



# Improved Second-Order Bounds for Prediction with Expert Advice

Nicolo Cesa-Bianchi, Yishay Mansour, Gilles Stoltz

## ► To cite this version:

Nicolo Cesa-Bianchi, Yishay Mansour, Gilles Stoltz. Improved Second-Order Bounds for Prediction with Expert Advice. 2006. <hal-00019799>

**HAL Id: hal-00019799**

**<https://hal.science/hal-00019799v1>**

Preprint submitted on 27 Feb 2006

**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.



HAL Authorization

# Improved Second-Order Bounds for Prediction with Expert Advice<sup>\*</sup>

Nicolò Cesa-Bianchi ([cesa-bianchi@dsi.unimi.it](mailto:cesa-bianchi@dsi.unimi.it))

*DSI, Università di Milano, via Comelico 39, 20135 Milano, Italy*

Yishay Mansour ([mansour@cs.tau.ac.il](mailto:mansour@cs.tau.ac.il))<sup>†</sup>

*School of computer Science, Tel-Aviv University, Tel Aviv, Israel*

Gilles Stoltz ([gilles.stoltz@ens.fr](mailto:gilles.stoltz@ens.fr))

*Département de Mathématiques et Applications, Ecole Normale Supérieure, 75005 Paris, France*

**Abstract.** This work studies external regret in sequential prediction games with both positive and negative payoffs. External regret measures the difference between the payoff obtained by the forecasting strategy and the payoff of the best action. In this setting, we derive new and sharper regret bounds for the well-known exponentially weighted average forecaster and for a new forecaster with a different multiplicative update rule. Our analysis has two main advantages: first, no preliminary knowledge about the payoff sequence is needed, not even its range; second, our bounds are expressed in terms of sums of squared payoffs, replacing larger first-order quantities appearing in previous bounds. In addition, our most refined bounds have the natural and desirable property of being stable under rescalings and general translations of the payoff sequence.

## 1. Introduction

The study of online forecasting strategies in adversarial settings has received considerable attention in the last few years. One of the goals of the research in this area is the design of randomized online algorithms that achieve a low external regret; i.e., algorithms able to minimize the difference between their expected cumulative payoff and the cumulative payoff achievable using the single best action (or, equivalently, the single best strategy in a given class).

If the payoffs are uniformly bounded, and there are finitely many actions, then there exist simple forecasting strategies whose external regret per time step vanishes irrespective to the choice of the payoff

---

<sup>\*</sup> An extended abstract appeared in the *Proceedings of the 18th Annual Conference on Learning Theory*, Springer, 2005. The work of all authors was supported in part by the IST Programme of the European Community, under the PASCAL Network of Excellence, IST-2002-506778.

<sup>†</sup> The work was done while the author was a fellow in the Institute of Advance studies, Hebrew University. His work was also supported by a grant no. 1079/04 from the Israel Science Foundation and an IBM faculty award.



sequence. In particular, under the assumption that all payoffs have the same sign (say positive), the best achieved rates for the regret are of the order of  $\sqrt{X^*}/n$ , where  $X^*/n$  is the highest average payoff among all actions after  $n$  time steps. If the payoffs were generated by an independent stochastic process, however, the tightest rate for the regret with respect to a fixed action should depend on the *variance* (rather than the average) of the observed payoffs for that action. Proving such a rate in a fully adversarial setting would be a fundamental result, and in this paper we propose new forecasting strategies that make a significant step towards this goal.

Generally speaking, one normally would expect any performance bound to be maintained under scaling and translation, since the units of measurement should not make a difference (for example, predicting the temperature should give similar performances irrespective to the scale, Celsius, Fahrenheit or Kelvin, on which the temperature is measured). However, in many computational settings this does not hold, for example in many domains there is a considerable difference between approximating a reward problem or its dual cost problem (although they have an identical optimal solution). Most of our bounds also assume no knowledge of the sequence of the ranges of the payoffs. For this reason it is important for us to stress that our bounds are stable under rescalings of the payoff sequence, even in the most general case of payoffs with arbitrary signs. The issues of invariance by translations and rescalings, discussed more in depth in Section 5.3, show that—in some sense—the bounds introduced in this paper are more “fundamental” than previous results. In order to describe our results we first set up our model and notations, and then we review previous related works.

In this paper we consider the following decision-theoretic variant proposed by Freund and Schapire (1997) of the framework of prediction with expert advice introduced by Littlestone and Warmuth (1994) and Vovk (1998). A forecaster repeatedly assigns probabilities to a fixed set of actions. After each assignment, the actual payoff associated to each action is revealed and new payoffs are set for the next round. The forecaster’s reward on each round is the average payoff of actions for that round, where the average is computed according to the forecaster’s current probability assignment. The goal of the forecaster is to achieve, on any sequence of payoffs, a cumulative reward close to  $X^*$ , the highest cumulative payoff among all actions. We call regret the difference between  $X^*$  and the cumulative reward achieved by the forecaster on the same payoff sequence.

In Section 2 we review the previously known bounds on the regret. The most basic one, obtained via the exponentially weighted average forecaster of Littlestone and Warmuth (1994) and Vovk (1998), bounds

the regret by a quantity of the order of  $M\sqrt{n \ln N}$ , where  $N$  is the number of actions and  $M$  is a known upper bound on the magnitude of payoffs.

In the special case of “one-sided games”, when all payoffs have the same sign (they are either always nonpositive or always nonnegative), Freund and Schapire (1997) showed that Littlestone and Warmuth’s weighted majority algorithm (1994) can be used to obtain a regret of the order of  $\sqrt{M|X^*| \ln N} + M \ln N$ . (If all payoffs are nonpositive, then the absolute value of each payoff is called *loss* and  $|X^*|$  is the cumulative loss of the best action.) By a simple rescaling and translation of payoffs, it is possible to reduce the more general “signed game”, in which each payoff might have an arbitrary sign, to either one of the one-sided games, and thus, bounds can be derived using this reduction. However the transformation also maps  $|X^*|$  to either  $Mn + X^*$  or  $Mn - X_n^*$ , thus significantly weakening the attractiveness of such a bound.

Recently, Allenberg-Neeman and Neeman (2004) proposed a direct analysis of the signed game avoiding this reduction. They proved that weighted majority (used in conjunction with a doubling trick) achieves the following: on any sequence of payoffs there exists an action  $j$  such that the regret is at most of order  $\sqrt{M(\ln N) \sum_{t=1}^n |x_{j,t}|}$ , where  $x_{j,t}$  is the payoff obtained by action  $j$  at round  $t$ , and  $M = \max_{i,t} |x_{i,t}|$  is a known upper bound on the magnitude of payoffs. Note that this bound does not relate the regret to the sum  $A_n^* = |x_{j^*,1}| + \dots + |x_{j^*,n}|$  of payoff magnitudes for the optimal action  $j^*$  (i.e., the one achieving  $X_n^*$ ). In particular, the bound of order  $\sqrt{MA_n^* \ln N} + M \ln N$  for one-sided games is only obtained if an estimate of  $A_n^*$  is available in advance.

