

# Sequential Learning of Principal Curves: Summarizing Data Streams on the Fly

Benjamin GUEDJ\* and Le LI†

May 22, 2018

## Abstract

When confronted with massive data streams, summarizing data with dimension reduction methods such as PCA raises theoretical and algorithmic pitfalls. Principal curves act as a nonlinear generalization of PCA and the present paper proposes a novel algorithm to automatically and sequentially learn principal curves from data streams. We show that our procedure is supported by regret bounds with optimal sublinear remainder terms. A greedy local search implementation that incorporates both sleeping experts and multi-armed bandit ingredients is presented, along with its regret bound and performance on a toy example and seismic data.

**keywords** sequential learning, principal curves, data streams, regret bounds, greedy algorithm, sleeping experts. MSC 2010: 68T10, 62L10, 62C99.

## 1 Introduction

Numerous methods have been proposed in the statistics and machine learning literature to sum up information and represent data by condensed and simpler to understand quantities. Among those methods, Principal Component Analysis (PCA) aims at identifying the maximal variance axes of data. This serves as a way to represent data in a more compact fashion and hopefully reveal as well as possible their variability. PCA has been introduced by [Pearson \(1901\)](#) and [Spearman \(1904\)](#) and further developed by [Hotelling \(1933\)](#). This is one of the most widely used procedures in multivariate exploratory analysis targeting dimension reduction or features extraction. Nonetheless, PCA is a linear procedure and the need for more sophisticated nonlinear techniques has led to the notion of principal curve. Principal curves may be seen as a nonlinear generalization of the first principal component. The goal is to obtain a curve which passes "in the middle" of data, as illustrated by [Figure 1](#). This notion has been at the heart of numerous applications in many different domains, such as physics ([Brunsdon, 2007](#); [Friedsam and Oren, 1989](#)), character and speech recognition ([Kégl and Krzyżak, 2002](#); [Reinhard and Niranjana, 1999](#)), mapping and geology ([Banfield and Raftery, 1992](#); [Brunsdon, 2007](#); [Stanford and Raftery, 2000](#)), to name but a few.

### 1.1 Earlier works on principal curves

The original definition of principal curve dates back to [Hastie and Stuetzle \(1989\)](#). A principal curve is a smooth ( $C^\infty$ ) parameterized curve  $\mathbf{f}(s) = (f_1(s), \dots, f_d(s))$  in  $\mathbb{R}^d$  which does not intersect itself, has finite length inside any bounded subset of  $\mathbb{R}^d$  and is self-consistent.

\*Inria, France. Web: <https://bguedj.github.io>, email [benjamin.guedj@inria.fr](mailto:benjamin.guedj@inria.fr), corresponding author.

†Université d'Angers & iAdvize, France. Email [le@iadvize.com](mailto:le@iadvize.com).

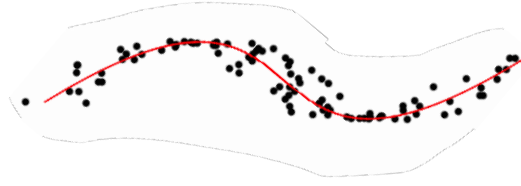


Figure 1: An example of principal curve.

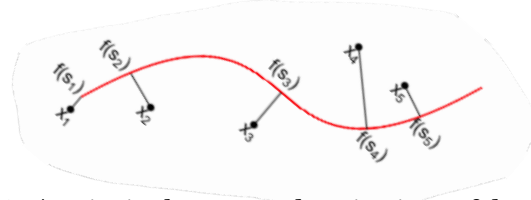


Figure 2: A principal curve and projections of data onto it.

This last requirement means that  $\mathbf{f}(s) = \mathbb{E}[X | s_{\mathbf{f}}(X) = s]$ , where  $X \in \mathbb{R}^d$  is a random vector and the so-called projection index  $s_{\mathbf{f}}(x)$  is the largest real number  $s$  minimizing the squared Euclidean distance between  $\mathbf{f}(s)$  and  $x$ , defined by

$$s_{\mathbf{f}}(x) = \sup \left\{ s : \|x - \mathbf{f}(s)\|_2^2 = \inf_{\tau} \|x - \mathbf{f}(\tau)\|_2^2 \right\}.$$

Self-consistency means that each point of  $\mathbf{f}$  is the average (under the distribution of  $X$ ) of all data points projected on  $\mathbf{f}$ , as illustrated by [Figure 2](#). However, an unfortunate consequence of this definition is that the existence is not guaranteed in general for a particular distribution, let alone for an online sequence for which no probabilistic assumption is made. [Kégl \(1999\)](#) proposed a new concept of principal curves which ensures the existence for a large class of distributions. Principal curves  $\mathbf{f}^*$  are defined as the curves minimizing the expected squared distance over a class  $\mathcal{F}_L$  of curves whose length is smaller than  $L > 0$ , namely,

$$\mathbf{f}^* \in \underset{\mathbf{f} \in \mathcal{F}_L}{\operatorname{arg\,inf}} \Delta(\mathbf{f}),$$

where

$$\Delta(\mathbf{f}) = \mathbb{E}[\Delta(\mathbf{f}, X)] = \mathbb{E} \left[ \inf_s \|\mathbf{f}(s) - X\|_2^2 \right].$$

If  $\mathbb{E}\|X\|_2^2 < \infty$ ,  $\mathbf{f}^*$  always exists but may not be unique. In practical situation where only i.i.d copies  $X_1, \dots, X_n$  of  $X$  are observed, [Kégl \(1999\)](#) considers classes  $\mathcal{F}_{k,L}$  of all polygonal lines with  $k$  segments and length not exceeding  $L$ , and chooses an estimator  $\hat{\mathbf{f}}_{k,n}$  of  $\mathbf{f}^*$  as the one within  $\mathcal{F}_{k,L}$  which minimizes the empirical counterpart

$$\Delta_n(\mathbf{f}) = \frac{1}{n} \sum_{i=1}^n \Delta(\mathbf{f}, X_i)$$

of  $\Delta(\mathbf{f})$ . It is proved in [Kégl et al. \(2000\)](#) that if  $X$  is almost surely bounded and  $k$  is proportional to  $n^{1/3}$ , then

$$\Delta(\hat{\mathbf{f}}_{k,n}) - \Delta(\mathbf{f}^*) = \mathcal{O}(n^{-1/3}).$$

As the task of finding a polygonal line with  $k$  segments and length at most  $L$  that minimizes  $\Delta_n(\mathbf{f})$  is computationally costly, [Kégl et al. \(2000\)](#) proposes the Polygonal Line algorithm. This iterative algorithm proceeds by fitting a polygonal line with  $k$  segments and considerably speeds up the exploration part by resorting to gradient descent. The two steps (projection and optimization) are similar to what is done by the  $k$ -means algorithm. However,

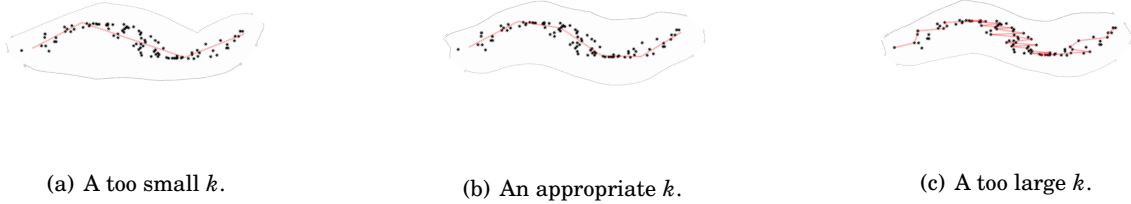


Figure 3: Principal curves with different number  $k$  of segments.

the Polygonal Line algorithm is not supported by theoretical bounds and leads to variable performance depending on the distribution of the observations.

As the number  $k$  of segments plays a crucial role (a too small  $k$  leads to a rough summary of data while a too large  $k$  yields overfitting, see Figure 3), Biau and Fischer (2012) aim to fill the gap by selecting an optimal  $k$  from both theoretical and practical perspectives. Their approach relies strongly on the theory of model selection by penalization introduced by Barron et al. (1999) and further developed by Birgé and Massart (2007). By considering countable classes  $\{\mathcal{F}_{k,\ell}\}_{k,\ell}$  of polygonal lines with  $k$  segments, total length  $\ell \leq L$  and whose vertices are on a lattice, the optimal  $(\hat{k}, \hat{\ell})$  is obtained as the minimizer of the criterion

$$\text{crit}(k, \ell) = \Delta_n(\hat{\mathbf{f}}_{k,\ell}) + \text{pen}(k, \ell),$$

where

$$\text{pen}(k, \ell) = c_0 \sqrt{\frac{k}{n}} + c_1 \frac{\ell}{n} + c_2 \frac{1}{\sqrt{n}} + \delta^2 \sqrt{\frac{w_{k,\ell}}{2n}}$$

is a penalty function where  $\delta$  stands for the diameter of observations,  $w_{k,\ell}$  denotes the weight attached to class  $\mathcal{F}_{k,\ell}$  and with constants  $c_0, c_1, c_2$  depending on  $\delta$ , the maximum length  $L$  and the dimension of observations. Biau and Fischer (2012) then prove that

$$\mathbb{E} \left[ \Delta(\hat{\mathbf{f}}_{\hat{k}, \hat{\ell}}) - \Delta(\mathbf{f}^*) \right] \leq \inf_{k,\ell} \left\{ \mathbb{E} \left[ \Delta(\hat{\mathbf{f}}_{k,\ell}) - \Delta(\mathbf{f}^*) \right] + \text{pen}(k, \ell) \right\} + \frac{\delta^2 \Sigma}{2^{3/2}} \sqrt{\frac{\pi}{n}}, \quad (1)$$

where  $\Sigma$  is a numerical constant. The expected loss of the final polygonal line  $\hat{\mathbf{f}}_{\hat{k}, \hat{\ell}}$  is close to the minimal loss achievable over  $\mathcal{F}_{k,\ell}$  up to a remainder term decaying as  $1/\sqrt{n}$ .

## 1.2 Motivation

The big data paradigm—where collecting, storing and analyzing massive amounts of large and complex data becomes the new standard—commands to revisit some of the classical statistical and machine learning techniques. The tremendous improvements of data acquisition infrastructures generates new continuous streams of data, rather than batch datasets. This has drawn a large interest to sequential learning. Existing theoretical works and practical implementations of principal curves are designed for the batch setting (Biau and Fischer, 2012; Kégl, 1999; Kégl and Krzyżak, 2002; Kégl et al., 2000; Sandilya and Kulkarni, 2002). To the best of our knowledge, very little effort has been put so far into extending principal curves algorithms to the sequential context (to the exception of Laparra and Malo, 2016, in a fairly different setting and with no theoretical results). The present paper aims at filling this gap: our goal is to propose an online perspective to principal curves by automatically and sequentially learning the best principal curve summarizing a data stream. This represents a clear improvement over the batch setting as when data streams are collected,

running a batch algorithm at each collect time is likely to be resources-demanding. Sequential learning takes advantage of the latest collected (set of) observations and therefore suffers a much smaller computational cost.

Sequential learning operates as follows: a blackbox reveals at each time  $t$  some deterministic value  $x_t, t = 1, 2, \dots$ , and a forecaster attempts to predict sequentially the next value based on past observations (and possibly other available information). The performance of the forecaster is no longer evaluated by its generalization error (as in the batch setting) but rather by a regret bound which quantifies the cumulative loss of a forecaster in the first  $T$  rounds with respect to some reference minimal loss. In sequential learning, the velocity of algorithms may be favored over statistical precision. An immediate use of aforementioned techniques (Biau and Fischer, 2012; Kégl et al., 2000; Sandilya and Kulkarni, 2002) at each time round  $t$  (treating data collected until  $t$  as a batch dataset) would result in a monumental algorithmic cost. Rather, we propose a novel algorithm which adapts to the sequential nature of data, *i.e.*, which takes advantage of previous computations. We refer the reader to the monograph Cesa-Bianchi and Lugosi (2006) for a thorough introduction to sequential learning.

