adding  $\tau_i^*$  to  $\tau$  (see Figure 4). Then  $\|\mu^* - \mu_\tau^*\|^2 \geq \|\mu^* - \mu_{\tau^\circ}^*\|^2$  because  $\tau^\circ$  is finer than  $\tau$ . Reducing  $\tau^\circ$  to a segmentation  $\tilde{\tau}^\circ$  of  $\{1, 2, \dots, \tau_i^*\}$  in  $k$  segments and  $\tau^*$  to a segmentation  $\tilde{\tau}^*$  of  $\{1, 2, \dots, \tau_i^*\}$  in  $i$  segments and defining  $\tilde{\mu}^* = (\mu_1^*, \dots, \mu_{\tau_i^*}^*) \in \mathcal{H}^i$ , we get back to a situation covered by the induction since  $i \leq D_{\tau^*} - 1$  and  $k < i$ . So,

$$\begin{aligned} \|\tilde{\mu}^* - \tilde{\mu}_{\tilde{\tau}^\circ}^*\|^2 &\geq \inf_{1 \leq j \leq i-1} \left\{ \frac{|\lambda_j^*| \cdot |\lambda_{j+1}^*|}{|\lambda_j^*| + |\lambda_{j+1}^*|} \cdot \left\| \tilde{\mu}_{\lambda_{j+1}^*}^* - \tilde{\mu}_{\lambda_j^*}^* \right\|_{\mathcal{H}}^2 \right\} \\ &\geq \inf_{1 \leq j \leq D_{\tau^*} - 1} \left\{ \frac{|\lambda_j^*| \cdot |\lambda_{j+1}^*|}{|\lambda_j^*| + |\lambda_{j+1}^*|} \cdot \left\| \mu_{\lambda_{j+1}^*}^* - \mu_{\lambda_j^*}^* \right\|_{\mathcal{H}}^2 \right\} \end{aligned}$$

and we get the result since  $\|\mu^* - \mu_{\tau^\circ}^*\|^2 \geq \|\tilde{\mu}^* - \tilde{\mu}_{\tilde{\tau}^\circ}^*\|^2$ . A symmetric reasoning can be applied if  $k \geq i$ , considering change-points to the right of  $\lambda_i^*$  and using that  $D - k + 1 < D_{\tau^*} - i + 1$  since  $D < D_{\tau^*}$ .

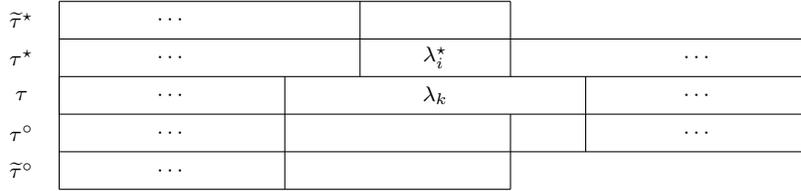


FIG 4. Proof of Lemma 3, Case (i):  $\lambda_i^*$  is a segment of  $\tau^*$  that is included in a segment of  $\tau$ . The segmentation  $\tau^\circ$  is obtained by joining  $\tau_i^*$  to the segmentation  $\tau$ .

**Case (ii)** Suppose that each inner segment of  $\tau^*$  contains a change-point of  $\tau$ . Since there are  $D_{\tau^*} - 2$  inner segments of  $\tau^*$  and  $D - 1 \leq D_{\tau^*} - 2$  change-points of  $\tau$ , there is at most (hence exactly) one change-point of  $\tau$  in each inner segment of  $\tau^*$ . Then  $D = D_{\tau^*} - 1$  and we are in the situation depicted in Figure 5.