In this paper we show new regret bounds for signed games. Our analysis has two main advantages: first, no preliminary knowledge about the payoff magnitude  $M$  or about the best cumulative payoff  $X^*$  is needed; second, our bounds are expressed in terms of sums of squared payoffs, such as  $x_{j,1}^2 + \dots + x_{j,n}^2$  and related forms. These quantities replace the larger terms  $M(|x_{j,1}| + \dots + |x_{j,n}|)$  appearing in the previous bounds. As an application of our results we obtain, without any preliminary knowledge on the payoff sequence, an improved regret bound for one-sided games of the order of  $\sqrt{(Mn - |X^*|)(|X^*|/n)(\ln N)}$ .

Some of our bounds are achieved using forecasters based on weighted majority run with a dynamic learning rate. However, we are able to obtain second-order bounds of a different flavor using a new forecaster that does not use the exponential probability assignments of weighted majority. In particular, unlike virtually all previously known forecasting schemes, the weights of this forecaster cannot be represented as the

gradient of an additive potential (see the monograph by Cesa-Bianchi and Lugosi, 2006 for an introduction to potential-based forecasters).

## 2. An overview of our results

We classify the existing regret bounds as zero-, first-, and second-order bounds. A zero-order regret bound depends only on the number of time steps and on upper bounds on the individual payoffs. A first-order bound has a main term that depends on a sum of payoffs, while the main term of a second order bound depends on a sum of squares of the payoffs. In this section we will also briefly discuss the information which the algorithms require in order to achieve the bounds.

We first introduce some notation and terminology. Our forecasting game is played in rounds. At each time step  $t = 1, 2, \dots$  the forecaster computes an assignment  $\mathbf{p}_t = (p_{1,t}, \dots, p_{N,t})$  of probabilities over the  $N$  actions. Then the payoff vector  $\mathbf{x}_t = (x_{1,t}, \dots, x_{N,t}) \in \mathbb{R}^N$  for time  $t$  is revealed and the forecaster's reward is  $\hat{x}_t = x_{1,t}p_{1,t} + \dots + x_{N,t}p_{N,t}$ . We define the cumulative reward of the forecaster by  $\hat{X}_n = \hat{x}_1 + \dots + \hat{x}_n$  and the cumulative payoff of action  $i$  by  $X_{i,n} = x_{i,1} + \dots + x_{i,n}$ . For all  $n$ , let  $X_n^* = \max_{i=1,\dots,N} X_{i,n}$  be the cumulative payoff of the best action up to time  $n$ . The forecaster's goal is to keep the *regret*  $X_n^* - \hat{X}_n$  as small as possible uniformly over  $n$ .

The one-sided games mentioned in the introduction are the *loss game*, where  $x_{i,t} \leq 0$  for all  $i$  and  $t$ , and the *gain game*, where  $x_{i,t} \geq 0$  for all  $i$  and  $t$ . We call *signed game* the setup in which no assumptions are made on the sign of the payoffs.

### 2.1. ZERO-ORDER BOUNDS

We say that a bound is of order zero whenever it only depends on bounds on the payoffs (or on the payoff ranges) and on the number of time steps  $n$ . The basic version of the exponentially weighted average forecaster of Littlestone and Warmuth (1994) ensures that the order of magnitude of the regret is  $M\sqrt{n \ln N}$  where  $M$  is a bound on the payoffs:  $|x_{i,t}| \leq M$  for all  $t \geq 1$  and  $i = 1, \dots, N$ . (Actually, the factor  $M$  may be replaced by a bound  $E$  on the *effective ranges* of the payoffs, defined by  $|x_{i,t} - x_{j,t}| \leq E$  for all  $t \geq 1$  and  $i, j = 1, \dots, N$ .) This basic version of this regret bound assumes that we have prior knowledge of both  $n$  and  $M$  (or  $E$ ).

In the case when  $n$  is not known in advance one can use a doubling trick (that is, restart the algorithm at times  $n = 2^k$  for  $k \geq \ln N$ ) and achieve a regret bound of the same order,  $M\sqrt{n \ln N}$  (only the

constant factor increases). Similarly, if  $M$  is not known in advance, one can restart the algorithm every time the maximum observed payoff exceeds the current estimate, and take the double of the old estimate as the new current estimate. Again, this influences the regret bound by only a constant factor. (The initial value of the estimate of  $M$  can be set to the maximal value in the first time step, see the techniques used in Section 3.)

A more elegant alternative, rather than the restarting the algorithm from scratch, is proposed by Auer, Cesa-Bianchi, and Gentile (2002) who consider a time-varying tuning parameter  $\eta_t \sim (1/M) \sqrt{(\ln N)/t}$ . They also derive a regret bound of the order of  $M \sqrt{n \ln N}$  uniformly over the number  $n$  of steps. Their method can be adapted along the lines of the techniques of Section 4.2 to deal with the case when  $M$  (or  $E$ ) is also unknown.

The results for the forecaster of Section 4 imply a zero-order bound sharper than  $E \sqrt{n \ln M}$ . This is presented in Corollary 1 and basically replaces  $E \sqrt{n}$  by  $\sqrt{E_1^2 + \dots + E_n^2}$ , where  $E_t$  is the *effective range* of the payoffs at round  $t$ ,

$$E_t = \max_{i=1,\dots,N} x_{i,t} - \min_{j=1,\dots,N} x_{j,t} . \quad (1)$$

## 2.2. ONE-SIDED GAMES: FIRST-ORDER REGRET BOUNDS

We say that a regret bound is first-order whenever its main term depends on a sum of payoffs. Since the payoff of any action is at most  $Mn$ , these bounds are usually sharper than zero-order bounds. More specifically, they have the potential of a huge improvement (when, for instance, the payoff of the best action is much smaller than  $Mn$ ) while they are at most worse by a constant factor with respect to their zero-order counterparts.

When all payoffs have the same sign Freund and Schapire (1997) first showed that Littlestone and Warmuth's weighted majority algorithm (1994) can be used as a basic ingredient to construct a forecasting strategy achieving a regret of order  $\sqrt{M|X_n^*| \ln N} + M \ln N$  where  $|X_n^*|$  is the absolute value of the cumulative payoff of the best action (i.e., the largest cumulative payoff in a gain game or the smallest cumulative loss in a loss game).

In order to achieve the above regret bound, the weighted majority algorithm needs prior knowledge of  $|X_n^*|$  (or a bound on it) and of the payoff magnitude  $M$ . As usual one can overcome this by a doubling trick. Doubling in this case is slightly more delicate, and would result in a bound of the order of  $\sqrt{M|X_n^*| \ln N} + M(\ln Mn) \ln N$ . Here

again, the techniques of Auer, Cesa-Bianchi, and Gentile (2002) could be adapted along the lines of the techniques of Section 4 to get a forecaster that, without restarting and without previous knowledge of  $M$  and  $X_n^*$ , achieves a regret bounded by a quantity of the order of  $\sqrt{M|X_n^*| \ln N} + M \ln N$ .

