



A new look at shifting regret

Nicolò Cesa-Bianchi, Pierre Gaillard, Gabor Lugosi, Gilles Stoltz

► To cite this version:

Nicolò Cesa-Bianchi, Pierre Gaillard, Gabor Lugosi, Gilles Stoltz. A new look at shifting regret. 2012.
hal-00670514v1

HAL Id: hal-00670514

<https://hal.science/hal-00670514v1>

Preprint submitted on 15 Feb 2012 (v1), last revised 27 Sep 2012 (v2)

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

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

A New Look at Shifting Regret

Nicolò Cesa-Bianchi

DSI, Università degli Studi di Milano

NICOLO.CESA-BIANCHI@UNIMI.IT

Pierre Gaillard

Ecole Normale Supérieure, Paris*

PIERRE.GAILLARD@ENS.FR

Gábor Lugosi

ICREA & Universitat Pompeu Fabra

GABOR.LUGOSI@UPF.EDU

Gilles Stoltz

Ecole Normale Supérieure†, Paris & HEC Paris, Jouy-en-Josas, France

GILLES.STOLTZ@ENS.FR

Abstract

We investigate extensions of well-known online learning algorithms such as fixed-share of [Herbster and Warmuth \(1998\)](#) or the methods proposed by [Bousquet and Warmuth \(2002\)](#). These algorithms use weight sharing schemes to perform as well as the best sequence of experts with a limited number of changes. Here we show, with a common, general, and simpler analysis, that weight sharing in fact achieves much more than what it was designed for. We use it to simultaneously prove new shifting regret bounds for online convex optimization on the simplex in terms of the total variation distance as well as new bounds for the related setting of adaptive regret. Finally, we exhibit the first logarithmic shifting bounds for exp-concave loss functions on the simplex.

Keywords: Prediction with expert advice, online convex optimization, tracking the best expert, shifting experts.

1. Introduction

Online convex optimization is a sequential prediction paradigm in which, at each time step, the learner chooses an element from a fixed convex set \mathcal{S} and then is given access to a convex loss function defined on the same set. The value of the function on the chosen element is the learner's loss. Many problems such as prediction with expert advice, sequential investment, and online regression/classification can be viewed as special cases of this general framework. Online learning algorithms are designed to minimize the regret. The standard notion of regret is the difference between the learner's cumulative loss and the cumulative loss of the single best element in \mathcal{S} . A much harder criterion to minimize is shifting regret, which is defined as the difference between the learner's cumulative loss and the cumulative loss of an arbitrary sequence of elements in \mathcal{S} . Shifting regret bounds are typically expressed in terms of the *shift*, a notion of regularity measuring the length of the trajectory in \mathcal{S} described by the comparison sequence (i.e., the sequence of elements against which the regret is evaluated).

In online convex optimization, shifting regret bounds for convex subsets $\mathcal{S} \subseteq \mathbb{R}^d$ are obtained for the online mirror descent (or follow-the-regularized-leader) algorithm. In this case the shift is

†. CNRS – Ecole normale supérieure, Paris – INRIA, within the project-team CLASSIC

†. CNRS – Ecole normale supérieure, Paris – INRIA, within the project-team CLASSIC

typically computed in terms of the p -norm of the difference of consecutive elements in the comparison sequence —see [Herbster and Warmuth \(2001\)](#) and [Cesa-Bianchi and Lugosi \(2006\)](#). In this paper we focus on the important special case when \mathcal{S} is the simplex, and investigate the online mirror descent with entropic regularizers. This family includes popular algorithms such as exponentially weighted average (EWA), Winnow, and exponentiated gradient. Proving general shifting bounds in this case is difficult due to the behavior of the regularizer at the boundary of the simplex. [Herbster and Warmuth \(2001\)](#) show shifting bounds for mirror descent with entropic regularizers using a 1-norm to measure the shift. In order to keep mirror descent from choosing points too close to the simplex boundary, they use a complex dynamic projection technique. When the comparison sequence is restricted to the corners of the simplex (which is the setting of prediction with expert advice), then the shift is naturally defined to be the number times the trajectory moves to a different corner. This problem is often called “tracking the best expert” —see, e.g., [Herbster and Warmuth \(1998\)](#); [Vovk \(1999\)](#); [Herbster and Warmuth \(2001\)](#); [Bousquet and Warmuth \(2002\)](#); [Györfy et al. \(2005\)](#), and it is well known that EWA with weight sharing, which corresponds to the fixed-share algorithm of [Herbster and Warmuth \(1998\)](#), achieves a good shifting bound in this setting. [Bousquet and Warmuth \(2002\)](#) introduce a generalization of the fixed-share algorithm, and prove various shifting bounds for any trajectory in the simplex. However, their bounds are expressed using a quantity that corresponds to a proper shift only for trajectories on the simplex corners.

Our analysis unifies, generalizes (and simplifies) the previously quite different proof techniques and algorithms used in [Herbster and Warmuth \(1998\)](#) and [Bousquet and Warmuth \(2002\)](#). Our bounds are expressed in terms of a notion of shift based on the total variation distance. The generalization of the “small expert set” result in [Bousquet and Warmuth \(2002\)](#) leads us to obtain better bounds when the sequence against which the regret is measured is sparse. When the trajectory is restricted to the corners of the simplex, we recover, and occasionally improve, the known shifting bounds for prediction with expert advice. Besides, our analysis also captures the setting of adaptive regret, a related notion of regret introduced by [Hazan and Seshadhri \(2009\)](#). It was known that shifting regret and adaptive regret had some connections but this connection is now seen to be even tighter, as both regrets can be viewed as instances of the same *alma mater* regret, which we minimize. Finally, we also show how to dynamically tune the parameters of our algorithms and review briefly the special case of exp-concave loss functions, exhibiting the first logarithmic shifting bounds for exp-concave loss functions on the simplex.

2. Preliminaries

We first define the sequential learning framework we work with. Even though our results hold in the general setting of online convex optimization, we present them in the, somewhat simpler, *online linear optimization* setup. We point out in Section 6 how these results may be generalized. Online linear optimization may be cast as a repeated game between the *forecaster* and the *environment* as follows. We use Δ_d to denote the simplex $\{\mathbf{q} \in [0, 1]^d : \|\mathbf{q}\|_1 = 1\}$.

Online linear optimization. For each round $t = 1, \dots, T$,

1. Forecaster chooses $\hat{\mathbf{p}}_t = (\hat{p}_{1,t}, \dots, \hat{p}_{d,t}) \in \Delta_d$;
2. Environment chooses a loss vector $\ell_t = (\ell_{1,t}, \dots, \ell_{d,t}) \in [0, 1]^d$;
3. Forecaster suffers loss $\hat{\mathbf{p}}_t^\top \ell_t$.

The goal of the forecaster is to minimize the accumulated loss $\widehat{L}_T = \sum_{t=1}^T \widehat{\mathbf{p}}_t^\top \ell_t$. In the now classical problem of prediction with expert advice, the goal of the forecaster is to compete with the best fixed component (often called “expert”) chosen in hindsight, that is, with $\min_{i=1,\dots,T} \sum_{t=1}^T \ell_{i,t}$. The focus of this paper is on more ambitious forecasters that compete with a richer class of sequences of components. Let $[d] = \{1, \dots, d\}$. We use $i_1^T = (i_1, \dots, i_T)$ to denote a sequence in $[d]^T$ and let $L_T(i_1^T) = \sum_{t=1}^T \ell_{i_t}$ be the cumulative linear loss of the sequence $i_1^T \in [d]^T$.

We start by introducing our main algorithmic tool, a generalized share algorithm. It is parametrized by the “mixing functions” $\psi_t : [0, 1]^{(t+1)d} \rightarrow \Delta_d$ for $t = 1, \dots, T$ that assign probabilities to past “pre-weights” as defined below. In all examples discussed in this paper, these mixing functions are quite simple but working with such a general model makes the main ideas more transparent. We then provide a simple lemma that serves as the starting point for analyzing different instances of the generalized share algorithm.

Algorithm 1: The generalized share algorithm.

