

Article

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

Le Li ¹  and Benjamin Guedj ^{2,*} 

¹ Department of Statistics, Central China Normal University, Wuhan 430079, China; leli@mail.ccnu.edu.cn

² Inria, Lille-Nord Europe Research Centre and Inria London, France and Centre for Artificial Intelligence, Department of Computer Science, University College London, London WC1V 6LJ, UK

* Correspondence: b.guedj@ucl.ac.uk

Abstract: When confronted with massive data streams, summarizing data with dimension reduction methods such as PCA raises theoretical and algorithmic pitfalls. A principal curve acts 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 (called *slpc*, for sequential learning principal curves) that incorporates both sleeping experts and multi-armed bandit ingredients is presented, along with its regret computation and performance on synthetic and real-life data.

Keywords: sequential learning; principal curves; data streams; regret bounds; greedy algorithm; sleeping experts



Citation: Li, L.; Guedj, B. Sequential Learning of Principal Curves: Summarizing Data Streams on the Fly. *Entropy* **2021**, *23*, 1534. <https://doi.org/10.3390/e23111534>

Academic Editor: Mohamed Medhat Gaber

Received: 22 August 2021
Accepted: 1 November 2021
Published: 18 November 2021

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

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 was introduced by [1,2] and further developed by [3]. This is one of the most widely used procedures in multivariate exploratory analysis targeting dimension reduction or feature 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 of skeletonization of data clouds has been at the heart of numerous applications in many different domains, such as physics [4,5], character and speech recognition [6,7], mapping and geology [5,8,9], to name but a few.

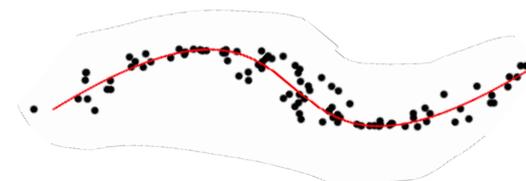


Figure 1. A principal curve.

1.1. Earlier Works on Principal Curves

The original definition of principal curve dates back to [10]. A principal curve is a smooth (C^∞) 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. 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.

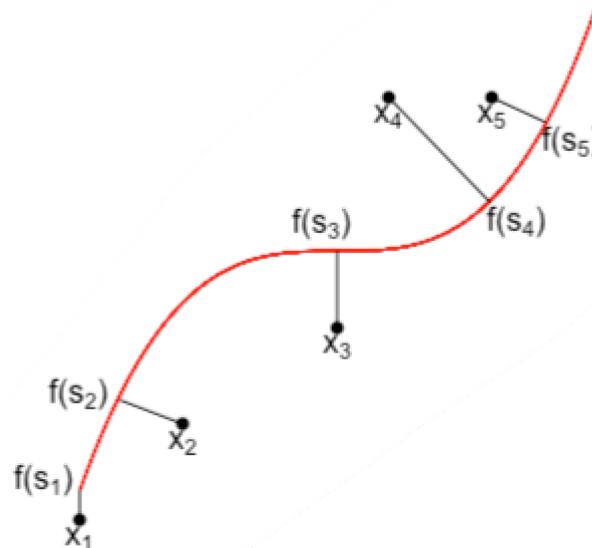


Figure 2. A principal curve and projections of data onto it.

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. In order to handle complex data structures, Ref. [11] proposed principal curves (PCOP) of principal oriented points (POPs) which are defined as the fixed points of an expectation function of points projected to a hyperplane minimising the total variance. To obtain POPs, a cluster analysis is performed on the hyperplane and only data in the local cluster are considered. Ref. [12] introduced the local principal curve (LPC), whose concept is similar to that of [11], but accelerates the computation of POPs by calculating local centers of mass instead of performing cluster analysis, and local principal component instead of principal direction. Later, Ref. [13] also considered LPC in data compression and regression to reduce the dimension of predictors space to low-dimension manifold. Ref. [14] extended the idea of localization to independent component analysis (ICA) by proposing a local-to-global non-linear ICA framework for visual and auditory signal. Ref. [15] considered principal curves from a different perspective: as the ridge of a smooth probability density function (PDF) generating dataset, where the ridges are collections of all points; the local gradient of a PDF is an eigenvector of the local Hessian, and the eigenvalues corresponding to the remaining eigenvectors are negative. To estimate principal curves based on this definition, the subspace constrained mean shift (SCMS) algorithm was proposed. All the local methods above require strong assumptions on the PDF, such as twice continuous differentiability, which may prove challenging to be satisfied in the settings of online sequential data. Ref. [16] proposed a new concept of principal curves which ensures its 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,

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

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 situations where only i.i.d. copies X_1, \dots, X_n of X are observed, the method of [16] 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 [17] that if X is almost surely bounded and $k \propto n^{1/3}$, 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 of at most L that minimizes $\Delta_n(\mathbf{f})$ is computationally costly, Ref. [17] proposed a 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, 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 of segments, k , plays a crucial role (a too small a k value leads to a poor summary of data, whereas a too-large k yields overfitting; see Figure 3), Ref. [18] aimed to fill the gap by selecting an optimal k from both theoretical and practical perspectives.

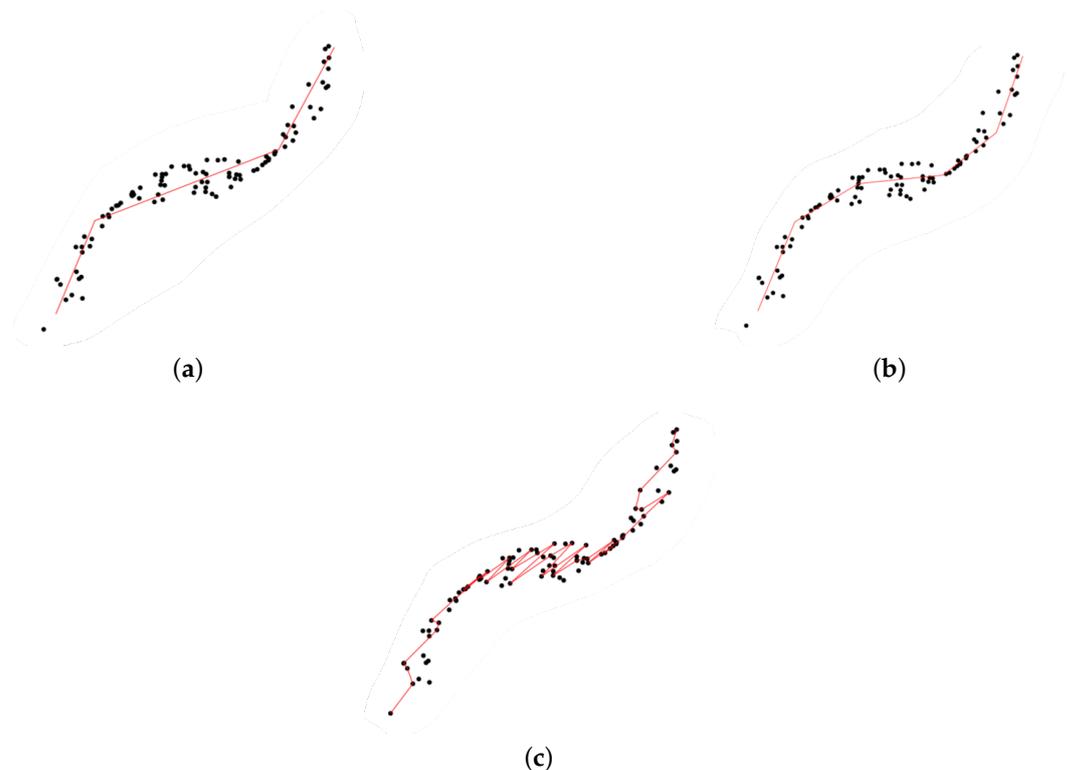


Figure 3. Principal curves with different numbers (k) of segments. (a) A too small k . (b) Right k . (c) A too large k .

Their approach relies strongly on the theory of model selection by penalization introduced by [19] and further developed by [20]. By considering countable classes $\{\mathcal{F}_{k,\ell}\}_{k,\ell}$ of polygonal lines with k segments and 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 δ stands for the diameter of observations and $w_{k,\ell}$ denotes the weight attached to class $\mathcal{F}_{k,\ell}$; and it has constants c_0, c_1, c_2 depending on δ , maximum length L and a certain number of dimensions of observations. Ref. [18] then proved 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 Σ 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 one 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 great interest to sequential learning. Extending the notion of principal curves to the sequential settings opens up immediate practical application possibilities. As an example, path planning for passengers' locations can help taxi companies to better optimize their fleet. Online algorithms that could yield instantaneous path summarization would be adapted to the sequential nature of geolocalized data. Existing theoretical works and practical implementations of principal curves are designed for the batch setting [7,16–18,21] and their adaptation to the sequential setting is not a smooth process. As an example, consider the algorithm in [18]. It is assumed that vertices of principal curves are located on a lattice, and its computational complexity is of order $\mathcal{O}(nN^p)$ where n is the number of observations, N the number of points on the lattice and p the maximum number of vertices. When p is large, running this algorithm at each epoch yields a monumental computational cost. In general, if data are not identically distributed or even adversary, algorithms that originally worked well in the batch setting may not be ideal when cast onto the online setting (see [22], Chapter 4). To the best of our knowledge, little effort has been put so far into extending principal curves algorithms to the sequential context.

