

Random-Walk Perturbations for Online Combinatorial Optimization

Luc Devroye, Gábor Lugosi and Gergely Neu

Abstract—We study online combinatorial optimization problems where a learner is interested in minimizing its cumulative regret in the presence of switching costs. To solve such problems, we propose a version of the follow-the-perturbed-leader algorithm in which the cumulative losses are perturbed by independent symmetric random walks. In the general setting, our forecaster is shown to enjoy near-optimal guarantees on both quantities of interest, making it the best known efficient algorithm for the studied problem. In the special case of prediction with expert advice, we show that the forecaster achieves an expected regret of the optimal order $O(\sqrt{n \log N})$ where n is the time horizon and N is the number of experts, while guaranteeing that the predictions are switched at most $O(\sqrt{n \log N})$ times, in expectation.

Index Terms—Online learning, Online combinatorial optimization, Follow the Perturbed Leader, Random walk

I. PRELIMINARIES

In this paper we study the problem of online prediction with expert advice (see [2]), and in particular, online linear optimization (see, e.g., [3], [4], [5], [6], [7], [8], [9], [10], [11], [12], [13]). The problem may be described as a repeated game between a *forecaster* and an adversary—the *environment*. At each time instant $t = 1, \dots, n$, the forecaster chooses one of the N available actions and suffers a loss corresponding to the chosen action i . Each action i is represented by a vector $v_i \in \mathbb{R}^d$, while the losses assigned by the environment at time t are described by the *loss vector* $\ell_t \in [0, 1]^d$. Thus, given the set of actions $\mathcal{S} = \{v_i : i = 1, 2, \dots, N\} \subseteq \mathbb{R}^d$ at every time instant t , the forecaster chooses, in a possibly randomized way, a vector $V_t \in \mathcal{S}$ and suffers loss $V_t^\top \ell_t$.

We consider the so-called *oblivious adversary* model in which the environment selects all losses before the prediction game starts and reveals the loss vector ℓ_t at time t after the forecaster has made its prediction. The losses are deterministic but the forecaster may randomize: at time t , the forecaster chooses a probability distribution p_t over the set of N actions and draws a random action V_t according to the distribution p_t . The prediction protocol is described in Figure 1.

The usual goal for the standard prediction problem is to devise an algorithm such that the cumulative loss $\widehat{L}_n = \sum_{t=1}^n V_t^\top \ell_t$ is as small as possible, in expectation and/or

L. Devroye is with the School of Computer Science, McGill University, Montreal, Canada H3A 2K6 (email: lucdevroye@gmail.com). G. Lugosi is with ICREA and the Department of Economics, Pompeu Fabra University, Ramon Trias Fargas 25–27, 08005 Barcelona, Spain (email: gabor.lugosi@gmail.com). G. Neu is with the SequeL team, INRIA Lille – Nord Europe, 40 avenue Halley, 59650 Villeneuve d’Ascq, France (email: gergely.neu@gmail.com).

A previous version of this paper was published in the Proceedings of the 26th Annual Conference on Learning Theory, 2013 as [1].

Parameters: set of actions $\mathcal{S} \subseteq \mathbb{R}^d$, number of rounds n ;
 The environment chooses the loss vector $\ell_t \in [0, 1]^d$ for all $t = 1, \dots, n$.
For all $t = 1, 2, \dots, n$, repeat

- 1) The forecaster chooses a probability distribution p_t over \mathcal{S} .
- 2) The forecaster draws an action V_t randomly according to p_t .
- 3) The environment reveals ℓ_t .
- 4) The forecaster suffers loss $V_t^\top \ell_t$.

Fig. 1. Online linear optimization.

with high probability (where probability is with respect to the forecaster’s randomization). Since we do not make any assumption on how the environment generates the losses ℓ_t , we cannot hope to minimize the above loss. Instead, a meaningful goal is to minimize the performance gap between our algorithm and the strategy that selects the best action chosen in hindsight. This performance gap is called the *regret* and is defined formally as

$$R_n = \max_{i \in \{1, 2, \dots, N\}} \sum_{t=1}^n (V_t - v)^{\top} \ell_t = \widehat{L}_n - L_n^*,$$

where we have also introduced the notation $L_n^* = \min_{v \in \mathcal{S}} v^{\top} \sum_{t=1}^n \ell_t$.

To gain simplicity in the presentation, we restrict our attention to the case of *online combinatorial optimization* in which $\mathcal{S} \subset \{0, 1\}^d$, that is, each action is represented as a binary vector. This special case arguably contains most important applications such as the *online shortest path* problem. In this example, a fixed directed acyclic graph of d edges is given with two distinguished vertices u and w . The forecaster, at every time instant t , chooses a directed path from u to w . Such a path is represented by its binary incidence vector $v \in \{0, 1\}^d$. The components of the loss vector $\ell_t \in [0, 1]^d$ represent losses assigned to the d edges and $v^{\top} \ell_t$ is the total loss assigned to the path v . Another (non-essential) simplifying assumption is that every action $v \in \mathcal{S}$ has the same number of 1’s: $\|v\|_1 = m$ for all $v \in \mathcal{S}$. The value of m plays an important role in the bounds presented in the paper.

A fundamental special case of the framework above is *prediction with expert advice*. In this setting, we have $m = 1$, $d = N$, and the learner has access to the unit vectors $\mathcal{S} = \{e_i\}_{i=1}^N$ as the decision set. Minimizing the regret in this setting is a well-studied problem (see the book of Cesa-Bianchi

and Lugosi [2]). It is known that no matter what algorithm the forecaster uses,

$$\liminf_{n, N \rightarrow \infty} \sup \frac{\mathbb{E}R_n}{\sqrt{(n/2) \ln N}} \geq 1,$$

where the supremum is taken with respect to all possible loss assignments with losses in $[0, 1]$. On the other hand, several prediction algorithms are known whose expected regret is of optimal order $O(\sqrt{n \log N})$ and many of them achieve a regret of this order with high probability. Perhaps the most popular one is the exponentially weighted average forecaster (a variant of weighted majority algorithm of Littlestone and Warmuth [14], and aggregating strategies of Vovk [15], also known as *Hedge* by Freund and Schapire [16]). The exponentially weighted average forecaster assigns probabilities to the actions that are inversely proportional to an exponential function of the loss accumulated by each action up to time t .

Another popular forecaster is the *follow the perturbed leader* (FPL) algorithm of Hannan [17]. Kalai and Vempala [7] showed that Hannan’s forecaster, when appropriately modified, indeed achieves an expected regret of optimal order. At time t , the FPL forecaster adds a random perturbation vector $\mathbf{Z}_t \in \mathbb{R}^d$ to the cumulative loss $\mathbf{L}_{t-1} = \sum_{s=1}^{t-1} \ell_s$ of each action and chooses an action \mathbf{v} that minimizes $\mathbf{v}^\top (\mathbf{L}_{t-1} + \mathbf{Z}_t)$. If the elements of the perturbation vector have joint density $(\eta/2)^N e^{-\eta \|\mathbf{z}\|_1}$ for $\eta \sim \sqrt{\log N/n}$, then the expected regret of the forecaster is of order $O(\sqrt{n \log N})$ ([7], see also [2], [18], [19]). This is true whether $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ are independent or not. If they are independent, then one may show that the regret is concentrated around its expectation. Another interesting choice is when $\mathbf{Z}_1 = \dots = \mathbf{Z}_n$, that is, the same perturbation is used over time. Even though this forecaster has an expected regret of optimal order, it may fail with reasonably high probability since its regret is much less concentrated.