The contributions of the present paper are twofold. We first propose a sequential principal curves algorithm, for which we derive regret bounds. We then move towards an implementation, illustrated on a toy dataset and a real-life dataset (seismic data). The sketch of our algorithm procedure is as follows. At each time round  $t$ , the number of segments of  $k_t$  is chosen automatically and the number of segments  $k_{t+1}$  in the next round is obtained by only using information about  $k_t$  and a small amount of past observations. The core of our procedure relies on computing a quantity which is linked to the mode of the so-called Gibbs quasi-posterior and is inspired by quasi-Bayesian learning. The use of quasi-Bayesian estimators is especially advocated by the PAC-Bayesian theory which originates in the machine learning community in the late 1990s, in the seminal works of Shawe-Taylor and Williamson (1997) and McAllester (1999a,b). In a recent preprint, Li et al. (2016) discuss the use of PAC-Bayesian tools for sequential learning.

The rest of the paper proceeds as follows. section 2 presents our notation and our online principal curve algorithm, for which we provide regret bounds with sublinear remainder terms in section 3. A practical implementation is proposed in section 4 and we illustrate its performance on a toy dataset and seismic data in section 5. Finally, we collect in section 6 proofs to original results claimed in the paper.

## 2 Notation

A parameterized curve in  $\mathbb{R}^d$  is a continuous function  $\mathbf{f}: I \rightarrow \mathbb{R}^d$  where  $I = [a, b]$  is a closed interval of the real line. The length of  $\mathbf{f}$  is given by

$$\mathcal{L}(\mathbf{f}) = \lim_{M \rightarrow \infty} \left\{ \sup_{a=s_0 < s_1 < \dots < s_M=b} \sum_{i=1}^M \|\mathbf{f}(s_i) - \mathbf{f}(s_{i-1})\|_2 \right\}.$$

Let  $x_1, x_2, \dots, x_T \in B(\mathbf{0}, \sqrt{d}R) \subset \mathbb{R}^d$  be a sequence of data, where  $B(\mathbf{c}, R)$  stands for the  $\ell_2$ -ball centered in  $\mathbf{c} \in \mathbb{R}^d$  with radius  $R > 0$ . Let  $\mathcal{Q}_\delta$  be a grid over  $B(\mathbf{0}, \sqrt{d}R)$ , *i.e.*,  $\mathcal{Q}_\delta = B(\mathbf{0}, \sqrt{d}R) \cap \Gamma_\delta$  where  $\Gamma_\delta$  is a lattice in  $\mathbb{R}^d$  with spacing  $\delta > 0$ . Let  $L > 0$  and define for each  $k \in \llbracket 1, p \rrbracket$  the collection  $\mathcal{F}_{k,L}$  of polygonal lines  $\mathbf{f}$  with  $k$  segments whose vertices are in  $\mathcal{Q}_\delta$  and such that  $\mathcal{L}(\mathbf{f}) \leq L$ . Denote by  $\mathcal{F}_p = \cup_{k=1}^p \mathcal{F}_{k,L}$  all polygonal lines with a number of segments  $\leq p$ , whose vertices are in  $\mathcal{Q}_\delta$  and whose length is at most  $L$ . Finally, let  $\mathcal{N}(\mathbf{f})$  denote the number of segments of  $\mathbf{f} \in \mathcal{F}_p$ . This strategy is illustrated by Figure 4.

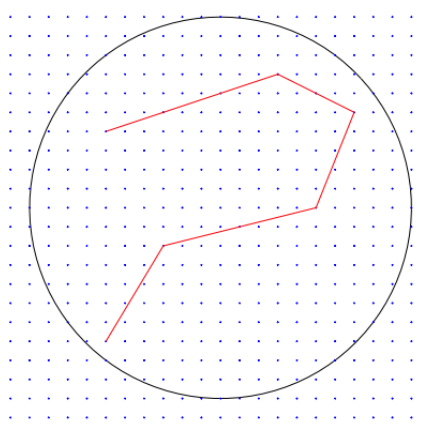


Figure 4: An example of a lattice  $\Gamma_\delta$  in  $\mathbb{R}^2$  with  $\delta = 1$  (spacing between blue points) and  $B(0, 10)$  (black circle). The red polygonal line is composed with vertices in  $\mathcal{Q} = B(0, 10) \cap \Gamma_\delta$ .

Our goal is to learn a time-dependent polygonal line which passes through the "middle" of data and gives a summary of all available observations  $x_1, \dots, x_{t-1}$  (denoted by  $(x_s)_{1:(t-1)}$  hereafter) before time  $t$ . Our output at time  $t$  is a polygonal line  $\hat{\mathbf{f}}_t \in \mathcal{F}_p$  depending on past information  $(x_s)_{1:(t-1)}$  and past predictions  $(\hat{\mathbf{f}}_s)_{1:(t-1)}$ . When  $x_t$  is revealed, the instantaneous loss at time  $t$  is computed as

$$\Delta(\hat{\mathbf{f}}_t, x_t) = \inf_{s \in I} \|\hat{\mathbf{f}}_t(s) - x_t\|_2^2. \quad (2)$$

In what follows, we investigate regret bounds for the cumulative loss based on (2). Given a measurable space  $\Theta$  (embedded with its Borel  $\sigma$ -algebra), we let  $\mathcal{P}(\Theta)$  denote the set of probability distributions on  $\Theta$ , and for some reference measure  $\pi$ , we let  $\mathcal{P}_\pi(\Theta)$  be the set of probability distributions absolutely continuous with respect to  $\pi$ .

For any  $k \in [1, p]$ , let  $\pi_k$  denote a probability distribution on  $\mathcal{F}_{k,L}$ . We define the *prior*  $\pi$  on  $\mathcal{F}_p = \cup_{k=1}^p \mathcal{F}_{k,L}$  as

$$\pi(\mathbf{f}) = \sum_{k \in [1, p]} w_k \pi_k(\mathbf{f}) \mathbb{1}_{\{\mathbf{f} \in \mathcal{F}_{k,L}\}}, \quad \mathbf{f} \in \mathcal{F}_p,$$

where  $w_1, \dots, w_p \geq 0$  and  $\sum_{k \in [1, p]} w_k = 1$ .

A quasi-Bayesian procedure would now consider the Gibbs quasi-posterior (note that this is not a proper posterior in all generality, hence the term "quasi")

$$\hat{\rho}_t(\cdot) \propto \exp(-\lambda S_t(\cdot)) \pi(\cdot),$$

where

$$S_t(\mathbf{f}) = S_{t-1}(\mathbf{f}) + \Delta(\mathbf{f}, x_t) + \frac{\lambda}{2} (\Delta(\mathbf{f}, x_t) - \Delta(\hat{\mathbf{f}}_t, x_t))^2,$$

as advocated by [Audibert \(2009\)](#) and [Li et al. \(2016\)](#) who then consider realisations from this quasi-posterior. In the present paper, we will rather consider a quantity linked to the mode of this quasi-posterior. Indeed, the mode of the quasi-posterior  $\hat{\rho}_{t+1}$  is

$$\arg \min_{\mathbf{f} \in \mathcal{F}_p} \left\{ \underbrace{\sum_{s=1}^t \Delta(\mathbf{f}, x_s)}_{(i)} + \underbrace{\frac{\lambda}{2} \sum_{s=1}^t (\Delta(\mathbf{f}, x_s) - \Delta(\hat{\mathbf{f}}_s, x_s))^2}_{(ii)} + \underbrace{\frac{\ln \pi(\mathbf{f})}{\lambda}}_{(iii)} \right\},$$

where (i) is a cumulative loss term, (ii) is a term controlling the variance of the prediction  $\mathbf{f}$  to past predictions  $\hat{\mathbf{f}}_s, s \leq t$ , and (iii) can be regarded as a penalty function on the complexity

of  $\mathbf{f}$  if  $\pi$  is well chosen. This mode hence has a similar flavor to follow the best expert or follow the perturbed leader in the setting of prediction with experts (see [Hutter and Poland, 2005](#) or [Cesa-Bianchi and Lugosi, 2006](#), Chapters 3 and 4) if we consider each  $\mathbf{f} \in \mathcal{F}_p$  as an expert which always delivers constant advice. These remarks yield [Algorithm 1](#).

---

**Algorithm 1** An algorithm to sequentially learn principal curves

---

- 1: **Input parameters:**  $p > 0, \eta > 0, \pi(z) = e^{-z} \mathbb{1}_{\{z > 0\}}$  and penalty function  $h : \mathcal{F}_p \rightarrow \mathbb{R}^+$
- 2: **Initialization:** For each  $\mathbf{f} \in \mathcal{F}_p$ , draw  $z_{\mathbf{f}} \sim \pi$  and  $\Delta_{\mathbf{f},0} = \frac{1}{\eta}(h(\mathbf{f}) - z_{\mathbf{f}})$
- 3: **For**  $t = 1, \dots, T$
- 4:     Get the data  $x_t$
- 5:     Obtain

$$\hat{\mathbf{f}}_t = \operatorname{arginf}_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{s=0}^{t-1} \Delta_{\mathbf{f},s} \right\},$$

where  $\Delta_{\mathbf{f},s} = \Delta(\mathbf{f}, x_s), s \geq 1$ .

- 6: **End for**
- 

### 3 Regret bounds for sequential learning of principal curves

We now present our main theoretical results.

**Theorem 1.** *For any sequence  $(x_t)_{1:T} \in B(0, \sqrt{d}R)$ ,  $R \geq 0$  and any penalty function  $h : \mathcal{F}_p \rightarrow \mathbb{R}^+$ , let  $\pi(z) = e^{-z} \mathbb{1}_{\{z > 0\}}$ . Let  $0 < \eta \leq \frac{1}{d(2R+\delta)^2}$ , then the procedure described in [Algorithm 1](#) satisfies*

$$\sum_{t=1}^T \mathbb{E}_{\pi} [\Delta(\hat{\mathbf{f}}_t, x_t)] \leq (1 + c_0(e-1)\eta) S_{T,h,\eta} + \frac{1}{\eta} \left( 1 + \ln \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} \right),$$

where  $c_0 = d(2R + \delta)^2$  and

$$S_{T,h,\eta} = \inf_{k \in [1,p]} \left\{ \inf_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{h(\mathbf{f})}{\eta} \right\} \right\}.$$

The expectation of the cumulative loss of polygonal lines  $\hat{\mathbf{f}}_1, \dots, \hat{\mathbf{f}}_T$  is upper-bounded by the smallest penalised cumulative loss over all  $k \in \{1, \dots, p\}$  up to a multiplicative term  $(1 + c_0(e-1)\eta)$  which can be made arbitrarily close to 1 by choosing a small enough  $\eta$ . However, this will lead to both a large  $h(\mathbf{f})/\eta$  in  $S_{T,h,\eta}$  and a large  $\frac{1}{\eta}(1 + \ln \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})})$ . In addition, another important issue is the choice of the penalty function  $h$ . For each  $\mathbf{f} \in \mathcal{F}_p$ ,  $h(\mathbf{f})$  should be large enough to ensure a small  $\sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})}$  while not too large to avoid overpenalization and a larger value for  $S_{T,h,\eta}$ . We therefore set

$$h(\mathbf{f}) \geq \ln(pe) + \ln \left| \{ \mathbf{f} \in \mathcal{F}_p, \mathcal{K}(\mathbf{f}) = k \} \right| \quad (3)$$

for each  $\mathbf{f}$  with  $k$  segments (where  $|M|$  denotes the cardinality of a set  $M$ ) since it leads to

$$\sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} = \sum_{k \in [1,p]} \sum_{\substack{\mathbf{f} \in \mathcal{F}_p \\ \mathcal{K}(\mathbf{f})=k}} e^{-h(\mathbf{f})} \leq \sum_{k \in [1,p]} \frac{1}{pe} \leq \frac{1}{e}.$$

The penalty function  $h(\mathbf{f}) = c_1 \mathcal{K}(\mathbf{f}) + c_2 L + c_3$  satisfies (3), where  $c_1, c_2, c_3$  are constants depending on  $R, d, \delta, p$  (this is proven in [Lemma 3](#)). We therefore obtain the following corollary.