Parameters: learning rate $\eta > 0$ and mixing functions ψ_t for $t = 1, \dots, T$

Initialization: $\widehat{\mathbf{p}}_1 = \mathbf{v}_1 = (1/d, \dots, 1/d)$

For each round $t = 1, \dots, T$,

1. Predict $\widehat{\mathbf{p}}_t$;
2. Observe loss $\ell_t \in [0, 1]^d$;
3. [loss update] **For** each $j = 1, \dots, d$ define

$$v_{j,t+1} = \frac{\widehat{p}_{j,t} e^{-\eta \ell_{j,t}}}{\sum_{i=1}^d \widehat{p}_{i,t} e^{-\eta \ell_{i,t}}} \quad \text{the current pre-weights,}$$

$$\underline{\mathbf{V}}_{t+1} = [v_{i,s}]_{i \in [d], 1 \leq s \leq t+1} \quad \text{the } d \times (t+1) \text{ matrix of all past and current pre-weights;}$$

4. [shared update] Define $\widehat{\mathbf{p}}_{t+1} = \psi_{t+1}(\underline{\mathbf{V}}_{t+1})$.
-

Lemma 1 For all $t \geq 1$ and for all $\mathbf{q}_t \in \Delta_d$, Algorithm 1 satisfies

$$(\widehat{\mathbf{p}}_t - \mathbf{q}_t)^\top \ell_t \leq \frac{1}{\eta} \sum_{i=1}^d q_{i,t} \ln \frac{v_{i,t+1}}{\widehat{p}_{i,t}} + \frac{\eta}{8}.$$

Proof By Hoeffding’s inequality,

$$\sum_{j=1}^d \widehat{p}_{j,t} \ell_{j,t} \leq -\frac{1}{\eta} \ln \left(\sum_{j=1}^d \widehat{p}_{j,t} e^{-\eta \ell_{j,t}} \right) + \frac{\eta}{8}. \quad (1)$$

By definition of $v_{i,t+1}$, for all $i = 1, \dots, d$ we then have $\sum_{j=1}^d \widehat{p}_{j,t} e^{-\eta \ell_{j,t}} = \frac{\widehat{p}_{i,t} e^{-\eta \ell_{i,t}}}{v_{i,t+1}}$,

which entails $\widehat{\mathbf{p}}_t^\top \ell_t \leq \ell_{i,t} + \frac{1}{\eta} \ln \frac{v_{i,t+1}}{\widehat{p}_{i,t}} + \frac{\eta}{8}$.

The proof is concluded by taking a convex aggregation with respect to \mathbf{q}_t . ■

3. Shifting bounds

In this section we prove shifting regret bounds for the generalized share algorithm. We compare the cumulative loss $\sum_{t=1}^T \hat{\mathbf{p}}_t^\top \ell_t$ of the forecaster with the loss of an arbitrary sequence of vectors $\mathbf{q}_1, \dots, \mathbf{q}_T$ in the simplex Δ_d , that is, with $\sum_{t=1}^T \mathbf{q}_t^\top \ell_t$. The bounds we obtain depend, of course, on the “regularity” of the comparison sequence. In the now classical results on tracking the best expert (as in [Herbster and Warmuth 1998](#); [Vovk 1999](#); [Herbster and Warmuth 2001](#); [Bousquet and Warmuth 2002](#)), this regularity is measured as the number of times $\mathbf{q}_t \neq \mathbf{q}_{t+1}$ (henceforth referred to as “hard shifts”). The main results of this paper show not only that these results may be generalized to obtain bounds in terms of “softer” regularity measures but that the same algorithms that were proposed with hard shift tracking in mind achieve such, perhaps surprisingly good, performance. Building on the general formulation introduced in Section 2, we derive such regret bounds for the fixed-share algorithm of [Herbster and Warmuth \(1998\)](#) and for the algorithms of [Bousquet and Warmuth \(2002\)](#).

In fact, it is advantageous to extend our analysis so that we not only compare the performance of the forecaster with sequences $\mathbf{q}_1, \dots, \mathbf{q}_T$ taking values in the simplex Δ_d of probability distributions but rather against arbitrary sequences $\mathbf{u}_1, \dots, \mathbf{u}_T \in \mathbb{R}_+^d$ of vectors with non-negative components. The loss of such a sequence is defined by $\sum_{t=1}^T \mathbf{u}_t^\top \ell_t$. For fair comparison, we measure the cumulative loss of the forecaster by $\sum_{t=1}^T \hat{\mathbf{p}}_t^\top \ell_t \|\mathbf{u}_t\|_1$. Of course, when $\mathbf{u}_t \in \Delta_d$, we recover the original notion of regret.

The norms $\|\mathbf{u}_1\|_1, \dots, \|\mathbf{u}_T\|_1$ may be viewed as a sequence of weights that give more or less importance to the instantaneous loss suffered at each step. Of particular interest is the case when $\|\mathbf{u}_t\|_1 \in [0, 1]$ which is the setting of “time selection functions” (see [Blum and Mansour 2007](#), Section 6). In particular, considering sequences $\|\mathbf{u}_t\|_1 \in \{0, 1\}$ that include the zero vector will provide us a simple way of deriving “adaptive” regret bounds, a notion introduced by [Hazan and Seshadhri \(2009\)](#).

The first regret bounds derived below measure the regularity of the sequence $\mathbf{u}_1^T = (\mathbf{u}_1, \dots, \mathbf{u}_T)$ in terms of the quantity

$$m(\mathbf{u}_1^T) = \sum_{t=1}^{T-1} D_{\text{TV}}(\mathbf{u}_{t+1}, \mathbf{u}_t) \quad (2)$$

where for $\mathbf{x} = (x_1, \dots, x_d), \mathbf{y} = (y_1, \dots, y_d) \in \mathbb{R}_+^d$, we define $D_{\text{TV}}(\mathbf{x}, \mathbf{y}) = \sum_{x_i \geq y_i} (x_i - y_i)$. Note that when $\mathbf{x}, \mathbf{y} \in \Delta_d$, we recover the total variation distance $D_{\text{TV}}(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_1$, while for general $\mathbf{x}, \mathbf{y} \in \mathbb{R}_+^d$, the quantity $D_{\text{TV}}(\mathbf{x}, \mathbf{y})$ is not necessarily symmetric and is always bounded by $\|\mathbf{x} - \mathbf{y}\|_1$. Note that when the vectors \mathbf{u}_t are incidence vectors $(0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^d$ of elements $i_t \in [d]$, then $m(\mathbf{u}_1^T)$ corresponds to the number of shifts of the sequence $i_1^T \in [d]^T$, and we recover from the results stated below the classical bounds for tracking the best expert.

3.1. Fixed-share update

We now analyze a specific instance of the generalized share algorithm corresponding to the update

$$\hat{p}_{j,t+1} = \sum_{i=1}^d \left(\frac{\alpha}{d} + (1 - \alpha) \mathbb{1}_{i=j} \right) v_{i,t+1} = \frac{\alpha}{d} + (1 - \alpha) v_{j,t+1}, \quad 0 \leq \alpha \leq 1. \quad (3)$$

Despite seemingly different statements, this update in Algorithm 1 can be seen to lead *exactly* to the fixed-share algorithm of [Herbster and Warmuth \(1998\)](#) for prediction with expert advice.

Proposition 2 *With the above update, for all $T \geq 1$, for all sequences ℓ_1, \dots, ℓ_T of loss vectors $\ell_t \in [0, 1]^d$, and for all $\mathbf{u}_1, \dots, \mathbf{u}_T \in \mathbb{R}_+^d$,*

$$\begin{aligned} \sum_{t=1}^T \|\mathbf{u}_t\|_1 \hat{\mathbf{p}}_t^\top \ell_t - \sum_{t=1}^T \mathbf{u}_t^\top \ell_t &\leq \frac{\|\mathbf{u}_1\|_1 \ln d}{\eta} + \frac{\eta}{8} \sum_{t=1}^T \|\mathbf{u}_t\|_1 \\ &\quad + \frac{m(\mathbf{u}_1^T)}{\eta} \ln \frac{d}{\alpha} + \frac{\sum_{t=1}^T \|\mathbf{u}_t\|_1 - m(\mathbf{u}_1^T) - 1}{\eta} \ln \frac{1}{1-\alpha}. \end{aligned}$$

We emphasize that the fixed-share forecaster does not need to “know” anything about the sequence of the norms $\|\mathbf{u}_t\|$. Of course, in order to minimize the obtained upper bound, the tuning parameters α, η need to be optimized and their values will depend on the maximal value of $m(\mathbf{u}_i^T)$ for the sequences one wishes to compete against. In particular, we obtain the following corollary, in which $h(x) = -x \ln x - (1-x) \ln(1-x)$ denotes the binary entropy function for $x \in [0, 1]$. We recall¹ that $h(x) \leq x \ln(e/x)$ for $x \in [0, 1]$.

Corollary 3 *Suppose Algorithm 1 is run with the update (3). Let $m_0 > 0$. For all $T \geq 1$, for all sequences ℓ_1, \dots, ℓ_T of loss vectors $\ell_t \in [0, 1]^d$, and for all $\mathbf{q}_1, \dots, \mathbf{q}_T \in \Delta_d$ with $m(\mathbf{q}_1^T) \leq m_0$,*

$$\sum_{t=1}^T \hat{\mathbf{p}}_t^\top \ell_t - \sum_{t=1}^T \mathbf{q}_t^\top \ell_t \leq \sqrt{\frac{T}{2} \left((m_0 + 1) \ln d + (T-1) h\left(\frac{m_0}{T-1}\right) \right)},$$

whenever η and α are optimally chosen in terms of m_0 and T .