Ref. [23] provided an incremental version of the SCMS algorithm [15] which is based on a definition of a principal curve as the ridge of a smooth probability density function generating observations. They applied the SCMS algorithm to the input points that are associated with the output points which are close to the new incoming sample and leave the remaining outputs unchanged. Hence, this algorithm can be used to deal with sequential data. As presented in the next section, our algorithm for sequentially updating principal curve vertices that are close to new data is similar in spirit to that of incremental SCMS. However, a difference is that our algorithm outputs polygonal lines. In addition, the computation complexity of our method is $\mathcal{O}(n^2)$, and incremental SCMS has $\mathcal{O}(n^3)$ complexity, where n is the number of observations. Ref. [24] considered sequential principal curves analysis in a fairly different setting in which the goal was to derive in an adaptive fashion a set of nonlinear sensors by using a set of preliminary principal curves. Unfolding sequentially principal curves and a sequential path for Jacobian integration were

considered. The “sequential” in this setting represented the generalization of principal curves to principal surfaces or even a principal manifold of higher dimensions. This way of sequentially exploiting principal curves was firstly proposed by [11] and later extended by [14,25,26] to give curvilinear representations using sequence of local-to-global curves. In addition, Refs. [15,27,28] presented, respectively, principal polynomial and non-parametric regressions to capture the nonlinear nature of data. However, these methods are not originally designed for treating sequential data. The present paper aims at filling this gap: our goal was to propose an online perspective to principal curves by automatically and sequentially learning the best principal curve summarizing a data stream. 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 [17,18,21] 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.

The contributions of the present paper are twofold. We first propose a sequential principal curve algorithm, for which we derived regret bounds. We then present an implementation, illustrated on a toy dataset and a real-life dataset (seismic data). The sketch of our algorithm’s 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 number 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 originated in the machine learning community in the late 1990s, in the seminal works of [29] and McAllester [30,31]. The PAC-Bayesian theory has been successfully adapted to sequential learning problems; see, for example, Ref. [32] for online clustering. We refer to [33,34] for a recent overview of the field.

The paper is organized 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 was proposed in Section 4, and we illustrate its performance on synthetic and real-life datasets in Section 5. Proofs of all original results claimed in the paper are collected in Section 6.

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 $x_1, x_2, \dots, x_T \in B(0, \sqrt{d}R) \subset \mathbb{R}^d$ be a sequence of data, where $B(\mathbf{c}, R)$ stands for the ℓ_2 -ball centered in $\mathbf{c} \in \mathbb{R}^d$ with radius $R > 0$. Let \mathcal{Q}_δ be a grid over $B(0, \sqrt{d}R)$, i.e., $\mathcal{Q}_\delta = B(0, \sqrt{d}R) \cap \Gamma_\delta$ where Γ_δ 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}_δ 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}_δ and whose length is at most L . Finally, let $\mathcal{K}(\mathbf{f})$ denote the number of segments of $\mathbf{f} \in \mathcal{F}_p$. This strategy is illustrated by Figure 4.

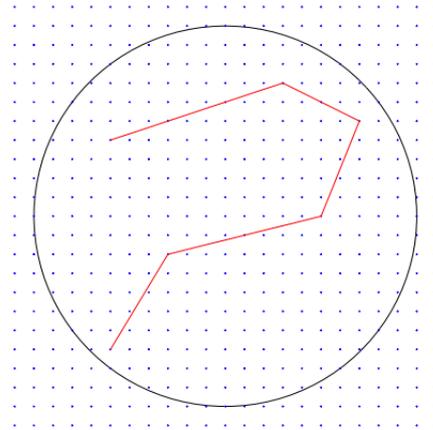


Figure 4. An example of a lattice Γ_δ in \mathbb{R}^2 with $\delta = 1$ (spacing between blue points) and $B(0,10)$ (black circle). The red polygonal line is composed of vertices in $\mathcal{Q}_\delta = B(0,10) \cap \Gamma_\delta$.

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

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

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

For any $k \in \llbracket 1, p \rrbracket$, let π_k denote a probability distribution on $\mathcal{F}_{k,L}$. We define the prior π on $\mathcal{F}_p = \cup_{k=1}^p \mathcal{F}_{k,L}$ as

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

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

We adopt a quasi-Bayesian-flavored procedure: 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 [32,35] who then considered realizations from this quasi-posterior. In the present paper, we will rather focus on a quantity linked to the mode of this quasi-posterior. Indeed, the mode of the quasi-posterior $\hat{\rho}_{t+1}$ is

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

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 π 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 [22,36], 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 Sequentially learning 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

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

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

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

The expectation of the cumulative loss of polygonal lines $\hat{\mathbf{f}}_1, \dots, \hat{\mathbf{f}}_T$ is upper-bounded by the smallest penalized 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 η . 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 $\sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})}$, but not too large to avoid overpenalization and a larger value for $S_{T,h,\eta}$. We therefore set

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

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

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

The penalty function $h(\mathbf{f}) = c_1\mathcal{K}(\mathbf{f}) + c_2L + c_3$ satisfies (3), where c_1, c_2, c_3 are constants depending on R, d, δ, p (this is proven in Lemma 3, in Section 6). We therefore obtain the following corollary.

Corollary 1. Under the assumptions of Theorem 1, let

$$\eta = \min \left\{ \frac{1}{d(2R + \delta)^2}, \sqrt{\frac{c_1p + c_2L + c_3}{c_0(e - 1) \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t)}} \right\}.$$

Then

$$\begin{aligned} \sum_{t=1}^T \mathbb{E} [\Delta(\hat{\mathbf{f}}_t, x_t)] &\leq \inf_{k \in \llbracket 1, p \rrbracket} \left\{ \inf_{\substack{\mathbf{f} \in \mathcal{F}_p \\ \mathcal{K}(\mathbf{f})=k}} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \sqrt{c_0(e - 1)r_{T,k,L}} \right\} \right\} \\ &\quad + \sqrt{c_0(e - 1)r_{T,p,L} + c_0(e - 1)(c_1p + c_2L + c_3)}, \end{aligned}$$

where $r_{T,k,L} = \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t)(c_1k + c_2L + c_3)$.

Proof. Note that

$$\sum_{t=1}^T \mathbb{E} [\Delta(\hat{\mathbf{f}}_t, x_t)] \leq S_{T,h,\eta} + \eta c_0(e - 1) \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + c_0(e - 1)(c_0p + c_2L + c_3),$$

and we conclude by setting

$$\eta = \sqrt{\frac{c_1p + c_2L + c_3}{c_0(e - 1) \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t)}}.$$

□

Sadly, Corollary 1 is not of much practical use since the optimal value for η 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 Sequentially and adaptively learning principal curves.

- 1: **Input parameters:** $p > 0, L > 0, \pi, h$ and $\eta_0 = \frac{\sqrt{c_1p + c_2L + c_3}}{c_0\sqrt{e-1}}$
- 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 $\hat{\mathbf{f}}_0 = \arg \inf_{\mathbf{f} \in \mathcal{F}_p} \Delta_{\mathbf{f},0}$
- 3: **For** $t = 1, \dots, T$
- 4: Compute $\eta_t = \frac{\sqrt{c_1p + c_2L + c_3}}{c_0\sqrt{(e-1)t}}$
- 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

$$\hat{\mathbf{f}}_t = \arg \inf_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{s=0}^{t-1} \Delta_{\mathbf{f},s} \right\}. \tag{4}$$

7: **End for**

Theorem 2. For any sequence $(x_t)_{1:T} \in B(0, \sqrt{d}R), R \geq 0$, let $h(\mathbf{f}) = c_1\mathcal{K}(\mathbf{f}) + c_2L + c_3$ 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

$$\eta_0 = \frac{\sqrt{c_1p + c_2L + c_3}}{c_0\sqrt{e-1}}, \quad \eta_t = \frac{\sqrt{c_1p + c_2L + c_3}}{c_0\sqrt{(e-1)t}},$$

where $t \geq 1$ and $c_0 = d(2R + \delta)^2$. Then the procedure described in Algorithm 2 satisfies

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

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 as (1), which is the most refined bound in the literature so far ([18] where the optimal values for k and L were 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 neighborhood $\mathcal{U}(\hat{\mathbf{f}}_t)$ (see the formal definition below) of $\hat{\mathbf{f}}_t$. This consideration is well suited for the principal curves 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 neighborhood 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 [37,38]. In this setting, defining the regret with respect to the best action in the whole set of actions in hindsight remains 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. Ref. [38] 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 [37] attains the optimal regret in the adversarial rewards case but has a runtime exponential in the number of all actions. Ref. [39] considered full and partial information with stochastic action availability and proposed an algorithm that runs in polynomial time. In what follows, we materialize our implementation by resorting to “sleeping experts”, i.e., a special set of available actions that adapts to the setting of principal curves.

Let σ denote an ordering of $|\mathcal{F}_p|$ actions, and \mathcal{A}_t a subset of the available 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 σ 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 $\max_\sigma \sum_{t=0}^T r_{\sigma(\mathcal{A}_t),t}$ (respectively, $\mathbb{E} \left[\max_\sigma \sum_{t=1}^T r_{\sigma(\mathcal{A}_t),t} \right]$ when \mathcal{A}_t is stochastic).