While the results presented above still hold in the general case where $m > 1$ when treating each $\mathbf{v} \in \mathcal{S}$ as a separate action, one may gain important computational advantage by taking the structure of the action set into account. In particular, as [7] emphasize, FPL-type forecasters may often be computed efficiently. Interestingly, this efficiency does not come at the price of inferior regret guarantees: as Neu and Bartók [20] have recently shown, an appropriately tuned version of FPL achieves the same regret of $O(m^{3/2} \sqrt{d \log d})$ as the straightforward extension of the exponentially weighted forecaster. The only known forecaster to achieve better performance than this is Component Hedge proposed by Koolen, Warmuth and Kivinen [11], guaranteeing a minimax optimal regret of $O(m \sqrt{n \log(d/m)})$. However, this forecaster can only be implemented efficiently for some special decision sets, and can still take $\Omega(d^6)$ time to run in the worst case (see [21]). In this paper, we propose an FPL-variant that retains the near-optimal regret guarantees of $O(m^{3/2} \sqrt{d \log d})$, while having nice additional properties discussed below.

Small regret is not the only desirable feature of an online forecasting algorithm. In many applications, one would like to define forecasters that do not change their prediction too often. For instance, consider a sequential routing problem on a computer network where predictions correspond to selecting

a path in a graph for each packet to traverse. In this situation, switching between routes might result in out-of-order delivery of packets due to changing delays, and eventually lead to decoding errors. Further examples of such problems include the online buffering problem described by Geulen, Voeking and Winkler [22] and the online lossy source coding problem of György and Neu [23]. A more abstract problem where the number of abrupt switches in the behavior is costly is the problem of online learning in Markovian decision processes, as described by Even-Dar, Kakade and Mansour [24] and Neu, György, Szepesvári, and Antos [25].

To be precise, define the number of action switches up to time n by

$$C_n = |\{1 < t \leq n : \mathbf{V}_{t-1} \neq \mathbf{V}_t\}| .$$

In particular, we are interested in defining randomized forecasters that achieve a regret R_n of near-optimal order while keeping the number of action switches C_n as small as possible. However, the usual forecasters with small regret—such as the exponentially weighted average forecaster or the FPL forecaster with i.i.d. perturbations—may switch actions a large number of times, typically $\Theta(n)$. Therefore, the design of special forecasters with small regret and small number of action switches is called for.

The first known algorithm to address this issue is the Follow the Lazy Leader (FLL) algorithm proposed by Kalai and Vempala [7]. This algorithm is designed to behave identically to their FPL algorithm *in expectation* (with an $O(m^{3/2} \sqrt{n \log d})$ bound on the regret, as shown by Neu and Bartók [20]), while guaranteeing that the expected number of action switches is $O(d \sqrt{(n/m) \log d})$. The “Shrinking Dartboard” algorithm proposed by Geulen, Voeking and Winkler [22] is based on a similar idea: this algorithm simulates the exponentially weighted forecaster in expectation, guaranteeing a regret of $O(m^{3/2} \sqrt{n \log d})$, while improving the upper bound on the expected number of switches to $O(\sqrt{mn \log d})$. In this paper, we propose a family of methods based on FPL in which perturbations are defined by independent symmetric random walks. We show that these intuitively appealing forecasters have similar regret and switch-number guarantees as Shrinking Dartboard and FLL.

In particular, we first propose an FPL-variant in which perturbations are generated by independent Gaussian random walks for each coordinate of the perturbation vector. We show that this algorithm guarantees a regret of $O(m^{3/2} \sqrt{n \log d})$, while keeping the number of switches bounded by $O(m \sqrt{n \log d})$. While this bound is inferior to that of the Shrinking Dartboard algorithm by a factor of $\sqrt{m \log d}$, that algorithm can only be efficiently implemented for some special decision sets \mathcal{S} —see [11] and [12] for some examples. On the other hand, our algorithm can be efficiently implemented whenever there exists an efficient implementation of the static optimization problem of finding $\arg \min_{\mathbf{v} \in \mathcal{S}} \mathbf{v}^\top \ell$ for any $\ell \in \mathbb{R}^d$. We compare our results to other results known in the literature in Table I.

Notice that our regret bound described above only guarantees that the number of switches is of $O(\sqrt{n \log N})$ in the setting of prediction with expert advice. In the second half of

	Regret	Switches	Efficient
Shrinking Dartboard [22]	$O(m^{3/2}\sqrt{n \log d})$	$O(\sqrt{mn \log d})$	sometimes
Follow the Lazy Leader [7]	$O(m^{3/2}\sqrt{n \log d})$	$O(d\sqrt{(n/m) \log d})$	always
Prediction by random-walk perturbations	$O(m^{3/2}\sqrt{n \log d})$	$O(m\sqrt{n \log d})$	always

TABLE I
THE RESULTS PRESENTED IN THIS PAPER VERSUS THE RESULTS OF [7] AND [22].

the paper, we show that this can be improved to $O(\sqrt{n \log N})$ by using symmetric *binary* random walks as perturbations. An interesting property of the resulting algorithm is that the expected regret can be directly upper bounded in terms of the expected number of switches.

We also note that a similar variant of the FPL forecaster was recently derived by Rakhlin, Shamir and Sridharan [26], who use perturbations of the form $Z_{i,t} = \sum_{s=t+1}^n X_{i,s}$, where $X_{i,s}$ are i.i.d. random variables with an arbitrary symmetric distribution. Rakhlin, Shamir and Sridharan exploit the fact that these perturbations can serve as a relaxation of the Rademacher complexity of the prediction game, and prove an $O(\sqrt{n \log N})$ bound on the expected regret of the resulting algorithm. While this approach cannot be used for analyzing our FPL-variant, our analysis can be directly applied to provide both regret and switch-number guarantees for their method. Additionally, note that our algorithm does not need to use prior knowledge of the number of rounds.

An interesting open question left for future research is characterizing the minimax rate for the quantity $R_n + C_n$ in terms of m . For all known algorithms (including ours), the upper bounds on this quantity depend on m as $m^{3/2}$, as contributed by the bound on R_n . However, this dependence is known to be suboptimal: the optimal linear dependence of the regret on m is achieved by the Component Hedge algorithm of [11]. We conjecture that the correct minimax rate of $R_n + C_n$ matches the minimax rate of R_n , which is of $\Theta(m\sqrt{n \log(d/m)})$. Such a bound could be achieved by a switch-constrained algorithm that implements Component Hedge on expectation in a similar way as the Shrinking Dartboard implements the exponentially weighted forecaster or FLL implements FPL. However, constructing such an implementation is far from trivial and is left here as an interesting open problem.

II. RANDOM-WALK PERTURBATIONS FOR ONLINE COMBINATORIAL OPTIMIZATION

To address the problem described in the previous section, we propose a variant of the Follow the Perturbed Leader (FPL) algorithm. The proposed forecaster perturbs the loss of each action at every time instant by a zero-mean Gaussian random variable with variance $\eta^2 > 0$ and chooses an action with minimal cumulative perturbed loss. More precisely, the algorithm draws independent random variables $X_{i,t} \sim \mathcal{N}(0, \eta^2)$ and the vector $\mathbf{X}_t = (X_{1,t}, \dots, X_{d,t})$ is added to the observed loss vector ℓ_{t-1} . At time t action $v \in \mathcal{S}$ is chosen that minimizes $\sum_{s=1}^t \mathbf{v}^\top (\ell_{s-1} + \mathbf{X}_s)$ (where we define ℓ_0 as the all-zero vector $\mathbf{0}$). Equivalently, the forecaster may be thought of as an FPL algorithm in which the cumulative losses \mathbf{L}_{t-1} are perturbed by the independent symmetric random walks

Algorithm 1 Online combinatorial optimization by random-walk perturbations.

Initialization: set $\mathbf{L}_0 = \mathbf{0}$ and $Z_0 = \mathbf{0}$.

For all $t = 1, 2, \dots, n$, **repeat**

- 1) Draw \mathbf{X}_t with i.i.d. Gaussian components

$$X_{i,t} \sim \mathcal{N}(0, \eta^2)$$

- 2) Let $\mathbf{Z}_t = \mathbf{Z}_{t-1} + \mathbf{X}_t$.
- 3) Choose action

$$\mathbf{V}_t = \arg \min_{v \in \mathcal{S}} \{v^\top (\mathbf{L}_{t-1} + \mathbf{Z}_t)\},$$

where ties are broken in favor of \mathbf{V}_{t-1} .