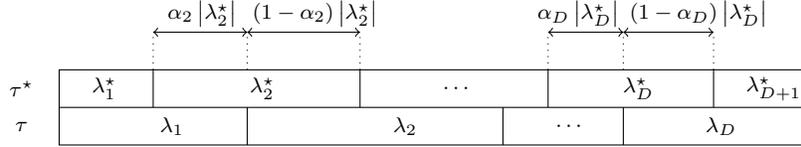


FIG 5. Proof of Lemma 3, Case (ii):  $D = D_{\tau^*} - 1$  and each inner segment of  $\tau^*$  contains exactly one change-point of  $\tau$ .

We can use Eq. (30) in Lemma 12 to lower bound the contribution of each  $\lambda \in \tau$  to  $\|\mu^* - \mu_\tau^*\|^2$ . For  $2 \leq i \leq D = D_{\tau^*} - 1$ , define  $\alpha_i := |\lambda_i^* \cap \lambda_{i-1}| / |\lambda_i^*|$ .

Then, we have

$$\begin{aligned}
 \|\mu^* - \mu_{\tau^*}^*\|^2 &\geq (1 \wedge \alpha_2) \frac{|\lambda_1^*| \cdot |\lambda_2^*|}{|\lambda_1^*| + |\lambda_2^*|} \cdot \left\| \mu_{\lambda_2^*}^* - \mu_{\lambda_1^*}^* \right\|_{\mathcal{H}}^2 \\
 &\quad + \sum_{j=2}^{D-1} \left( [(1 - \alpha_j) \wedge \alpha_{j+1}] \cdot \frac{|\lambda_j^*| \cdot |\lambda_{j+1}^*|}{|\lambda_j^*| + |\lambda_{j+1}^*|} \cdot \left\| \mu_{\lambda_{j+1}^*}^* - \mu_{\lambda_j^*}^* \right\|_{\mathcal{H}}^2 \right) \\
 &\quad + [(1 - \alpha_D) \wedge 1] \frac{|\lambda_D^*| \cdot |\lambda_{D+1}^*|}{|\lambda_D^*| + |\lambda_{D+1}^*|} \cdot \left\| \mu_{\lambda_{D+1}^*}^* - \mu_{\lambda_D^*}^* \right\|_{\mathcal{H}}^2 \\
 &\geq [1 \wedge \alpha_2 + (1 - \alpha_2) \wedge \alpha_3 + \dots + (1 - \alpha_{D-1}) \wedge \alpha_D + (1 - \alpha_D) \wedge 1] \\
 &\quad \times \inf_{1 \leq j \leq D_{\tau^*} - 1} \left\{ \frac{|\lambda_j^*| \cdot |\lambda_{j+1}^*|}{|\lambda_j^*| + |\lambda_{j+1}^*|} \cdot \left\| \mu_{\lambda_{j+1}^*}^* - \mu_{\lambda_j^*}^* \right\|_{\mathcal{H}}^2 \right\}.
 \end{aligned}$$

Since  $\alpha_i \geq 0$  for any  $2 \leq i \leq D_{\tau^*} - 1$ , it is straightforward to show that

$$\alpha_2 + (1 - \alpha_2) \wedge \alpha_3 + \dots + (1 - \alpha_D) \geq 1,$$

which concludes the proof.  $\square$

### A.3.3. Proof of Lemma 4

Let us define  $\delta := \min\{n_{\underline{\Delta}_{\tau^*}}, d_{\infty}^{(1)}(\tau^*, \tau)\}$ . If  $\delta = 0$ , then Eq. (18) holds true. We assume from now on that  $\delta > 0$ .