Our procedure starts with a **partition** step which aims at identifying the “relevant” neighborhood of an observation $x \in \mathbb{R}^d$ with respect to a given polygonal line, and then proceeds with the definition of the **neighborhood** of an action \mathbf{f} . We then provide the full implementation and prove a regret bound.

Partition. For any polygonal line \mathbf{f} with k segments, we denote by $\vec{\mathbf{V}} = (v_1, \dots, v_{k+1})$ 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 $\mathbf{f}(\vec{\mathbf{V}})$ to represent the polygonal line formed by connecting consecutive vertices in $\vec{\mathbf{V}}$ if no confusion arises. Let $V_i, i = 1, \dots, k + 1$ and $S_i, i = 1, \dots, k$ be the Voronoi partitions of \mathbb{R}^d with respect to \mathbf{f} , i.e., regions consisting of all points closer to vertex v_i or segment s_i . Figure 5 shows an example of Voronoi partition with respect to \mathbf{f} with three segments.

Neighborhood. For any $x \in \mathbb{R}^d$, we define the neighborhood $\mathcal{N}(x)$ with respect to \mathbf{f} as the union of all Voronoi partitions whose closure intersects with two vertices connecting the projection $\mathbf{f}(s_f(x))$ of x to \mathbf{f} . For example, for the point x in Figure 5, its neighborhood $\mathcal{N}(x)$ is the union of S_2, V_3, S_3 and V_4 . In addition, let $\mathcal{N}_t(x) = \{x_s \in \mathcal{N}(x), s = 1, \dots, t\}$ be the set of observations $x_{1:t}$ belonging to $\mathcal{N}(x)$ and $\bar{\mathcal{N}}_t(x)$ be its average. Let $\mathcal{D}(M) = \sup_{x,y \in M} \|x - y\|_2$ denote the diameter of set $M \subset \mathbb{R}^d$. We finally define the local grid $\mathcal{Q}_{\delta,t}(x)$ of $x \in \mathbb{R}^d$ at time t as

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

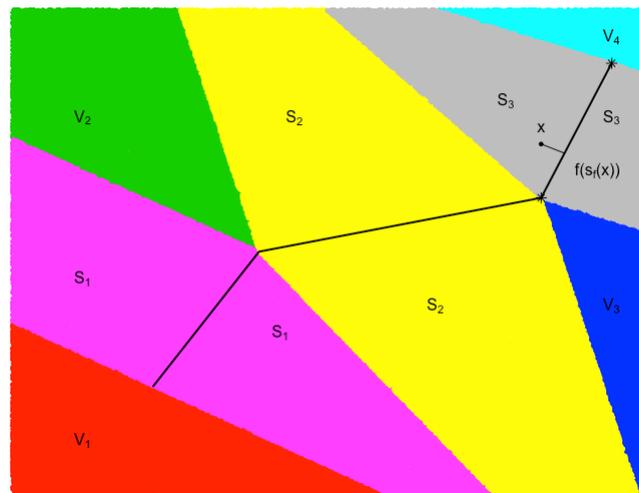


Figure 5. An example of a Voronoi partition.

We can finally proceed to the definition of the neighborhood $\mathcal{U}(\hat{\mathbf{f}}_t)$ of $\hat{\mathbf{f}}_t$. Assume $\hat{\mathbf{f}}_t$ has $k_t + 1$ vertices $\vec{\mathbf{V}} = (\underbrace{v_{1:i_t-1}}_{(i)}, \underbrace{v_{i_t:j_t-1}}_{(ii)}, \underbrace{v_{j_t:k_t+1}}_{(iii)})$, where vertices of (ii) belong to $\mathcal{Q}_{\delta,t}(x_t)$ while those of (i) and (iii) do not. The neighborhood $\mathcal{U}(\hat{\mathbf{f}}_t)$ consists of \mathbf{f} sharing vertices (i) and (iii) with $\hat{\mathbf{f}}_t$, but can be equipped with different vertices (ii) in $\mathcal{Q}_{\delta,t}(x_t)$; i.e.,

$$\mathcal{U}(\hat{\mathbf{f}}_t) = \left\{ \mathbf{f}(\vec{\mathbf{V}}), \vec{\mathbf{V}} = (v_{1:i_t-1}, v_{1:m}, v_{j_t:k_t+1}) \right\},$$

where $v_{1:m} \in \mathcal{Q}_{\delta,t}(x_t)$ and m is given by

$$m = \begin{cases} j_t - i_t - 1 & \text{reduce segments by 1 unit,} \\ j_t - i_t & \text{same number of segments,} \\ j_t - i_t + 1 & \text{increase segments by 1 unit.} \end{cases}$$

In Algorithm 3, we initiate the principal curve $\hat{\mathbf{f}}_1$ as the first component line segment whose vertices are the two farthest projections of data $x_{1:t_0}$ (t_0 can be set to 20 in practice) on the first component line. The reward of \mathbf{f} at round t in this setting is therefore $r_{\mathbf{f},t} = c_0 - \Delta(\mathbf{f}, x_{t_0+t})$. 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 and rewards of others are estimated by a constant, and we update our action from the neighborhood $\mathcal{U}(\hat{\mathbf{f}}_{t-1})$ of the previous action $\hat{\mathbf{f}}_{t-1}$. This local update (or search) greatly reduces computation complexity since $|\mathcal{U}(\hat{\mathbf{f}}_{t-1})| \ll |\mathcal{F}_p|$ when p is large. 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 β needs to be carefully calibrated since it should not be too large to ensure that the condition $cond(t)$ is non-empty; otherwise, all rewards are estimated by the same constant and thus lead 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 neighborhood of $\hat{\mathbf{f}}_t$ even if we are in the exploration phase at time $t + 1$. Conversely, very small β could result in large bias for the estimation $\frac{r_{\mathbf{f},t}}{\mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t)}$ of $r_{\mathbf{f},t}$. Note that the exploitation phase is close yet different to the label efficient prediction ([40], Remark 1.1) since we allow an action at time t to be different from the previous one. Ref. [41] 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 does not have an explicit form. However, due to the simple exponential distribution of z_t chosen in our case, an explicit form of $\mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t)$ is straightforward.

Algorithm 3 A locally greedy algorithm for sequentially learning 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 $z_t \sim \pi$.
- 5: Let

$$\hat{\sigma}_t = \text{sort} \left(\mathbf{f}, \sum_{s=1}^{t-1} \hat{r}_{\mathbf{f},s} - \frac{1}{\eta_{t-1}} h(\mathbf{f}) + \frac{1}{\eta_{t-1}} z_t \right),$$

i.e., sorting all $\mathbf{f} \in \mathcal{F}_p$ in descending order 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:

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

- 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:

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

where \mathcal{H}_t denotes all the randomness before time t and $cond(t) = \{ \mathbf{f} \in \mathcal{F}_p : \mathbb{P}(\hat{\mathbf{f}}_t = \mathbf{f} | \mathcal{H}_t) > \beta \}$. 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**
-

Theorem 3. Assume that $p > 6$, $T \geq 2|\mathcal{F}_p|^2$ and let $\beta = |\mathcal{F}_p|^{-\frac{1}{2}}T^{-\frac{1}{4}}$, $\alpha = \frac{c_0}{\beta}$, $\hat{c}_0 = \frac{2c_0}{\beta}$, $\epsilon = 1 - |\mathcal{F}_p|^{\frac{1}{2}-\frac{3}{p}}T^{-\frac{1}{4}}$ and

$$\eta_1 = \eta_2 = \dots = \eta_T = \frac{\sqrt{c_1p + c_2L + c_3}}{\sqrt{T(e-1)}\hat{c}_0}.$$

Then the procedure described in Algorithm 3 satisfies the regret bound

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

The proof of Theorem 3 is presented in Section 6. 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 the neighborhood 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 vanish; hence, the upper bound reduces to Theorem 2. In addition, our algorithm achieves an order that is smaller (from the perspective of both the number $|\mathcal{F}_p|$ of all actions and the total rounds T) than [39] since at each time, the availability of actions for our algorithm can be either the whole action set or a neighborhood of the previous action while [39] consider at each time only partial and independent stochastic available set of actions generated from a predefined distribution.

5. Numerical Experiments

We illustrate the performance of Algorithm 3 on synthetic and real-life data. Our implementation (hereafter denoted by `s1pc`—Sequential Learning of Principal Curves) is conducted with the R language and thus our most natural competitors are the R package `princurve`, which is the algorithm from [10], and `incremental`, which is the algorithm from SCMS [23]. We let $p = 50$, $R = \max_{t=1, \dots, T} \|x_t\|_2 / \sqrt{d}$, $L = 0.1p\sqrt{d}R$. The spacing δ of the lattice is adjusted with respect to data scale.

Synthetic data We generate a dataset $\{x_t \in \mathbb{R}^2, t = 1, \dots, 500\}$ uniformly along the curve $y = 0.05 \times (x - 5)^3$, $x \in [0, 10]$. Table 1 shows the regret (first row) for

- the ground truth (sum of squared distances of all points to the true curve),
- `princurve` and `incremental` SCMS (sum of squared distances between observation x_{t+1} and fitted `princurve` on observations $x_{1:t}$),
- `s1pc` (regret being equal to $\sum_{t=0}^{T-1} \mathbb{E}[\Delta(\hat{\mathbf{f}}_{t+1}, x_{t+1})]$ in both cases).