**Corollary 1.** Under the assumptions of [Theorem 1](#), let

$$\eta = \min \left\{ \frac{1}{d(2R + \delta)^2}, \sqrt{\frac{c_1 p + c_2 L + c_3}{c_0(e-1) \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t)}} \right\}.$$

Then

$$\sum_{t=1}^T \mathbb{E} [\Delta(\hat{\mathbf{f}}_t, x_t)] \leq \inf_{k \in [1, p]} \left\{ \inf_{\substack{\mathbf{f} \in \mathcal{F}_p \\ \mathcal{K}(\mathbf{f})=k}} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \sqrt{c_0(e-1)r_{T,k,L}} \right\} \right\} \\ + \sqrt{c_0(e-1)r_{T,p,L} + c_0(e-1)(c_1 p + c_2 L + c_3)},$$

where  $r_{T,k,L} = \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t)(c_1 k + c_2 L + c_3)$ .

*Proof.* Note that

$$\sum_{t=1}^T \mathbb{E} [\Delta(\hat{\mathbf{f}}_t, x_t)] \leq S_{T,h,\eta} + \eta c_0(e-1) \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + c_0(e-1)(c_1 p + c_2 L + c_3),$$

and we conclude by setting

$$\eta = \sqrt{\frac{c_1 p + c_2 L + c_3}{c_0(e-1) \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t)}}.$$

□

However [Corollary 1](#) is not usable in practice since the optimal value for  $\eta$  depends on  $\inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t)$  which is obviously unknown, even more so at time  $t = 0$ . We therefore provide an adaptive refinement of [Algorithm 1](#) in the following [Algorithm 2](#).

---

**Algorithm 2** An adaptive algorithm to sequentially learn principal curves

---

- 1: **Input parameters:**  $p > 0$ ,  $L > 0$ ,  $\pi$ ,  $h$  and  $\eta_0 = \frac{\sqrt{c_1 p + c_2 L + c_3}}{c_0 \sqrt{e-1}}$
- 2: **Initialization:** For each  $\mathbf{f} \in \mathcal{F}_p$ , draw  $z_{\mathbf{f}} \sim \pi$ ,  $\Delta_{\mathbf{f},0} = \frac{1}{\eta_0}(h(\mathbf{f}) - z_{\mathbf{f}})$  and  $\hat{\mathbf{f}}_0 = \operatorname{arginf}_{\mathbf{f} \in \mathcal{F}_p} \Delta_{\mathbf{f},0}$
- 3: **For**  $t = 1, \dots, T$
- 4:     Compute  $\eta_t = \frac{\sqrt{c_1 p + c_2 L + c_3}}{c_0 \sqrt{(e-1)t}}$
- 5:     Get data  $x_t$  and compute  $\Delta_{\mathbf{f},t} = \Delta(\mathbf{f}, x_t) + \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}}\right)(h(\mathbf{f}) - z_{\mathbf{f}})$
- 6:     Obtain

$$\hat{\mathbf{f}}_t = \operatorname{arginf}_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{s=0}^{t-1} \Delta_{\mathbf{f},s} \right\}. \quad (4)$$

- 7: **End for**
- 

**Theorem 2.** For any sequence  $(x_t)_{1:T} \in B(0, \sqrt{d}R)$ ,  $R \geq 0$ , let  $h(\mathbf{f}) = c_1 \mathcal{K}(\mathbf{f}) + c_2 L + c_3$  where  $c_1, c_2, c_3$  are constants depending on  $R, d, \delta, \ln p$ . Let  $\pi(z) = e^{-z} \mathbb{1}_{\{z > 0\}}$  and

$$\eta_0 = \frac{\sqrt{c_1 p + c_2 L + c_3}}{c_0 \sqrt{e-1}}, \quad \eta_t = \frac{\sqrt{c_1 p + c_2 L + c_3}}{c_0 \sqrt{(e-1)t}}, \quad t \geq 1,$$

where  $c_0 = d(2R + \delta)^2$ . Then the procedure described in [Algorithm 2](#) satisfies

$$\sum_{t=1}^T \mathbb{E} [\Delta(\hat{\mathbf{f}}_t, x_t)] \leq \inf_{k \in [1, p]} \left\{ \inf_{\substack{\mathbf{f} \in \mathcal{F}_p \\ \mathcal{K}(\mathbf{f})=k}} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + c_0 \sqrt{(e-1)T(c_1 k + c_2 L + c_3)} \right\} \right\} \\ + 2c_0 \sqrt{(e-1)T(c_1 p + c_2 L + c_3)}.$$

The message of this regret bound is that the expected cumulative loss of polygonal lines  $\hat{\mathbf{f}}_1, \dots, \hat{\mathbf{f}}_T$  is upper-bounded by the minimal cumulative loss over all  $k \in \{1, \dots, p\}$ , up to an additive term which is sublinear in  $T$ . The actual magnitude of this remainder term is  $\sqrt{kT}$ . When  $L$  is fixed, the number  $k$  of segments is a measure of complexity of the retained polygonal line. This bound therefore yields the same magnitude than (1) which is the most refined bound in the literature so far (Biau and Fischer, 2012, where the optimal values for  $k$  and  $L$  are obtained in a model selection fashion).

## 4 Implementation

The argument of the infimum in Algorithm 2 is taken over  $\mathcal{F}_p = \cup_{k=1}^p \mathcal{F}_{k,L}$  which has a cardinality of order  $|\mathcal{Q}_\delta|^p$ , making any greedy search largely time-consuming. We instead turn to the following strategy: given a polygonal line  $\hat{\mathbf{f}}_t \in \mathcal{F}_{k_t,L}$  with  $k_t$  segments, we consider, with a certain proportion, the availability of  $\hat{\mathbf{f}}_{t+1}$  within a neighbourhood  $\mathcal{U}(\hat{\mathbf{f}}_t)$  (see the formal definition below) of  $\hat{\mathbf{f}}_t$ . This consideration is suited for principal curve setting since if observation  $x_t$  is close to  $\hat{\mathbf{f}}_t$ , one can expect that the polygonal line which well fits observations  $x_s, s = 1, \dots, t$  lies in a neighbourhood of  $\hat{\mathbf{f}}_t$ . In addition, if each polygonal line  $\mathbf{f}$  is regarded as an action, we no longer assume that all actions are available at all times, and allow the set of available actions to vary at each time. This is a model known as "sleeping experts (or actions)" in prior work (Auer et al., 2003; Kleinberg et al., 2008). In this setting, defining the regret with respect to the best action in the whole set of actions in hindsight remains to be difficult since that action might sometimes be unavailable. Hence it is natural to define the regret with respect to the best ranking of all actions in the hindsight according to their losses or rewards, and at each round one chooses among the available actions by selecting the one which ranks the highest in this ordering. Kleinberg et al. (2008) introduced this notion of regret and studied both the full-information (best action) and partial-information (multi-armed bandit) settings with stochastic and adversarial rewards and adversarial action availability. They pointed out that the EXP4 algorithm (Auer et al., 2003) attains the optimal regret in adversarial rewards case but runs in time exponential in the number of all actions. Kanade et al. (2009) considered full and partial information with stochastic action availability and proposed an algorithm that runs in polygonal time. In what follows, we will realise our implementation by resorting to "sleeping experts" *i.e.*, a special available set of actions that adapts to the setting of principal curves.

Let  $\sigma$  denote an ordering of  $|\mathcal{F}_p|$  actions, and  $\mathcal{A}_t$  an available subset of the actions at round  $t$ . We let  $\sigma(\mathcal{A}_t)$  denote the highest ranked action in  $\mathcal{A}_t$ . In addition, for any action  $\mathbf{f} \in \mathcal{F}_p$  we define the reward  $r_{\mathbf{f},t}$  of  $\mathbf{f}$  at round  $t, t \geq 0$  by

$$r_{\mathbf{f},t} = c_0 - \Delta(\mathbf{f}, x_t).$$

It is clear that  $r_{\mathbf{f},t} \in (0, c_0)$ . The convention from losses to gains is done in order to facilitate the subsequent performance analysis. The reward of an ordering  $\sigma$  is the cumulative reward of the selected action at each time

$$\sum_{t=1}^T r_{\sigma(\mathcal{A}_t),t},$$

and the reward of the best ordering is  $\max_{\sigma} \sum_{t=0}^T r_{\sigma(\mathcal{A}_t),t}$  ( $\mathbb{E}[\max_{\sigma} \sum_{t=1}^T r_{\sigma(\mathcal{A}_t),t}]$  when  $\mathcal{A}_t$  is stochastic).

Our procedure starts with a **partition** step which aims at identifying the "relevant" neighbourhood of an observation  $x \in \mathbb{R}^d$  with respect to a given polygonal line, and then proceeds with the definition of **neighbourhood** of an action  $\mathbf{f}$ , and finally we give details of implementation and show its regret bound.



**Partition** For any polygonal line  $\mathbf{f}$  with  $k$  segments, we denote by  $\bar{\mathbf{V}} = (v_1, \dots, v_{k+1})$  its vertices and by  $s_i, i = 1, \dots, k$  the line segments connecting  $v_i$  and  $v_{i+1}$ . In the sequel, we use  $\mathbf{f}(\bar{\mathbf{V}})$  to represent the polygonal line formed by connecting consecutive vertices in  $\bar{\mathbf{V}}$  if no confusion is possible. In addition, for  $x \in \mathbb{R}^d$ , let  $\Delta(x, v_i) = \|x - v_i\|_2^2$  be the squared distance between  $x$  and  $v_i$ , and let  $\Delta(x, s_i)$  be the squared distance between  $x$  and line segment  $s_i$ , *i.e.*,

$$\Delta(x, s_i) = \begin{cases} \|x - v_i\|_2^2 & \text{if } s_{s_i}(x) = v_i, \\ \|x - v_{i+1}\|_2^2 & \text{if } s_{s_i}(x) = v_{i+1}, \\ \|x - v_i\|_2^2 - \left( (x - v_i)^T \frac{v_{i+1} - v_i}{\|v_{i+1} - v_i\|_2} \right)^2 & \text{otherwise,} \end{cases}$$

where  $s_{s_i}(x)$  is the projection index of  $x$  to line segment  $s_i$ . In this step,  $\mathbb{R}^d$  can be partitioned into at most  $(2k + 1)$  disjoint sets  $V_i, i = 1, \dots, k + 1$  and  $S_i, i = 1, \dots, k$  defined as

$$V_i = \left\{ x \in \mathbb{R}^d : \Delta(x, v_i) = \Delta(\mathbf{f}, x), x \notin \cup_{s=1}^{i-1} V_s \right\},$$

$$S_i = \left\{ x \in \mathbb{R}^d : \Delta(x, s_i) = \Delta(\mathbf{f}, x), x \notin \cup_{m=1}^{i-1} S_m \cup_{i=1}^{k+1} V_i \right\},$$

where  $V_i, i = 1, \dots, k + 1$  is the set of all  $x \in \mathbb{R}^d$  whose closest vertex of  $\mathbf{f}$  is  $v_i$ , and  $S_i, i = 1, \dots, k$  is the set of all  $x \in \mathbb{R}^d$  whose closest segment of  $\mathbf{f}$  is  $s_i$ . In the sequel, we use  $V_{i:j}, i \leq j$  (resp.,  $S_{i:j}$  and  $v_{i:j}$ ) to abbreviate the sequence of partitions  $V_i, \dots, V_j$  (resp.,  $S_i, \dots, S_j$  and  $v_i, \dots, v_j$ ).

For any  $x \in \mathbb{R}^d$ , let us denote by  $\mathcal{V}(x)$  the set of vertices of  $\mathbf{f}$  connecting to the projection  $\mathbf{f}(s_{\mathbf{f}}(x))$  and denote by  $\mathcal{N}(x)$  the neighbourhood  $x$  with respect to  $\mathbf{f}$ , *i.e.*,

$$\mathcal{V}(x) = \begin{cases} v_{1:2} & \text{if } x \in V_1 \cup S_1, \\ v_{1:3} & \text{if } x \in V_2, \\ v_{i:i+1} & \text{if } x \in S_i, i = 2, \dots, k-1, \\ v_{i-1:i+1} & \text{if } x \in V_i, i = 3, \dots, k-1, \\ v_{k-1:k+1} & \text{if } x \in V_k, \\ v_{k:k+1} & \text{if } x \in S_k \cup V_{k+1}, \end{cases} \quad \mathcal{N}(x) = \begin{cases} V_{1:2} \cup S_{1:2} & \text{if } x \in V_1 \cup S_1, \\ V_{1:3} \cup S_{1:3} & \text{if } x \in V_2, \\ V_{i:i+1} \cup S_{i-1:i+1} & \text{if } x \in S_i, i = 2, \dots, k-1, \\ V_{i-1:i+1} \cup S_{i-2:i+1} & \text{if } x \in V_i, i = 3, \dots, k-1, \\ V_{k-1:k+1} \cup S_{k-2:k} & \text{if } x \in V_k, \\ V_{k:k+1} \cup S_{k-1:k} & \text{if } x \in S_k \cup V_{k+1}. \end{cases}$$

Note that those definitions of  $\mathcal{V}(x)$  and  $\mathcal{N}(x)$  hold whenever  $k \geq 4$ . If  $1 \leq k < 4$ ,  $\mathcal{V}(x)$  and  $\mathcal{N}(x)$  are identical. Next, let

$$\mathcal{N}_t(x) = \left\{ x_{s,,}, s = 1, \dots, t, x_s \in \mathcal{N}(x) \right\}$$

be the set of observations  $x_{1:t}$  belonging to  $\mathcal{N}(x)$ . We finally let  $\bar{\mathcal{N}}_t(x)$  stand for the average of all observations in  $\mathcal{N}_t(x)$  and  $\mathcal{D}(M) = \sup_{x,y \in M} \|x - y\|_2$  denotes the diameter of set  $M \in \mathbb{R}^d$ .