- 4) Observe the loss vector ℓ_t , suffer loss $\mathbf{V}_t^\top \ell_t$.
 - 5) Set $\mathbf{L}_t = \mathbf{L}_{t-1} + \ell_t$.
-

$\mathbf{Z}_t = \sum_{s=1}^t \mathbf{X}_s$. This is the way the algorithm is presented in Algorithm 1.

Conceptually, the difference between standard FPL and the proposed version is the way the perturbations are generated: while common versions of FPL use perturbations that are generated in an i.i.d. fashion, the perturbations of the algorithm proposed here are dependent. This will enable us to control the number of action switches during the learning process. Note that the standard deviation of these perturbations at time t is still of order \sqrt{t} just like for the standard FPL forecaster with optimal parameter settings.

To obtain intuition why this approach will solve our problem, first consider a problem with action set $\mathcal{S} = \{(1, 0)^\top, (0, 1)^\top\}$ and an environment that generates equal losses, say $\ell_{i,t} = 0$ for all i and t . When using i.i.d. perturbations, FPL switches actions with probability $1/2$ in each round, thus yielding $C_t = t/2 + O(\sqrt{t})$ with overwhelming probability. The same holds for the exponentially weighted average forecaster. On the other hand, when using the random-walk perturbations described above, we only switch between the actions when the leading random walk is changed, that is, when the difference of the two random walks—which is also a symmetric random walk—hits zero. This distribution is well understood and the probability that this occurs more than $x\sqrt{n}$ times during the first n steps is roughly $2\mathbb{P}\{N > 2x\} \leq 2e^{-2x^2}$ where N is a standard normal random variable (see [27, Section III.4]). Thus, in this case we see that the number of switches is bounded by $O\left(\sqrt{n \log(1/\delta)}\right)$, with probability at least $1 - \delta$. As we show below, assuming all-zero losses is the worst case for the number of switches.

Even though we only prove bounds for the expected regret and the expected number of switches, the above example gives

some intuition about the upper tail probabilities. While the above idea can be extended to the case of non-zero loss sequences to obtain high-confidence switch-number guarantees, proving similar results for the general setting is a highly nontrivial problem. We note that by our Lemma 2 (stated later in Section III), the regret of our algorithm can be directly bounded in terms of the number of switches, thus we can guarantee upper bounds of $O(\sqrt{n})$ on both C_n and R_n with high probability. We are not aware of any other algorithm that provides high-confidence guarantees on both quantities of interest even in this simple special case.

The next theorem bounds the performance of the proposed forecaster. We are not only interested in the regret but also the number of switches $C_n = \sum_{t=1}^n \mathbb{1}_{\{\mathbf{V}_{t+1} \neq \mathbf{V}_t\}}$. The regret is of similar order as that of the standard FPL forecaster, up to an additive logarithmic factor. Moreover, the expected number of switches is $O(m\sqrt{n} \log d)$. Remarkably, the dependence on d is only logarithmic and it is the weight m of the actions that plays an important role.

Theorem 1: The expected regret and the expected number of action switches satisfy (under the oblivious adversary model),

$$\mathbb{E}R_n \leq m\sqrt{2n \log d} \left(\frac{2m}{\eta} + \eta \right) + \frac{m^2(\log n + 1)}{\eta^2}$$

and

$$\begin{aligned} \mathbb{E}C_n &\leq \sum_{t=1}^n \frac{m(1 + \eta\sqrt{2 \log d})\sqrt{2 \log d}}{\eta\sqrt{t}} \\ &+ \sum_{t=1}^n \frac{m(1 + 2\eta\sqrt{2 \log d} + \eta^2(2 \log d + \sqrt{2 \log d} + 1))}{4\eta^2 t}. \end{aligned}$$

In particular, setting $\eta = \sqrt{2m}$ yields

$$\mathbb{E}R_n \leq 4m^{3/2}\sqrt{n \log d} + \frac{m(\log n + 1)}{2}.$$

and

$$\mathbb{E}C_n = O(m\sqrt{n} \log d).$$

The proof of the regret bound follows the steps of the proof of Theorem 1 in [20], and is deferred to the appendix. The more interesting part is the bound for the expected number of action switches $\mathbb{E}C_n = \sum_{t=1}^n \mathbb{P}[\mathbf{V}_{t+1} \neq \mathbf{V}_t]$. For proving a bound on this quantity, we study the evolution of the *lead pack* A_t defined as the set of actions with near-optimal loss given by

$$\begin{aligned} A_t &= \{ \mathbf{w} \in \mathcal{S} : (\mathbf{w} - \mathbf{V}_t)^\top (\mathbf{L}_{t-1} + \mathbf{Z}_t) \\ &\leq \|\mathbf{w} - \mathbf{V}_t\|_1 \cdot \|\boldsymbol{\ell}_t + \mathbf{X}_{t+1}\|_\infty \}. \end{aligned} \quad (1)$$

We sometimes refer to $\|\boldsymbol{\ell}_t + \mathbf{X}_{t+1}\|_\infty$ as the *diameter* of the lead pack. Observe that no action outside A_t can take the lead at time $t+1$, since if $\mathbf{w} \notin A_t$, then

$$(\mathbf{w} - \mathbf{V}_t)^\top (\mathbf{L}_{t-1} + \mathbf{Z}_t) > |(\mathbf{w} - \mathbf{V}_t)^\top (\boldsymbol{\ell}_t + \mathbf{X}_{t+1})|$$

so $\mathbf{w}^\top (\mathbf{L}_t + \mathbf{Z}_{t+1}) > \mathbf{V}_t^\top (\mathbf{L}_t + \mathbf{Z}_{t+1})$ and \mathbf{w} cannot be the new leader. It follows that we can upper bound the probability of switching as

$$\mathbb{P}[\mathbf{V}_{t+1} \neq \mathbf{V}_t] \leq \mathbb{P}[|A_t| > 1],$$

which leaves us with the problem of upper bounding $\mathbb{P}[|A_t| > 1]$. The following lemma gives a bound of this quantity. Putting this statement together with the well-known facts that $\mathbb{E}[\|\mathbf{X}_1\|_\infty] \leq \eta\sqrt{2 \log d}$ and $\mathbb{E}[\|\mathbf{X}_1\|_\infty^2] \leq \eta^2(2 \log d + \sqrt{2 \log d} + 1)$ (see, e.g., [28]) proves the second statement of the theorem.

Lemma 1: For each $t = 1, 2, \dots, n$,

$$\begin{aligned} \mathbb{P}[|A_t| > 1 | \mathbf{X}_{t+1}] &\leq \frac{m \|\boldsymbol{\ell}_t + \mathbf{X}_{t+1}\|_\infty \sqrt{2 \log d}}{\eta\sqrt{t}} \\ &+ \frac{m \|\boldsymbol{\ell}_t + \mathbf{X}_{t+1}\|_\infty^2}{2\eta^2 t}. \end{aligned}$$

Proof: We use the notation $\mathbb{P}_t[\cdot] = \mathbb{P}[\cdot | \mathbf{X}_{t+1}]$ and $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot | \mathbf{X}_{t+1}]$. Also, let

$$\mathbf{h}_t = \boldsymbol{\ell}_t + \mathbf{X}_{t+1} \quad \text{and} \quad \mathbf{H}_t = \sum_{s=0}^{t-1} \mathbf{h}_s.$$

Furthermore, we use the shorthand notation $c = \|\mathbf{h}_t\|_\infty$.