The mean computation time with different values for the time horizons T are also reported.

Table 1. The first line is the regret (cumulative loss) on synthetic data (average over 10 trials, with standard deviation in brackets). Second and third lines are the average computation time for two values of the time horizon T . `princurve` and `incremental` SCMS are deterministic, hence the zero standard deviation for regret.

Ground Truth	Princurve	Incremental SCMS	s1pc
2.48 (0)	26.02 (0)	19.09 (0)	20.83 (3.23)
T = 500	0.029 s (0.0001 s)	18.79 s (0.007 s)	1.44 s (0.030 s)
T = 5000	0.35 s (0.006 s)	>60 s (NA)	4.13 s (0.807 s)

Table 1 demonstrates the advantages of our method `s1pc`, as it achieved the optimal tradeoff between performance (in terms of regret) and runtime. Although `princurve` outperformed the other two algorithms in terms of computation time, it yielded the largest

regret, since it outputs a curve which does not pass in “the middle of data” but rather bends towards the curvature of the data cloud, as shown in Figure 6 where the predicted principal curves \hat{f}_{t+1} for `princurve`, `incremental SCMS` and `s1pc` are presented. `incremental SCMS` and `s1pc` both yielded satisfactory results, although the mean computation time of `s1pc` was significantly smaller than that of `incremental SCMS` (the reason being that eigenvectors of the Hessian of PDF need to be computed in `incremental SCMS`). Figure 7 showed, respectively, the estimation of the regret of `s1pc` and its per-round value (i.e., the cumulative loss divided by the number of rounds) both with respect to the round t . The jumps in the per-round curve occurred at the beginning, due to the initialization from a first principal component and to the collection of new data. When data accumulates, the vanishing pattern of the per-round curve illustrates that the regret is sublinear in t , which matches our aforementioned theoretical results.

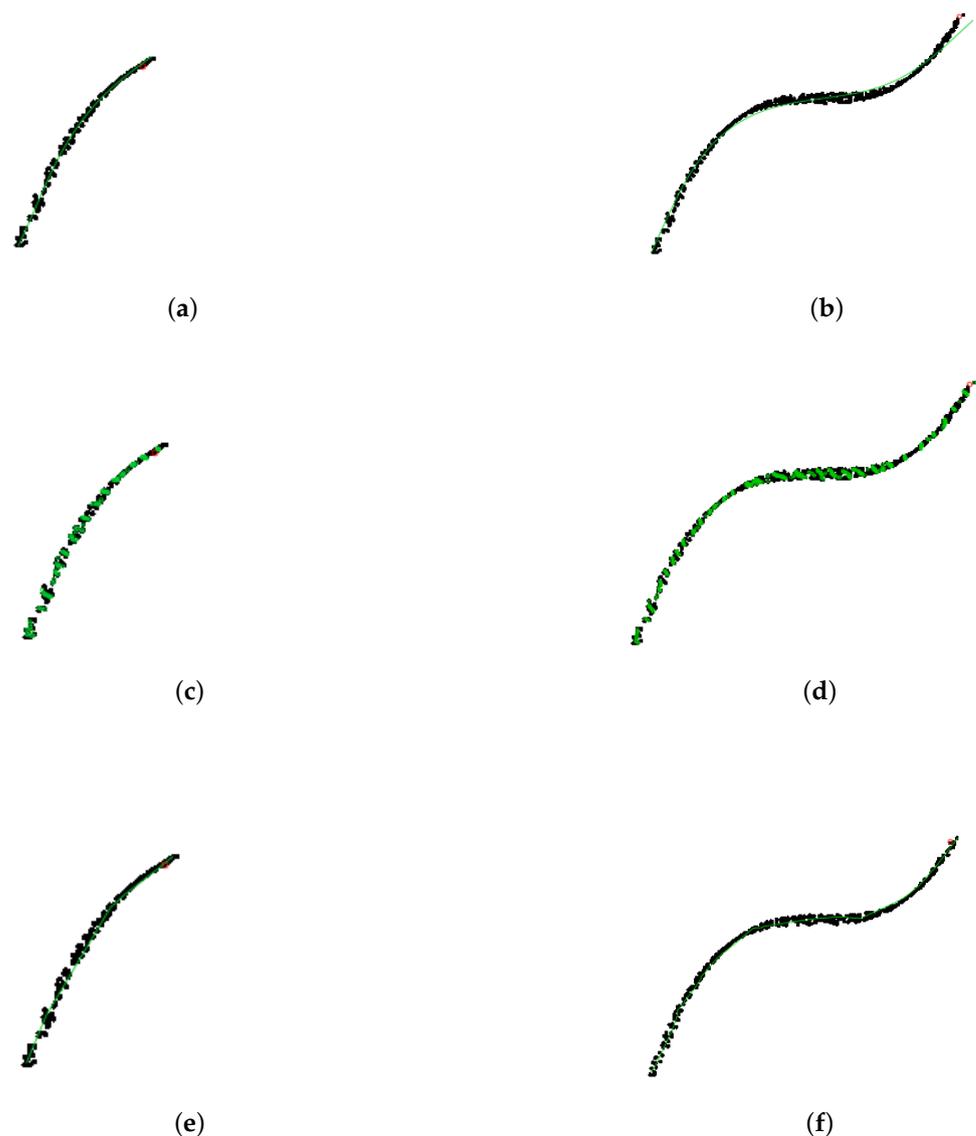


Figure 6. Synthetic data. Black dots represent data $x_{1:t}$. The red point is the new observation x_{t+1} . `princurve` (solid red) and `s1pc` (solid green). (a) $t = 150$, `princurve`. (b) $t = 450$, `princurve`. (c) $t = 150$, `incremental SCMS`. (d) $t = 450$, `incremental SCMS`. (e) $t = 150$, `s1pc`. (f) $t = 450$, `s1pc`.

In addition, to better illustrate the way `s1pc` works between two epochs, Figure 8 focuses on the impact of collecting a new data point on the principal curve. We see that

only a local vertex is impacted, whereas the rest of the principal curve remains unaltered. This cutdown in algorithmic complexity is one the key assets of s1pc.

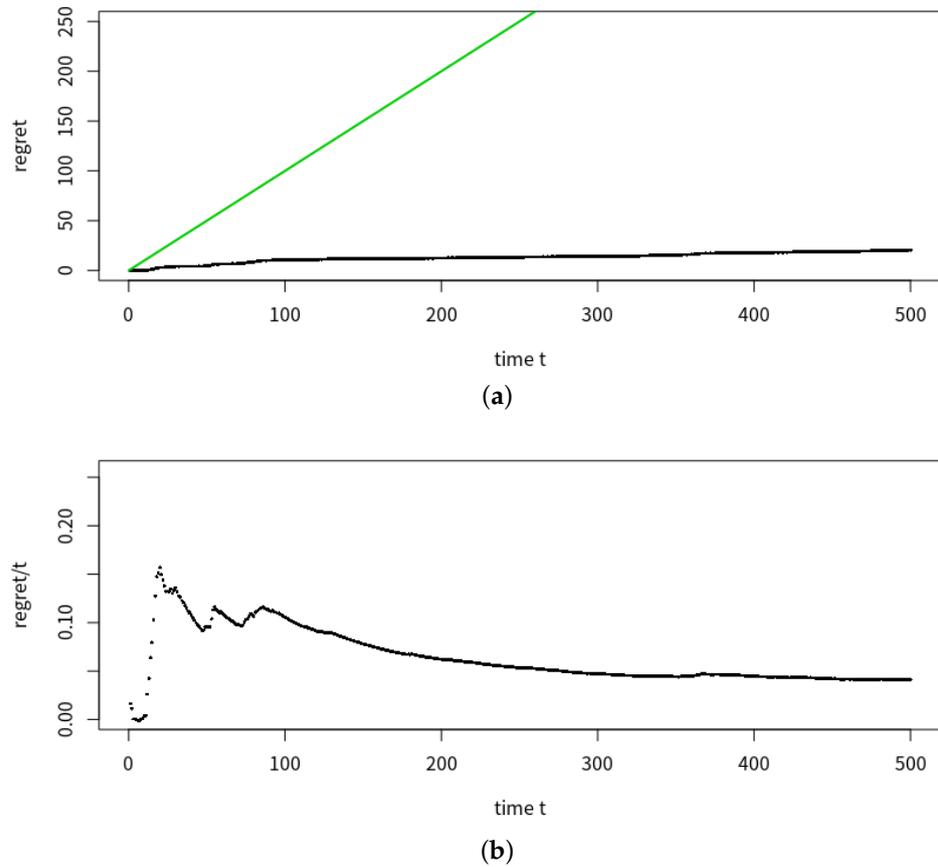


Figure 7. Mean estimation of regret and per-round regret of s1pc with respect to time round t , for the horizon $T = 500$. (a) Mean estimation of the regret of s1pc over 20 trials (black line) and a bisection line (green) with respect to time round t . (b) Per-round of estimated regret of s1pc with respect to t .