We initiate the principal curve  $\hat{\mathbf{f}}_1$  as the first component line segment whose vertices are the two farthest projections of data  $x_{1:t_0}$  on the first component line. Given  $\hat{\mathbf{f}}_t, t \geq 1$  with  $k_t$  segments, we obtain firstly the Voronoi-like partition  $V_{1:k_t+1}$  and  $S_{1:k_t}$ , then identify the adjacent vertices set  $\mathcal{V}(x_{t+t_0}) = \{v_{i_i:j_i}\}$  and set  $\mathcal{N}_{t+t_0}(x_{t+t_0})$  for observations  $x_{1:t+t_0}$  in  $\mathcal{N}(x_{t+t_0})$ .

**Neighbourhood** For each  $t$ , let us first define the local vertices set  $\mathcal{Q}_{\delta,t}(x)$  around  $x \in \mathbb{R}^d$  as

$$\mathcal{Q}_{\delta,t}(x) = B(\bar{\mathcal{N}}_t(x), \mathcal{D}(\mathcal{N}_t(x))) \cap \mathcal{Q}_{\delta},$$

where  $B(x, r)$  is a  $\ell_2$ -ball in  $\mathbb{R}^d$ , centered in  $x$  with radius  $r > 0$ . Then for  $k \in \{k_t - 1, k_t, k_t + 1\}$ , a candidate  $\bar{\mathbf{V}}$  (equivalently a candidate  $\bar{\mathbf{f}}(\bar{\mathbf{V}})$ ) with  $k + 1$  vertices is formed as

$$\bar{\mathbf{V}} = (v_{1:i_t-1}, v_{1:m}, v_{j_t+1:k_t+1}),$$

where  $m = k + 1 - k_t + j_t - i_t$  and  $v_{1:m}$  are distinct  $m$  vertices chosen from the local vertices set  $\mathcal{Q}_{\delta, t+t_0}(x_{t+t_0})$  around  $x_{t+t_0}$ . In other words, a candidate  $\bar{\mathbf{V}}$  is built by keeping all the vertices of  $\hat{\mathbf{f}}_t$  that do not belong to  $\mathcal{V}(x_{t+t_0})$  and by updating  $m$  vertices drawn from  $\mathcal{Q}_{\delta, t+t_0}(x_{t+t_0})$ . We define the neighbourhood  $\mathcal{U}(\hat{\mathbf{f}}_t)$  of  $\hat{\mathbf{f}}_t$  by

$$\mathcal{U}(\hat{\mathbf{f}}_t) = \left\{ \bar{\mathbf{f}}(\bar{\mathbf{V}}), \bar{\mathbf{V}} \text{ such that } v_{1:m} \in \mathcal{Q}_{\delta, t+t_0}(x_{t+t_0}) \text{ for } k \in \{k_t - 1, k_t, k_t + 1\} \right\}. \quad (5)$$

The cardinality  $|\mathcal{U}(\hat{\mathbf{f}}_t)|$  is upper bounded by  $|\mathcal{F}_p|^{\frac{3}{p}}$  since all elements in  $|\mathcal{U}(\hat{\mathbf{f}}_t)|$  differ with  $\hat{\mathbf{f}}_t$  up to at most three vertices.

Finally, we give our implementation with action availability that may vary at each time in [Algorithm 3](#).

---

**Algorithm 3** A locally greedy algorithm to sequentially learn principal curves

---

1: **Input parameters:**  $p > 0, R > 0, L > 0, \epsilon > 0, \alpha > 0, 1 > \beta > 0$  and any penalty function  $h$

2: **Initialization:** Given  $(x_t)_{1:t_0}$ , obtain  $\hat{\mathbf{f}}_1$  as the first principal component

3: **For**  $t = 2, \dots, T$

4: Draw  $I_t \sim \text{Bernoulli}(\epsilon)$  and  $z_{\mathbf{f}} \sim \pi$ .

5: Let  $\hat{\sigma}_t = \text{sort}\left(\mathbf{f}, \sum_{s=1}^{t-1} \hat{r}_{\mathbf{f},s} - \frac{1}{\eta_{t-1}} h(\mathbf{f}) + \frac{1}{\eta_{t-1}} z_{\mathbf{f}}\right)$ , i.e., descending sorting all  $\mathbf{f} \in \mathcal{F}_p$  according to their perturbed cumulative reward till  $t - 1$ .

6: If  $I_t = 1$ , set  $\mathcal{A}_t = \mathcal{F}_p$  and  $\hat{\mathbf{f}}_t = \hat{\sigma}^t(\mathcal{A}_t)$  and observe  $r_{\hat{\mathbf{f}}_t, t}$

7:

$$\hat{r}_{\mathbf{f}, t} = r_{\mathbf{f}, t} \text{ for } \mathbf{f} \in \mathcal{F}_p.$$

8: If  $I_t = 0$ , set  $\mathcal{A}_t = \mathcal{U}(\hat{\mathbf{f}}_{t-1})$ ,  $\hat{\mathbf{f}}_t = \hat{\sigma}^t(\mathcal{A}_t)$  and observe  $r_{\hat{\mathbf{f}}_t, t}$

9:

$$\hat{r}_{\mathbf{f}, t} = \begin{cases} \frac{r_{\mathbf{f}, t}}{\mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t)} & \text{if } \mathbf{f} \in \mathcal{U}(\hat{\mathbf{f}}_{t-1}) \cap \text{cond}(t) \text{ and } \hat{\mathbf{f}}_t = \mathbf{f}, \\ \alpha & \text{otherwise,} \end{cases}$$

where  $\mathcal{H}_t$  denotes all the randomness before time  $t$  and  $\text{cond}(t) = \{\mathbf{f} \in \mathcal{F}_p : \mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t) > \beta\}$ . In particular, when  $t = 1$ , we set  $\hat{r}_{\mathbf{f}, 1} = r_{\mathbf{f}, 1}$  for all  $\mathbf{f} \in \mathcal{F}_p$ ,  $\mathcal{U}(\hat{\mathbf{f}}_0) = \emptyset$  and  $\hat{r}_{\hat{\sigma}^1(\mathcal{U}(\hat{\mathbf{f}}_0)), 1} \equiv 0$ .

10: **End for**

---

The [Algorithm 3](#) has an exploration phase (when  $I_t = 1$ ) and an exploitation phase ( $I_t = 0$ ). In the exploration phase, it is allowed to observe rewards of all actions and to choose an optimal perturbed action from the set  $\mathcal{F}_p$  of all actions. In the exploitation phase, only rewards of a part of actions can be accessed to and rewards of others are estimated by a constant, and we update our action from the neighbourhood  $\mathcal{U}(\hat{\mathbf{f}}_{t-1})$  of the previous action

$\hat{\mathbf{f}}_{t-1}$ . This local update (or search) can hence reduce computation complexity. In addition, this local search will be enough to account for the case when  $x_t$  locates in  $\mathcal{U}(\hat{\mathbf{f}}_{t-1})$ . The parameter  $\beta$  needs to be carefully calibrated since on the one hand,  $\beta$  should not be too big to make sure that the condition  $\text{cond}(t)$  is non-empty, otherwise, all rewards are estimated by the same constant and thus leading to the same descending ordering of tuples for both  $(\sum_{s=1}^{t-1} \hat{r}_{\mathbf{f},s}, \mathbf{f} \in \mathcal{F}_p)$  and  $(\sum_{s=1}^t \hat{r}_{\mathbf{f},s}, \mathbf{f} \in \mathcal{F}_p)$ . Therefore, we may face the risk of having  $\hat{\mathbf{f}}_{t+1}$  in the neighbourhood of  $\hat{\mathbf{f}}_t$  even if we are in the exploration phase at time  $t+1$ ; on the other hand, very small  $\beta$  could result in large bias for estimation  $\frac{r_{\mathbf{f},t}}{\mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t)}$  of  $r_{\mathbf{f},t}$ . In addition, the exploitation phase is similar but different to the label efficient prediction (Cesa-Bianchi et al., 2005, Remark 1.1) since we allow an action at time  $t$  to be different from the previous one. Moreover, Neu and Bartók (2013) have proposed the *Geometric Resampling* method to estimate the conditional probability  $\mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t)$  since this quantity often doesn't have an explicit form. However, due to the simple exponential distribution of  $z_{\mathbf{f}}$  chosen in our case, an explicit form of  $\mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t)$  is straightforward.

**Theorem 3.** Assume that  $p > 6$ ,  $T \geq 2|\mathcal{F}_p|^2$  and let

$$\begin{aligned} \beta &= |\mathcal{F}_p|^{-\frac{1}{2}} T^{-\frac{1}{4}}, & \alpha &= \frac{c_0}{\beta}, & \hat{c}_0 &= \frac{2c_0}{\beta}, \\ \eta_1 = \eta_2 = \dots = \eta_T &= \frac{\sqrt{c_1 p + c_2 L + c_3}}{\sqrt{T(e-1)\hat{c}_0}}, & \epsilon &= 1 - |\mathcal{F}_p|^{\frac{1}{2} - \frac{3}{p}} T^{-\frac{1}{4}}. \end{aligned}$$

Then the procedure described in Algorithm 3 satisfies the regret bound

$$\sum_{t=1}^T \mathbb{E}[\Delta(\hat{\mathbf{f}}_t, x_t)] \leq \inf_{\mathbf{f} \in \mathcal{F}_p} \mathbb{E} \left[ \sum_{t=1}^T \Delta(\mathbf{f}, t) \right] + \mathcal{O}(T^{\frac{3}{4}}).$$

*Proof.* With the given value of parameters, the assumption of Lemma 4, Lemma 5 and Lemma 6 are satisfied; Combining their inequalities lead to the following

$$\begin{aligned} \sum_{t=1}^T \mathbb{E} \left[ r_{\hat{\mathbf{f}}_t, t} \right] &\geq \mathbb{E} \left[ \max_{\sigma} \left\{ \sum_{t=1}^T r_{\sigma(\mathcal{A}_t), t} - \frac{1}{\eta} h(\sigma(\mathcal{A}_t)) \right\} \right] - 2\alpha\beta(1-\epsilon) \sum_{t=1}^T |\mathcal{U}(\hat{\mathbf{f}}_{t-1})| \\ &\quad - \hat{c}_0^2(e-1)\eta T - \hat{c}_0(e-1)(c_1 p + c_2 L + c_3) \\ &\quad - (1 - |\mathcal{F}_p|\beta) \sqrt{2T \left[ \frac{c_0^2}{\beta} + \alpha^2(1-\beta) + (c_0 + 2\alpha)^2 \right] \ln \left( \frac{1}{\beta} \right)} - |\mathcal{F}_p|\beta c_0 T \\ &\geq \mathbb{E} \left[ \max_{\sigma} \left\{ \sum_{t=1}^T r_{\sigma(\mathcal{A}_t), t} - \frac{1}{\eta} h(\sigma(\mathcal{A}_t)) \right\} \right] - (1-\epsilon) |\mathcal{F}_p|^{\frac{3}{p}} c_0 T \\ &\quad - \hat{c}_0^2(e-1)\eta T - \hat{c}_0(e-1)(c_1 p + c_2 L + c_3) \\ &\quad - (1 - |\mathcal{F}_p|\beta) \sqrt{2T \left[ \frac{c_0^2}{\beta} + \alpha^2(1-\beta) + (c_0 + 2\alpha)^2 \right] \ln \left( \frac{1}{\beta} \right)} - |\mathcal{F}_p|\beta c_0 T \\ &\geq \mathbb{E} \left[ \max_{\sigma} \left\{ \sum_{t=1}^T r_{\sigma(\mathcal{A}_t), t} - \frac{1}{\eta} h(\sigma(\mathcal{A}_t)) \right\} \right] - \mathcal{O} \left( |\mathcal{F}_p|^{\frac{1}{2}} T^{\frac{3}{4}} \right), \end{aligned}$$

where the second inequality is due to the fact that the cardinality  $|\mathcal{U}(\hat{\mathbf{f}}_{t-1})|$  is upper bounded by  $|\mathcal{F}_p|^{\frac{3}{p}}$  for  $t \geq 1$ . In addition, using the definition of  $r_{\mathbf{f},t}$  that  $r_{\mathbf{f},t} = c_0 - \Delta(\mathbf{f}, x_t)$  terminates the proof of Theorem 3.  $\square$