If we only consider vectors of the form $\mathbf{q}_t = (0, \dots, 0, 1, 0, \dots, 0)$ then $m(\mathbf{q}_1^T)$ corresponds to the number of times $\mathbf{q}_{t+1} \neq \mathbf{q}_t$ in the sequence \mathbf{q}_1^T . We thus recover [Herbster and Warmuth \(1998, Theorem 1\)](#) and [Bousquet and Warmuth \(2002, Lemma 6\)](#) from the much more general Proposition 2.

Proof of Proposition 2 Applying Lemma 1 with $\mathbf{q}_t = \mathbf{u}_t / \|\mathbf{u}_t\|_1$, and multiplying by $\|\mathbf{u}_t\|_1$, we get for all $t \geq 1$ and $\mathbf{u}_t \in \mathbb{R}_+^d$

$$\|\mathbf{u}_t\|_1 \hat{\mathbf{p}}_t^\top \ell_t - \mathbf{u}_t^\top \ell_t \leq \frac{1}{\eta} \sum_{i=1}^d u_{i,t} \ln \frac{v_{i,t+1}}{\hat{p}_{i,t}} + \frac{\eta}{8} \|\mathbf{u}_t\|_1. \quad (4)$$

We now examine

$$\sum_{i=1}^d u_{i,t} \ln \frac{v_{i,t+1}}{\hat{p}_{i,t}} = \sum_{i=1}^d \left(u_{i,t} \ln \frac{1}{\hat{p}_{i,t}} - u_{i,t-1} \ln \frac{1}{v_{i,t}} \right) + \sum_{i=1}^d \left(u_{i,t-1} \ln \frac{1}{v_{i,t}} - u_{i,t} \ln \frac{1}{v_{i,t+1}} \right). \quad (5)$$

For the first term on the right-hand side, we have

$$\sum_{i=1}^d \left(u_{i,t} \ln \frac{1}{\hat{p}_{i,t}} - u_{i,t-1} \ln \frac{1}{v_{i,t}} \right) = \sum_{i: u_{i,t} \geq u_{i,t-1}} \left((u_{i,t} - u_{i,t-1}) \ln \frac{1}{\hat{p}_{i,t}} + u_{i,t-1} \ln \frac{v_{i,t}}{\hat{p}_{i,t}} \right)$$

1. As can be seen by noting that $\ln(1/(1-x)) < x/(1-x)$

$$+ \sum_{i: u_{i,t} < u_{i,t-1}} \underbrace{\left((u_{i,t} - u_{i,t-1}) \ln \frac{1}{v_{i,t}} + u_{i,t} \ln \frac{v_{i,t}}{\widehat{p}_{i,t}} \right)}_{\leq 0}. \quad (6)$$

In view of the update (3), we have $1/\widehat{p}_{i,t} \leq d/\alpha$ and $v_{i,t}/\widehat{p}_{i,t} \leq 1/(1-\alpha)$. Substituting in (6), we get

$$\begin{aligned} & \sum_{i=1}^d \left(u_{i,t} \ln \frac{1}{\widehat{p}_{i,t}} - u_{i,t-1} \ln \frac{1}{v_{i,t}} \right) \\ & \leq \sum_{i: u_{i,t} \geq u_{i,t-1}} (u_{i,t} - u_{i,t-1}) \ln \frac{d}{\alpha} + \left(\sum_{i: u_{i,t} \geq u_{i,t-1}} u_{i,t-1} + \sum_{i: u_{i,t} < u_{i,t-1}} u_{i,t} \right) \ln \frac{1}{1-\alpha} \\ & = D_{TV}(\mathbf{u}_t, \mathbf{u}_{t-1}) \ln \frac{d}{\alpha} + \underbrace{\left(\sum_{i=1}^d u_{i,t} - \sum_{i: u_{i,t} \geq u_{i,t-1}} (u_{i,t} - u_{i,t-1}) \right)}_{=\|\mathbf{u}_t\|_1 - D_{TV}(\mathbf{u}_t, \mathbf{u}_{t-1})} \ln \frac{1}{1-\alpha}. \end{aligned}$$

The sum of the second term in (5) telescopes. Substituting the obtained bounds in the first sum of the right-hand side in (5), and summing over $t = 2, \dots, T$, leads to

$$\begin{aligned} \sum_{t=2}^T \sum_{i=1}^d u_{i,t} \ln \frac{v_{i,t+1}}{\widehat{p}_{i,t}} & \leq m(\mathbf{u}_1^T) \ln \frac{d}{\alpha} + \left(\sum_{t=2}^T \|\mathbf{u}_t\|_1 - 1 - m(\mathbf{u}_1^T) \right) \ln \frac{1}{1-\alpha} \\ & \quad + \sum_{i=1}^d u_{i,1} \ln \frac{1}{v_{i,2}} - \underbrace{u_{i,T} \ln \frac{1}{v_{i,T+1}}}_{\leq 0}. \end{aligned}$$

We hence get from (4), which we use in particular for $t = 1$,

$$\begin{aligned} \sum_{t=1}^T \|\mathbf{u}_t\|_1 \widehat{\mathbf{p}}_t^\top \boldsymbol{\ell}_t - \mathbf{u}_t^\top \boldsymbol{\ell}_t & \leq \frac{1}{\eta} \sum_{i=1}^d u_{i,1} \ln \frac{1}{\widehat{p}_{i,1}} + \frac{\eta}{8} \sum_{t=1}^T \|\mathbf{u}_t\|_1 \\ & \quad + \frac{m(\mathbf{u}_1^T)}{\eta} \ln \frac{d}{\alpha} + \frac{\sum_{t=1}^T \|\mathbf{u}_t\|_1 - 1 - m(\mathbf{u}_1^T)}{\eta} \ln \frac{1}{1-\alpha}. \quad \blacksquare \end{aligned}$$

3.2. Sparse sequences: Bousquet-Warmuth updates

Bousquet and Warmuth (2002) proposed forecasters that are able to efficiently compete with the best sequence of experts among all those sequences that only switch a bounded number of times and also take a small number of different values. Such “sparse” sequences of experts appear naturally in many applications. In this section we show that their algorithms in fact work very well in comparison with a much larger class of sequences $\mathbf{u}_1, \dots, \mathbf{u}_T$ that are “regular”—that is, $m(\mathbf{u}_1^T)$, defined in (2) is small—and “sparse” in the sense that the quantity

$$n(\mathbf{u}_1^T) = \sum_{i=1}^d \max_{t=1, \dots, T} u_{i,t}$$

is small. Note that when $\mathbf{q}_t \in \Delta_d$ for all t , then two interesting upper bounds can be provided. First, denoting the union of the supports of these convex combinations by $S \subseteq [d]$, we have $n(\mathbf{q}_1^T) \leq |S|$, the cardinality of S . Also,

$$n(\mathbf{q}_1^T) \leq \left| \{ \mathbf{q}_t, \ t = 1, \dots, T \} \right|,$$

the cardinality of the pool of convex combinations. Thus, $n(\mathbf{u}_1^T)$ generalizes the notion of sparsity of Bousquet and Warmuth (2002).

Here we consider a family of shared updates of the form

$$\hat{p}_{j,t} = (1 - \alpha)v_{j,t} + \alpha \frac{w_{j,t}}{Z_t}, \quad 0 \leq \alpha \leq 1, \quad (7)$$

where the $w_{j,t}$ are nonnegative weights that may depend on past and current pre-weights and $Z_t = \sum_{i=1}^d w_{i,t}$ is a normalization constant. Shared updates of this form were proposed by Bousquet and Warmuth (2002, Sections 3 and 5.2).

Apart from generalizing the regret bounds of Bousquet and Warmuth (2002), we believe that the analysis given below is significantly simpler and more transparent. We are also able to slightly improve their original bounds.

We focus on choices of the weights $w_{j,t}$ that satisfy the following conditions: there exists a constant $C \geq 1$ such that for all $j = 1, \dots, d$ and $t = 1, \dots, T$,

$$v_{j,t} \leq w_{j,t} \leq 1 \quad \text{and} \quad C w_{j,t+1} \geq w_{j,t}. \quad (8)$$

The next result improves on Proposition 2 when $T \ll d$ and $n(\mathbf{u}_1^T) \ll m(\mathbf{u}_1^T)$, that is, when the dimension (or number of experts) d is large but the sequence \mathbf{u}_1^T is sparse.