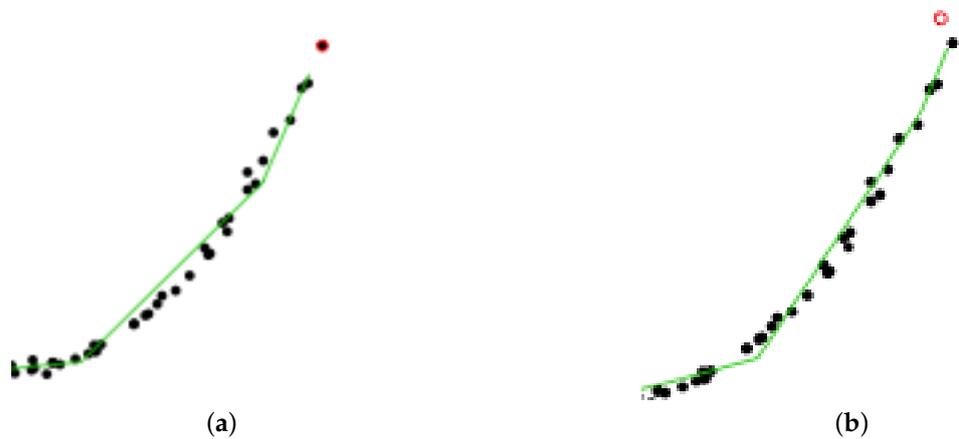


Figure 8. Synthetic data. Zooming in: how a new data point impacts the principal curve only locally. (a) At time $t = 97$. (b) And at time $t = 98$.

Synthetic data in high dimension. We also apply our algorithm on a dataset $\{x_t \in \mathbb{R}^6, t = 1, 2, \dots, 200\}$ in higher dimension. It is generated uniformly along a parametric curve whose coordinates are

$$\begin{pmatrix} 0.5t \cos(t) \\ 0.5t \sin(t) \\ 0.5t \\ -t \\ \sqrt{t} \\ 2\ln(t+1) \end{pmatrix}$$

where t takes 100 equidistant values in $[0, 2\pi]$. To the best of our knowledge, [10,16,18] only tested their algorithm on 2-dimensional data. This example aims at illustrating that our algorithm also works on higher dimensional data. Table 2 shows the regret for the ground truth, princurve and s1pc.

Table 2. Regret (cumulative loss) on synthetic high dimensional data in (average over 10 trials, with standard deviation in brackets). princurve and incremental SCMS are deterministic, hence the zero standard deviation.

Ground Truth	Princurve	Incremental SCMS	s1pc
3.290 (0)	14.204 (0)	5.38 (0)	6.797 (0.409)

In addition, Figure 9 shows the behaviour of s1pc (green) on each dimension.

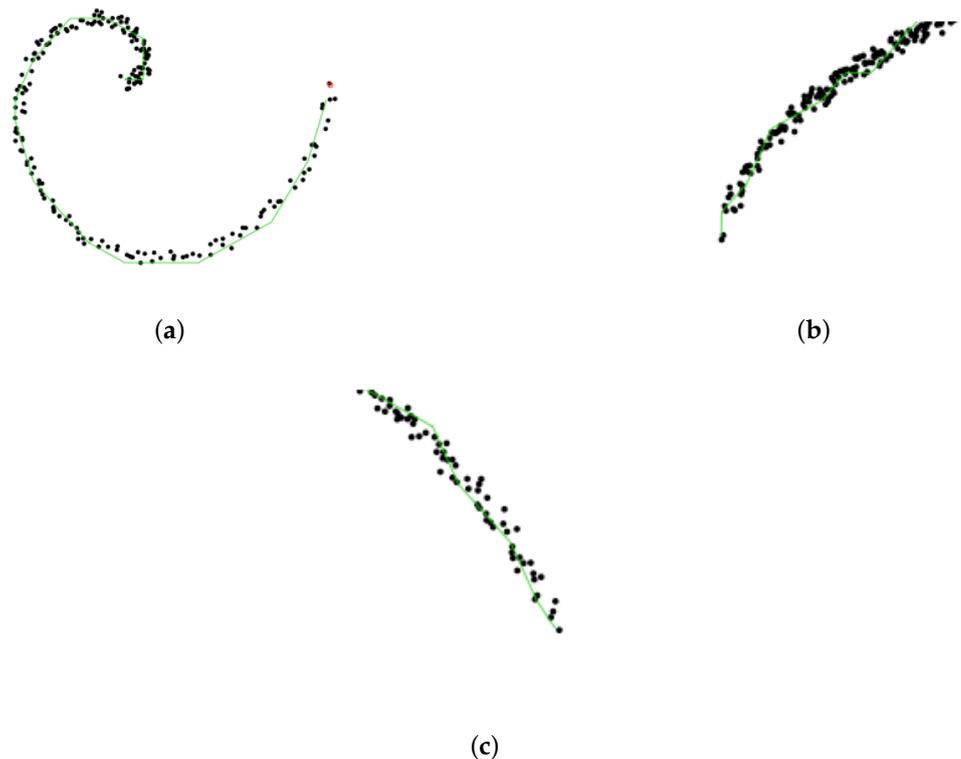


Figure 9. s1pc (green line) on synthetic high dimensional data from different perspectives. Black dots represent recordings $x_{1..99}$; the red dot is the new recording x_{200} . (a) s1pc, $t = 199$, 1st and 2nd coordinates. (b) s1pc, $t = 199$, 3th and 5th coordinates. (c) s1pc, $t = 199$, 4th and 6th coordinates.

Seismic data. Seismic data spanning long periods of time are essential for a thorough understanding of earthquakes. The “Centennial Earthquake Catalog” [42] 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. We focus on a particularly representative seismic active zone (a lithospheric border close to Australia) whose longitude is between $E130^\circ$ to $E180^\circ$ and latitude between $S70^\circ$ to $N30^\circ$, with $T = 218$ seismic recordings. As shown in Figure 10, `s1pc` recovers nicely the tectonic plate boundary, but both `princurve` and `incremental SCMS` with well-calibrated bandwidth fail to do so.

Lastly, since no ground truth is available, we used the R^2 coefficient to assess the performance (residuals are replaced by the squared distance between data points and their projections onto the principal curve). The average over 10 trials was 0.990.

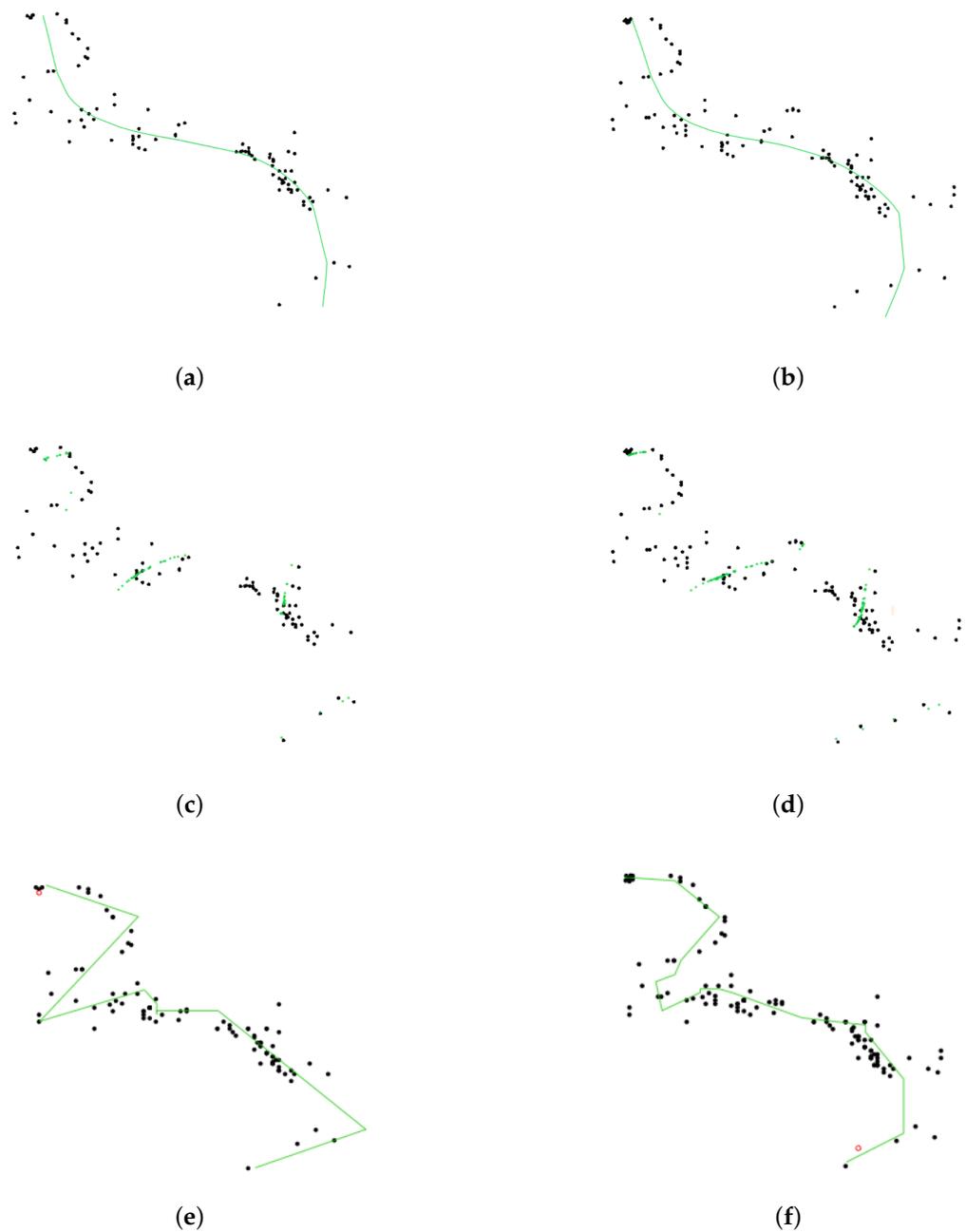


Figure 10. Seismic data. Black dots represent seismic recordings $x_{1:t}$; the red dot is the new recording x_{t+1} . (a) `princurve`, $t = 100$. (b) `princurve`, $t = 125$. (c) `incremental SCMS`, $t = 100$. (d) `incremental SCMS`, $t = 125$. (e) `s1pc`, $t = 100$. (f) `s1pc`, $t = 125$.