### 2.3. SIGNED GAMES: FIRST-ORDER REGRET BOUNDS

As mentioned in the introduction, one can translate a signed game to a one-sided game as follows. Consider a signed game with payoffs  $x_{i,t} \in [-M, M]$ . Provided that  $M$  is known to the forecaster, he may use the translation  $x'_{i,t} = x_{i,t} + M$  to convert the signed game into a gain game. For the resulting gain game, by using the techniques described above, one can derive a regret bound of the order of

$$\sqrt{(\ln N) (Mn + X_n^*)} + M \ln N . \quad (2)$$

Similarly, using the translation  $x'_{i,t} = x_{i,t} - M$ , we get a loss game, for which one can derive the similar regret bound

$$\sqrt{(\ln N) (Mn - X_n^*)} + M \ln N . \quad (3)$$

The main weakness of the transformation is that the bounds (2) and (3) are essentially zero-order bounds, though this depends on the precise value of  $X_n^*$ . (Note that when  $M$  is unknown, or to get tighter bounds, one may use the translation  $x'_{i,t} = x_{i,t} - \min_{j=1,\dots,N} x_{j,t}$  from signed games to gain games, or the translation  $x'_{i,t} = x_{i,t} - \max_{j=1,\dots,N} x_{j,t}$  from signed games to loss games.)

Recently, Allenberg-Neeman and Neeman (2004) proposed a direct analysis of the signed game avoiding this reduction. They give a simple algorithm whose regret is of the order of  $\sqrt{MA_n^* \ln N} + M \ln N$  where  $A_n^* = |x_{k_n^*,1}| + \dots + |x_{k_n^*,n}|$  is the sum of the absolute values of the payoffs of the best expert  $k_n^*$  for the rounds  $1, \dots, n$ . Since  $A_n^* = |X_n^*|$  in case of a one-sided game, this is indeed a generalization to signed games of Freund and Schapire's first-order bound for one-sided games. Though Allenberg-Neeman and Neeman need prior knowledge of both  $M$  and  $A_n^*$  to tune the parameters of the algorithm, a direct extension of their results along the lines of Section 3.1 gives the first-order bound

$$\begin{aligned} & \sqrt{M(\ln N) \max_{t=1,\dots,n} A_t^*} + M \ln N \\ &= \sqrt{M(\ln N) \max_{t=1,\dots,n} \sum_{s=1}^t |x_{k_t^*,s}|} + M \ln N \end{aligned} \quad (4)$$

which holds when only  $M$  is known.

#### 2.4. SECOND-ORDER BOUNDS ON THE REGRET

A regret bound is second-order whenever its main term is a function of a sum of squared payoffs (or on a quantity that is homogeneous in such a sum). Ideally, they are a function of

$$Q_n^* = \sum_{t=1}^n x_{k_n^*, t}^2 .$$

Expressions involving squared payoffs are at the core of many analyses in the framework of prediction with expert advice, especially in the presence of limited feedback. (See, for instance, the bandit problem, studied by Auer et al., 2002, and more generally prediction under partial monitoring and the work of Cesa-Bianchi, Lugosi, and Stoltz, 2005, Cesa-Bianchi, Lugosi, and Stoltz, 2004, Piccolboni and Schindelhauer, 2001.) However, to the best of our knowledge, the bounds presented here are the first ones to explicitly include second-order information extracted from the payoff sequence.

In Section 3 we give a very simple algorithm whose regret is of the order of  $\sqrt{Q_n^* \ln N} + M \ln N$ . Since  $Q_n^* \leq MA_n^*$ , this bound improves on the first-order bounds. Even though our basic algorithm needs prior knowledge of both  $M$  and  $Q_n^*$  to tune its parameters, we are able to extend it (essentially by using various doubling tricks) and achieve a bound of the order of

$$\sqrt{(\ln N) \max_{t=1, \dots, n} Q_t^*} + M \ln N = \sqrt{(\ln N) \max_{t=1, \dots, n} \sum_{s=1}^t x_{k_t^*, s}^2} + M \ln N \quad (5)$$

without using any prior knowledge about  $Q_n^*$ . (The extension is not as straightforward as one would expect, since the quantities  $Q_t^*$  are not necessarily monotone over time.)

Note that this bound is less sensitive to extreme values. For instance, in case of a loss game (i.e., all payoffs are nonpositive),  $Q_t^* \leq ML_t^*$ , where  $L_t^*$  is the cumulative loss of the best action up to time  $t$ . Therefore,  $\max_{s \leq n} Q_s^* \leq ML_n^*$  and the bound (5) is at least as good as the family of bounds called “improvements for small losses” (or first-order bounds) presented in Section 2.2. However, it is easy to exhibit examples where the new bound is far better by considering sequences of outcomes where there are some “outliers” among the  $x_{i,t}$ . These outliers may raise the maximum  $M$  significantly, whereas they have only little impact on the  $\max_{s \leq n} Q_s^*$ .



We also analyze the weighted majority algorithm in Section 4, and show how exponential weights with a time varying parameter can be used to derive a regret bound of the order of  $\sqrt{V_n \ln N} + E \ln N$  where  $V_n$  is the cumulative variance of the forecaster's rewards on the given sequence and  $E$  is the range of the payoffs. (Again, we derive first the bound in the case where the payoff range is known, and then extend it to the case where the payoff range is unknown.) The above bound is somewhat different from standard regret bounds because it depends on the predictions of the forecaster. In Sections 4.4 and 5 we show how one can use such a bound to derive regret bounds which only depend on the sequence of payoffs.

### 3. A new algorithm for sequential prediction

We introduce a new forecasting strategy for the signed game. In Theorem 4, the main result of this section, we show that, without any preliminary knowledge of the sequence of payoffs, the regret of a variant of this strategy is bounded by a quantity defined in terms of the sums  $Q_{i,n} = x_{i,1}^2 + \dots + x_{i,n}^2$ . Since  $Q_{i,n} \leq M(|x_{i,1}| + \dots + |x_{i,n}|)$ , such second-order bounds are generally better than all previously known bounds (see Section 2).

Our basic forecasting strategy, which we call  $\text{prod}(\eta)$ , has an input parameter  $\eta > 0$  and maintains a set of  $N$  weights. At time  $t = 1$  the weights are initialized with  $w_{i,1} = 1$  for  $i = 1, \dots, N$ . At each time  $t = 1, 2, \dots$ ,  $\text{prod}(\eta)$  computes the probability assignment  $\mathbf{p}_t = (p_{1,t}, \dots, p_{N,t})$ , where  $p_{i,t} = w_{i,t}/W_t$  and  $W_t = w_{1,t} + \dots + w_{N,t}$ . After the payoff vector  $\mathbf{x}_t$  is revealed, the weights are updated using the rule  $w_{i,t+1} = w_{i,t}(1 + \eta x_{i,t})$ . The following simple fact plays a key role in our analysis.