Proposition 4 *Suppose Algorithm 1 is run with the shared update (7) with weights satisfying the conditions (8). Then for all $T \geq 1$, for all sequences ℓ_1, \dots, ℓ_T of loss vectors $\ell_t \in [0, 1]^d$, and for all sequences $\mathbf{u}_1, \dots, \mathbf{u}_T \in \mathbb{R}_+^d$,*

$$\begin{aligned} \sum_{t=1}^T \|\mathbf{u}_t\|_1 \hat{\mathbf{p}}_t^\top \ell_t - \sum_{t=1}^T \mathbf{u}_t^\top \ell_t &\leq \frac{n(\mathbf{u}_1^T) \ln d}{\eta} + \frac{n(\mathbf{u}_1^T) T \ln C}{\eta} + \frac{\eta}{8} \sum_{t=1}^T \|\mathbf{u}_t\|_1 \\ &\quad + \frac{m(\mathbf{u}_1^T)}{\eta} \ln \frac{\max_{t \leq T} Z_t}{\alpha} + \frac{\sum_{t=1}^T \|\mathbf{u}_t\|_1 - m(\mathbf{u}_1^T) - 1}{\eta} \ln \frac{1}{1 - \alpha}. \end{aligned}$$

Proof The beginning and the end of the proof are similar to the one of Proposition 2, as they do not depend on the specific weight update. In particular, inequalities (4) and (5) remain the same. The proof is modified after (6), which this time we upper bound using the first condition in (8),

$$\begin{aligned} \sum_{i=1}^d \left(u_{i,t} \ln \frac{1}{\hat{p}_{i,t}} - u_{i,t-1} \ln \frac{1}{v_{i,t}} \right) &= \sum_{i: u_{i,t} \geq u_{i,t-1}} (u_{i,t} - u_{i,t-1}) \ln \frac{1}{\hat{p}_{i,t}} + u_{i,t-1} \ln \frac{v_{i,t}}{\hat{p}_{i,t}} \\ &\quad + \sum_{i: u_{i,t} < u_{i,t-1}} \underbrace{(u_{i,t} - u_{i,t-1})}_{\leq 0} \underbrace{\ln \frac{1}{v_{i,t}}}_{\geq \ln(1/w_{i,t})} + u_{i,t} \ln \frac{v_{i,t}}{\hat{p}_{i,t}}. \quad (9) \end{aligned}$$

By definition of the shared update (7), we have $1/\widehat{p}_{i,t} \leq Z_t/(\alpha w_{i,t})$ and $v_{i,t}/\widehat{p}_{i,t} \leq 1/(1-\alpha)$. We then upper bound the quantity at hand in (9) by

$$\begin{aligned} & \sum_{i: u_{i,t} \geq u_{i,t-1}} (u_{i,t} - u_{i,t-1}) \ln \left(\frac{Z_t}{\alpha w_{i,t}} \right) + \left(\sum_{i: u_{i,t} \geq u_{i,t-1}} u_{i,t-1} + \sum_{i: u_{i,t} < u_{i,t-1}} u_{i,t} \right) \ln \frac{1}{1-\alpha} \\ & + \sum_{i: u_{i,t} < u_{i,t-1}} (u_{i,t} - u_{i,t-1}) \ln \frac{1}{w_{i,t}} \\ = & D_{\text{TV}}(\mathbf{u}_t, \mathbf{u}_{t-1}) \ln \frac{Z_t}{\alpha} + (\|\mathbf{u}_t\|_1 - D_{\text{TV}}(\mathbf{u}_t, \mathbf{u}_{t-1})) \ln \frac{1}{1-\alpha} + \sum_{i=1}^d (u_{i,t} - u_{i,t-1}) \ln \frac{1}{w_{i,t}}. \end{aligned}$$

Proceeding as in the end of the proof of Proposition 2, we then get the claimed bound, provided that we can show that

$$\sum_{t=2}^T \sum_{i=1}^d (u_{i,t} - u_{i,t-1}) \ln \frac{1}{w_{i,t}} \leq n(\mathbf{u}_1^T) (\ln d + T \ln C) - \|\mathbf{u}_1\|_1 \ln d,$$

which we do next. Indeed, the left-hand side can be rewritten as

$$\begin{aligned} & \sum_{t=2}^T \sum_{i=1}^d \left(u_{i,t} \ln \frac{1}{w_{i,t}} - u_{i,t} \ln \frac{1}{w_{i,t+1}} \right) + \sum_{t=2}^T \sum_{i=1}^d \left(u_{i,t} \ln \frac{1}{w_{i,t+1}} - u_{i,t-1} \ln \frac{1}{w_{i,t}} \right) \\ \leq & \left(\sum_{t=2}^T \sum_{i=1}^d u_{i,t} \ln \frac{C w_{i,t+1}}{w_{i,t}} \right) + \left(\sum_{i=1}^d u_{i,T} \ln \frac{1}{w_{i,T+1}} - \sum_{i=1}^d u_{i,1} \ln \frac{1}{w_{i,2}} \right) \\ \leq & \left(\sum_{i=1}^d \left(\max_{t=1,\dots,T} u_{i,t} \right) \sum_{t=2}^T \ln \frac{C w_{i,t+1}}{w_{i,t}} \right) + \left(\sum_{i=1}^d \left(\max_{t=1,\dots,T} u_{i,t} \right) \ln \frac{1}{w_{i,T+1}} - \sum_{i=1}^d u_{i,1} \ln \frac{1}{w_{i,2}} \right) \\ = & \sum_{i=1}^d \left(\max_{t=1,\dots,T} u_{i,t} \right) \left((T-1) \ln C + \ln \frac{1}{w_{i,2}} \right) - \sum_{i=1}^d u_{i,1} \ln \frac{1}{w_{i,2}}, \end{aligned}$$

where we used $C \geq 1$ for the first inequality and the second condition in (8) for the second inequality. The proof is concluded by noting that (8) entails $w_{i,2} \geq (1/C)w_{i,1} \geq (1/C)v_{i,1} = 1/(dC)$ and that the coefficient $\max_{t=1,\dots,T} u_{i,t} - u_{i,1}$ in front of $\ln(1/w_{i,2})$ is nonnegative. \blacksquare

We now generalize Corollaries 8 and 9 of Bousquet and Warmuth (2002) by showing two specific instances of the generic update (7) that satisfy (8). The first update uses $w_{j,t} = \max_{s \leq t} v_{j,s}$. Then (8) is satisfied with $C = 1$. Moreover, since a sum of maxima of nonnegative elements is smaller than the sum of the sums, $Z_t \leq \min\{d, t\} \leq T$. This immediately gives the following result.

Corollary 5 Suppose Algorithm 1 is run with the update (7) with $w_{j,t} = \max_{s \leq t} v_{j,s}$. For all $T \geq 1$, for all sequences ℓ_1, \dots, ℓ_T of loss vectors $\ell_t \in [0, 1]^d$, and for all $\mathbf{q}_1, \dots, \mathbf{q}_T \in \Delta_d$,

$$\sum_{t=1}^T \widehat{\mathbf{p}}_t^\top \ell_t - \sum_{t=1}^T \mathbf{q}_t^\top \ell_t \leq \frac{n(\mathbf{q}_1^T) \ln d}{\eta} + \frac{\eta}{8} T + \frac{m(\mathbf{q}_1^T)}{\eta} \ln \frac{T}{\alpha} + \frac{T - m(\mathbf{q}_1^T) - 1}{\eta} \ln \frac{1}{1-\alpha}.$$

The second update we discuss uses $w_{j,t} = \max_{s \leq t} e^{\gamma(s-t)} v_{j,s}$ in (7) for some $\gamma > 0$. Both conditions in (8) are satisfied with $C = e^\gamma$. One also has that

$$Z_t \leq d \quad \text{and} \quad Z_t \leq \sum_{\tau \geq 0} e^{-\gamma\tau} = \frac{1}{1 - e^{-\gamma}} \leq \frac{1}{\gamma}$$

as $e^x \geq 1 + x$ for all real x . The bound of Proposition 4 then instantiates as

$$\frac{n(\mathbf{q}_1^T) \ln d}{\eta} + \frac{n(\mathbf{q}_1^T) T \gamma}{\eta} + \frac{\eta T}{8} + \frac{m(\mathbf{q}_1^T)}{\eta} \ln \frac{\min\{d, 1/\gamma\}}{\alpha} + \frac{T - m(\mathbf{q}_1^T) - 1}{\eta} \ln \frac{1}{1 - \alpha}$$

when sequences $\mathbf{u}_t = \mathbf{q}_t \in \Delta_d$ are considered. This bound is best understood when γ is tuned optimally based on T and on two bounds m_0 and n_0 over the quantities $m(\mathbf{q}_1^T)$ and $n(\mathbf{q}_1^T)$. Indeed, by optimizing $n_0 T \gamma + m_0 \ln(1/\gamma)$, i.e., by choosing $\gamma = m_0/(n_0 T)$, one gets a bound that improves on the one of the previous corollary:

Corollary 6 *Let $m_0, n_0 > 0$. Suppose Algorithm 1 is run with the update $w_{j,t} = \max_{s \leq t} e^{\gamma(s-t)} v_{j,s}$ where $\gamma = m_0/(n_0 T)$. For all $T \geq 1$, for all sequences ℓ_1, \dots, ℓ_T of loss vectors $\ell_t \in [0, 1]^d$, and for all $\mathbf{q}_1, \dots, \mathbf{q}_T \in \Delta_d$ such that $m(\mathbf{q}_1^T) \leq m_0$ and $n(\mathbf{q}_1^T) \leq n_0$, we have*

$$\begin{aligned} \sum_{t=1}^T \hat{\mathbf{p}}_t^\top \ell_t - \sum_{t=1}^T \mathbf{q}_t^\top \ell_t &\leq \frac{n_0 \ln d}{\eta} + \frac{m_0}{\eta} \left(1 + \ln \min \left\{ d, \frac{n_0 T}{m_0} \right\} \right) \\ &\quad + \frac{\eta T}{8} + \frac{m_0}{\eta} \ln \frac{1}{\alpha} + \frac{T - m_0 - 1}{\eta} \ln \frac{1}{1 - \alpha}. \end{aligned}$$

As the factors $e^{-\gamma t}$ cancel out in the numerator and denominator of the ratio in (7), there is a straightforward implementation of the algorithm (not requiring the knowledge of T) that needs to maintain only d weights.