Back to Seismic Data. Figure 11 was taken from the USGS website (<https://earthquake.usgs.gov/data/centennial/>) and gives the global locations of earthquakes for the period 1900–1999. The seismic data (latitude, longitude, magnitude of earthquakes, etc.) used in the present paper may be downloaded from this website.

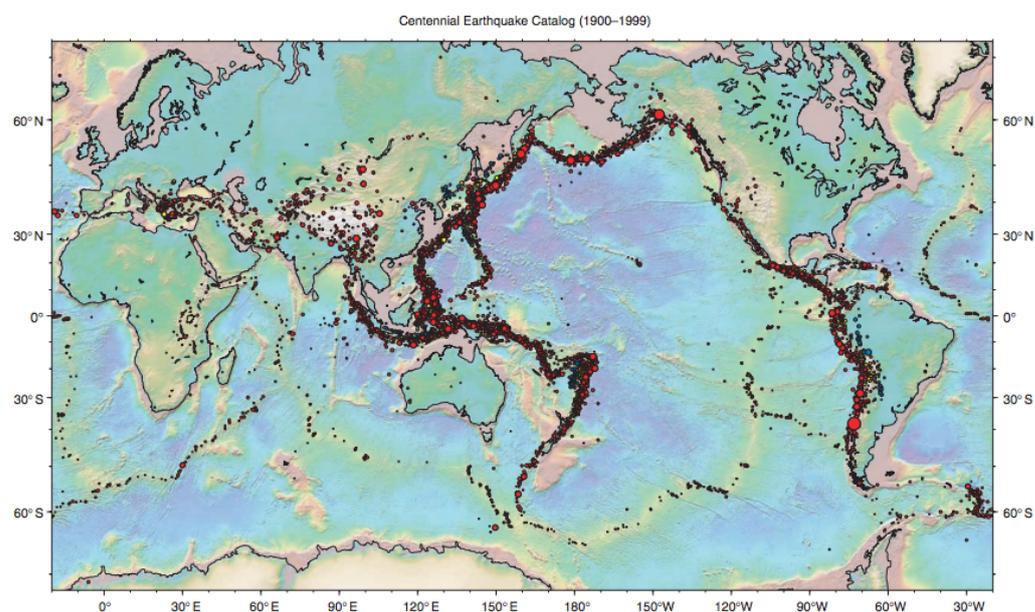


Figure 11. Seismic data from <https://earthquake.usgs.gov/data/centennial/>.

Daily Commute Data. The identification of segments of personal daily commuting trajectories can help taxi or bus companies to optimize their fleets and increase frequencies on segments with high commuting activity. Sequential principal curves appear to be an ideal tool to address this learning problem: we tested our algorithm on trajectory data from the University of Illinois at Chicago (https://www.cs.uic.edu/~boxu/mp2p/gps_data.html). The data were obtained from the GPS reading systems carried by two of the laboratory members during their daily commute for 6 months in the Cook county and the Dupage county of Illinois. Figure 12 presents the learning curves yielded by `princurve` and `s1pc` on geolocalization data for the first person, on May 30. A particularly remarkable asset of `s1pc` is that abrupt curvature in the data sequence was perfectly captured, whereas `princurve` does not enjoy the same flexibility. Again, we used the R^2 coefficient to assess the performance (where residuals are replaced by the squared distances between data points and their projections onto the principal curve). The average over 10 trials was 0.998.

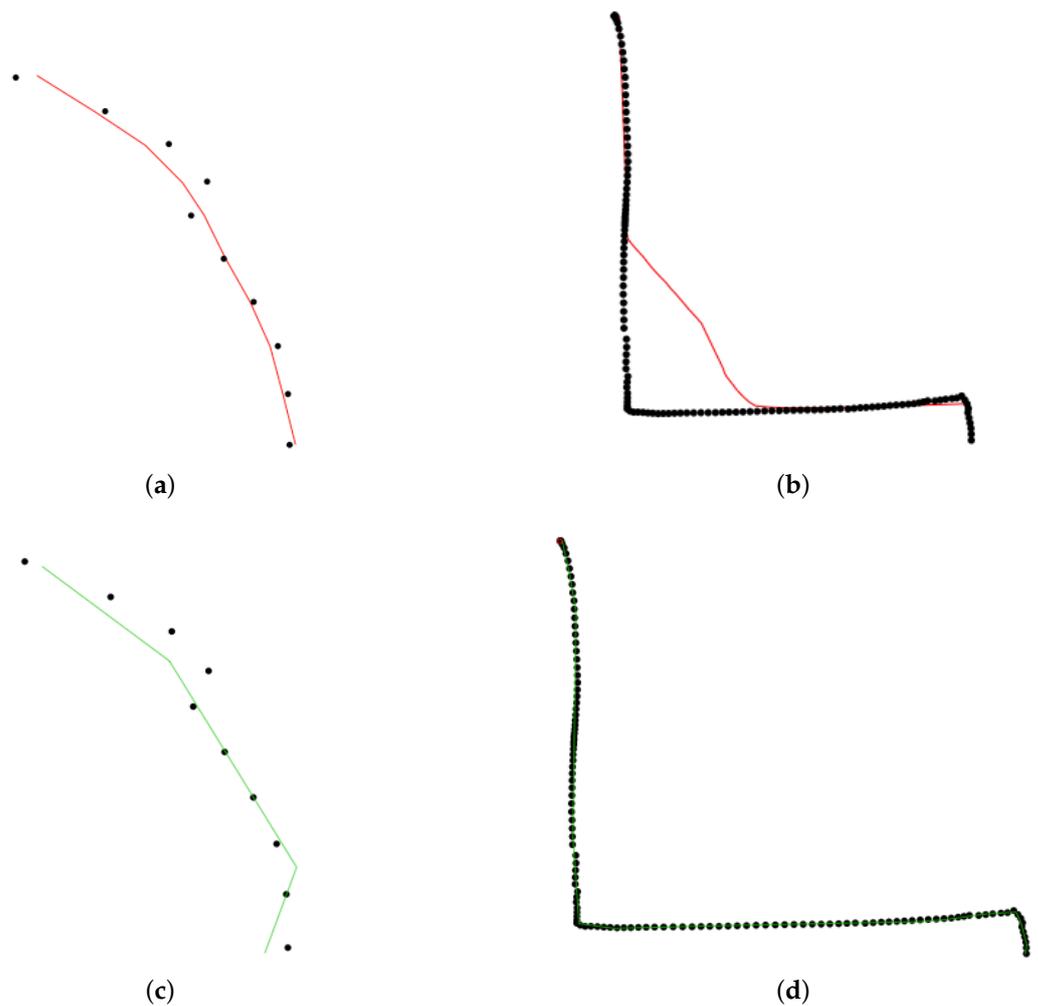


Figure 12. Daily commute data. Black dots represent collected locations $x_{1:t}$. The red point is the new observation x_{t+1} . princurve (solid red) and slpc (solid green). (a) $t = 10$, princurve. (b) $t = 127$, princurve. (c) $t = 10$, slpc. (d) $t = 127$, slpc.