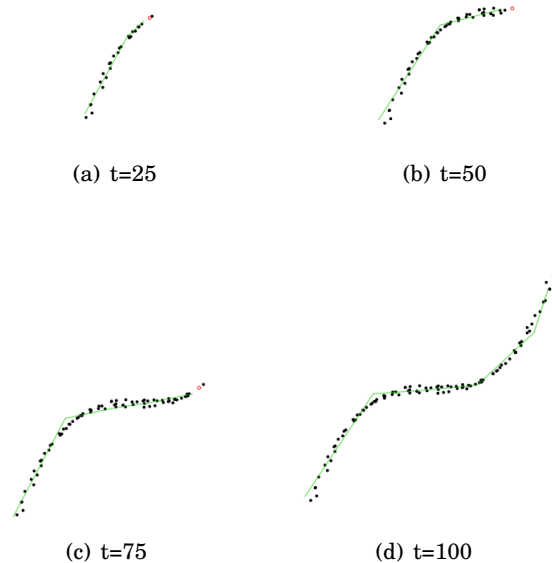


Figure 5: A principal curve (green) sequentially learned (black dots represent data  $x_{1:t}$ , red is the new observation  $x_{t+1}$ ).

One can see that the regret is upper bounded by a term of order  $\left(|\mathcal{F}_p|^{\frac{1}{2}} T^{\frac{3}{4}}\right)$ , sublinear in  $T$ .

The term  $(1-\epsilon)c_0T = c_0|\mathcal{F}_p|^{\frac{1}{2}} T^{\frac{3}{4}}$  is the price to pay for the local search (with a proportion  $1-\epsilon$ ) of polygonal line  $\hat{\mathbf{f}}_t$  in neighbourhood of the previous  $\hat{\mathbf{f}}_{t-1}$ . If  $\epsilon = 1$ , we would have that  $\hat{c}_0 = c_0$  and the last two terms in the first inequality of [Theorem 3](#) would disappear, hence the upper bound reduces to the one in [Theorem 2](#). In addition, our algorithm achieves an order that is smaller (from the point of view of both the number  $|\mathcal{F}_p|$  of all actions and the total rounds  $T$ ) than [Kanade et al. \(2009\)](#) since at each time, the availability of actions for our algorithm can be either the whole action set or a neighbourhood of the previous action while [Kanade et al. \(2009\)](#) considers at each time only partial and independent stochastic available set of actions generated from a predefined distribution.

## 5 Numerical experiments

Finally, we illustrate the performance of [Algorithm 3](#) on a toy example and seismic data (from [Engdahl and Villaseñor, 2002](#)).

**A toy example** The data set

$$\{x_t \in \mathbb{R}^2, t = 1, \dots, 100\}$$

attached to this example is generated uniformly along the curve  $y = -0.05 \times (x - 5)^3$ ,  $x \in [0, 10]$ . We set  $p = 20$ ,  $\delta = 0.5$ ,  $R = 10$ ,  $d = 2$ ,  $L = 0.01 \times p \times R \times \sqrt{d}$ . [Figure 5](#) presents the sequentially learned principal curve  $\hat{\mathbf{f}}_{t+1}$  at times  $t = 25, 50, 75, 100$ . A zoom on what happens between two consecutive time stamps is presented in [Figure 6](#). We see that collecting a new data point only impacts a local vertex, whereas the rest of the principal curve remains unchanged. The implementation is conducted in the R language and the code is available upon request.

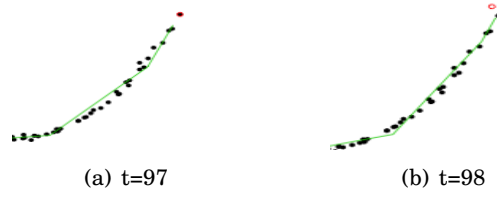


Figure 6: Zooming in: how a new data point changes the principal curve.

**Seismic data** Seismic data spanning long periods of time are essential for a thorough understanding of earthquakes. The “Centennial Earthquake Catalog” (Engdahl and Villaseñor, 2002) aims at providing a realistic picture of the seismicity distribution on Earth. It consists in a global catalog of locations and magnitudes of instrumentally recorded earthquakes from 1900 to 2008. Figure 7 is taken from the USGS website<sup>1</sup> and gives the global earthquakes locations on the period 1900–1999. The seismic data (latitude, longitude, magnitude of earthquakes, *etc.*) used in the present paper may be downloaded from this website. We focus on a particularly representative seismic active zone (a lithospheric border close to Australia) whose longitude is between E130° to E180° and latitude between S70° to N30°. The number of seismic recordings inside this area is  $T = 218$ . Figure 8 presents the learned principal curve: our procedure recovers nicely the tectonic plate boundary.

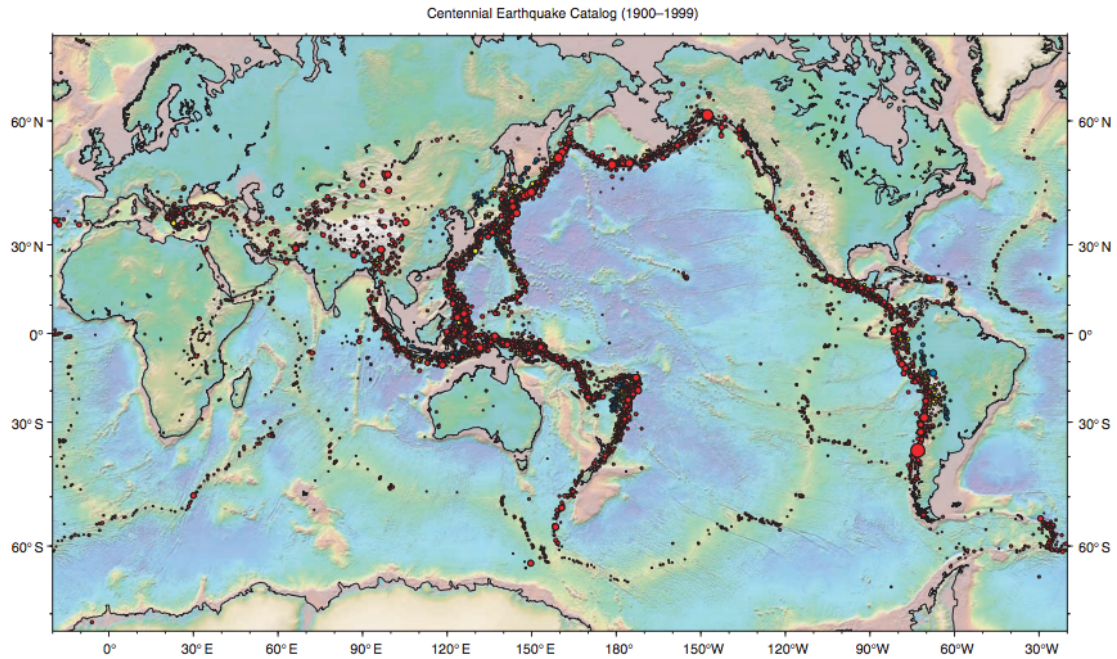


Figure 7: Seismic data from <https://earthquake.usgs.gov/data/centennial/>

## 6 Proofs

This section contains the proof of Theorem 2 (note that Theorem 1 is a straightforward consequence, with  $\eta_t = \eta$ ,  $t = 0, \dots, T$ ) and necessary lemmas for Theorem 3. Let us first

<sup>1</sup><https://earthquake.usgs.gov/data/centennial/>

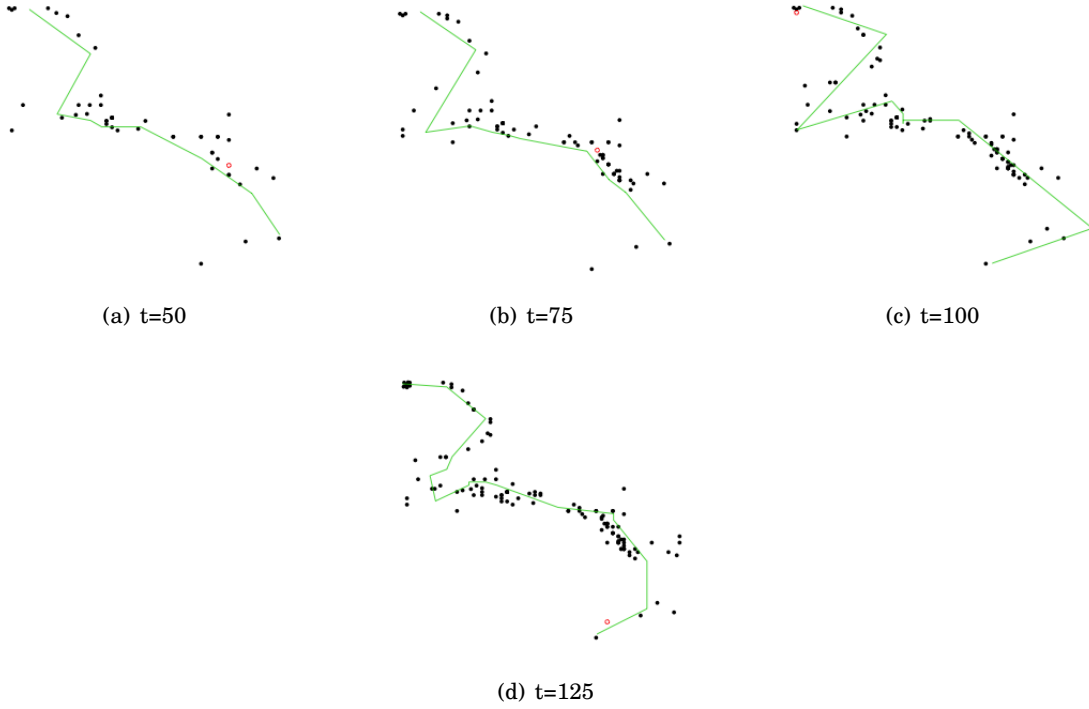


Figure 8: Seismic data: the learned principal curve (green) recovers the tectonic plate boundary (black dots represent data  $x_{1:t}$ , red is the new observation  $x_{t+1}$ ).

define for each  $t = 0, \dots, T$  the following forecaster sequence  $(\hat{\mathbf{f}}_t^*)_t$

$$\hat{\mathbf{f}}_0^* = \operatorname{arg\,inf}_{\mathbf{f} \in \mathcal{F}_p} \{\Delta_{\mathbf{f},0}\} = \operatorname{arg\,inf}_{\mathbf{f} \in \mathcal{F}_p} \left\{ \frac{1}{\eta_0} h(\mathbf{f}) - \frac{1}{\eta_0} z_{\mathbf{f}} \right\},$$

$$\hat{\mathbf{f}}_t^* = \operatorname{arg\,inf}_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{s=0}^t \Delta_{\mathbf{f},s} \right\} = \operatorname{arg\,inf}_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{s=1}^t \Delta(\mathbf{f}, x_s) + \frac{1}{\eta_{t-1}} h(\mathbf{f}) - \frac{1}{\eta_{t-1}} z_{\mathbf{f}} \right\}, \quad t \geq 1.$$

Note that  $\hat{\mathbf{f}}_t^*$  is an "illegal" forecaster since it peeks into the future. In addition, denote by

$$\mathbf{f}^* = \operatorname{arg\,inf}_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{1}{\eta_T} h(\mathbf{f}) \right\}$$

the polygonal line in  $\mathcal{F}_p$  which minimizes the cumulative loss in the first  $T$  rounds plus a penalty term.  $\mathbf{f}^*$  is deterministic while  $\hat{\mathbf{f}}_t^*$  is a random quantity (since it depends on  $z_{\mathbf{f}}$ ,  $\mathbf{f} \in \mathcal{F}_p$  drawn from  $\pi$ ). If several  $\mathbf{f}$  attain the infimum, we choose  $\mathbf{f}_T^*$  as the one having the smallest complexity. We now enunciate the first (out of three) intermediary technical result.