*Lemma 1.* For all  $z \geq -1/2$ ,  $\ln(1+z) \geq z - z^2$ .

*Proof.* Let  $f(z) = \ln(1+z) - z + z^2$ . Note that

$$f'(z) = \frac{1}{1+z} - 1 + 2z = \frac{z(1+2z)}{1+z}$$

so that  $f'(z) \leq 0$  for  $-1/2 \leq z \leq 0$  and  $f'(z) \geq 0$  for  $z \geq 0$ . Hence the minimum of  $f$  is achieved in 0 and equals 0, concluding the proof.  $\square$

We are now ready to state a lower bound on the cumulative reward of  $\text{prod}(\eta)$  in terms of the quantities  $Q_{k,n}$ .

*Lemma 2.* Assume there exists  $M > 0$  such that the payoffs satisfy  $x_{i,t} \geq -M$  for  $t = 1, \dots, n$  and  $i = 1, \dots, N$ . For any sequence of payoffs, for any action  $k$ , for any  $\eta \leq 1/(2M)$ , and for any  $n \geq 1$ , the cumulative reward of  $\text{prod}(\eta)$  is lower bounded as

$$\hat{X}_n \geq X_{k,n} - \frac{\ln N}{\eta} - \eta Q_{k,n} .$$

*Proof.* For any  $k = 1, \dots, N$ , note that  $x_{k,t} \geq -M$  and  $\eta \leq 1/(2M)$  imply  $\eta x_{k,t} \geq -1/2$ . Hence, we can apply Lemma 1 to  $\eta x_{k,t}$  and get

$$\begin{aligned} \ln \frac{W_{n+1}}{W_1} &\geq \ln \frac{W_{k,n+1}}{W_1} = -\ln N + \ln \prod_{t=1}^n (1 + \eta x_{k,t}) = -\ln N + \sum_{t=1}^n \ln(1 + \eta x_{k,t}) \\ &\geq -\ln N + \sum_{t=1}^n (\eta x_{k,t} - \eta^2 x_{k,t}^2) = -\ln N + \eta X_{k,n} - \eta^2 Q_{k,n} . \end{aligned} \quad (6)$$

On the other hand,

$$\begin{aligned} \ln \frac{W_{n+1}}{W_1} &= \sum_{t=1}^n \ln \frac{W_{t+1}}{W_t} = \sum_{t=1}^n \ln \left( \sum_{i=1}^N p_{i,t} (1 + \eta x_{i,t}) \right) \\ &= \sum_{t=1}^n \ln \left( 1 + \eta \sum_{i=1}^N x_{i,t} p_{i,t} \right) \leq \eta \hat{X}_n \end{aligned} \quad (7)$$

where in the last step we used  $\ln(1+z_t) \leq z_t$  for all  $z_t = \eta \sum_{i=1}^N x_{i,t} p_{i,t} \geq -1/2$ . Combining (6) and (7), and dividing by  $\eta > 0$ , we get

$$\hat{X}_n \geq -\frac{\ln N}{\eta} + X_{k,n} - \eta Q_{k,n} .$$

Our choice of  $\eta$  gives the claimed bound.  $\square$

By choosing  $\eta$  appropriately, we can optimize the bound as follows.

*Theorem 1.* Assume there exists  $M > 0$  such that the payoffs satisfy  $x_{i,t} \geq -M$  for  $t = 1, \dots, n$  and  $i = 1, \dots, N$ . For any  $Q > 0$ , if  $\text{prod}(\eta)$  is run with

$$\eta = \min \left\{ 1/(2M), \sqrt{(\ln N)/Q} \right\} \quad (8)$$

then for any sequence of payoffs, for any action  $k$ , and for any  $n \geq 1$  such that  $Q_{k,n} \leq Q$ ,

$$\hat{X}_n \geq X_{k,n} - \max \left\{ 2\sqrt{Q \ln N}, 4M \ln N \right\} .$$

### 3.1. UNKNOWN BOUND ON QUADRATIC VARIATION (Q)

To achieve the bound stated in Theorem 1, the parameter  $\eta$  must be tuned using preliminary knowledge of a lower bound on the payoffs and an upper bound on the quantities  $Q_{k,n}$ . In this and the following sections we remove these requirements one by one. We start by introducing a new algorithm that, using a doubling trick over **prod**, avoids any preliminary knowledge of an upper bound on the  $Q_{k,n}$ .

Let  $k_t^*$  be the index of the best action up to time  $t$ ; that is,  $k_t^* \in \operatorname{argmax}_k X_{k,t}$  (ties are broken by choosing the action  $k$  with minimal associated  $Q_{k,t}$ ). We denote the associated quadratic penalty by

$$Q_t^* = Q_{k_t^*}^* = \sum_{s=1}^t x_{k_t^*,s}^2 .$$

Ideally, our regret bound should depend on  $Q_n^*$  and be of the form  $\sqrt{Q_n^* \ln N} + M \ln N$ . However, note that the sequence  $Q_1^*, Q_2^*, \dots$  is not necessarily monotone, since if at time  $t+1$  the best action changes, then  $Q_t^*$  and  $Q_{t+1}^*$  are not related. Therefore, we cannot use a straightforward doubling trick, as this only applies to monotone sequences. Our solution is to express the bound in terms of the smallest nondecreasing sequence that upper bounds the original sequence  $(Q_t^*)_{t \geq 1}$ . This is a general trick to handle situations where the penalty terms are not monotone.

Let **prod-Q**( $M$ ) be the prediction algorithm that receives a quantity  $M > 0$  as input parameter and repeatedly runs **prod**( $\eta_r$ ), where  $\eta_r$  is defined below. The parameter  $M$  is a bound on the payoffs, such that for all  $i = 1, \dots, N$  and  $t = 1, \dots, n$ , we have  $|x_{i,t}| \leq M$ . The  $r$ -th parameter  $\eta_r$  corresponds to the parameter  $\eta$  defined in (8) for  $M$  and  $Q = 4^r M^2$ . Namely, we choose

$$\eta_r = \min \left\{ 1/(2M), \sqrt{\ln N}/(2^r M) \right\} .$$

We call epoch  $r$ ,  $r = 0, 1, \dots$ , the sequence of time steps when **prod-Q** is running **prod**( $\eta_r$ ). The last step of epoch  $r \geq 0$  is the time step  $t = t_r$  when  $Q_t^* > 4^r M^2$  happens for the first time. When a new epoch  $r+1$  begins, **prod** is restarted with parameter  $\eta_{r+1}$ .