We start by analyzing $\mathbb{P}_t[|A_t| = 1]$:

$$\begin{aligned} \mathbb{P}_t[|A_t| = 1] &= \\ &= \sum_{\mathbf{v} \in \mathcal{S}} \mathbb{P}_t[\forall \mathbf{w} \neq \mathbf{v} : (\mathbf{w} - \mathbf{v})^\top \mathbf{H}_t > \|\mathbf{w} - \mathbf{v}\|_1 c] \\ &= \sum_{\mathbf{v} \in \mathcal{S}} \int_{y \in \mathbb{R}} f_{\mathbf{v}}(y) \mathbb{P}_t[\forall \mathbf{w} \neq \mathbf{v} : \\ &\quad \mathbf{w}^\top \mathbf{H}_t > y + \|\mathbf{w} - \mathbf{v}\|_1 c | \mathbf{v}^\top \mathbf{H}_t = y] dy, \end{aligned} \quad (2)$$

where $f_{\mathbf{v}}$ is the density of $\mathbf{v}^\top \mathbf{H}_t$.

Next, we crucially use the fact that the conditional distributions of correlated Gaussian random variables are also Gaussian. In particular, let us fix an arbitrary $\mathbf{v} \in \mathcal{S}$ and define $k_{\mathbf{w}} = (m - \|\mathbf{w} - \mathbf{v}\|_1)$ for all $\mathbf{w} \in \mathcal{S}$. Then, the covariance between the perturbed loss of \mathbf{w} and \mathbf{v} is given as

$$\text{cov}(\mathbf{w}^\top \mathbf{H}_t, \mathbf{v}^\top \mathbf{H}_t) = \eta^2 (m - \|\mathbf{w} - \mathbf{v}\|_1) t = \eta^2 k_{\mathbf{w}} t.$$

Organizing all actions $\mathbf{w} \in \mathcal{S} \setminus \mathbf{v}$ into a matrix $\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{N-1})$, the conditional distribution of $\mathbf{W}^\top \mathbf{H}_t$ is an $(N-1)$ -variate Gaussian distribution with mean vector

$$\mu_{\mathbf{v}}(y) = \left(\mathbf{w}_1^\top \mathbf{L}_{t-1} + y \frac{k_{\mathbf{w}_1}}{m}, \dots, \mathbf{w}_{N-1}^\top \mathbf{L}_{t-1} + y \frac{k_{\mathbf{w}_{N-1}}}{m} \right)^\top$$

and covariance matrix $\Sigma_{\mathbf{v}}$, given that $\mathbf{v}^\top \mathbf{H}_t = y$. Defining $\mathbf{K} = (k_{\mathbf{w}_1}, \dots, k_{\mathbf{w}_{N-1}})^\top$ and using the notation

$$\varphi_{\mathbf{v}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^{N-1} |\Sigma_{\mathbf{v}}|}} \exp\left(-\frac{\mathbf{x}^\top \Sigma_{\mathbf{v}}^{-1} \mathbf{x}}{2}\right),$$

we get that

$$\begin{aligned}
& \mathbb{P}_t [\forall \mathbf{w} \neq \mathbf{v} : \mathbf{w}^\top \mathbf{H}_t > y + \|\mathbf{w} - \mathbf{v}\|_1 c | \mathbf{v}^\top \mathbf{H}_t = y] \\
&= \int_{z_i=y+(m-k_{w_i})c}^{\infty} \cdots \int \varphi_{\mathbf{v}}(\mathbf{z} - \mu_{\mathbf{v}}(y)) d\mathbf{z} \\
&= \int_{z_i=y+mc}^{\infty} \cdots \int \varphi_{\mathbf{v}}(\mathbf{z} - \mu_{\mathbf{v}}(y) - c\mathbf{K}) d\mathbf{z} \\
&= \int_{z_i=y+mc}^{\infty} \cdots \int \varphi_{\mathbf{v}}(\mathbf{z} - \mu_{\mathbf{v}}(y + mc)) d\mathbf{z} \\
&= \mathbb{P}_t [\forall \mathbf{w} \neq \mathbf{v} : \mathbf{w}^\top \mathbf{H}_t > y + mc | \mathbf{v}^\top \mathbf{H}_t = y + mc],
\end{aligned}$$

where we used $\mu_{\mathbf{v}}(y + mc) = \mu_{\mathbf{v}}(y) + c\mathbf{K}$. Using this, we rewrite (2) as

$$\begin{aligned}
& \mathbb{P}_t [|A_t| = 1] \\
&= \sum_{\mathbf{v} \in \mathcal{S}} \int_{y \in \mathbb{R}} f_{\mathbf{v}}(y) \mathbb{P}_t [\forall \mathbf{w} \neq \mathbf{v} : \mathbf{w}^\top \mathbf{H}_t > y | \mathbf{v}^\top \mathbf{H}_t = y] dy \\
&\quad - \sum_{\mathbf{v} \in \mathcal{S}} \int_{y \in \mathbb{R}} (f_{\mathbf{v}}(y) - f_{\mathbf{v}}(y - mc)) \cdot \\
&\quad \cdot \mathbb{P}_t [\forall \mathbf{w} \neq \mathbf{v} : \mathbf{w}^\top \mathbf{H}_t > y | \mathbf{v}^\top \mathbf{H}_t = y] dy
\end{aligned}$$

Observing that the first term above corresponds to the probability that the leader is unique and that this event holds almost surely under our perturbation scheme, we get that the last term equals $\mathbb{P}_t [|A_t| > 1]$. To treat this term, we use that $\mathbf{v}^\top \mathbf{H}_t$ is Gaussian with mean $\mathbf{v}^\top \mathbf{L}_{t-1}$ and variance $\eta^2 m t$ to obtain

$$\begin{aligned}
f_{\mathbf{v}}(y) - f_{\mathbf{v}}(y - mc) &= f_{\mathbf{v}}(y) \left(1 - \frac{f_{\mathbf{v}}(y - mc)}{f_{\mathbf{v}}(y)} \right) \\
&\leq f_{\mathbf{v}}(y) \left(\frac{mc^2}{2\eta^2 t} - \frac{c(y - \mathbf{v}^\top \mathbf{L}_{t-1})}{\eta^2 t} \right).
\end{aligned}$$

Also observe that the density of $\mathbf{V}_t^\top \mathbf{H}_t$ at $y \in \mathbb{R}$ is exactly

$$\sum_{\mathbf{v} \in \mathcal{S}} f_{\mathbf{v}}(y) \mathbb{P}_t [\forall \mathbf{w} \neq \mathbf{v} : \mathbf{w}^\top \mathbf{H}_t > y | \mathbf{v}^\top \mathbf{H}_t = y].$$

Thus,

$$\begin{aligned}
\mathbb{P}_t [|A_t| > 1] &\leq \sum_{\mathbf{v} \in \mathcal{S}} \int_{y \in \mathbb{R}} (f_{\mathbf{v}}(y) - f_{\mathbf{v}}(y - mc)) \cdot \\
&\quad \cdot \mathbb{P}_t [\forall \mathbf{w} \neq \mathbf{v} : \mathbf{w}^\top \mathbf{H}_t > y | \mathbf{v}^\top \mathbf{H}_t = y] dy \\
&\leq \sum_{\mathbf{v} \in \mathcal{S}} \int_{y \in \mathbb{R}} f_{\mathbf{v}}(y) \mathbb{P}_t [\forall \mathbf{w} \neq \mathbf{v} : \mathbf{w}^\top \mathbf{H}_t > y | \mathbf{v}^\top \mathbf{H}_t = y] \cdot \\
&\quad \cdot \left(\frac{mc^2}{2\eta^2 t} - \frac{c(y - \mathbf{v}^\top \mathbf{L}_{t-1})}{\eta^2 t} \right) dy \\
&= \frac{mc^2}{2\eta^2 t} - \frac{c\mathbb{E}[\mathbf{V}_t^\top \mathbf{Z}_t]}{\eta^2 t} \leq \frac{mc^2}{2\eta^2 t} + \frac{mc\mathbb{E}[\|\mathbf{Z}_t\|_{\infty}]}{\eta^2 t}.
\end{aligned}$$

Using the definition of c and $\mathbb{E}[\|\mathbf{Z}_t\|_{\infty}] \leq \eta\sqrt{2t \log d}$ gives the result. ■

III. RANDOM-WALK PERTURBATIONS FOR PREDICTION WITH EXPERT ADVICE

In this section we refine our method to obtain bounds of optimal order in the special case of prediction with expert advice, that is, when the set of actions is the set of $d = N$ unit vectors. In this case, we use I_t to denote the unique coordinate i of \mathbf{V}_t with $V_{i,t} = 1$. While the straightforward application of Algorithm 1 and Theorem 1 guarantees a regret of optimal order in this case, the switch-number guarantees are of a suboptimal $O(\sqrt{n} \log N)$. In this section, we propose a variant of our algorithm that achieves order-optimal regret guarantees while switching its predictions only $O(\sqrt{n} \log N)$ times, similarly to the algorithms of [7] and [22].