In contrast, the corresponding algorithm of Bousquet and Warmuth (2002), using the updates $\hat{p}_{j,t} = (1 - \alpha)v_{j,t} + \alpha S_t^{-1} \sum_{s \leq t-1} (s-t)^{-1} v_{j,s}$ or $\hat{p}_{j,t} = (1 - \alpha)v_{j,t} + \alpha S_t^{-1} \max_{s \leq t-1} (s-t)^{-1} v_{j,s}$, where S_t denote normalization factors, needs to maintain $O(dT)$ weights with a naive implementation, and $O(d \ln T)$ weights with a more sophisticated one. In addition, the obtained bounds are slightly worse than the one stated above in Corollary 6 as an additional factor of $m_0 \ln(1 + \ln T)$ is present in Bousquet and Warmuth (2002, Corollary 9).

4. Adaptive regret

Next we show how the results of the previous section, e.g., Proposition 2, imply guarantees in terms of adaptive regret—a notion introduced by Hazan and Seshadhri (2009) as follows. For $\tau_0 \in \{1, \dots, T\}$, the τ_0 -adaptive regret of a forecaster is defined by

$$\mathcal{R}_T^{\tau_0\text{-adapt}} = \max_{\substack{[r, s] \subset [1, T] \\ s+1-r \leq \tau_0}} \left\{ \sum_{t=r}^s \hat{\mathbf{p}}_t^\top \ell_t - \min_{\mathbf{q} \in \Delta_d} \sum_{t=r}^s \mathbf{q}^\top \ell_t \right\}. \quad (10)$$

Adaptive regret is an alternative way to measure the performance of a forecaster against a changing environment. It is a straightforward observation that adaptive regret bounds also lead to shifting

regret bounds (in terms of hard shifts). Here we show that these two notions of regret share an even tighter connection, as they can be both viewed as instances of the same *alma mater* bound, e.g., Proposition 2.

Hazan and Seshadhri (2009) essentially considered the case of online convex optimization with exp-concave loss function (see Section 6 below). In case of general convex functions, they also mentioned that the greedy projection forecaster of Zinkevich (2003) —i.e., mirror descent with a quadratic regularizer— enjoys adaptive regret guarantees. This forecaster can be implemented on the simplex in time $O(d)$ —see, e.g., Duchi et al. (2008). We now show that the simpler fixed-share algorithm has a similar adaptive regret bound.

Proposition 7 *Suppose that Algorithm 1 is run with the shared update (3). Then for all $T \geq 1$, for all sequences ℓ_1, \dots, ℓ_T of loss vectors $\ell_t \in [0, 1]^d$, and for all $\tau_0 \in \{1, \dots, T\}$,*

$$\mathcal{R}_T^{\tau_0\text{-adapt}} \leq \frac{1}{\eta} \ln \frac{d}{\alpha} + \frac{\tau_0 - 1}{\eta} \ln \frac{1}{1 - \alpha} + \frac{\eta}{8} \tau_0.$$

In particular, when η and α are chosen optimally (depending on τ_0 and T)

$$\mathcal{R}_T^{\tau_0\text{-adapt}} \leq \sqrt{\frac{\tau_0}{2} \left(\tau_0 h\left(\frac{1}{\tau_0}\right) + \ln d \right)} \leq \sqrt{\frac{\tau_0}{2} \ln(ed\tau_0)}.$$

Proof For $1 \leq r \leq s \leq T$ and $\mathbf{q} \in \Delta_d$, the regret in the right-hand side of (10) equals the regret considered in Proposition 2 against the sequence \mathbf{u}_1^T defined as $\mathbf{u}_t = \mathbf{q}$ for $t = r, \dots, s$ and $\mathbf{0} = (0, \dots, 0)$ for the remaining t . When $r \geq 2$, this sequence is such that $D_{TV}(\mathbf{u}_r, \mathbf{u}_{r-1}) = D_{TV}(\mathbf{q}, \mathbf{0}) = 1$ and $D_{TV}(\mathbf{u}_{s+1}, \mathbf{u}_s) = D_{TV}(\mathbf{0}, \mathbf{q}) = 0$ so that $m(\mathbf{u}_1^T) = 1$, while $\|\mathbf{u}_1\|_1 = 0$. When $r = 1$, we have $\|\mathbf{u}_1\|_1 = 1$ and $m(\mathbf{u}_1^T) = 0$. In all cases, $m(\mathbf{u}_1^T) + \|\mathbf{u}_1\|_1 = 1$. Specializing the bound of Proposition 2 to the thus defined sequence \mathbf{u}_1^T gives the result. ■

5. Online tuning of the parameters

The forecasters studied above need their parameters η and α to be tuned according to various quantities, including the time horizon T . We show here how the trick of Auer et al. (2002) of having these parameters vary over time can be extended to our setting. For the sake of concreteness we focus on the fixed-share update, i.e., Algorithm 1 run with the update (3). We respectively replace steps 3 and 4 of its description by the loss and shared updates

$$v_{j,t+1} = \frac{\hat{p}_{j,t}^{\frac{\eta_t}{\eta_{t-1}}} e^{-\eta_t \ell_{j,t}}}{\sum_{i=1}^d \hat{p}_{i,t}^{\frac{\eta_t}{\eta_{t-1}}} e^{-\eta_t \ell_{i,t}}} \quad \text{and} \quad p_{j,t+1} = \frac{\alpha_t}{d} + (1 - \alpha_t) v_{j,t+1}, \quad (11)$$

for all $t \geq 1$ and all $j \in [d]$, where (η_t) and (α_t) are two sequences of positive numbers, indexed by $\tau \geq 1$. We also conventionally define $\eta_0 = \eta_1$. Proposition 2 is then adapted in the following way (when $\eta_t \equiv \eta$ and $\alpha_t \equiv \alpha$, Proposition 2 is exactly recovered).

Proposition 8 *The forecaster based on the above updates (11) is such that whenever $\eta_t \leq \eta_{t-1}$ and $\alpha_t \leq \alpha_{t-1}$ for all $t \geq 1$, the following performance bound is achieved. For all $T \geq 1$, for all sequences ℓ_1, \dots, ℓ_T of loss vectors $\ell_t \in [0, 1]^d$, and for all $\mathbf{u}_1, \dots, \mathbf{u}_T \in \mathbb{R}_+^d$,*

$$\begin{aligned} \sum_{t=1}^T \|\mathbf{u}_t\|_1 \hat{\mathbf{p}}_t^\top \ell_t - \sum_{t=1}^T \mathbf{u}_t^\top \ell_t &\leq \left(\frac{\|\mathbf{u}_1\|_1}{\eta_1} + \sum_{t=2}^T \|\mathbf{u}_t\|_1 \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \right) \ln d \\ &\quad + \frac{m(\mathbf{u}_1^T)}{\eta_T} \ln \frac{d(1 - \alpha_T)}{\alpha_T} + \sum_{t=2}^T \frac{\|\mathbf{u}_t\|_1}{\eta_{t-1}} \ln \frac{1}{1 - \alpha_t} + \sum_{t=1}^T \frac{\eta_{t-1}}{8} \|\mathbf{u}_t\|_1. \end{aligned}$$

Due to space constraints, we only instantiate the obtained bound to the case of T -adaptive regret guarantees, when T is unknown and/or can increase without bounds.