*Theorem 2.* Given  $M > 0$ , for all  $n \geq 1$  and all sequences of payoffs bounded by  $M$ , i.e.,  $\max_{1 \leq i \leq N} \max_{1 \leq t \leq n} |x_{i,t}| \leq M$ , the cumulative reward of algorithm **prod-Q**( $M$ ) satisfies

$$\begin{aligned} \hat{X}_n &\geq X_n^* - 8\sqrt{(\ln N) \max_{s \leq n} Q_s^*} \\ &\quad - 2M \left( 1 + \log_4 n + 2(1 + \lfloor (\log_2 \ln N)/2 \rfloor) \ln N \right) \\ &= X_n^* - O \left( \sqrt{(\ln N) \max_{s \leq n} Q_s^*} + M \ln n + M \ln N \ln \ln N \right) . \end{aligned}$$

*Proof.* We denote by  $R$  the index of the last epoch and let  $t_R = n$ . If we have only one epoch, then the theorem follows from Theorem 1 applied with a bound of  $Q = M^2$  on the squared payoffs of the best expert. Therefore, for the rest of the proof we assume  $R \geq 1$ . Let

$$X_k^{(r)} = \sum_{s=t_{r-1}+1}^{t_r-1} x_{k,s}, \quad Q_k^{(r)} = \sum_{s=t_{r-1}+1}^{t_r-1} x_{k,s}^2, \quad \hat{X}^{(r)} = \sum_{s=t_{r-1}+1}^{t_r-1} \hat{x}_s$$

where the sums are over all the time steps  $s$  in epoch  $r$  except the last one,  $t_r$ . (Here  $t_{-1}$  is conventionally set to 0.) We also denote  $k_r = k_{t_r-1}^*$  the index of the best overall expert up to time  $t_r - 1$  (one time step before the end of epoch  $r$ ). We have that  $Q_{k_r}^{(r)} \leq Q_{k_r, t_r-1} = Q_{t_r-1}^*$ . Now, by definition of the algorithm,  $Q_{t_r-1}^* \leq 4^r M^2$ . Theorem 1 (applied to time steps  $t_{r-1} + 1, \dots, t_r - 1$ ) shows that

$$\hat{X}^{(r)} \geq X_{k_r}^{(r)} - \max \left\{ 2\sqrt{4^r M^2 \ln N}, 4M \ln N \right\}.$$

The maximum in the right-hand side equals  $2^{r+1} M \sqrt{\ln N}$  when  $r > r_0 = 1 + \lfloor (\log_2 \ln N)/2 \rfloor$ . Summing over  $r = 0, \dots, R$  we get

$$\begin{aligned} \hat{X}_n &= \sum_{r=0}^R \left( \hat{X}^{(r)} + \hat{x}_{k_r, t_r} \right) \\ &\geq \sum_{r=0}^R \left( \hat{x}_{k_r, t_r} + X_{k_r}^{(r)} \right) - 4(1 + r_0)M \ln N - \sum_{r=r_0+1}^R 2\sqrt{4^r M^2 \ln N} \\ &\geq \sum_{r=0}^R \left( \hat{x}_{k_r, t_r} + X_{k_r}^{(r)} \right) - 4(1 + r_0)M \ln N - 2^{R+2} M \sqrt{\ln N} \\ &\geq \sum_{r=0}^R X_{k_r}^{(r)} - (R + 1)M - 4(1 + r_0)M \ln N - 2^{R+2} M \sqrt{\ln N}. \quad (9) \end{aligned}$$

Now, since  $k_0$  is the index of the expert with largest payoff up to time  $t_0 - 1$ , we have that  $X_{k_1, t_1-1} = X_{k_1}^{(0)} + x_{k_1, t_0} + X_{k_1}^{(1)} \leq X_{k_0}^{(0)} + X_{k_1}^{(1)} + M$ . By a simple induction, we in fact get

$$X_{k_R, t_R-1} \leq \sum_{r=0}^{R-1} \left( X_{k_r}^{(r)} + M \right) + X_{k_R}^{(R)}. \quad (10)$$

As, in addition,  $X_{k_R, t_R-1} = X_{k_{n-1}^*, n-1}$  and  $X_{k_n^*, n}$  may only differ by at most  $M$ , combining (9) and (10) we have indeed proven that

$$\hat{X}_n \geq X_{k_n^*, n} - \left( 2(R + 1)M + 4M(1 + r_0) \ln N + 2^{R+2} M \sqrt{\ln N} \right).$$

The proof is concluded by noting first, that  $R \leq \log_4 n$ , and second that, as  $R \geq 1$ ,  $\max_{s \leq n} Q_s^* \geq 4^{R-1} M^2$  by definition of the algorithm.  $\square$

### 3.2. UNKNOWN BOUND ON PAYOFFS (M)

In this section we show how one can overcome the case when there is no a priori bound on the payoffs. In the next section we combine the techniques of this section and Section 3.1 to deal with the case when both parameters are unknown

Let  $\text{prod-M}(Q)$  be the prediction algorithm that receives a number  $Q > 0$  as input parameter and repeatedly runs  $\text{prod}(\eta_r)$ , where the  $\eta_r$ ,  $r = 0, 1, \dots$ , are defined below. We call epoch  $r$  the sequence of time steps when  $\text{prod-M}$  is running  $\text{prod}(\eta_r)$ . At the beginning,  $r = 0$  and  $\text{prod-M}(Q)$  runs  $\text{prod}(\eta_0)$ , where

$$M_0 = \sqrt{Q/(4 \ln N)} \quad \text{and} \quad \eta_0 = 1/(2M_0) = \sqrt{(\ln N)/Q}.$$

For all  $t \geq 1$ , we denote

$$M_t = \max_{s=1, \dots, t} \max_{i=1, \dots, N} 2^{\lceil \log_2 |x_{i,s}| \rceil}.$$

The last step of epoch  $r \geq 0$  is the time step  $t = t_r$  when  $M_t > M_{t_{r-1}}$  happens for the first time (conventionally, we set  $M_{t_{-1}} = M_0$ ). When a new epoch  $r + 1$  begins,  $\text{prod}$  is restarted with parameter  $\eta_{r+1} = 1/(2M_{t_r})$ .

Note that  $\eta_0 = 1/(2M_0)$  in round 0 and  $\eta_r = 1/(2M_{t_{r-1}})$  in any round  $r \geq 1$ , where  $M_{t_0} > M_0$  and  $M_{t_r} \geq 2M_{t_{r-1}}$  for each  $r \geq 1$ .

*Theorem 3.* For any sequence of payoffs, for any action  $k$ , and for any  $n \geq 1$  such that  $Q_{k,n} \leq Q$ , the cumulative reward of algorithm  $\text{prod-M}(Q)$  is lower bounded as

$$\hat{X}_n \geq X_{k,n} - 2\sqrt{Q \ln N} - 12 M (1 + \ln N)$$

where  $M = \max_{1 \leq i \leq N} \max_{1 \leq t \leq n} |x_{i,t}|$ .