The algorithm—presented as Algorithm 2—is obtained by replacing the Gaussian increments in Algorithm 1 by independent random variables that take values $\pm 1/2$ with equal probabilities. The benefit of using this perturbation scheme is that the diameter of the lead packs defined in (1) can be upper bounded by a constant, and thus we can eliminate higher moments of $\|\mathbf{X}_{t+1}\|_{\infty}$ in the upper bound on $\mathbb{P}[|A_t| > 1]$. To gain further intuition, notice that choosing any *fixed* (i.e., non-random) lead-pack diameter in the proof of Lemma 1 would allow experts from outside the lead pack to take the lead with positive probability. While this probability can be decreased at the expense of slightly expanding the diameter, its rate of decay is not sufficiently fast for improving our previously presented results. In fact, balancing the two terms constituting the probability of switching gives the exact same result. On the other hand, when using $\pm 1/2$ -valued random increments, it is possible to set a fixed diameter of 1 that ensures that the new leader comes from the lead pack with probability 1, and thus the extra term vanishes.

Algorithm 2 Random-walk perturbations for prediction with expert advice.

Initialization: set $L_{i,0} = 0$ and $Z_{i,0} = 0$ for all $i = 1, 2, \dots, N$.

For all $t = 1, 2, \dots, n$, **repeat**

- 1) Draw $X_{i,t}$ for all $i = 1, 2, \dots, N$ such that

$$X_{i,t} = \begin{cases} \frac{1}{2} & \text{with probability } \frac{1}{2} \\ -\frac{1}{2} & \text{with probability } \frac{1}{2}. \end{cases}$$

- 2) Let $Z_{i,t} = Z_{i,t-1} + X_{i,t}$ for all $i = 1, 2, \dots, N$.
- 3) Choose action

$$I_t = \arg \min_i (L_{i,t-1} + Z_{i,t}),$$

where ties are broken in favor of I_{t-1} .

- 4) Observe losses $\ell_{i,t}$ for all $i = 1, 2, \dots, N$, suffer loss $\ell_{I_t,t}$.
 - 5) Set $L_{i,t} = L_{i,t-1} + \ell_{i,t}$ for all $i = 1, 2, \dots, N$.
-

The next theorem summarizes our performance bounds for the proposed forecaster.

Theorem 2: The expected regret and expected number of switches of actions of the forecaster of Algorithm 2 satisfy, for all possible loss sequences (under the oblivious-adversary

model),

$$\mathbb{E}R_n \leq 2\mathbb{E}C_n \leq 8\sqrt{2n \log N} + 16 \log n + 16 .$$

While it is possible to perform the analysis of Algorithm 2 similarly to that of Algorithm 1, we take a different path: The proof we present below is based on the observation that the regret of Algorithm 2 can be bounded in terms of the number of action switches. The next simple lemma formalizes this statement.

Lemma 2: Fix any $i \in \{1, 2, \dots, N\}$. Then

$$\widehat{L}_n - L_{i,n} \leq 2C_n + Z_{i,n+1} - \sum_{t=1}^{n+1} X_{I_{t-1},t} .$$

Proof: We apply Lemma 3.1 of [2] (sometimes referred to as the “be-the-leader” lemma) for the sequence $(\ell_{\cdot,t-1} + X_{\cdot,t})_{t=1}^{\infty}$ with $\ell_{j,0} = 0$ for all $j \in \{1, 2, \dots, N\}$, obtaining

$$\begin{aligned} \sum_{t=1}^{n+1} (\ell_{I_t,t-1} + X_{I_t,t}) &\leq \sum_{t=1}^{n+1} (\ell_{i,t-1} + X_{i,t}) \\ &= L_{i,n} + Z_{i,n+1} . \end{aligned}$$

Reordering terms, we get

$$\sum_{t=1}^n \ell_{I_t,t} \leq L_{i,n} + \sum_{t=1}^{n+1} (\ell_{I_{t-1},t-1} - \ell_{I_t,t-1}) + Z_{i,n} - \sum_{t=1}^{n+1} X_{I_t,t} . \quad (3)$$

The last term can be rewritten as

$$- \sum_{t=1}^{n+1} X_{I_t,t} = - \sum_{t=1}^{n+1} X_{I_{t-1},t} + \sum_{t=1}^{n+1} (X_{I_{t-1},t} - X_{I_t,t}) .$$

Now notice that $X_{I_{t-1},t} - X_{I_t,t}$ and $\ell_{I_{t-1},t-1} - \ell_{I_t,t-1}$ are both zero when $I_t = I_{t-1}$ and are upper bounded by 1 otherwise. That is, we get that

$$\begin{aligned} \sum_{t=1}^{n+1} (\ell_{I_{t-1},t-1} - \ell_{I_t,t-1}) + \sum_{t=1}^{n+1} (X_{I_{t-1},t} - X_{I_t,t}) \\ \leq 2 \sum_{t=1}^{n+1} \mathbb{1}_{\{I_{t-1} \neq I_t\}} = 2C_n . \end{aligned}$$

Putting everything together gives the statement of the lemma. \blacksquare

Next we analyze the number of switches C_n . Similarly to the analysis of Algorithm 1, we study the *lead pack* A_t defined as

$$A_t = \{i \in \{1, \dots, N\} : L_{i,t-1} + Z_{i,t} < L_{I_t,t-1} + Z_{I_t,t} + 2\} .$$

Once again, observe that no action from outside the lead pack has a positive probability of taking the lead at time $t+1$. We bound the probability of lead change as

$$\mathbb{P}[I_t \neq I_{t+1}] \leq \frac{1}{2} \mathbb{P}[|A_t| > 1] .$$

The key to the proof of the theorem is the following lemma that gives an upper bound for the probability that the lead pack contains more than one action. It implies, in particular, that

$$\mathbb{E}[C_n] \leq 4\sqrt{2n \log N} + 4 \log n + 4 ,$$

which is what we need to prove the bounds of Theorem 2.