Corollary 9 *The forecaster based on the above updates with $\eta_t = \sqrt{(\ln(dt))/t}$ for $t \geq 3$ and $\eta_0 = \eta_1 = \eta_2 = \eta_3$ on the one hand, $\alpha_t = 1/t$ on the other hand, is such that for all $T \geq 3$ and for all sequences ℓ_1, \dots, ℓ_T of loss vectors $\ell_t \in [0, 1]^d$,*

$$\max_{[r,s] \subset [1,T]} \left\{ \sum_{t=r}^s \hat{\mathbf{p}}_t^\top \ell_t - \min_{\mathbf{q} \in \Delta_d} \sum_{t=r}^s \mathbf{q}^\top \ell_t \right\} \leq \sqrt{2T \ln(dT)} + \sqrt{3 \ln(3d)}.$$

Proof The sequence $n \mapsto \ln(n)/n$ is only non-increasing after round $n \geq 3$, so that the defined sequences of (α_t) and (η_t) are non-increasing, as desired. For a given pair (r, s) and a given $\mathbf{q} \in \Delta_d$, we consider the sequence ν_1^T defined in the proof of Proposition 7; it satisfies that $m(\mathbf{u}_1^T) \leq 1$ and $\|\mathbf{u}_t\|_1 \leq 1$ for all $t \geq 1$. Therefore, Proposition 8 ensures that

$$\begin{aligned} \sum_{t=r}^s \hat{\mathbf{p}}_t^\top \ell_t - \min_{\mathbf{q} \in \Delta_d} \sum_{t=r}^s \mathbf{q}^\top \ell_t &\leq \frac{\ln d}{\eta_T} + \frac{1}{\eta_T} \ln \underbrace{\frac{d(1 - \alpha_T)}{\alpha_T}}_{\leq dT} + \underbrace{\sum_{t=2}^T \frac{1}{\eta_{t-1}} \ln \frac{1}{1 - \alpha_t}}_{\leq (1/\eta_T) \sum_{t=2}^T \ln(t/(t-1)) = (\ln T)/\eta_T} + \sum_{t=1}^T \frac{\eta_{t-1}}{8}. \end{aligned}$$

It only remains to substitute the proposed values of η_t and to note that

$$\sum_{t=1}^T \eta_{t-1} \leq 3\eta_3 + \sum_{t=3}^{T-1} \frac{1}{\sqrt{t}} \sqrt{\ln(dT)} \leq 3\sqrt{\frac{\ln(3d)}{3}} + 2\sqrt{T} \sqrt{\ln(dT)}.$$

■

6. Online convex optimization and exp-concave loss functions

By using a standard reduction, the results of the previous sections can be applied to online convex optimization on the simplex. In this setting, at each step t the forecaster chooses $\hat{\mathbf{p}}_t \in \Delta_d$ and then is given access to a convex loss $\ell_t : \Delta_d \rightarrow [0, 1]$. Now, using Algorithm 1 with the loss vector $\ell_t \in \partial \ell_t(\hat{\mathbf{p}}_t)$ given by a subgradient of ℓ_t leads to the desired bounds. Indeed, by the convexity of ℓ_t , the regret at each time t with respect to any vector $\mathbf{u}_t \in \mathbb{R}_+^d$ with $\|\mathbf{u}_t\|_1 > 0$ is then bounded as

$$\|\mathbf{u}_t\|_1 \left(\ell_t(\hat{\mathbf{p}}_t) - \ell_t\left(\frac{\mathbf{u}_t}{\|\mathbf{u}_t\|_1}\right) \right) \leq (\|\mathbf{u}_t\|_1 \hat{\mathbf{p}}_t - \mathbf{u}_t)^\top \ell_t.$$

6.1. Exp-concave loss functions

Recall that a loss function ℓ_t is called η_0 -exp-concave if $e^{-\eta_0 \ell_t}$ is concave. (In particular, exp-concavity implies convexity.) [Bousquet and Warmuth \(2002\)](#) study shifting regret for exp-concave loss functions. However, they define the regret of an element \mathbf{q}_1^T of the comparison class (a sequence of elements in Δ_d) by

$$\sum_{t=1}^T \left(\ell_t(\hat{\mathbf{p}}_t) - \mathbf{q}_1^T \ell_t \right) \quad (12)$$

where $\ell_t = (\ell_t(e_1), \dots, \ell_t(e_d))$ and e_1, \dots, e_d are the elements of the canonical basis of \mathbb{R}^d . This corresponds to the linear optimization case studied in the previous sections. However, due to exp-concavity, (1) can be replaced by an application of Jensen's inequality, namely,

$$\ell_t(\hat{\mathbf{p}}_t) \leq -\frac{1}{\eta_0} \ln \left(\sum_{j=1}^d \hat{p}_{j,t} e^{-\eta_0 \ell_t(e_j)} \right).$$

Hence the various propositions and corollaries of Sections 3 and 4 still hold true for the regret (12) up to some modifications (deletion of the terms linear in η , assumption of exp-concavity, boundedness no longer needed). For the sake of concreteness, we illustrate the required modifications on Proposition 4.

Proposition 10 *Suppose Algorithm 1 is run with the shared update (7) with weights satisfying the conditions (8) and for the choice $\eta = \eta_0$. Then for all $T \geq 1$, for all sequences ℓ_1, \dots, ℓ_T of η_0 -exp-concave loss functions, and for all sequences $\mathbf{u}_1, \dots, \mathbf{u}_T \in \mathbb{R}_+^d$,*

$$\begin{aligned} \sum_{t=1}^T \|\mathbf{u}_t\|_1 \ell_t(\hat{\mathbf{p}}_t) - \sum_{t=1}^T \mathbf{u}_t^T \ell_t &\leq \frac{(\mathbf{u}_1^T) \ln d}{\eta_0} + \frac{n(\mathbf{u}_1^T) T \ln C}{\eta_0} \\ &\quad + \frac{m(\mathbf{u}_1^T)}{\eta_0} \ln \frac{\max_{t \leq T} Z_t}{\alpha} + \frac{\sum_{t=1}^T \|\mathbf{u}_t\|_1 - m(\mathbf{u}_1^T) - 1}{\eta_0} \ln \frac{1}{1 - \alpha}. \end{aligned}$$

We now turn to the more ambitious goal of controlling regrets of the form $\sum_{t=1}^T (\ell_t(\hat{\mathbf{p}}_t) - \ell_t(\mathbf{q}_t))$ where losses ℓ_t are exp-concave. [Hazan and Seshadhri \(2009\)](#) constructed algorithms with T -adaptive regret of the order of $O(\ln^2 T)$ and running in time $\text{poly}(d, \log T)$. They also constructed different algorithms with T -adaptive regret bounded by $O(\ln T)$ and running time $\text{poly}(d, T)$.

Next, we show the first logarithmic shifting bounds for exp-concave loss functions. However, we only do so against sequences \mathbf{q}_1^T of elements in Δ_d , i.e., we offer here no general bound in terms of linear vectors \mathbf{u}_1^T that would unify here as well the view between tracking bounds and adaptive regret bounds. Besides, we get shifting bounds only in terms of hard shifts

$$s(\mathbf{q}_1^T) = |\{t = 2, \dots, T : \mathbf{q}_t \neq \mathbf{q}_{t-1}\}|.$$

Obviously, getting unifying bounds in terms of soft shifts of sequences \mathbf{u}_1^T of linear vectors is an important open question, which we leave for future research. To get our bound, we mix ideas of [Herbster and Warmuth \(1998\)](#) and [Blum and Kalai \(1997\)](#). We define a prior over the sequences of convex weight vectors as the distribution of the following homogeneous Markov chain $\mathbf{Q}_1, \mathbf{Q}_2, \dots$: The starting vector \mathbf{Q}_1 is drawn at random according to the uniform distribution μ over Δ_d . Then,

given \mathbf{Q}_{t-1} , the next element \mathbf{Q}_t is equal to \mathbf{Q}_{t-1} with probability $1 - \alpha$ and with probability α is drawn at random according to μ . In the sequel, all probabilities \mathbb{P} and expectations \mathbb{E} will be with respect to this Markov chain. Now, the convex weight vector used at time $t \geq 1$ by the forecaster is

$$\hat{\mathbf{p}}_t = \frac{\mathbb{E}[\mathbf{Q}_t e^{-\eta_0 L_{t-1}(\mathbf{Q}_1^{t-1})}]}{\mathbb{E}[e^{-\eta_0 L_{t-1}(\mathbf{Q}_1^{t-1})}]}, \quad \text{where} \quad L_{t-1}(\mathbf{Q}_1^{t-1}) = \sum_{s=1}^{t-1} \ell_s(\mathbf{Q}_s) \quad (13)$$

(with the convention that an empty sum is null). For this forecaster, we get the following performance bound, whose proof can be found in appendix.