*Proof.* As in the proof of Theorem 2, we denote by  $R$  the index of the last epoch and let  $t_R = n$ . We assume  $R \geq 1$  (otherwise, the theorem follows directly from Theorem 1 applied with a lower bound of  $-M_0$  on the payoffs). Note that at time  $n$  we have either  $M_n \leq M_{t_{R-1}}$ , implying  $M_n = M_{t_R} = M_{t_{R-1}}$ , or  $M_n > M_{t_{R-1}}$ , implying  $M_n = M_{t_R} = 2M_{t_{R-1}}$ . In both cases,  $M_{t_R} \geq M_{t_{R-1}}$ . In addition, since  $R \geq 1$ , we also have  $M_{t_R} \leq 2M$ .

Similarly to the proof of Theorem 2, for all epochs  $r$  and actions  $k$  introduce

$$X_k^{(r)} = \sum_{s=t_{r-1}+1}^{t_r-1} x_{k,s}, \quad Q_k^{(r)} = \sum_{s=t_{r-1}+1}^{t_r-1} x_{k,s}^2, \quad \hat{X}^{(r)} = \sum_{s=t_{r-1}+1}^{t_r-1} \hat{x}_s$$

where, as before, we set  $t_{-1} = 0$ . Applying Lemma 2 to each epoch  $r = 0, \dots, R$  we get that  $\hat{X}_n - X_{k,n}$  is equal to

$$\begin{aligned} \hat{X}_n - X_{k,n} &= \sum_{r=0}^R \left( \hat{X}^{(r)} - X_k^{(r)} \right) + \sum_{r=0}^R (\hat{x}_{t_r} - x_{k,t_r}) \\ &\geq - \sum_{r=0}^R \frac{\ln N}{\eta_r} - \sum_{r=0}^R \eta_r Q_k^{(r)} + \sum_{r=0}^R (\hat{x}_{t_r} - x_{k,t_r}) . \end{aligned}$$

We bound each sum separately. For the first sum, since  $M_{t_s} \geq 2^{s-r} M_{t_r}$  for each  $0 \leq r \leq s \leq R-1$ , we have for  $s \leq R-1$ ,

$$\sum_{r=0}^s M_{t_r} \leq \sum_{r=0}^s 2^{r-s} M_{t_s} \leq 2 M_{t_s} . \quad (11)$$

Thus,

$$\sum_{r=0}^R \frac{\ln N}{\eta_r} = \sum_{r=0}^R 2 M_{t_{r-1}} \ln N \leq 2 (M_{t_{-1}} + 2 M_{t_{R-1}}) \ln N \leq 6 M_{t_R} \ln N$$

where we used (11) and  $M_{t_{-1}} = M_0 \leq M_{t_{R-1}} \leq M_{t_R}$ . For the second sum, using the fact that  $\eta_r$  decreases with  $r$ , we have

$$\sum_{r=0}^R \eta_r Q_k^{(r)} \leq \eta_0 \sum_{r=0}^R Q_k^{(r)} \leq \eta_0 Q_{k,n} \leq \sqrt{\frac{\ln N}{Q}} Q = \sqrt{Q \ln N} .$$

Finally, using (11) again,

$$\sum_{r=0}^R |\hat{x}_{t_r} - x_{k,t_r}| \leq \sum_{r=0}^R 2 M_{t_r} \leq 2 (2 M_{t_{R-1}} + M_{t_R}) \leq 6 M_{t_R} .$$

The resulting lower bound  $6 M_{t_R} (1 + \ln N) + \sqrt{Q \ln N}$  implies the one stated in the theorem by recalling that, when  $R \geq 1$ ,  $M_{t_R} \leq 2 M$ .  $\square$

### 3.3. UNKNOWN BOUNDS ON BOTH PAYOFFS (M) AND QUADRATIC VARIATION (Q)

We now show a regret bound for the case when  $M$  and the  $Q_{k,n}$  are both unknown. We consider again the notation of the beginning of Section 3.1. The quantities of interest for the doubling trick of Section 3.1 were the homogeneous quantities  $(1/M^2) \max_{s \leq t} Q_s^*$ . Here we assume no knowledge of  $M$ . We propose a doubling trick on the only homogeneous quantities we have access to, that is,  $\max_{s \leq t} (Q_s^*/M_s^2)$ , where  $M_t$  is defined in Section 3.2 and the maximum is needed for the same reasons of monotonicity as in Section 3.1.

We define the new (parameterless) prediction algorithm **prod-MQ**. Intuitively, the algorithm can be thought as running, at the low level, the algorithm **prod-Q**( $M_t$ ). When the value of  $M_t$  changes, we restart **prod-Q**( $M_t$ ), with the new value but keep track of  $Q_t^*$ .

Formally, we define the prediction algorithm **prod-MQ** in the following way. Epochs are indexed by pairs  $(r, s)$ . At the beginning of each epoch  $(r, s)$ , the algorithm takes a fresh start and runs **prod**( $\eta_{(r,s)}$ ), where  $\eta_{(r,s)}$ , for  $r = 0, 1, \dots$  and  $s = 0, 1, \dots$ , is defined by

$$\eta_{(r,s)} = \min \left\{ 1 / \left( 2M^{(r)} \right), \sqrt{\ln N} / \left( 2^{S_{r-1}+s} M^{(r)} \right) \right\}$$

and  $M^{(r)}$ ,  $S_r$  are defined below.

At the beginning,  $r = 0$ ,  $s = 0$ , and since **prod**( $\eta$ ) always sets  $\mathbf{p}_1$  to be the uniform distribution irrespective to the choice of  $\eta$ , without loss of generality we assume that **prod** is started at epoch  $(0, 0)$  with  $M^{(0)} = M_1$  and  $S_{-1} = 0$ .

The last step of epoch  $(r, s)$  is the time step  $t = t_{(r,s)}$  when either:

$$(C1) \quad Q_t^* > 4^{S_{r-1}+s} M_t^2 \text{ happens for the first time}$$

or

$$(C2) \quad M_t > M^{(r)} \text{ happens for the first time.}$$

If epoch  $(r, s)$  ends because of (C1), the next epoch is  $(r, s + 1)$ , and the value of  $M^{(r)}$  is unchanged. If epoch  $(r, s)$  ends because of (C2), the next epoch is  $(r + 1, 0)$ ,  $S_r = S_{r-1} + s$ , and  $M^{(r+1)} = M_t$ .

Note that within epochs indexed by the same  $r$ , the payoffs in all steps but the last one are bounded by  $M^{(r)}$ . Note also that the quantities  $S_r$  count the number of times an epoch ended because of (C1). Finally, note that there are  $S_r - S_{r-1} + 1$  epochs  $(r, s)$  for a given  $r \geq 0$ , indexed by  $s = 0, \dots, S_r - S_{r-1}$ .

*Theorem 4.* For any sequence of payoffs and for any  $n \geq 1$ , the cumulative reward of algorithm **prod-MQ** satisfies

$$\begin{aligned}\hat{X}_n &\geq X_n^* - 32M\sqrt{q\ln N} \\ &\quad - 22M(1 + \ln N) - 2M\log_2 n - 4M\lceil(\log_2 \ln N)/2\rceil \\ &= X_n^* - O\left(M\sqrt{q\ln N} + M\ln n + M\ln N\right)\end{aligned}$$

where  $M = \max_{1 \leq i \leq N} \max_{1 \leq t \leq n} |x_{i,t}|$  and  $q = \max \left\{ 1, \max_{s \leq n} \frac{Q_s^*}{M_s^2} \right\}$ .