Lemma 3:

$$\mathbb{P}[|A_t| > 1] \leq 4\sqrt{2\frac{\log N}{t}} + \frac{8}{t} .$$

Proof: Define $p_t(k) = \mathbb{P}[Z_{i,t} = \frac{k}{2}]$ for all $k = -t, \dots, t$, and we let S_t denote the set of leaders at time t :

$$S_t = \left\{ j \in \{1, \dots, N\} : L_{j,t-1} + Z_{j,t} = \min_i \{L_{i,t-1} + Z_{i,t}\} \right\} .$$

The forecaster picks $I_t \in S_t$ arbitrarily when $I_{t-1} \notin S_t$, otherwise it stays with $I_t = I_{t-1}$. Let us start with analyzing the probability of the event $\{A_t = \{I_t\}\} = \{|A_t| = 1\}$:

$$\begin{aligned} \mathbb{P}[|A_t| = 1] &= \sum_{k=-t}^t \sum_{j=1}^N p_t(k) \mathbb{P} \left[\min_{i \neq j} \{L_{i,t-1} + Z_{i,t}\} \geq L_{j,t-1} + \frac{k}{2} + 2 \right] \\ &\geq \sum_{k=-t+4}^t \sum_{j=1}^N p_t(k-4) \mathbb{P} \left[\min_{i \neq j} \{L_{i,t-1} + Z_{i,t}\} \geq L_{j,t-1} + \frac{k}{2} \right] \\ &= \sum_{k=-t+4}^t \sum_{j=1}^N p_t(k) \mathbb{P} \left[\min_{i \neq j} \{L_{i,t-1} + Z_{i,t}\} \geq L_{j,t-1} + \frac{k}{2} \right] \frac{p_t(k-4)}{p_t(k)} . \end{aligned}$$

Before proceeding, we need to make two observations. First of all,

$$\begin{aligned} \sum_{j=1}^N p_t(k) \mathbb{P} \left[\min_{i \in \{1, 2, \dots, N\} \setminus j} \{L_{i,t-1} + Z_{i,t}\} \geq L_{j,t-1} + \frac{k}{2} \right] \\ \geq \mathbb{P} \left[\exists j \in S_t : Z_{j,t} = \frac{k}{2} \right] \geq \mathbb{P} \left[\min_{j \in S_t} Z_{j,t} = \frac{k}{2} \right] , \end{aligned}$$

where the first inequality follows from the union bound and the second from the fact that the latter event implies the former. Also notice that $\frac{t+Z_{1,t}}{2}$ is binomially distributed with parameters t and $1/2$ and therefore $p_t(k) = \binom{t}{\frac{t+k}{2}} \frac{1}{2^t}$. Hence,

$$\begin{aligned} \frac{p_t(k-4)}{p_t(k)} &= \frac{\binom{t+k}{2}! \binom{t-k}{2}!}{\binom{t+k}{2}! \binom{t-k}{2}!} \\ &= 1 + \frac{4(t+1)(k-2)}{(t-k+2)(t-k+4)} . \end{aligned}$$

It can be easily verified that

$$\frac{4(t+1)(k-2)}{(t-k+2)(t-k+4)} \geq \frac{4(t+1)(k-2)}{(t+2)(t+4)}$$

holds for all $k \in [-t, t]$. Using our first observation and the bound on $\mathbb{P}[|A_t| = 1]$, we get

$$\mathbb{P}[|A_t| = 1] \geq \sum_{k=-t+4}^t \mathbb{P} \left[\min_{j \in S_t} Z_{j,t} = \frac{k}{2} \right] \frac{p_t(k-4)}{p_t(k)} .$$

Along with our second observation, this implies

$$\begin{aligned}
\mathbb{P}[|A_t| > 1] &\leq 1 - \sum_{k=-t+4}^t \mathbb{P}\left[\min_{j \in S_t} Z_{j,t} = \frac{k}{2}\right] \frac{p_t(k-4)}{p_t(k)} \\
&\leq 1 - \sum_{k=-t+4}^t \mathbb{P}\left[\min_{j \in S_t} Z_{j,t} = \frac{k}{2}\right] \left(1 + \frac{4(t+1)(k-2)}{(t+2)(t+4)}\right) \\
&\leq \sum_{k=-t}^t \mathbb{P}\left[\min_{j \in S_t} Z_{j,t} = \frac{k}{2}\right] \left(\frac{4(2-k)(t+1)}{(t+2)(t+4)}\right) \\
&= \frac{8(t+1)}{(t+2)(t+4)} - 8 \frac{t+1}{(t+2)(t+4)} \mathbb{E}\left[\min_{j \in S_t} Z_{j,t}\right] \\
&\leq \frac{8}{t} + \frac{8}{t} \mathbb{E}\left[\max_{j \in \{1,2,\dots,N\}} Z_{j,t}\right].
\end{aligned}$$

Now using $\mathbb{E}[\max_j Z_{j,t}] \leq \sqrt{\frac{t \log N}{2}}$ implies

$$\mathbb{P}[|A_t| > 1] \leq 4\sqrt{\frac{2 \log N}{t}} + \frac{8}{t}$$

as desired. \blacksquare

APPENDIX

Proof of the first statement of Theorem 1: The proof is based on the proof of Theorem 4.2 of [2] and Theorem 3 of [13], with key insights taken from Neu and Bartók [20]. The main difference from those proofs is that the standard deviation of our perturbations changes over time, however, this issue is easy to treat. First, we define an infeasible “forecaster” that peeks one step into the future and uses perturbation $\widehat{\mathbf{Z}}_t = \sqrt{t}\mathbf{X}_1$:

$$\widehat{\mathbf{V}}_t = \arg \min_{\mathbf{w} \in \mathcal{S}} \mathbf{w}^\top (\mathbf{L}_t + \widehat{\mathbf{Z}}_t).$$

Now fix any $\mathbf{v} \in \mathcal{S}$. Applying Lemma 3.1 of [2] for the sequence $(\ell_{t-1} + \widehat{\mathbf{Z}}_t - \widehat{\mathbf{Z}}_{t-1})_{t=1}^\infty$ with $\ell_0 = \mathbf{0}$, we get

$$\sum_{t=1}^n \widehat{\mathbf{V}}_t^\top (\ell_t + (\widehat{\mathbf{Z}}_t - \widehat{\mathbf{Z}}_{t-1})) \leq \mathbf{v}^\top (\mathbf{L}_n + \widehat{\mathbf{Z}}_n).$$

After reordering, we obtain

$$\begin{aligned}
\sum_{t=1}^n \mathbf{V}_t^\top \ell_t &\leq \mathbf{v}^\top \mathbf{L}_n + \mathbf{v}^\top \widehat{\mathbf{Z}}_n + \sum_{t=1}^n (\mathbf{V}_t - \widehat{\mathbf{V}}_t)^\top \ell_t \\
&\quad - \sum_{t=1}^n \widehat{\mathbf{V}}_t^\top (\widehat{\mathbf{Z}}_t - \widehat{\mathbf{Z}}_{t-1}) \\
&= \mathbf{v}^\top \mathbf{L}_n + \mathbf{v}^\top \widehat{\mathbf{Z}}_n + \sum_{t=1}^n (\mathbf{V}_t - \widehat{\mathbf{V}}_t)^\top \ell_t \\
&\quad + \sum_{t=1}^n (\sqrt{t-1} - \sqrt{t}) \widehat{\mathbf{V}}_t^\top \mathbf{X}_1
\end{aligned}$$

The last term can be bounded as

$$\begin{aligned}
\sum_{t=1}^n (\sqrt{t-1} - \sqrt{t}) \widehat{\mathbf{V}}_t^\top \mathbf{X}_1 &\leq \sum_{t=1}^n (\sqrt{t} - \sqrt{t-1}) \left| \widehat{\mathbf{V}}_t^\top \mathbf{X}_1 \right| \\
&\leq m \sum_{t=1}^n (\sqrt{t} - \sqrt{t-1}) \|\mathbf{X}_1\|_\infty \\
&\leq m\sqrt{n} \|\mathbf{X}_1\|_\infty.
\end{aligned}$$

Taking expectations, we obtain the bound

$$\mathbb{E} \widehat{\mathbf{L}}_n - \mathbf{v}^\top \mathbf{L}_n \leq \sum_{t=1}^n \mathbb{E} \left[(\mathbf{V}_t - \widehat{\mathbf{V}}_t)^\top \ell_t \right] + \eta m \sqrt{2n \log d},$$

where we used $\mathbb{E}[\|\mathbf{X}_1\|_\infty] \leq \eta \sqrt{2 \log d}$.