Proposition 11 *For all $T \geq 1$, for all sequences ℓ_1, \dots, ℓ_T of η_0 -exp-concave loss functions taking values in $[0, L]$, the cumulative loss of the above forecaster is bounded for all sequences $\mathbf{q}_1, \dots, \mathbf{q}_T \in \Delta_d$ by*

$$\begin{aligned} \sum_{t=1}^T \ell_t(\hat{\mathbf{p}}_t) - \sum_{t=1}^T \ell_t(\mathbf{q}_t) &\leq \frac{(s(\mathbf{q}_1^T) + 1)(d - 1)}{\eta} \max \left\{ 1, \ln \frac{e \eta L T}{(s(\mathbf{q}_1^T) + 1)(d - 1)} \right\} \\ &\quad + \frac{s(\mathbf{q}_1^T)}{\eta} \ln \frac{1}{\alpha} + \frac{T - s(\mathbf{q}_1^T) - 1}{\eta} \ln \frac{1}{1 - \alpha}. \end{aligned}$$

Under the imposition of a bound s_0 on the numbers of hard shifts $s(\mathbf{q}_1^T)$ and up to a tuning of α in terms of s_0 and T , the last two terms of the bound are smaller than $T h(s_0/T) \leq s_0 \ln(es_0/T)$ and therefore, the whole regret bound is $O((ds_0/\eta_0) \ln T)$.

References

- P. Auer, N. Cesa-Bianchi, and C. Gentile. Adaptive and self-confident on-line learning algorithms. *Journal of Computer and System Sciences*, 64:48–75, 2002.
- A. Blum and A. Kalai. Universal portfolios with and without transaction costs. In *Proceedings of the 10th Annual Conference on Learning Theory (COLT)*, pages 309–313. Springer, 1997.
- A. Blum and Y. Mansour. From external to internal regret. *Journal of Machine Learning Research*, 8:1307–1324, 2007.
- O. Bousquet and M.K. Warmuth. Tracking a small set of experts by mixing past posteriors. *Journal of Machine Learning Research*, 3:363–396, 2002.
- N. Cesa-Bianchi and G. Lugosi. *Prediction, learning, and games*. Cambridge University Press, 2006.
- J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra. Efficient projections onto the ℓ_1 -ball for learning in high dimensions. In *Proceedings of the 25th International Conference on Machine Learning, ICML 2008*, 2008.
- A. György, T. Linder, and G. Lugosi. Tracking the best of many experts. In *Proceedings of the 18th Annual Conference on Learning Theory (COLT)*, pages 204–216, Bertinoro, Italy, Jun. 2005. Springer.
- E. Hazan and C. Seshadhri. Efficient learning algorithms for changing environment. *Proceedings of the 26th International Conference of Machine Learning (ICML)*, 2009.
- M. Herbster and M. Warmuth. Tracking the best expert. *Machine Learning*, 32:151–178, 1998.
- M. Herbster and M. Warmuth. Tracking the best linear predictor. *Journal of Machine Learning Research*, 1:281–309, 2001.
- V. Vovk. Derandomizing stochastic prediction strategies. *Machine Learning*, 35(3):247–282, Jun. 1999.
- M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th International Conference on Machine Learning, ICML 2003*, 2003.

Appendix A. Proof of Proposition 8

We first adapt Lemma 1.

Lemma 12 *The forecaster based on the loss and shared updates (11) satisfies, for all $t \geq 1$ and for all $\mathbf{q}_t \in \Delta_d$,*

$$(\hat{\mathbf{p}}_t - \mathbf{q}_t)^\top \boldsymbol{\ell}_t \leq \sum_{i=1}^d q_{i,t} \left(\frac{1}{\eta_{t-1}} \ln \frac{1}{\hat{p}_{i,t}} - \frac{1}{\eta_t} \ln \frac{1}{v_{i,t+1}} \right) + \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \ln d + \frac{\eta_{t-1}}{8},$$

whenever $\eta_t \leq \eta_{t-1}$.

Proof By Hoeffding's inequality,

$$\sum_{j=1}^d \hat{p}_{j,t} \ell_{j,t} \leq -\frac{1}{\eta_{t-1}} \ln \left(\sum_{j=1}^d \hat{p}_{j,t} e^{-\eta_{t-1} \ell_{j,t}} \right) + \frac{\eta_{t-1}}{8}.$$

By Jensen's inequality, since $\eta_t \leq \eta_{t-1}$ and thus $x \mapsto x^{\frac{\eta_{t-1}}{\eta_t}}$ is convex,

$$\frac{1}{d} \sum_{j=1}^d \hat{p}_{j,t} e^{-\eta_{t-1} \ell_{j,t}} = \frac{1}{d} \sum_{j=1}^d \left(\hat{p}_{j,t}^{\frac{\eta_t}{\eta_{t-1}}} e^{-\eta_t \ell_{j,t}} \right)^{\frac{\eta_{t-1}}{\eta_t}} \geq \left(\frac{1}{d} \sum_{j=1}^d \hat{p}_{j,t}^{\frac{\eta_t}{\eta_{t-1}}} e^{-\eta_t \ell_{j,t}} \right)^{\frac{\eta_{t-1}}{\eta_t}}.$$

Substituting in Hoeffding's bound we get

$$\hat{\mathbf{p}}_t^\top \boldsymbol{\ell}_t \leq -\frac{1}{\eta_t} \ln \left(\sum_{j=1}^d \hat{p}_{j,t}^{\frac{\eta_t}{\eta_{t-1}}} e^{-\eta_t \ell_{j,t}} \right) + \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \ln d + \frac{\eta_{t-1}}{8}.$$

Now, by definition of the loss update in (11), for all $i \in [d]$,

$$\sum_{j=1}^d \hat{p}_{j,t}^{\frac{\eta_t}{\eta_{t-1}}} e^{-\eta_t \ell_{j,t}} = \frac{1}{v_{i,t+1}} \hat{p}_{i,t}^{\frac{\eta_t}{\eta_{t-1}}} e^{-\eta_t \ell_{i,t}},$$

which, after substitution in the previous bound leads to the inequality

$$\hat{\mathbf{p}}_t^\top \boldsymbol{\ell}_t \leq \ell_{i,t} + \frac{1}{\eta_{t-1}} \ln \frac{1}{\hat{p}_{i,t}} - \frac{1}{\eta_t} \ln \frac{1}{v_{i,t+1}} + \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \ln d + \frac{\eta_{t-1}}{8},$$

valid for all $i \in [d]$. The proof is concluded by taking a convex aggregation over i with respect to \mathbf{q}_t . \blacksquare

The proof of Proposition 8 follows the steps of the one of Proposition 2; we sketch it below.

Proof of Proposition 8 Applying Lemma 12 with $\mathbf{q}_t = \mathbf{u}_t / \|\mathbf{u}_t\|_1$, and multiplying by $\|\mathbf{u}_t\|_1$, we get for all $t \geq 1$ and $\mathbf{u}_t \in \mathbb{R}_+^d$,

$$\begin{aligned} \|\mathbf{u}_t\|_1 \widehat{\mathbf{p}}_t^\top \boldsymbol{\ell}_t - \mathbf{u}_t^\top \boldsymbol{\ell}_t &\leq \frac{1}{\eta_{t-1}} \sum_{i=1}^d u_{i,t} \ln \frac{1}{\widehat{p}_{i,t}} - \frac{1}{\eta_t} \sum_{i=1}^d u_{i,t} \ln \frac{1}{v_{i,t+1}} \\ &\quad + \|\mathbf{u}_t\|_1 \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \ln d + \frac{\eta_{t-1}}{8} \|\mathbf{u}_t\|_1. \end{aligned} \quad (14)$$

We will sum these bounds over $t \geq 1$ to get the desired result but need to perform first some additional boundings for $t \geq 2$; in particular, we examine

$$\begin{aligned} &\frac{1}{\eta_{t-1}} \sum_{i=1}^d u_{i,t} \ln \frac{1}{\widehat{p}_{i,t}} - \frac{1}{\eta_t} \sum_{i=1}^d u_{i,t} \ln \frac{1}{v_{i,t+1}} \\ &= \frac{1}{\eta_{t-1}} \sum_{i=1}^d \left(u_{i,t} \ln \frac{1}{\widehat{p}_{i,t}} - u_{i,t-1} \ln \frac{1}{v_{i,t}} \right) + \sum_{i=1}^d \left(\frac{u_{i,t-1}}{\eta_{t-1}} \ln \frac{1}{v_{i,t}} - \frac{u_{i,t}}{\eta_t} \ln \frac{1}{v_{i,t+1}} \right), \end{aligned} \quad (15)$$

where the first difference in the right-hand side can be bounded as in (6) by

$$\begin{aligned} &\sum_{i=1}^d \left(u_{i,t} \ln \frac{1}{\widehat{p}_{i,t}} - u_{i,t-1} \ln \frac{1}{v_{i,t}} \right) \\ &\leq \sum_{i: u_{i,t} \geq u_{i,t-1}} \left((u_{i,t} - u_{i,t-1}) \ln \frac{1}{\widehat{p}_{i,t}} + u_{i,t-1} \ln \frac{v_{i,t}}{\widehat{p}_{i,t}} \right) + \sum_{i: u_{i,t} < u_{i,t-1}} u_{i,t} \ln \frac{v_{i,t}}{\widehat{p}_{i,t}} \\ &\leq D_{TV}(\mathbf{u}_t, \mathbf{u}_{t-1}) \ln \frac{d}{\alpha_t} + (\|\mathbf{u}_t\|_1 - D_{TV}(\mathbf{u}_t, \mathbf{u}_{t-1})) \ln \frac{1}{1 - \alpha_t} \\ &\leq D_{TV}(\mathbf{u}_t, \mathbf{u}_{t-1}) \ln \frac{d(1 - \alpha_T)}{\alpha_T} + \|\mathbf{u}_t\|_1 \ln \frac{1}{1 - \alpha_t}, \end{aligned} \quad (16)$$