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 the proof of Theorem 3 (which involves intermediary lemmas). Let us first define for each $t = 0, \dots, T$ the following forecaster sequence $(\hat{f}_t^*)_t$

$$\hat{f}_0^* = \arg \inf_{f \in \mathcal{F}_p} \{\Delta_{f,0}\} = \arg \inf_{f \in \mathcal{F}_p} \left\{ \frac{1}{\eta_0} h(f) - \frac{1}{\eta_0} z_f \right\},$$

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

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

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

the polygonal line in \mathcal{F}_p which minimizes the cumulative loss in the first T rounds plus a penalty term. f^* is deterministic, and \hat{f}_t^* is a random quantity (since it depends on z_f ,

$\mathbf{f} \in \mathcal{F}_p$ drawn from π). If several \mathbf{f} attain the infimum, we chose \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 $B(0, \sqrt{d}R)$,

$$\sum_{t=0}^T \Delta_{\hat{\mathbf{f}}_t^*, t} \leq \sum_{t=0}^T \Delta_{\hat{\mathbf{f}}_T^*, t}, \quad \pi\text{-almost surely.} \tag{5}$$

Proof. Proof by induction on T . Clearly (5) holds for $T = 0$. Assume that (5) 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 (5) and the definition of $\hat{\mathbf{f}}_T^*$, for $k \geq 1$, we have π -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}_T^* . Hence

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

where the second inequality is due to $\mathbb{E}[Z_{\mathbf{f}_T^*}] = 0$ and $\left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) > 0$ for $t = 0, 1, \dots, T$ since η_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 $y \geq 0$

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

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

$$\mathbb{E} \left[\sum_{t=1}^T \Delta(\hat{\mathbf{f}}_t^*, x_t) \right] \leq \inf_{\mathbf{f} \in \mathcal{F}_p} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{1}{\eta_T} h(\mathbf{f}) \right\} + \frac{1}{\eta_T} \left(1 + \ln \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} \right). \tag{6}$$

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 $(x_t) \in B(0, \sqrt{d}R)$, $t = 1, \dots, T$,

$$\sum_{t=1}^T \mathbb{E} \left[\Delta(\hat{\mathbf{f}}_t, x_t) \right] \leq \sum_{t=1}^T (1 + \eta_{t-1} c_0 (e - 1)) \mathbb{E} \left[\Delta(\hat{\mathbf{f}}_t^*, x_t) \right]. \tag{7}$$

Proof. Let us denote by

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

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

$$\begin{aligned} \ln|\{\mathbf{f} \in \mathcal{F}_p, \mathcal{K}(\mathbf{f}) = k\}| &\leq (\ln(8peV_d) + 3d^{\frac{3}{2}} - d)k + \left(\frac{\ln 2}{\delta\sqrt{d}} + \frac{d}{\delta}\right)L + d \ln\left(\frac{\sqrt{d}(2R + \delta)}{\delta}\right) \\ &\stackrel{\Delta}{=} c_1k + c_2L + c_3, \end{aligned}$$

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}_δ . 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

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

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 [16]. \square

We now have all the ingredients to prove Theorem 1 and Theorem 2.

First, combining (6) and (7) yields that

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

Assume that $\eta_t = \eta$, $t = 0, \dots, T$ and $h(\mathbf{f}) = c_1\mathcal{K}(\mathbf{f}) + c_2L + c_3$ 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

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[\Delta(\hat{\mathbf{f}}_t, x_t)] &\leq S_{T,h,\eta} + \frac{1}{\eta} \left(\frac{1}{2} + \ln \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} \right) + c_0(e-1)\eta \sum_{t=1}^T \mathbb{E}[\Delta(\hat{\mathbf{f}}_t^*, x_t)] \\ &\leq S_{T,h,\eta} + c_0(e-1)\eta S_{T,h,\eta} \\ &\leq S_{T,h,\eta} + \eta c_0(e-1) \inf_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + c_0(e-1)(c_1p + c_2L + c_3), \end{aligned}$$

where

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

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

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

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

$$\begin{aligned} \sum_{t=1}^T \mathbb{E} \left[\Delta(\hat{\mathbf{f}}_t, x_t) \right] &\leq \inf_{k \in \llbracket 1,p \rrbracket} \left\{ \inf_{\substack{\mathbf{f} \in \mathcal{F}_p \\ \mathcal{X}(\mathbf{f})=k}} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + \frac{h(\mathbf{f})}{\eta_T} \right\} \right\} + \frac{1}{\eta_T} \left(1 + \ln \sum_{\mathbf{f} \in \mathcal{F}_p} e^{-h(\mathbf{f})} \right) \\ &\quad + c_0^2(e-1) \sum_{t=1}^T \eta_{t-1} \\ &\leq \inf_{k \in \llbracket 1,p \rrbracket} \left\{ \inf_{\substack{\mathbf{f} \in \mathcal{F}_p \\ \mathcal{X}(\mathbf{f})=k}} \left\{ \sum_{t=1}^T \Delta(\mathbf{f}, x_t) + c_0 \sqrt{(e-1)T(c_0 k + c_2 L + c_3)} \right\} \right\} \\ &\quad + 2c_0 \sqrt{(e-1)T(c_0 p + c_2 L + c_3)}, \end{aligned}$$

which concludes the proof of Theorem 2.

Lemma 4. Using Algorithm 3, if $0 < \epsilon \leq 1, 0 < \beta < 1, \alpha \geq \frac{(1-\beta)c_0}{\beta}$ 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

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

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

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

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 $\mathbb{P} \left(\hat{\sigma}^t(\mathcal{A}_t) = \mathbf{f} \middle| \mathcal{H}_t \right) \geq \beta$. For $t = 1$, the above inequality is trivial since $\hat{r}_{\hat{\sigma}^1(\mathcal{A}_1), 1} \equiv 0$ by its definition. Hence, for $t \geq 1$, one has

$$\begin{aligned}
 \mathbb{E} \left[r_{\hat{\mathbf{f}}_t, t} \middle| \mathcal{H}_t \right] &= \epsilon \mathbb{E} \left[r_{\hat{\sigma}^t(\mathcal{F}_p), t} \middle| \mathcal{H}_t, I_t = 1 \right] + (1 - \epsilon) \mathbb{E} \left[r_{\hat{\sigma}^t(\mathcal{A}_t), t} \middle| \mathcal{H}_t, I_t = 0 \right] \\
 &\geq \mathbb{E} \left[\hat{r}_{\hat{\mathbf{f}}_t, t} \middle| \mathcal{H}_t \right] - 2\alpha\beta|\mathcal{A}_t|.
 \end{aligned} \tag{8}$$

Summing on both sides of inequality (8) 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

$$\mathbb{E} \left[\max_{\hat{\sigma}} \left\{ \sum_{t=1}^T \hat{r}_{\hat{\sigma}(\mathcal{A}_t), t} - \frac{1}{\eta} h(\hat{\sigma}(\mathcal{A}_t)) \right\} \right] - \sum_{t=1}^T \mathbb{E} \left[\hat{r}_{\hat{\sigma}^t(\mathcal{A}_t), t} \right] \leq \hat{c}_0^2(e - 1)\eta T + \hat{c}_0(e - 1)(c_1p + c_2L + c_3).$$

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

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

where in the second inequality we use that $r_{\mathbf{f}, t} \leq c_0$ for all \mathbf{f} and t , and that $\mathbb{P} \left(\hat{\mathbf{f}}_t = \mathbf{f} \middle| \mathcal{H}_t \right) \geq \beta$ when $\mathbf{f} \in \mathcal{U} \left(\hat{\mathbf{f}}_{t-1} \right) \cap \text{cond}(t)$. The rest of the proof is similar to those of Lemmas 1 and 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 Lemmas 1 and 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 the 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 [43] (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

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

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

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

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

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 $\max_{\mathbf{f} \in \mathcal{F}_p} \sum_{t=1}^T (r_{\mathbf{f},t} - \hat{r}_{\mathbf{f},t}) \leq c_0T$ almost surely, we terminate the proof of Lemma 6. \square

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

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

With those values, the assumptions of Lemmas 4, 5 and 6 are satisfied. Combining their results lead to the following

$$\begin{aligned}
\sum_{t=1}^T \mathbb{E} [r_{\hat{\mathbf{f}}_t, t}] &\geq \mathbb{E} \left[\max_{\sigma} \left\{ \sum_{t=1}^T r_{\sigma(\mathcal{A}_t), t} - \frac{1}{\eta} h(\sigma(\mathcal{A}_t)) \right\} \right] - 2\alpha\beta(1-\epsilon) \sum_{t=1}^T |\mathcal{U}(\hat{\mathbf{f}}_{t-1})| \\
&\quad - \hat{c}_0^2(e-1)\eta T - \hat{c}_0(e-1)(c_1p + c_2L + c_3) \\
&\quad - (1 - |\mathcal{F}_p|\beta) \sqrt{2T \left[\frac{c_0^2}{\beta} + \alpha^2(1-\beta) + (c_0 + 2\alpha)^2 \right]} \ln\left(\frac{1}{\beta}\right) - |\mathcal{F}_p|\beta c_0 T \\
&\geq \mathbb{E} \left[\max_{\sigma} \left\{ \sum_{t=1}^T r_{\sigma(\mathcal{A}_t), t} - \frac{1}{\eta} h(\sigma(\mathcal{A}_t)) \right\} \right] - (1-\epsilon) |\mathcal{F}_p|^{\frac{3}{p}} c_0 T \\
&\quad - \hat{c}_0^2(e-1)\eta T - \hat{c}_0(e-1)(c_1p + c_2L + c_3) \\
&\quad - (1 - |\mathcal{F}_p|\beta) \sqrt{2T \left[\frac{c_0^2}{\beta} + \alpha^2(1-\beta) + (c_0 + 2\alpha)^2 \right]} \ln\left(\frac{1}{\beta}\right) - |\mathcal{F}_p|\beta c_0 T \\
&\geq \mathbb{E} \left[\max_{\sigma} \left\{ \sum_{t=1}^T r_{\sigma(\mathcal{A}_t), t} - \frac{1}{\eta} h(\sigma(\mathcal{A}_t)) \right\} \right] - \mathcal{O}\left(|\mathcal{F}_p|^{\frac{1}{2}} T^{\frac{3}{4}}\right),
\end{aligned}$$

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

Author Contributions: Conceptualization, L.L. and B.G.; Formal analysis, L.L. and B.G.; Methodology, B.G.; Project administration, B.G.; Software, L.L.; Supervision, B.G.; Writing—original draft, L.L. and B.G.; Writing—review and editing, L.L. and B.G. All authors have read and agreed to the published version of the manuscript.