Thus, we are left with the problem of bounding $\mathbb{E}[(\mathbf{V}_t - \widehat{\mathbf{V}}_t)^\top \ell_t]$ for each $t \geq 1$. Similarly to [20], we do this by introducing

$$p_t(\mathbf{u}) = \mathbb{P}[\mathbf{V}_t = \mathbf{u}] \quad \text{and} \quad \widehat{p}_t(\mathbf{u}) = \mathbb{P}[\widehat{\mathbf{V}}_t = \mathbf{u}]$$

for all $\mathbf{u} \in \mathcal{S}$ and studying the relationship between the distributions p_t and \widehat{p}_t . To this end, let us fix an arbitrary $\mathbf{u} \in \mathcal{S}$ and define the “sparse loss vector” $\tilde{\ell}_t(\mathbf{u})$ with its k -th component being $\tilde{\ell}_{k,t}(\mathbf{u}) = u_k \ell_{k,t}$. Let

$$\widetilde{\mathbf{V}}_t(\mathbf{u}) = \arg \min_{\mathbf{w} \in \mathcal{S}} \mathbf{w}^\top (\mathbf{L}_{t-1} + \tilde{\ell}_t(\mathbf{u}) + \widehat{\mathbf{Z}}_t)$$

and

$$\tilde{p}_t(\mathbf{u}) = \mathbb{P}[\widetilde{\mathbf{V}}_t(\mathbf{u}) = \mathbf{u}].$$

As shown in Lemma 2 of [20], $\tilde{p}_t(\mathbf{u}) \leq \widehat{p}_t(\mathbf{u})$ holds independently of the distribution of the perturbations, given that all components of the loss vector are nonnegative. Now define

$$\mathbf{w}_t(\mathbf{z}) = \arg \min_{\mathbf{w} \in \mathcal{S}} \mathbf{w}^\top (\mathbf{L}_{t-1} + \mathbf{z})$$

for all $\mathbf{z} \in \mathbb{R}^d$ and let $f_t(\mathbf{z})$ be the density of \mathbf{Z}_t (which coincides with the density of $\widehat{\mathbf{Z}}_t$). For all $\mathbf{u} \in \mathcal{S}$, we have

$$\begin{aligned}
p_t(\mathbf{u}) &= \mathbb{E}[\mathbb{1}_{\{\mathbf{w}_t(\mathbf{Z}_t) = \mathbf{u}\}}] \\
&= \int_{\mathbf{z} \in \mathbb{R}^d} f_t(\mathbf{z}) \mathbb{1}_{\{\mathbf{w}_t(\mathbf{z}) = \mathbf{u}\}} d\mathbf{z} \\
&= \int_{\mathbf{z} \in \mathbb{R}^d} f_t(\mathbf{z} + \tilde{\ell}_t(\mathbf{u})) \mathbb{1}_{\{\mathbf{w}_t(\mathbf{z} + \tilde{\ell}_t(\mathbf{u})) = \mathbf{u}\}} d\mathbf{z} \\
&= \mathbb{E}[\mathbb{1}_{\{\mathbf{w}_t(\widehat{\mathbf{Z}}_t + \tilde{\ell}_t(\mathbf{u})) = \mathbf{u}\}}] \\
&\quad + \int_{\mathbf{z} \in \mathbb{R}^d} (f_t(\mathbf{z} + \tilde{\ell}_t(\mathbf{u})) - f_t(\mathbf{z})) \mathbb{1}_{\{\mathbf{w}_t(\mathbf{z}) = \mathbf{u}\}} d\mathbf{z} \\
&= \tilde{p}_t(\mathbf{u}) + \int_{\mathbf{z} \in \mathbb{R}^d} (f_t(\mathbf{z} + \tilde{\ell}_t(\mathbf{u})) - f_t(\mathbf{z})) \mathbb{1}_{\{\mathbf{w}_t(\mathbf{z}) = \mathbf{u}\}} d\mathbf{z}.
\end{aligned}$$

While the first term is upper bounded by $\widehat{p}_t(\mathbf{u})$, the last one can be upper bounded as

$$\begin{aligned}
&\int_{\mathbf{z} \in \mathbb{R}^d} f_t(\mathbf{z}) \left(1 - \exp\left(\frac{(\mathbf{z} - \tilde{\ell}_t(\mathbf{u}))^\top \tilde{\ell}_t(\mathbf{u})}{\eta^2 t}\right)\right) \mathbb{1}_{\{\mathbf{w}_t(\mathbf{z}) = \mathbf{u}\}} d\mathbf{z} \\
&\leq - \int_{\mathbf{z} \in \mathbb{R}^d} f_t(\mathbf{z}) \left(\frac{(\mathbf{z} - \tilde{\ell}_t(\mathbf{u}))^\top \tilde{\ell}_t(\mathbf{u})}{\eta^2 t}\right) \mathbb{1}_{\{\mathbf{w}_t(\mathbf{z}) = \mathbf{u}\}} d\mathbf{z} \\
&\leq \frac{p_t(\mathbf{u}) \|\tilde{\ell}_t(\mathbf{u})\|_2^2}{\eta^2 t} + \frac{1}{\eta^2 t} \int_{\mathbf{z} \in \mathbb{R}^d} f_t(\mathbf{z}) \|\mathbf{z}^\top \tilde{\ell}_t(\mathbf{u})\| \mathbb{1}_{\{\mathbf{w}_t(\mathbf{z}) = \mathbf{u}\}} d\mathbf{z} \\
&\leq \frac{p_t(\mathbf{u}) m}{\eta^2 t} + \frac{m}{\eta^2 t} \int_{\mathbf{z} \in \mathbb{R}^d} f_t(\mathbf{z}) \|\mathbf{z}\|_\infty \mathbb{1}_{\{\mathbf{w}_t(\mathbf{z}) = \mathbf{u}\}} d\mathbf{z},
\end{aligned}$$

where we have used that $\|\tilde{\ell}_t(\mathbf{u})\|_2^2 \leq m$ and $\|\tilde{\ell}_t(\mathbf{u})\|_1 \leq m$ hold by the definition of $\tilde{\ell}_t(\mathbf{u})$. Using that $\mathbb{P}[\hat{\mathbf{V}}_t = \mathbf{u}] = \hat{p}_t(\mathbf{u})$, we obtain

$$\begin{aligned} \mathbb{E}[\mathbf{V}_t^\top \ell_t] &= \sum_{\mathbf{u} \in \mathcal{S}} p_t(\mathbf{u}) \mathbf{u}^\top \ell_t \leq \sum_{\mathbf{u} \in \mathcal{S}} \hat{p}_t(\mathbf{u}) \mathbf{u}^\top \ell_t \\ &+ \sum_{\mathbf{u} \in \mathcal{S}} \left(\frac{p_t(\mathbf{u})m}{\eta^2 t} + \frac{m}{\eta^2 t} \int_{\mathbf{z} \in \mathbb{R}^d} f_t(\mathbf{z}) \|\mathbf{z}\|_\infty \mathbb{1}_{\{\mathbf{w}_t(\mathbf{z})=\mathbf{u}\}} d\mathbf{z} \right) \mathbf{u}^\top \ell_t \\ &\leq \mathbb{E}[\hat{\mathbf{V}}_t^\top \ell_t] + \frac{m^2}{\eta^2 t} + \frac{m^2}{\eta^2 t} \int_{\mathbf{z} \in \mathbb{R}^d} f_t(\mathbf{z}) \|\mathbf{z}\|_\infty \sum_{\mathbf{u} \in \mathcal{S}} \mathbb{1}_{\{\mathbf{w}_t(\mathbf{z})=\mathbf{u}\}} d\mathbf{z} \\ &= \mathbb{E}[\hat{\mathbf{V}}_t^\top \ell_t] + \frac{m^2}{\eta^2 t} + \frac{m^2}{\eta^2 t} \mathbb{E}[\|\mathbf{Z}_t\|_\infty] \\ &\leq \mathbb{E}[\hat{\mathbf{V}}_t^\top \ell_t] + \frac{m^2}{\eta^2 t} + \frac{m^2}{\eta} \sqrt{\frac{2 \log d}{t}}, \end{aligned}$$

where we used $\mathbb{E}[\|\mathbf{Z}_t\|_\infty] \leq \eta \sqrt{2t \log d}$ in the last step.