where we used for the second inequality that the shared update in (11) is such that $1/\widehat{p}_{i,t} \leq d/\alpha_t$ and $v_{i,t}/\widehat{p}_{i,t} \leq 1/(1 - \alpha_t)$, and for the third inequality, that $\alpha_t \geq \alpha_T$ and $x \mapsto (1 - x)/x$ is increasing on $(0, 1]$. Summing (15) over $t = 2, \dots, T$ using (16) and the fact that $\eta_t \geq \eta_T$, we get

$$\begin{aligned} &\sum_{t=2}^T \left(\frac{1}{\eta_{t-1}} \sum_{i=1}^d u_{i,t} \ln \frac{1}{\widehat{p}_{i,t}} - \frac{1}{\eta_t} \sum_{i=1}^d u_{i,t} \ln \frac{1}{v_{i,t+1}} \right) \\ &\leq \frac{m(\mathbf{u}_1^T)}{\eta_T} \ln \frac{d(1 - \alpha_T)}{\alpha_T} + \sum_{t=2}^T \frac{\|\mathbf{u}_t\|_1}{\eta_{t-1}} \ln \frac{1}{1 - \alpha_t} + \underbrace{\sum_{i=1}^d \left(\frac{u_{i,1}}{\eta_1} \ln \frac{1}{v_{i,2}} - \frac{u_{i,T}}{\eta_T} \ln \frac{1}{v_{i,T+1}} \right)}_{\geq 0}. \end{aligned}$$

An application of (14) —including for $t = 1$, for which we recall that $\widehat{p}_{i,1} = 1/d$ and $\eta_1 = \eta_0$ by convention— concludes the proof. \blacksquare

Appendix B. Proof of Proposition 11

Proof By the definition of exp-concavity and by application of Jensen's inequality to the distribution \mathbb{P}_t over $(\Delta_d)^t$ with density

$$\mathbf{r}_1^t \mapsto \frac{1}{\mathbb{E}\left[e^{-\eta_0 L_{t-1}(\mathbf{r}_1^{t-1})}\right]} e^{-\eta_0 L_{t-1}(\mathbf{r}_1^{t-1})} \times 1$$

with respect to the marginal distribution of \mathbb{P} over $(\Delta_d)^t$, we have that

$$\exp(-\eta_0 \ell_t(\hat{\mathbf{p}}_t)) = \exp\left(-\eta_0 \ell_t(\mathbb{E}_t[\mathbf{Q}_t])\right) \geq \mathbb{E}_t\left[\exp(-\eta_0 \ell_t(\mathbf{Q}_t))\right] = \frac{\mathbb{E}\left[e^{-\eta_0 L_t(\mathbf{Q}_1^t)}\right]}{\mathbb{E}\left[e^{-\eta_0 L_{t-1}(\mathbf{Q}_1^{t-1})}\right]}.$$

Thus, a telescoping sum appears,

$$\sum_{t=1}^T \ell_t(\hat{\mathbf{p}}_t) = \sum_{t=1}^T -\frac{1}{\eta_0} \ln e^{-\eta_0 \ell_t(\hat{\mathbf{p}}_t)} \leq -\frac{1}{\eta_0} \ln \mathbb{E}\left[e^{-\eta_0 L_T(\mathbf{Q}_1^T)}\right].$$

It suffices to lower bound the expectation. To do so, we define for all sequences \mathbf{r}_1^k the set of the sequences of k weight vectors that only shift when \mathbf{r}_1^k does and that at each such shift are ε -close to the corresponding values of the \mathbf{r}_t :

$$S_{\varepsilon, \mathbf{r}_1^k} = \left\{ \mathbf{s}_1^k \in \mathcal{X}^k : \forall t \in \{2, \dots, k\}, \mathbf{s}_t \neq \mathbf{s}_{t-1} \Rightarrow \mathbf{r}_t \neq \mathbf{r}_{t-1} \right. \\ \left. \text{and } \forall t \in \{1, \dots, k\}, \mathbf{s}_t = (1 - \varepsilon)\mathbf{r}_t + \varepsilon \mathbf{w}_t \text{ for some } \mathbf{w}_t \in \mathcal{X} \right\}.$$

Note that the second defining constraint is equivalent to the same constraint only at the shifting times of \mathbf{r}_1^k , in view of the first constraint. Since exp-concave loss functions are in particular convex, we get that for all $\mathbf{s}_1^T \in S_{\varepsilon, \mathbf{q}_1^T}$,

$$\sum_{t=1}^T \ell_t(\mathbf{s}_t) \leq (1 - \varepsilon) \sum_{t=1}^T \ell_t(\mathbf{q}_t) + \varepsilon \sum_{t=1}^T \ell_t(\mathbf{w}_t) \leq \sum_{t=1}^T \ell_t(\mathbf{q}_t) + \varepsilon LT.$$

Thus,

$$-\frac{1}{\eta_0} \ln \mathbb{E}\left[e^{-\eta_0 L_T(\mathbf{Q}_1^T)}\right] \leq -\frac{1}{\eta_0} \ln \mathbb{E}\left[e^{-\eta_0 L_T(\mathbf{Q}_1^T)} \mathbb{I}_{\{\mathbf{Q}_1^T \in S_{\varepsilon, \mathbf{q}_1^T}\}}\right] \\ \leq \sum_{t=1}^T \ell_t(\mathbf{q}_t) + \varepsilon LT - \frac{1}{\eta_0} \ln \mathbb{P}(S_{\varepsilon, \mathbf{q}_1^T}).$$

Furthermore, we show by induction on t that for all $t \geq 1$,

$$\mathbb{P}(S_{\varepsilon, \mathbf{q}_1^t}) \geq \varepsilon^{d-1} (1 - \alpha)^{t-s(\mathbf{q}_1^T)-1} \left(\alpha \varepsilon^{d-1}\right)^{s(\mathbf{q}_1^T)}.$$

This is true for $t = 1$ as $S_{\varepsilon, \mathbf{q}_1} = (1 - \varepsilon)\mathbf{q}_1 + \varepsilon\mathcal{X}$ has a \mathbb{P} -probability given by its μ -probability, which is equal to $\varepsilon^{d-1} \mu(\mathcal{X}) = \varepsilon^{d-1}$, and as by convention, $s(\mathbf{q}_1) = 0$. Besides, when $t \geq 2$, we have by definition of \mathbb{P} (cf. its defining transition probability distributions) and $S_{\varepsilon, \mathbf{q}_1^t}$ (cf. the \mathbf{s}_1^t can only shift when the \mathbf{q}_1^t do) that

$$\mathbb{P}(S_{\varepsilon, \mathbf{q}_1^t}) \geq \begin{cases} (1 - \alpha) \mathbb{P}(S_{\varepsilon, \mathbf{q}_1^{t-1}}) & \text{when } \mathbf{q}_t = \mathbf{q}_{t-1} \\ \alpha \mathbb{P}(S_{\varepsilon, \mathbf{q}_1^{t-1}}) \mu(S_{\varepsilon, \mathbf{r}_t}) = \alpha \varepsilon^{d-1} \mathbb{P}(S_{\varepsilon, \mathbf{q}_1^{t-1}}) & \text{when } \mathbf{q}_t \neq \mathbf{q}_{t-1}, \end{cases}$$

which concludes the induction.

Substituting the obtained bound, we have proved so far that

$$\sum_{t=1}^T \ell_t(\hat{\mathbf{p}}_t) - \sum_{t=1}^T \ell_t(\mathbf{q}_t) \leq \varepsilon LT - \frac{1}{\eta_0} \ln \left(\varepsilon^{d-1} (1 - \alpha)^{t-s(\mathbf{q}_1^T)-1} \left(\alpha \varepsilon^{d-1} \right)^{s(\mathbf{q}_1^T)} \right).$$

$\varepsilon \in [0, 1]$ is a parameter of the analysis, it can be optimized to minimize

$$\varepsilon LT + \frac{(s(\mathbf{q}_1^T) + 1)(d - 1)}{\eta_0} \ln \frac{1}{\varepsilon}$$

and get the claimed bound. This is achieved by choosing

$$\varepsilon = \min \left\{ 1, \frac{(s(\mathbf{q}_1^T) + 1)(d - 1)}{\eta_0 LT} \right\}.$$

■