The proof is in the Appendix.

#### 4. Second-order bounds for weighted majority

In this section we derive new regret bounds for the weighted majority forecaster of Littlestone and Warmuth (1994) using a time-varying learning rate. This allows us to avoid the doubling tricks of Section 3 and keep the assumption that no knowledge on the payoff sequence is available to the forecaster beforehand.

Similarly to the results of Section 3, the main term in the new bounds depends on second-order quantities associated to the sequence of payoffs. However, the precise definition of these quantities makes the bounds of this section generally not comparable to the bounds obtained in Section 3.

The weighted majority forecaster using the sequence  $\eta_2, \eta_3, \dots > 0$  of learning rates assigns at time  $t$  a probability distribution  $\mathbf{p}_t$  over the  $N$  experts defined by  $\mathbf{p}_1 = (1/N, \dots, 1/N)$  and

$$p_{i,t} = \frac{e^{\eta_t X_{i,t-1}}}{\sum_{j=1}^N e^{\eta_t X_{j,t-1}}} \quad \text{for } i = 1, \dots, N \text{ and } t \geq 2. \quad (12)$$

Note that the quantities  $\eta_t > 0$  may depend on the past payoffs  $x_{i,s}$ ,  $i = 1, \dots, N$  and  $s = 1, \dots, t-1$ . The analysis of Auer, Cesa-Bianchi, and Gentile (2002), for a related variant of weighted majority, is at the core of the proof of the following lemma (proof in Appendix).

*Lemma 3.* Consider any nonincreasing sequence  $\eta_2, \eta_3, \dots$  of positive learning rates and any sequence  $\mathbf{x}_1, \mathbf{x}_2, \dots \in \mathbb{R}^N$  of payoff vectors. Define the nonnegative function  $\Phi$  by

$$\begin{aligned}\Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) &= -\sum_{i=1}^N p_{i,t} x_{i,t} + \frac{1}{\eta_t} \ln \sum_{i=1}^N p_{i,t} e^{\eta_t x_{i,t}} \\ &= \frac{1}{\eta_t} \ln \left( \sum_{i=1}^N p_{i,t} e^{\eta_t (x_{i,t} - \hat{x}_t)} \right) .\end{aligned}$$



Then the weighted majority forecaster (12) run with the sequence  $\eta_2, \eta_3, \dots$  satisfies, for any  $n \geq 1$  and for any  $\eta_1 \geq \eta_2$ ,

$$\hat{X}_n - X_n^* \geq - \left( \frac{2}{\eta_{n+1}} - \frac{1}{\eta_1} \right) \ln N - \sum_{t=1}^n \Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) .$$

Let  $Z_t$  be the random variable with range  $\{x_{1,t}, \dots, x_{N,t}\}$  and distribution  $\mathbf{p}_t$ . Note that  $\mathbb{E}Z_t$  is the expected payoff  $\hat{x}_t$  of the forecaster using distribution  $\mathbf{p}_t$  at time  $t$ . Introduce

$$\text{Var } Z_t = \mathbb{E}Z_t^2 - \mathbb{E}^2 Z_t = \sum_{i=1}^N p_{i,t} x_{i,t}^2 - \left( \sum_{i=1}^N p_{i,t} x_{i,t} \right)^2 .$$

Hence  $\text{Var } Z_t$  is the variance of the payoffs at time  $t$  under the distribution  $\mathbf{p}_t$  and the cumulative variance  $V_n = \text{Var } Z_1 + \dots + \text{Var } Z_n$  is the main second-order quantity used in this section. The next result bounds  $\Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t)$  in terms of  $\text{Var } Z_t$ .

*Lemma 4.* For all payoff vectors  $\mathbf{x}_t = (x_{1,t}, \dots, x_{N,t})$ , all probability distributions  $\mathbf{p}_t = (p_{1,t}, \dots, p_{N,t})$ , and all learning rates  $\eta_t \geq 0$ , we have

$$\Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) \leq E$$

where  $E$  is such that  $|x_{i,t} - x_{j,t}| \leq E$  for all  $i, j = 1, \dots, N$ . If, in addition,  $0 \leq \eta_t |x_{i,t} - x_{j,t}| \leq 1$  for all  $i, j = 1, \dots, N$ , then

$$\Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) \leq (e - 2)\eta_t \text{Var } Z_t .$$

*Proof.* The first inequality is straightforward. To prove the second one we use  $e^a \leq 1 + a + (e - 2)a^2$  for  $|a| \leq 1$ . Consequently, noting that  $\eta_t |x_{i,t} - \hat{x}_t| \leq 1$  for all  $i$  by assumption, we have that

$$\begin{aligned} \Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) &= \frac{1}{\eta_t} \ln \left( \sum_{i=1}^N p_{i,t} e^{\eta_t (x_{i,t} - \hat{x}_t)} \right) \\ &\leq \frac{1}{\eta_t} \ln \left( \sum_{i=1}^N p_{i,t} \left( 1 + \eta_t (x_{i,t} - \hat{x}_t) + (e - 2)\eta_t^2 (x_{i,t} - \hat{x}_t)^2 \right) \right) . \end{aligned}$$

Using  $\ln(1 + a) \leq a$  for all  $a > -1$  and some simple algebra concludes the proof of the second inequality.  $\square$

In Auer et al. (2002) a very similar result is proven, except that there the variance is further bounded (up to a multiplicative factor) by the expectation  $\hat{x}_t$  of  $Z_t$ .

#### 4.1. KNOWN BOUND ON THE PAYOFF RANGES (E)

We now introduce a time-varying learning rate based on  $V_n$ . For simplicity, we assume in a first time that a bound  $E$  on the payoff ranges  $E_t$ , defined in (1), is known beforehand and turn back to the general case in Theorem 6. The sequence  $\eta_2, \eta_3, \dots$  is defined as

$$\eta_t = \min \left\{ \frac{1}{E}, C \sqrt{\frac{\ln N}{V_{t-1}}} \right\} \quad (13)$$

for  $t \geq 2$ , with  $C = \sqrt{2(\sqrt{2} - 1)/(e - 2)} \approx 1.07$ .

Note that  $\eta_t$  depends on the forecaster's past predictions. This is in the same spirit as the self-confident learning rates considered in Auer, Cesa-Bianchi, and Gentile (2002).