Because  $n_{\underline{\Delta}_{\tau^*}} \geq \delta$ , for any  $1 \leq i \leq D_{\tau^*} - 1$ , we can write  $|\tau_{i+1}^* - \tau_i^*| \geq \delta$ . On the other hand, because  $d_{\infty}^{(1)}(\tau^*, \tau) \geq \delta$ , there exists  $i \in \{1, \dots, D_{\tau^*} - 1\}$  such that, for any  $j \in \{1, \dots, D - 1\}$ ,  $|\tau_i^* - \tau_j| \geq \delta$ . Since  $\delta \leq n_{\underline{\Delta}_{\tau^*}}$ , this also holds true for  $j = 0$  and  $j = D$ . Let us define, as illustrated by Figure 6,

$$\lambda^{\circ} := \{\tau_i^* - \delta + 1, \dots, \tau_i^*, \tau_i^* + 1, \dots, \tau_{i+1}^* + \delta\}.$$

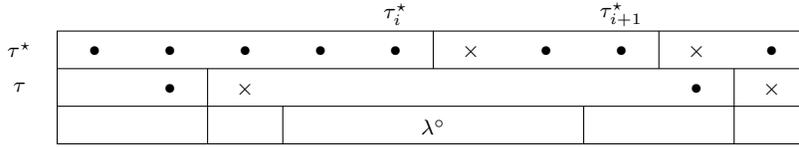


FIG 6. Construction of  $\lambda^{\circ}$  in the proof of Lemma 4. In this case,  $\delta = 2$  since  $\underline{\Delta}_{\tau^*} = 2/10$  (the rightmost segment of  $\tau^*$  is of size 2) and  $d_{\infty}(\tau^*, \tau) = 3$  (achieved in  $\tau_i^*$ ).

Since  $\lambda^{\circ}$  is included in a segment of  $\tau$ ,

$$\|\mu^* - \mu_{\tau^*}^*\|^2 \geq \sum_{j \in \lambda^{\circ}} \|\mu_j^* - (\mu_{\tau^*}^*)_j\|_{\mathcal{H}}^2 \geq \sum_{j \in \lambda^{\circ}} \|\mu_j^* - \mu_{\lambda^{\circ}}^*\|_{\mathcal{H}}^2.$$

Because of the hypothesis we made,  $\lambda^\circ$  only intersects  $\lambda_i^*$  and  $\lambda_{i+1}^*$  among the segments of  $\tau^*$ , so Eq. (29) in Lemma 12 shows that

$$\sum_{j \in \lambda^\circ} \|\mu_j^* - \mu_{\lambda^\circ}^*\|_{\mathcal{H}}^2 = \frac{|\lambda^\circ \cap \lambda_i^*| \cdot |\lambda^\circ \cap \lambda_{i+1}^*|}{|\lambda^\circ \cap \lambda_i^*| + |\lambda^\circ \cap \lambda_{i+1}^*|} \left\| \mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^* \right\|_{\mathcal{H}}^2 = \frac{\delta}{2} \left\| \mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^* \right\|_{\mathcal{H}}^2 \geq \frac{\delta}{2} \underline{\Delta}^2,$$

hence the result.  $\square$

#### A.4. Proof of Lemma 7

In this proof, since  $\tau$  is fixed, we denote by  $\lambda_1, \dots, \lambda_D$  the segments of  $\tau$ , that is,  $\lambda_i = \{\tau_{i-1} + 1, \dots, \tau_i\}$ .

First, notice that

$$L_\tau = \langle \mu^* - \mu_{\tau^*}^*, \varepsilon \rangle = \sum_{i=1}^{D_{\tau^*}} \left\langle \mu_{\lambda_i^*}^*, \sum_{j \in \lambda_i^*} \varepsilon_j \right\rangle_{\mathcal{H}} - \sum_{i=1}^{D_\tau} \left\langle \mu_{\lambda_i}^*, \sum_{j \in \lambda_i} \varepsilon_j \right\rangle_{\mathcal{H}}. \quad (33)$$