**Lemma 1.** For any sequence  $x_1, \dots, x_T$  in  $B(0, \sqrt{d}R)$ ,

$$\sum_{t=0}^T \Delta_{\hat{\mathbf{f}}_t^*, t} \leq \sum_{t=0}^T \Delta_{\mathbf{f}_T^*, t}, \quad \pi\text{-almost surely.} \quad (6)$$

*Proof.* Proof by induction on  $T$ . Clearly (6) holds for  $T = 0$ . Assume that (6) holds for  $T - 1$ :

$$\sum_{t=0}^{T-1} \Delta_{\hat{\mathbf{f}}_t^*, t} \leq \sum_{t=0}^{T-1} \Delta_{\hat{\mathbf{f}}_{T-1}^*, t}.$$

Adding  $\Delta_{\hat{\mathbf{f}}_T^*, T}$  to both sides of the above inequality concludes the proof.  $\square$

By (6) and the definition of  $\hat{\mathbf{f}}_T^*$ , for  $k \geq 1$ , we have  $\pi$ -almost surely that

$$\begin{aligned} \sum_{t=1}^T \Delta(\hat{\mathbf{f}}_t^*, x_t) &\leq \sum_{t=1}^T \Delta(\hat{\mathbf{f}}_T^*, x_t) + \frac{1}{\eta_T} h(\hat{\mathbf{f}}_T^*) - \frac{1}{\eta_T} Z_{\hat{\mathbf{f}}_T^*} + \sum_{t=0}^T \left( \frac{1}{\eta_{t-1}} - \frac{1}{\eta_t} \right) (h(\hat{\mathbf{f}}_t^*) - Z_{\hat{\mathbf{f}}_t^*}) \\ &\leq \sum_{t=1}^T \Delta(\mathbf{f}^*, x_t) + \frac{1}{\eta_T} h(\mathbf{f}^*) - \frac{1}{\eta_T} Z_{\mathbf{f}^*} + \sum_{t=0}^T \left( \frac{1}{\eta_{t-1}} - \frac{1}{\eta_t} \right) (h(\hat{\mathbf{f}}_t^*) - Z_{\hat{\mathbf{f}}_t^*}) \\ &= \inf_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{1}{\eta_T} h(\mathbf{f}) \right\} - \frac{1}{\eta_T} Z_{\mathbf{f}^*} + \sum_{t=0}^T \left( \frac{1}{\eta_{t-1}} - \frac{1}{\eta_t} \right) (h(\hat{\mathbf{f}}_t^*) - Z_{\hat{\mathbf{f}}_t^*}), \end{aligned}$$

where  $1/\eta_{-1} = 0$  by convention. The second and third inequality is due to respectively the definition of  $\hat{\mathbf{f}}_T^*$  and  $\mathbf{f}^*$ . Hence

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T \Delta(\hat{\mathbf{f}}_t^*, x_t) \right] &\leq \inf_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{1}{\eta_T} h(\mathbf{f}) \right\} - \frac{1}{\eta_T} \mathbb{E}[Z_{\mathbf{f}^*}] + \sum_{t=0}^T \mathbb{E} \left[ \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) (-h(\hat{\mathbf{f}}_t^*) + Z_{\hat{\mathbf{f}}_t^*}) \right] \\ &\leq \inf_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{1}{\eta_T} h(\mathbf{f}) \right\} + \sum_{t=1}^T \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \mathbb{E} \left[ \sup_{\mathbf{f} \in \mathcal{F}_p} (-h(\mathbf{f}) + Z_{\mathbf{f}}) \right] \\ &= \inf_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{1}{\eta_T} h(\mathbf{f}) \right\} + \frac{1}{\eta_T} \mathbb{E} \left[ \sup_{\mathbf{f} \in \mathcal{F}_p} (-h(\mathbf{f}) + Z_{\mathbf{f}}) \right], \end{aligned}$$

where the second inequality is due to  $\mathbb{E}[Z_{\mathbf{f}^*}] = 0$  and  $\left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) > 0$  for  $t = 0, 1, \dots, T$  since  $\eta_t$  is decreasing in  $t$  in [Theorem 2](#). In addition, for  $y \geq 0$ , one has

$$\mathbb{P}(-h(\mathbf{f}) + Z_{\mathbf{f}} > y) = e^{-h(\mathbf{f})-y}.$$

Hence, for any  $y \geq 0$

$$\mathbb{P} \left( \sup_{\mathbf{f} \in \mathcal{F}_p} (-h(\mathbf{f}) + Z_{\mathbf{f}}) > y \right) \leq \sum_{\mathbf{f} \in \mathcal{F}_p} \mathbb{P}(Z_{\mathbf{f}} \geq h(\mathbf{f}) + y) = \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} e^{-y} = u e^{-y},$$

where  $u = \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})}$ . Therefore, we have

$$\begin{aligned} \mathbb{E} \left[ \sup_{\mathbf{f} \in \mathcal{F}_p} (-h(\mathbf{f}) + Z_{\mathbf{f}}) - \ln u \right] &\leq \mathbb{E} \left[ \max \left( 0, \sup_{\mathbf{f} \in \mathcal{F}_p} (-h(\mathbf{f}) + Z_{\mathbf{f}} - \ln u) \right) \right] \\ &\leq \int_0^\infty \mathbb{P} \left( \max \left( 0, \sup_{\mathbf{f} \in \mathcal{F}_p} (-h(\mathbf{f}) + Z_{\mathbf{f}} - \ln u) \right) > y \right) dy \\ &\leq \int_0^\infty \mathbb{P} \left( \sup_{\mathbf{f} \in \mathcal{F}_p} (-h(\mathbf{f}) + Z_{\mathbf{f}}) > y + \ln u \right) dy \\ &\leq \int_0^\infty u e^{-(y+\ln u)} dy = 1. \end{aligned}$$

We thus obtain

$$\mathbb{E} \left[ \sum_{t=1}^T \Delta(\hat{\mathbf{f}}_t^*, x_t) \right] \leq \inf_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{1}{\eta_T} h(\mathbf{f}) \right\} + \frac{1}{\eta_T} \left( 1 + \ln \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} \right). \quad (7)$$

Next, we control the regret of [Algorithm 2](#).

**Lemma 2.** Assume that  $z_{\mathbf{f}}$  is sampled from the symmetric exponential distribution in  $\mathbb{R}$ , i.e.,  $\pi(z) = e^{-z} \mathbb{1}_{\{z>0\}}$ . Assume that  $\sup_{t=1, \dots, T} \eta_{t-1} \leq \frac{1}{d(2R+\delta)^2}$ , and define  $c_0 = d(2R+\delta)^2$ . Then for any sequence  $(x_t) \in B(0, \sqrt{d}R)$ ,  $t = 1, \dots, T$ ,

$$\sum_{t=1}^T \mathbb{E}[\Delta(\hat{\mathbf{f}}_t, x_t)] \leq \sum_{t=1}^T (1 + \eta_{t-1} c_0 (e-1)) \mathbb{E}[\Delta(\hat{\mathbf{f}}_t^*, x_t)]. \quad (8)$$

*Proof.* Let us denote by

$$F_t(\mathbf{Z}_{\mathbf{f}}) = \Delta(\hat{\mathbf{f}}_t, x_t) = \Delta\left(\operatorname{arg\,inf}_{\mathbf{f} \in \mathcal{F}} \left( \sum_{s=1}^{t-1} \Delta(\mathbf{f}, x_s) + \frac{1}{\eta_{t-1}} h(\mathbf{f}) - \frac{1}{\eta_{t-1}} \mathbf{Z}_{\mathbf{f}} \right), x_t\right)$$

the instantaneous loss suffered by the polygonal line  $\hat{\mathbf{f}}_t$  when  $x_t$  is obtained. We have

$$\begin{aligned} \mathbb{E}[\Delta(\hat{\mathbf{f}}_t^*, x_t)] &= \int F_t(z - \eta_{t-1} \Delta(\mathbf{f}, x_t)) \pi(z) dz \\ &= \int F_t(z) \pi(z + \eta_{t-1} \Delta(\mathbf{f}, x_t)) dz \\ &= \int F_t(z) e^{-(z + \eta_{t-1} \Delta(\mathbf{f}, x_t))} dz \\ &\geq e^{-\eta_{t-1} d(2R+\delta)^2} \int F_t(z) e^{-z} dz \\ &= e^{-\eta_{t-1} d(2R+\delta)^2} \mathbb{E}[\Delta(\hat{\mathbf{f}}_t, x_t)], \end{aligned}$$

where the inequality is due to the fact that  $\Delta(\mathbf{f}, x) \leq d(2R+\delta)^2$  holds uniformly for any  $\mathbf{f} \in \mathcal{F}_p$  and  $x \in B(0, \sqrt{d}R)$ . Finally, summing on  $t$  on both sides and using the elementary inequality  $e^x \leq 1 + (e-1)x$  if  $x \in (0, 1)$  concludes the proof.  $\square$

**Lemma 3.** For  $k \in \llbracket 1, p \rrbracket$ , we control the cardinality of set  $\{\mathbf{f} \in \mathcal{F}_p, \mathcal{K}(\mathbf{f}) = k\}$  as

$$\begin{aligned} \ln|\{\mathbf{f} \in \mathcal{F}_p, \mathcal{K}(\mathbf{f}) = k\}| &\leq \left( \ln(8peV_d) + 3d^{\frac{3}{2}} - d \right) k + \left( \frac{\ln 2}{\delta \sqrt{d}} + \frac{d}{\delta} \right) L + d \ln \left( \frac{\sqrt{d}(2R+\delta)}{\delta} \right) \\ &\triangleq c_1 k + c_2 L + c_3, \end{aligned}$$

where  $V_d$  denotes the volume of the unit ball in  $\mathbb{R}^d$ .

*Proof.* First, let  $N_{k,\delta}$  denote the set of polygonal lines with  $k$  segments and whose vertices are in  $\mathcal{Q}_\delta$ . Notice that  $N_{k,\delta}$  is different from  $\{\mathbf{f} \in \mathcal{F}_p, \mathcal{K}(\mathbf{f}) = k\}$  and that

$$|\{\mathbf{f} \in \mathcal{F}_p, \mathcal{K}(\mathbf{f}) = k\}| \leq \binom{p}{k} |N_{k,\delta}|.$$

Hence

$$\begin{aligned} \ln|\{\mathbf{f} \in \mathcal{F}_p, \mathcal{K}(\mathbf{f}) = k\}| &\leq \ln \binom{p}{k} + \ln |N_{k,\delta}| \\ &\leq k \ln \frac{pe}{k} + k \left( \ln 8V_d + 3d^{\frac{3}{2}} - d \right) + \left( \frac{\ln 2}{\sqrt{d}\delta} + \frac{d}{\delta} \right) L + d \ln \left( \frac{\sqrt{d}(2R+\delta)}{\delta} \right) \\ &\leq k \ln(pe) + k \left( \ln 8V_d + 3d^{\frac{3}{2}} - d \right) + \left( \frac{\ln 2}{\sqrt{d}\delta} + \frac{d}{\delta} \right) L + d \ln \left( \frac{\sqrt{d}(2R+\delta)}{\delta} \right), \end{aligned}$$

where the second inequality is a consequence to the elementary inequality  $\binom{p}{k} \leq \left(\frac{pe}{k}\right)^k$  combined with Lemma 2 in [Kégl \(1999\)](#).  $\square$



We now have all the ingredients to prove [Theorem 1](#) and [Theorem 2](#).