Putting everything together, we obtain

$$\begin{aligned} \mathbb{E} \hat{L}_n - \mathbf{v}^\top \mathbf{L}_n &\leq \sum_{t=1}^n \frac{m^2}{\eta^2 t} + \sum_{t=1}^n \frac{m^2}{\eta} \sqrt{\frac{2 \log d}{t}} + \eta m \sqrt{2n \log d} \\ &\leq \frac{2m^2 \sqrt{2n \log d}}{\eta} + \eta m \sqrt{2n \log d} + \frac{m^2 (\log n + 1)}{\eta^2}, \end{aligned}$$

concluding the proof of the statement. \blacksquare

ACKNOWLEDGMENT

This research was supported in part by the Natural Sciences and Engineering Research Council (NSERC) of Canada, the Spanish Ministry of Science and Technology grant MTM2012-37195, INRIA, the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreement 270327 (project ComplACS), the Ministry of Higher Education and Research, and by FUI project Hermès. The authors thank László and János Györfi, as well as László Németh for organizing the workshop PICSA 2012 in Jásd, during which some major details of the paper were worked out.

REFERENCES

- [1] L. Devroye, G. Lugosi, and G. Neu, “Prediction by random-walk perturbation,” in *Proceedings of the 25th Annual Conference on Learning Theory* (S. Shalev-Shwartz and I. Steinwart, eds.), pp. 460–473, 2013.
- [2] N. Cesa-Bianchi and G. Lugosi, *Prediction, Learning, and Games*. New York, NY, USA: Cambridge University Press, 2006.
- [3] C. Gentile and M. Warmuth, “Linear hinge loss and average margin,” in *Advances in Neural Information Processing Systems (NIPS)*, pp. 225–231, 1998.
- [4] J. Kivinen and M. Warmuth, “Relative loss bounds for multidimensional regression problems,” *Machine Learning*, vol. 45, pp. 301–329, 2001.
- [5] A. Grove, N. Littlestone, and D. Schuurmans, “General convergence results for linear discriminant updates,” *Machine Learning*, vol. 43, pp. 173–210, 2001.
- [6] E. Takimoto and M. Warmuth, “Paths kernels and multiplicative updates,” *Journal of Machine Learning Research*, vol. 4, pp. 773–818, 2003.
- [7] A. Kalai and S. Vempala, “Efficient algorithms for online decision problems,” *Journal of Computer and System Sciences*, vol. 71, pp. 291–307, 2005.
- [8] M. Warmuth and D. Kuzmin, “Randomized online PCA algorithms with regret bounds that are logarithmic in the dimension,” *Journal of Machine Learning Research*, vol. 9, pp. 2287–2320, 2008.

- [9] D. P. Helmbold and M. Warmuth, “Learning permutations with exponential weights,” *Journal of Machine Learning Research*, vol. 10, pp. 1705–1736, 2009.
- [10] E. Hazan, S. Kale, and M. Warmuth, “Learning rotations with little regret,” in *Proceedings of the 23rd Annual Conference on Learning Theory (COLT)*, pp. 144–154, 2010.
- [11] W. Koolen, M. Warmuth, and J. Kivinen, “Hedging structured concepts,” in *Proceedings of the 23rd Annual Conference on Learning Theory (COLT)*, pp. 93–105, 2010.
- [12] N. Cesa-Bianchi and G. Lugosi, “Combinatorial bandits,” *Journal of Computer and System Sciences*, vol. 78, pp. 1404–1422, 2012.
- [13] J. Y. Audibert, S. Bubeck, and G. Lugosi, “Regret in online combinatorial optimization,” *Mathematics of Operations Research*, vol. 39, pp. 31–45, 2014.
- [14] N. Littlestone and M. Warmuth, “The weighted majority algorithm,” *Information and Computation*, vol. 108, pp. 212–261, 1994.
- [15] V. Vovk, “Aggregating strategies,” in *Proceedings of the third annual workshop on Computational learning theory (COLT)*, pp. 371–386, 1990.
- [16] Y. Freund and R. E. Schapire, “A decision-theoretic generalization of on-line learning and an application to boosting,” *Journal of Computer and System Sciences*, vol. 55, pp. 119–139, 1997.
- [17] J. Hannan, “Approximation to Bayes risk in repeated play,” *Contributions to the theory of games*, vol. 3, pp. 97–139, 1957.
- [18] M. Hutter and J. Poland, “Prediction with expert advice by following the perturbed leader for general weights,” in *ALT*, pp. 279–293, 2004.
- [19] J. Poland, “FPL analysis for adaptive bandits,” in *In 3rd Symposium on Stochastic Algorithms, Foundations and Applications (SAGA’05)*, pp. 58–69, 2005.
- [20] G. Neu and G. Bartók, “An efficient algorithm for learning with semi-bandit feedback,” in *Proceedings of the 24th International Conference on Algorithmic Learning Theory* (S. Jain, R. Munos, F. Stephan, and T. Zeugmann, eds.), vol. 8139 of *Lecture Notes in Computer Science*, pp. 234–248, Springer, 2013.
- [21] D. Suehiro, K. Hatano, S. Kijima, E. Takimoto, and K. Nagano, “Online prediction under submodular constraints,” in *Algorithmic Learning Theory*, vol. 7568 of *Lecture Notes in Computer Science*, pp. 260–274, Springer Berlin Heidelberg, 2012.
- [22] S. Geulen, B. Voeccking, and M. Winkler, “Regret minimization for online buffering problems using the weighted majority algorithm,” in *Proceedings of the 23rd Annual Conference on Learning Theory (COLT 2010)* (A. Kalai and M. Mohri, eds.), pp. 132–143, 2010.
- [23] A. Györfi and G. Neu, “Near-optimal rates for limited-delay universal lossy source coding,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT 2011)*, 2011.
- [24] E. Even-Dar, S. M. Kakade, and Y. Mansour, “Online Markov decision processes,” *Mathematics of Operations Research*, vol. 34, no. 3, pp. 726–736, 2009.
- [25] G. Neu, A. Györfi, Cs. Szepesvári, and A. Antos, “Online Markov decision processes under bandit feedback,” in *Advances in Neural Information Processing Systems 23* (J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, eds.), pp. 1804–1812, 2011.
- [26] S. Rakhlin, O. Shamir, and K. Sridharan, “Relax and randomize : From value to algorithms,” in *Advances in Neural Information Processing Systems 25*, pp. 2150–2158, 2012.
- [27] W. Feller, *An Introduction to Probability Theory and its Applications, Vol. 1*. New York: John Wiley, 1968.
- [28] S. Boucheron, G. Lugosi, and P. Massart, *Concentration inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2013.

Luc Devroye (b. Belgium) obtained his Ph.D. from the University of Texas in 1976, and joined the School of Computer Science at McGill University in Montreal, Canada, in 1977. His research interests include probability theory as applied to the analysis of algorithms, mathematical statistics, machine learning, pattern recognition, and random number generation.

Gábor Lugosi graduated in electrical engineering at the Technical University of Budapest in 1987, and received his Ph.D. from the Hungarian Academy of Sciences in 1991. Since 1996, he has been at the Department of Economics, Pompeu Fabra University. In 2006 he became an ICREA research professor. His research interests include learning theory, nonparametric statistics, inequalities in probability, random structures, and information theory.

Gergely Neu received his M.Sc. degree in Electrical Engineering and his Ph.D. degree in Technical Informatics from the Budapest University of Technology and Economics (Hungary) in 2008 and 2013, respectively. Since 2013, he is a postdoctoral fellow at the SequeL team of INRIA Lille – Nord Europe. His research interests include reinforcement learning, online learning, and bandit problems.