Now, if  $D_\tau < D_{\tau^*}$  we arbitrarily define  $\lambda_{D_\tau+1} = \dots = \lambda_{D_{\tau^*}} = \emptyset$ , so that  $\sum_{j \in \lambda_i} \varepsilon_j = 0$  for every  $i \in \{D_\tau + 1, \dots, D_{\tau^*}\}$ . Similarly, if  $D_{\tau^*} < D_\tau$ , we define  $\lambda_{D_{\tau^*}+1}^* = \dots = \lambda_{D_\tau}^* = \emptyset$ . We also define  $\mu_\emptyset^* = \mu_n^*$  by convention. Then, defining  $D^+ := \max\{D_{\tau^*}, D_\tau\}$ , we can rewrite Eq. (33) as follows:

$$\begin{aligned} L_\tau &= \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i^*}^*, \sum_{j \in \lambda_i^*} \varepsilon_j \right\rangle_{\mathcal{H}} - \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i}^*, \sum_{j \in \lambda_i} \varepsilon_j \right\rangle_{\mathcal{H}} \\ &= \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i^*}^* - \mu_{\lambda_i}^*, \sum_{j \in \lambda_i^*} \varepsilon_j \right\rangle_{\mathcal{H}} + \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i}^*, \sum_{j \in \lambda_i^*} \varepsilon_j - \sum_{j \in \lambda_i} \varepsilon_j \right\rangle_{\mathcal{H}} \\ &= \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i^*}^* - \mu_{\lambda_i}^*, \sum_{j \in \lambda_i^*} \varepsilon_j \right\rangle_{\mathcal{H}} + \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i}^* - \mu_n^*, \sum_{j \in \lambda_i^*} \varepsilon_j - \sum_{j \in \lambda_i} \varepsilon_j \right\rangle_{\mathcal{H}}, \end{aligned}$$

since

$$\sum_{i=1}^{D^+} \left( \sum_{j \in \lambda_i^*} \varepsilon_j - \sum_{j \in \lambda_i} \varepsilon_j \right) = 0.$$

Then, by the triangle inequality and Cauchy-Schwarz inequality,

$$\begin{aligned} |L_\tau| &\leq \sum_{i=1}^{D^+} \left\| \mu_{\lambda_i^*}^* - \mu_{\lambda_i}^* \right\|_{\mathcal{H}} \left\| \sum_{j \in \lambda_i^*} \varepsilon_j \right\|_{\mathcal{H}} + \sum_{i=1}^{D^+} \left\| \mu_{\lambda_i}^* - \mu_n^* \right\|_{\mathcal{H}} \left\| \sum_{j \in \lambda_i^*} \varepsilon_j - \sum_{j \in \lambda_i} \varepsilon_j \right\|_{\mathcal{H}} \\ &\leq \text{diam conv} \{ \mu_j^* / j \in \{1, \dots, n\} \} \left[ \sum_{i=1}^{D^+} \left\| \sum_{j \in \lambda_i^*} \varepsilon_j \right\|_{\mathcal{H}} + \sum_{i=1}^{D^+} \left( \left\| \sum_{j \in \lambda_i^*} \varepsilon_j \right\|_{\mathcal{H}} + \left\| \sum_{j \in \lambda_i} \varepsilon_j \right\|_{\mathcal{H}} \right) \right] \\ &\leq 3 \text{diam conv} \{ \mu_j^* / j \in \{1, \dots, n\} \} D^+ \sup_{1 \leq a < b \leq n} \left\| \sum_{j=a}^b \varepsilon_j \right\|_{\mathcal{H}} \end{aligned}$$

where we used that  $\mu_\lambda^* \in \text{conv}\{\mu_j^* / j \in \{1, \dots, n\}\}$  for any segment  $\lambda$ . Since the diameter of the convex hull of a finite set of points is equal to the diameter of the set, we have

$$\begin{aligned} \text{diam conv}\{\mu_j^* / j \in \{1, \dots, n\}\} &= \text{diam}\{\mu_j^* / j \in \{1, \dots, n\}\} \\ &\leq (D_{\tau^*} - 1)\bar{\Delta} < D_{\tau^*}\bar{\Delta}. \end{aligned}$$

Using also Lemma 5, we get the result. □