First, combining (7) and (8) yields that

$$\begin{aligned} \sum_{t=1}^T \mathbb{E} [\Delta(\hat{\mathbf{f}}_t, x_t)] &\leq \inf_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{1}{\eta_T} h(\mathbf{f}) \right\} + \frac{1}{\eta_T} \left( \frac{1}{2} + \ln \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} \right) \\ &\quad + c_0(e-1) \sum_{t=1}^T \eta_{t-1} \mathbb{E} [\Delta(\hat{\mathbf{f}}_t^*, x_t)] \\ &\leq \inf_{k \in [1, p]} \left\{ \inf_{\substack{\mathbf{f} \in \mathcal{F}_p \\ \mathcal{K}(\mathbf{f})=k}} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{h(\mathbf{f})}{\eta_T} \right\} \right\} + \frac{1}{\eta_T} \left( \frac{1}{2} + \ln \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} \right) \\ &\quad + c_0(e-1) \sum_{t=1}^T \eta_{t-1} \mathbb{E} [\Delta(\hat{\mathbf{f}}_t^*, x_t)]. \end{aligned}$$

Assume that  $\eta_t = \eta$ ,  $t = 0, \dots, T$  and  $h(\mathbf{f}) = c_1 \mathcal{K}(\mathbf{f}) + c_2 L + c_3$  for  $\mathbf{f} \in \mathcal{F}_p$ , then  $\left( \frac{1}{2} + \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} \right) \leq 0$  and moreover

$$\begin{aligned} \sum_{t=1}^T \mathbb{E} [\Delta(\hat{\mathbf{f}}_t, x_t)] &\leq S_{T, h, \eta} + \frac{1}{\eta} \left( \frac{1}{2} + \ln \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} \right) + c_0(e-1)\eta \sum_{t=1}^T \mathbb{E} [\Delta(\hat{\mathbf{f}}_t^*, x_t)] \\ &\leq S_{T, h, \eta} + c_0(e-1)\eta S_{T, h, \eta} \\ &\leq S_{T, h, \eta} + \eta c_0(e-1) \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + c_0(e-1)(c_1 p + c_2 L + c_3), \end{aligned}$$

where

$$S_{T, h, \eta} = \inf_{k \in [1, p]} \left\{ \inf_{\substack{\mathbf{f} \in \mathcal{F}_p \\ \mathcal{K}(\mathbf{f})=k}} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{h(\mathbf{f})}{\eta} \right\} \right\}$$

and the second inequality is obtained with [Lemma 1](#). By setting

$$\eta = \sqrt{\frac{c_1 p + c_2 L + c_3}{c_0(e-1) \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t)}}$$

we obtain

$$\begin{aligned} \sum_{t=1}^T \mathbb{E} [\Delta(\hat{\mathbf{f}}_t, x_t)] &\leq \inf_{k \in [1, p]} \left\{ \inf_{\substack{\mathbf{f} \in \mathcal{F}_p \\ \mathcal{K}(\mathbf{f})=k}} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \sqrt{c_0(e-1)r_{T, k, L}} \right\} \right\} \\ &\quad + \sqrt{c_0(e-1)L_{T, p, L} + c_0(e-1)c_1 p + c_2 L + c_3}, \end{aligned}$$

where  $r_{T, k, L} = \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t)(c_1 k + c_2 L + c_3)$ . This proves [Theorem 1](#).

Finally, assume that

$$\eta_0 = \frac{\sqrt{c_1 p + c_2 L + c_3}}{c_0 \sqrt{e-1}} \quad \text{and} \quad \eta_t = \frac{\sqrt{c_1 p + c_2 L + c_3}}{c_0 \sqrt{e-1} t}, \quad t = 1, \dots, T.$$

Since  $\mathbb{E} [\Delta(\hat{\mathbf{f}}_t^*, x_t)] \leq c_0$  for any  $t = 1, \dots, T$ , we have

$$\sum_{t=1}^T \mathbb{E} [\Delta(\hat{\mathbf{f}}_t, x_t)] \leq \inf_{k \in [1, p]} \left\{ \inf_{\substack{\mathbf{f} \in \mathcal{F}_p \\ \mathcal{K}(\mathbf{f})=k}} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{h(\mathbf{f})}{\eta_T} \right\} \right\} + \frac{1}{\eta_T} \left( 1 + \ln \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} \right) + c_0^2(e-1) \sum_{t=1}^T \eta_{t-1}$$

$$\leq \inf_{k \in [1, p]} \left\{ \inf_{\substack{\mathbf{f} \in \mathcal{F}_p \\ \mathcal{K}(\mathbf{f})=k}} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + c_0 \sqrt{(e-1)T(c_0 k + c_2 L + c_3)} \right\} \right\} \\ + 2c_0 \sqrt{(e-1)T(c_0 p + c_2 L + c_3)},$$

which concludes the proof of [Theorem 2](#).

**Lemma 4.** Under the [Algorithm 3](#), if  $1 \geq \epsilon > 0$ ,  $1 > \beta > 0$ ,  $\alpha \geq \frac{(1-\beta)c_0}{\beta}$  and  $|\mathcal{U}(\hat{\mathbf{f}}_{t-1})| \geq 2$  for all  $t \geq 2$ , where  $|\mathcal{U}(\hat{\mathbf{f}}_{t-1})|$  is the cardinality of  $\mathcal{U}(\hat{\mathbf{f}}_{t-1})$ , then we have

$$\sum_{t=1}^T \mathbb{E} \left[ r_{\hat{\mathbf{f}}_t, t} \right] \geq \sum_{t=1}^T \mathbb{E} \left[ \hat{r}_{\hat{\sigma}^t(\mathcal{A}_t), t} \right] - 2(1-\epsilon)\alpha\beta \sum_{t=1}^T |\mathcal{U}(\hat{\mathbf{f}}_{t-1})|.$$

*Proof.* First notice that  $\mathcal{A}_t = \mathcal{U}(\hat{\mathbf{f}}_{t-1})$  if  $I_t = 0$ , and that for  $t \geq 2$

$$\begin{aligned} \mathbb{E} \left[ r_{\hat{\mathbf{f}}_t, t} \middle| \mathcal{H}_t, I_t = 0 \right] &= \mathbb{E} \left[ r_{\hat{\sigma}^t(\mathcal{A}_t), t} \middle| \mathcal{H}_t, I_t = 0 \right] \\ &= \sum_{\mathbf{f} \in \mathcal{A}_t \cap \text{cond}(t)} r_{\mathbf{f}, t} \mathbb{P} \left( \hat{\sigma}^t(\mathcal{A}_t) = \mathbf{f} \middle| \mathcal{H}_t \right) + \sum_{\mathbf{f} \in \mathcal{A}_t \cap \text{cond}(t)^c} r_{\mathbf{f}, t} \mathbb{P} \left( \hat{\sigma}^t(\mathcal{A}_t) = \mathbf{f} \middle| \mathcal{H}_t \right) \\ &\geq \sum_{\mathbf{f} \in \mathcal{A}_t \cap \text{cond}(t)} r_{\mathbf{f}, t} + \sum_{\mathbf{f} \in \mathcal{A}_t \cap \text{cond}(t)^c} \alpha \mathbb{P} \left( \hat{\sigma}^t(\mathcal{A}_t) = \mathbf{f} \middle| \mathcal{H}_t \right) \\ &\quad - (1-\beta) \sum_{\mathbf{f} \in \mathcal{A}_t \cap \text{cond}(t)} r_{\mathbf{f}, t} - \sum_{\mathbf{f} \in \mathcal{A}_t \cap \text{cond}(t)^c} (\alpha - r_{\mathbf{f}, t}) \mathbb{P} \left( \hat{\sigma}^t(\mathcal{A}_t) = \mathbf{f} \middle| \mathcal{H}_t \right) \\ &= \mathbb{E} \left[ \hat{r}_{\hat{\sigma}^t(\mathcal{A}_t), t} \middle| \mathcal{H}_t, I_t = 0 \right] - (1-\beta) \sum_{\mathbf{f} \in \mathcal{A}_t \cap \text{cond}(t)} r_{\mathbf{f}, t} - \sum_{\mathbf{f} \in \mathcal{A}_t \cap \text{cond}(t)^c} (\alpha - r_{\mathbf{f}, t}) \mathbb{P} \left( \hat{\sigma}^t(\mathcal{A}_t) = \mathbf{f} \middle| \mathcal{H}_t \right) \\ &\geq \mathbb{E} \left[ \hat{r}_{\hat{\sigma}^t(\mathcal{A}_t), t} \middle| \mathcal{H}_t, I_t = 0 \right] - (1-\beta)c_0 |\mathcal{A}_t| - \alpha\beta |\mathcal{A}_t| \\ &\geq \mathbb{E} \left[ \hat{r}_{\hat{\sigma}^t(\mathcal{A}_t), t} \middle| \mathcal{H}_t, I_t = 0 \right] - 2\alpha\beta |\mathcal{A}_t|, \end{aligned}$$

where  $\text{cond}(t)^c$  denotes the complement of set  $\text{cond}(t)$ ; the first inequality above is due to the assumption that for all  $\mathbf{f} \in \mathcal{A}_t \cap \text{cond}(t)$ , we have  $\mathbb{P}(\hat{\sigma}^t(\mathcal{A}_t) = \mathbf{f} | \mathcal{H}_t) \geq \beta$ . For  $t = 1$ , the above inequality is trivial since  $\hat{r}_{\hat{\sigma}^1(\mathcal{U}(\hat{\mathbf{f}}_0)), 1} \equiv 0$  by its definition. Hence, for  $t \geq 1$ , one has

$$\begin{aligned} \mathbb{E} \left[ r_{\hat{\mathbf{f}}_t, t} \middle| \mathcal{H}_t \right] &= \epsilon \mathbb{E} \left[ r_{\hat{\sigma}^t(\mathcal{F}_p), t} \middle| \mathcal{H}_t, I_t = 1 \right] + (1-\epsilon) \mathbb{E} \left[ r_{\hat{\sigma}^t(\mathcal{A}_t), t} \middle| \mathcal{H}_t, I_t = 0 \right] \\ &\geq \mathbb{E} \left[ \hat{r}_{\hat{\mathbf{f}}_t, t} \middle| \mathcal{H}_t \right] - 2\alpha\beta |\mathcal{A}_t|. \end{aligned} \tag{9}$$

Summing on both side of inequality (9) over  $t$  terminates the proof of [Lemma 4](#).  $\square$

**Lemma 5.** Let  $\hat{c}_0 = \frac{c_0}{\beta} + \alpha$ . If  $0 < \eta_1 = \eta_2 = \dots = \eta_T = \eta < \frac{1}{\hat{c}_0}$ , then we have

$$\mathbb{E} \left[ \max_{\hat{\sigma}} \left\{ \sum_{t=1}^T \hat{r}_{\hat{\sigma}^t(\mathcal{A}_t), t} - \frac{1}{\eta} h(\hat{\sigma}(\mathcal{A}_t)) \right\} \right] - \sum_{t=1}^T \mathbb{E} \left[ \hat{r}_{\hat{\sigma}^t(\mathcal{A}_t), t} \right] \leq \hat{c}_0^2 (e-1)\eta T + \hat{c}_0 (e-1)(c_1 p + c_2 L + c_3).$$

*Proof.* By the definition of  $\hat{r}_{\mathbf{f}, t}$  in [Algorithm 3](#), for any  $\mathbf{f} \in \mathcal{F}_p$  and  $t \geq 1$ , we have uniformly

$$\hat{r}_{\mathbf{f}, t} \leq \max \left\{ \frac{r_{\mathbf{f}, t}}{\mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t)}, \alpha, r_{\mathbf{f}, t} \right\} \leq \max \left\{ \frac{c_0}{\beta}, \alpha \right\} \leq \hat{c}_0,$$

where in the second inequality we use that  $r_{\mathbf{f},t} \leq c_0$  for all  $\mathbf{f}$  and  $t$ , and that  $\mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t) \geq \beta$  when  $\mathbf{f} \in \mathcal{U}(\hat{\mathbf{f}}_{t-1}) \cap \text{cond}(t)$ . The rest of the proof is similar to that of [Lemma 1](#) and [Lemma 2](#). In fact, if we define by  $\hat{\Delta}(\mathbf{f}, x_t) = \hat{c}_0 - \hat{r}_{\mathbf{f},t}$ , then one can easily observe the following relation when  $I_t = 1$  (similar relation in the case that  $I_t = 0$ )

$$\begin{aligned} \hat{\mathbf{f}}_t &= \hat{\sigma}^t(\mathcal{F}_p) = \arg \max_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{s=1}^{t-1} \hat{r}_{\mathbf{f},s} + \frac{1}{\eta} (z_{\mathbf{f}} - h(\mathbf{f})) \right\} \\ &= \arg \min_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{s=1}^{t-1} \hat{\Delta}(\mathbf{f}, x_s) + \frac{1}{\eta} (h(\mathbf{f}) - z_{\mathbf{f}}) \right\}. \end{aligned}$$