Funding: LL is funded and supported by the Fundamental Research Funds for the Central Universities (Grand No. 30106210158) and National Natural Science Foundation of China (Grant No. 61877023), the Fundamental Research Funds for the Central Universities (CCNU19TD009). BG is supported in part by the U.S. Army Research Laboratory and the U. S. Army Research Office, and by the U.K. Ministry of Defence and the U.K. Engineering and Physical Sciences Research Council (EPSRC) under grant number EP/R013616/1. BG acknowledges partial support from the French National Agency for Research, grants ANR-18-CE40-0016-01 and ANR-18-CE23-0015-02.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Pearson, K. On lines and planes of closest fit to systems of point in space. *Philos. Mag.* **1901**, *2*, 559–572. [\[CrossRef\]](#)
2. Spearman, C. “General Intelligence”, Objectively Determined and Measured. *Am. J. Psychol.* **1904**, *15*, 201–292. [\[CrossRef\]](#)
3. Hotelling, H. Analysis of a complex of statistical variables into principal components. *J. Educ. Psychol.* **1933**, *24*, 417–441. [\[CrossRef\]](#)
4. Friedsam, H.; Oren, W.A. The application of the principal curve analysis technique to smooth beamlines. In Proceedings of the 1st International Workshop on Accelerator Alignment, Stanford, CA, USA, 31 July–2 August 1989.
5. Brunson, C. Path estimation from GPS tracks. In Proceedings of the 9th International Conference on GeoComputation, Maynooth, Ireland, 3–5 September 2007.
6. Reinhard, K.; Niranjan, M. Parametric Subspace Modeling Of Speech Transitions. *Speech Commun.* **1999**, *27*, 19–42. [\[CrossRef\]](#)
7. Kégl, B.; Krzyżak, A. Piecewise linear skeletonization using principal curves. *IEEE Trans. Pattern Anal. Mach. Intell.* **2002**, *24*, 59–74. [\[CrossRef\]](#)
8. Banfield, J.D.; Raftery, A.E. Ice floe identification in satellite images using mathematical morphology and clustering about principal curves. *J. Am. Stat. Assoc.* **1992**, *87*, 7–16. [\[CrossRef\]](#)

9. Stanford, D.C.; Raftery, A.E. Finding curvilinear features in spatial point patterns: Principal curve clustering with noise. *IEEE Trans. Pattern Anal. Mach. Intell.* **2000**, *22*, 601–609. [[CrossRef](#)]
10. Hastie, T.; Stuetzle, W. Principal curves. *J. Am. Stat. Assoc.* **1989**, *84*, 502–516. [[CrossRef](#)]
11. Delicado, P. Another Look at Principal Curves and Surfaces. *J. Multivar. Anal.* **2001**, *77*, 84–116. [[CrossRef](#)]
12. Einbeck, J.; Tutz, G.; Evers, L. Local principal curves. *Stat. Comput.* **2005**, *15*, 301–313. [[CrossRef](#)]
13. Einbeck, J.; Tutz, G.; Evers, L. Data Compression and Regression through Local Principal Curves and Surfaces. *Int. J. Neural Syst.* **2010**, *20*, 177–192. [[CrossRef](#)]
14. Malo, J.; Gutiérrez, J. V1 non-linear properties emerge from local-to-global non-linear ICA. *Netw. Comput. Neural Syst.* **2006**, *17*, 85–102. [[CrossRef](#)]
15. Ozertem, U.; Erdogmus, D. Locally Defined Principal Curves and Surfaces. *J. Mach. Learn. Res.* **2011**, *12*, 1249–1286.
16. Kégl, B. Principal Curves: Learning, Design, and Applications. Ph.D. Thesis, Concordia University, Montreal, QC, Canada, 1999.
17. Kégl, B.; Krzyżak, A.; Linder, T.; Zeger, K. Learning and design of principal curves. *IEEE Trans. Pattern Anal. Mach. Intell.* **2000**, *22*, 281–297. [[CrossRef](#)]
18. Biau, G.; Fischer, A. Parameter selection for principal curves. *IEEE Trans. Inf. Theory* **2012**, *58*, 1924–1939. [[CrossRef](#)]
19. Barron, A.; Birgé, L.; Massart, P. Risk bounds for model selection via penalization. *Probab. Theory Relat. Fields* **1999**, *113*, 301–413. [[CrossRef](#)]
20. Birgé, L.; Massart, P. Minimal penalties for Gaussian model selection. *Probab. Theory Relat. Fields* **2007**, *183*, 33–73. [[CrossRef](#)]
21. Sandilya, S.; Kulkarni, S.R. Principal curves with bounded turn. *IEEE Trans. Inf. Theory* **2002**, *48*, 2789–2793. [[CrossRef](#)]
22. Cesa-Bianchi, N.; Lugosi, G. *Prediction, Learning and Games*; Cambridge University Press: New York, NY, USA, 2006.
23. Rudzicz, F.; Ghassabeh, Y.A. Incremental algorithm for finding principal curves. *IET Signal Process.* **2015**, *9*, 521–528.
24. Laparra, V.; Malo, J. Sequential Principal Curves Analysis. *arXiv* **2016**, arXiv:1606.00856.
25. Laparra, V.; Jiménez, S.; Camps-Valls, G.; Malo, J. Nonlinearities and Adaptation of Color Vision from Sequential Principal Curves Analysis. *Neural Comput.* **2012**, *24*, 2751–2788. [[CrossRef](#)]
26. Laparra, V.; Malo, J. Visual Aftereffects and Sensory Nonlinearities from a single Statistical Framework. *Front. Hum. Neurosci.* **2015**, *9*. [[CrossRef](#)]
27. Laparra, V.; Jiménez, S.; Tuia, D.; Camps-Valls, G.; Malo, J. Principal Polynomial Analysis. *Int. J. Neural Syst.* **2014**, *24*, 1440007. [[CrossRef](#)]
28. Laparra, V.; Malo, J.; Camps-Valls, G. Dimensionality Reduction via Regression in Hyperspectral Imagery. *IEEE J. Sel. Top. Signal Process.* **2015**, *9*, 1026–1036. [[CrossRef](#)]
29. Shawe-Taylor, J.; Williamson, R.C. A PAC analysis of a Bayes estimator. In Proceedings of the 10th annual conference on Computational Learning Theory, Nashville, TN, USA, 6–9 July 1997; pp. 2–9. [[CrossRef](#)]
30. McAllester, D.A. Some PAC-Bayesian Theorems. *Mach. Learn.* **1999**, *37*, 355–363. [[CrossRef](#)]
31. McAllester, D.A. PAC-Bayesian Model Averaging. In Proceedings of the 12th Annual Conference on Computational Learning Theory, Santa Cruz, CA, USA, 7–9 July 1999; pp. 164–170.
32. Li, L.; Guedj, B.; Loustau, S. A quasi-Bayesian perspective to online clustering. *Electron. J. Stat.* **2018**, *12*, 3071–3113. [[CrossRef](#)]
33. Guedj, B. A Primer on PAC-Bayesian Learning. In Proceedings of the Second Congress of the French Mathematical Society, Long Beach, CA, USA, 10 June 2019; pp. 391–414.
34. Alquier, P. User-friendly introduction to PAC-Bayes bounds. *arXiv* **2021**, arXiv:2110.11216.
35. Audibert, J.Y. Fast Learning Rates in Statistical Inference through Aggregation. *Ann. Stat.* **2009**, *37*, 1591–1646. [[CrossRef](#)]
36. Hutter, M.; Poland, J. Adaptive Online Prediction by Following the Perturbed Leader. *J. Mach. Learn. Res.* **2005**, *6*, 639–660.
37. Auer, P.; Cesa-Bianchi, N.; Freund, Y.; Schapire, R.E. The Nonstochastic multiarmed Bandit problem. *SIAM J. Comput.* **2003**, *32*, 48–77. [[CrossRef](#)]
38. Kleinberg, R.D.; Niculescu-Mizil, A.; Sharma, Y. Regret Bounds for Sleeping Experts and Bandits. In *COLT*; Springer: Berlin/Heidelberg, Germany, 2008.
39. Kanade, V.; McMahan, B.; Bryan, B. Sleeping Experts and Bandits with Stochastic Action Availability and Adversarial Rewards. *Artif. Intell. Stat.* **2009**, *3*, 1137–1155.
40. Cesa-Bianchi, N.; Lugosi, G.; Stoltz, G. Minimizing regret with label-efficient prediction. *IEEE Trans. Inf. Theory* **2005**, *51*, 2152–2162. [[CrossRef](#)]
41. Neu, G.; Bartók, G. *An Efficient Algorithm for Learning with Semi-Bandit Feedback*; Lecture Notes in Computer Science; Springer: Berlin/Heidelberg, Germany, 2013; Volume 8139, pp. 234–248.
42. Engdahl, E.R.; Villaseñor, A. 41 Global seismicity: 1900–1999. *Int. Geophys.* **2002**, *81*, 665–690.
43. Chung, F.; Lu, L. Concentration Inequalities and Martingale Inequalities: A Survey. *Internet Math.* **2006**, *3*, 79–127. [[CrossRef](#)]