Then applying [Lemma 1](#) and [Lemma 2](#) on this newly defined sequence  $\hat{\Delta}(\hat{\mathbf{f}}_t, x_t), t = 1, \dots, T$  leads to the result of [Lemma 5](#).  $\square$

The proof of upcoming [Lemma 6](#) requires the following submartingale inequality: let  $Y_0, \dots, Y_T$  be a sequence of random variable adapted to random events  $\mathcal{H}_0, \dots, \mathcal{H}_T$  such that for  $1 \leq t \leq T$ , the following three conditions hold

$$\mathbb{E}[Y_t | \mathcal{H}_t] \leq 0, \quad \text{Var}(Y_t | \mathcal{H}_t) \leq a^2, \quad Y_t - \mathbb{E}[Y_t | \mathcal{H}_t] \leq b.$$

Then for any  $\lambda > 0$ ,

$$\mathbb{P}\left(\sum_{t=1}^T Y_t > Y_0 + \lambda\right) \leq \exp\left(-\frac{\lambda^2}{2T(a^2 + b^2)}\right).$$

The proof can be found in [Chung and Lu \(2006, Theorem 7.3\)](#).

**Lemma 6.** Assume that  $0 < \beta < \frac{1}{|\mathcal{F}_p|}$ ,  $\alpha \geq \frac{c_0}{\beta}$  and  $\eta > 0$ , then we have

$$\begin{aligned} \mathbb{E}\left[\max_{\sigma} \left\{ \sum_{t=1}^T r_{\sigma(\mathcal{A}_t),t} - \frac{1}{\eta} h(\sigma(\mathcal{A}_t)) \right\}\right] - \mathbb{E}\left[\max_{\hat{\sigma}} \left\{ \sum_{t=1}^T \hat{r}_{\hat{\sigma}(\mathcal{A}_t),t} - \frac{1}{\eta} h(\hat{\sigma}(\mathcal{A}_t)) \right\}\right] \\ \leq (1 - |\mathcal{F}_p| \beta) \sqrt{2T \left[ \frac{c_0^2}{\beta} + \alpha^2(1 - \beta) + (c_0 + 2\alpha)^2 \right] \ln\left(\frac{1}{\beta}\right)} + |\mathcal{F}_p| \beta c_0 T. \end{aligned}$$

*Proof.* First, we have almost surely that

$$\max_{\sigma} \left\{ \sum_{t=1}^T r_{\sigma(\mathcal{A}_t),t} - \frac{1}{\eta} h(\sigma(\mathcal{A}_t)) \right\} - \max_{\hat{\sigma}} \left\{ \sum_{t=1}^T \hat{r}_{\hat{\sigma}(\mathcal{A}_t),t} - \frac{1}{\eta} h(\hat{\sigma}(\mathcal{A}_t)) \right\} \leq \max_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T (r_{\mathbf{f},t} - \hat{r}_{\mathbf{f},t}).$$

Denote by  $Y_{\mathbf{f},t} = r_{\mathbf{f},t} - \hat{r}_{\mathbf{f},t}$ . Since

$$\mathbb{E}\left[\hat{r}_{\mathbf{f},t} | \mathcal{H}_t\right] = \begin{cases} r_{\mathbf{f},t} + (1 - \epsilon)\alpha(1 - \mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t)) & \text{if } \mathbf{f} \in \mathcal{U}(\hat{\mathbf{f}}_{t-1}) \cap \text{cond}(t), \\ \epsilon r_{\mathbf{f},t} + (1 - \epsilon)\alpha & \text{otherwise,} \end{cases}$$

and  $\alpha > c_0 \geq r_{\mathbf{f},t}$  uniformly for any  $\mathbf{f}$  and  $t$ , then we have uniformly that  $\mathbb{E}[Y_t | \mathcal{H}_t] \leq 0$ , hence satisfying the first condition.

For the second condition, if  $\mathbf{f} \in \mathcal{U}(\hat{\mathbf{f}}_{t-1}) \cap \text{cond}(t)$ , then

$$\text{Var}(Y_t | \mathcal{H}_t) = \mathbb{E}\left[\hat{r}_{\mathbf{f},t}^2 | \mathcal{H}_t\right] - (\mathbb{E}[\hat{r}_{\mathbf{f},t} | \mathcal{H}_t])^2$$

$$\begin{aligned}
&\leq \epsilon r_{\mathbf{f},t}^2 + (1-\epsilon) \left[ \frac{r_{\mathbf{f},t}^2}{\mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t)} + \alpha (1 - \mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t)) \right] \\
&\quad - [r_{\mathbf{f},t} + (1-\epsilon)\alpha (1 - \mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t))]^2 \\
&\leq \frac{r_{\mathbf{f},t}^2}{\beta} + \alpha^2(1-\beta) \leq \frac{c_0^2}{\beta} + \alpha^2(1-\beta).
\end{aligned}$$

Similarly, for  $\mathbf{f} \notin \mathcal{U}(\hat{\mathbf{f}}_{t-1}) \cap \text{cond}(t)$ , one can have  $\text{Var}(Y_t | \mathcal{H}_t) \leq \alpha^2$ .

Moreover, for the third condition, since

$$\mathbb{E}[Y_{\mathbf{f},t} | \mathcal{H}_t] \geq -2\alpha,$$

then

$$Y_{\mathbf{f},t} - \mathbb{E}[Y_{\mathbf{f},t} | \mathcal{H}_t] \leq r_{\mathbf{f},t} + 2\alpha \leq c_0 + 2\alpha.$$

Setting  $\lambda = \sqrt{2T \left[ \frac{c_0^2}{\beta} + \alpha^2(1-\beta) + (c_0 + 2\alpha)^2 \right] \ln\left(\frac{1}{\beta}\right)}$  leads to

$$\mathbb{P}\left(\sum_{t=1}^T Y_{\mathbf{f},t} \geq \lambda\right) \leq \beta.$$

Hence the following inequality holds with probability  $1 - |\mathcal{F}_p| \beta$

$$\max_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T (r_{\mathbf{f},t} - \hat{r}_{\mathbf{f},t}) \leq \sqrt{2T \left[ \frac{c_0^2}{\beta} + \alpha^2(1-\beta) + (c_0 + 2\alpha)^2 \right] \ln\left(\frac{1}{\beta}\right)}.$$

Finally, noticing that  $\max_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T (r_{\mathbf{f},t} - \hat{r}_{\mathbf{f},t}) \leq c_0 T$  almost surely, we terminate the proof of [Lemma 6](#). □

## References

- J.-Y. Audibert. Fast learning rates in statistical inference through aggregation. *The Annals of Statistics*, 37(4):1591–1646, 2009. [5](#)
- P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal of Computing*, 32(1):48–77, 2003. [8](#)
- J. D. Banfield and A. E. Raftery. Ice floe identification in satellite images using mathematical morphology and clustering about principal curves. *Journal of the American Statistical Association*, 87(417):7–16, 1992. [1](#)
- A. Barron, L. Birgé, and P. Massart. Risk bounds for model selection via penalization. *Probability Theory and Related Fields*, 113:301–413, 1999. [3](#)
- G. Biau and A. Fischer. Parameter selection for principal curves. *IEEE Transactions on Information Theory*, 58(3):1924–1939, 2012. [3](#), [4](#), [8](#)
- L. Birgé and P. Massart. Minimal penalties for gaussian model selection. *Probability Theory and Related Fields*, 183:33–73, 2007. [3](#)

- C. Brunson. Path estimation from GPS tracks. In *Proceedings of the 9th International Conference on GeoComputation*, National Centre for Geocomputation, National University of Ireland, Maynooth, Eire, 2007. 1
- N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning and Games*. Cambridge University Press, New York, 2006. 4, 6
- N. Cesa-Bianchi, G. Lugosi, and G. Stoltz. Minimizing regret with label-efficient prediction. *IEEE Transactions on Information Theory*, 51:2152–2162, 2005. 11
- F. Chung and L. Lu. Concentration inequalities and martingale inequalities: A survey. *Internet Mathematics*, 3:79–127, 2006. 19
- E. R. Engdahl and A. Villaseñor. 41 global seismicity: 1900–1999. *International Geophysics*, 81:665–690, 2002. 12, 13
- H. Friedsam and W. A. Oren. The application of the principal curve analysis technique to smooth beamlines. In *Proceedings of the 1st International Workshop on Accelerator Alignment*, 1989. 1
- T. Hastie and W. Stuetzle. Principal curves. *Journal of the American Statistical Association*, 84:502–516, 1989. 1
- H. Hotelling. Analysis of a complex of statistical variables into principal components. *Journal of educational psychology*, 24(6):417, 1933. 1
- M. Hutter and J. Poland. Adaptive online prediction by following the perturbed leader. *Journal of Machine Learning Research*, 6:639–660, 2005. 6
- V. Kanade, B. McMahan, and B. Bryan. Sleeping experts and bandits with stochastic action availability and adversarial rewards. *AISTATS*, 3:1137–1155, 2009. 8, 12
- B. Kégl. *Principal curves: learning, design, and applications*. PhD thesis, Concordia University Montreal, Quebec, 1999. 2, 3, 16
- B. Kégl and A. Krzyżak. Piecewise linear skeletonization using principal curves. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(1):59–74, 2002. 1, 3
- B. Kégl, A. Krzyżak, T. Linder, and K. Zeger. Learning and design of principal curves. *IEEE transactions on pattern analysis and machine intelligence*, 22(3):281–297, 2000. 2, 3, 4
- R. D. Kleinberg, A. Niculescu-Mizil, and Y. Sharma. *Regret Bounds for Sleeping Experts and Bandits*. In COLT. Springer, 2008. 8
- Valero Laparra and Jesús Malo. Sequential principal curves analysis. arXiv preprint, 2016. URL <https://arxiv.org/abs/1606.00856>. 3
- L. Li, B. Guedj, and S. Loustau. A Quasi-Bayesian Perspective to Online Clustering. arXiv preprint, 2016. URL <https://arxiv.org/abs/1602.00522>. 4, 5
- D. A. McAllester. Some PAC-Bayesian theorems. *Machine Learning*, 37(3):355–363, 1999a. 4
- D. A. McAllester. PAC-Bayesian model averaging. In *Proceedings of the 12th annual conference on Computational Learning Theory*, pages 164–170. ACM, 1999b. 4

- G. Neu and G. Bartók. An efficient algorithm for learning with semi-bandit feedback. In *Lecture Notes in Computer Science*, volume 8139, pages 234–248. Springer, Berlin, Heidelberg, 2013. [11](#)
- K Pearson. On lines and planes of closest fit to systems of point in space. *Philosophical Magazine*, 2(11):559–572, 1901. [1](#)
- K. Reinhard and M. Niranjan. Parametric subspace modeling of speech transitions. *Speech Communication*, 27:19–42, 1999. [1](#)
- S. Sandilya and S. R. Kulkarni. Principal curves with bounded turn. *IEEE Transactions on Information Theory*, 48:2789–2793, 2002. [3](#), [4](#)
- J. Shawe-Taylor and R. C. Williamson. A PAC analysis of a Bayes estimator. In *Proceedings of the 10th annual conference on Computational Learning Theory*, pages 2–9. ACM, 1997. doi: 10.1145/267460.267466. [4](#)
- C. Spearman. "General Intelligence", Objectively Determined and Measured. *The American Journal of Psychology*, 15(2):201–292, 1904. [1](#)
- D. C. Stanford and A. E. Raftery. Finding curvilinear features in spatial point patterns: principal curve clustering with noise. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(6):601–609, 2000. [1](#)