*Theorem 5.* Provided a bound  $E$  on the payoff ranges is known beforehand, i.e.,  $\max_{t=1, \dots, n} \max_{i,j=1, \dots, N} |x_{i,t} - x_{j,t}| \leq E$ , the weighted majority forecaster using the time-varying learning rate (13) achieves, for all sequences of payoffs and for all  $n \geq 1$ ,

$$\hat{X}_n - X_n^* \geq -4\sqrt{V_n \ln N} - 2E \ln N - E/2 .$$

*Proof.* We start by applying Lemma 3 using the learning rate (13), and setting  $\eta_1 = \eta_2$  for the analysis,

$$\begin{aligned} \hat{X}_n - X_n^* &\geq - \left( \frac{2}{\eta_{n+1}} - \frac{1}{\eta_1} \right) \ln N - \sum_{t=1}^n \Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) \\ &\geq -2 \max \left\{ E \ln N, (1/C) \sqrt{V_n \ln N} \right\} - (e - 2) \sum_{t=1}^n \eta_t \text{Var } Z_t \end{aligned}$$

where  $C$  is defined in (13) and the second inequality follows from the second bound of Lemma 4. We now denote by  $T$  the first time step  $t$  when  $V_t > E^2/4$ . Using that  $\eta_t \leq 1/E$  for all  $t$  and  $V_T \leq E^2/2$ , we get

$$\sum_{t=1}^n \eta_t \text{Var } Z_t \leq \frac{E}{2} + \sum_{t=T+1}^n \eta_t \text{Var } Z_t . \quad (14)$$

We bound the last sum using  $\eta_t \leq C \sqrt{(\ln N)/V_{t-1}}$  for  $t \geq T + 1$  (note that, for  $t \geq T + 1$ ,  $V_{t-1} \geq V_T > E^2/4 > 0$ ). This yields

$$\sum_{t=T+1}^n \eta_t \text{Var } Z_t \leq C \sqrt{\ln N} \sum_{t=T+1}^n \frac{V_t - V_{t-1}}{\sqrt{V_{t-1}}} .$$

Since  $V_t \leq V_{t-1} + E^2/4$  and  $V_{t-1} \geq E^2/4$  for  $t \geq T + 1$ , we have

$$\begin{aligned} \frac{V_t - V_{t-1}}{\sqrt{V_{t-1}}} &= \frac{\sqrt{V_t} + \sqrt{V_{t-1}}}{\sqrt{V_{t-1}}} (\sqrt{V_t} - \sqrt{V_{t-1}}) \\ &\leq (\sqrt{2} + 1) (\sqrt{V_t} - \sqrt{V_{t-1}}) = \frac{\sqrt{V_t} - \sqrt{V_{t-1}}}{\sqrt{2} - 1}. \end{aligned}$$

Therefore, by a telescoping argument,

$$\begin{aligned} \sum_{t=T+1}^n \eta_t \text{Var } Z_t &\leq \frac{C\sqrt{\ln N}}{\sqrt{2} - 1} (\sqrt{V_n} - \sqrt{V_T}) \\ &\leq \frac{C}{\sqrt{2} - 1} \sqrt{V_n \ln N}. \end{aligned} \quad (15)$$

Putting things together, we have already proved that

$$\begin{aligned} \hat{X}_n - X_n^* &\geq -2 \max \left\{ E \ln N, (1/C) \sqrt{V_n \ln N} \right\} \\ &\quad - \frac{e-2}{2} E - \frac{C(e-2)}{\sqrt{2}-1} \sqrt{V_n \ln N}. \end{aligned}$$

In the case when  $\sqrt{V_n} \geq CE\sqrt{\ln N}$ , the regret  $\hat{X}_n - X_n^*$  is bounded from below by

$$- \left( \frac{2}{C} + \frac{C(e-2)}{\sqrt{2}-1} \right) \sqrt{V_n \ln N} - \frac{e-2}{2} E \geq -4\sqrt{V_n \ln N} - E/2,$$

where we substituted the value of  $C$  and obtained a constant for the leading term equal to  $2\sqrt{2(e-2)}/\sqrt{\sqrt{2}-1} \leq 3.75$ . When  $\sqrt{V_n} \leq CE\sqrt{\ln N}$ , the lower bound is more than

$$\begin{aligned} -2E \ln N - \frac{C(e-2)}{\sqrt{2}-1} \sqrt{V_n \ln N} - \frac{e-2}{2} E \\ \geq -2E \ln N - 2\sqrt{V_n \ln N} - E/2. \end{aligned}$$

This concludes the proof.  $\square$

#### 4.2. UNKNOWN BOUND ON THE PAYOFF RANGES (E)

We present the adaptation needed when no bound on the real-valued payoff range is known beforehand. For any sequence of payoff vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots$  and for all  $t = 1, 2, \dots$ , we define, similarly to Section 3.2, a quantity that keeps track of the payoff ranges seen so far. More precisely,  $E_t = 2^k$ , where  $k \in \mathbb{Z}$  is the smallest integer such that

$\max_{s=1,\dots,t} \max_{i,j=1,\dots,N} |x_{i,s} - x_{j,s}| \leq 2^k$ . Now let the sequence  $\eta_2, \eta_3, \dots$  be defined as

$$\eta_t = \min \left\{ \frac{1}{E_{t-1}}, C \sqrt{\frac{\ln N}{V_{t-1}}} \right\} \quad (16)$$

for  $t \geq 2$ , with  $C = \sqrt{2(\sqrt{2} - 1)/(e - 2)}$ .

We are now ready to state and prove the main result of this section, which bounds the regret in terms of the variance of the predictions. We show in the next section how this bound leads to more intrinsic bounds on the regret.

*Theorem 6.* Consider the weighted majority forecaster using the time varying learning rate (16). Then, for all sequences of payoffs and for all  $n \geq 1$ ,

$$\hat{X}_n - X_n^* \geq -4\sqrt{V_n \ln N} - 4E \ln N - 6E$$

where  $E = \max_{t=1,\dots,n} \max_{i,j=1,\dots,N} |x_{i,t} - x_{j,t}|$ .

*Proof.* The proof is similar to the one of Theorem 5, we only have to deal with the estimation of the payoff ranges. We apply again Lemma 3,

$$\begin{aligned} \hat{X}_n - X_n^* &\geq -\left(\frac{2}{\eta_{n+1}} - \frac{1}{\eta_1}\right) \ln N - \sum_{t=1}^n \Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) \\ &\geq -2 \max \left\{ E_n \ln N, (1/C) \sqrt{V_n \ln N} \right\} - \sum_{t=1}^n \Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) \\ &= -2 \max \left\{ E_n \ln N, (1/C) \sqrt{V_n \ln N} \right\} \\ &\quad - \sum_{t \in \mathcal{T}} \Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) - \sum_{t \notin \mathcal{T}} \Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) \end{aligned}$$

where  $C$  is defined in (16), and  $\mathcal{T}$  is the set of time steps  $t \geq 2$  when  $E_t = E_{t-1}$  (note that  $1 \notin \mathcal{T}$  by definition). Thus  $\mathcal{T}$  is a finite union of intervals of integers,  $\mathcal{T} = \llbracket 1, n \rrbracket \setminus \{t_1, \dots, t_R\}$ , where we denote  $t_1 = 1$  and let  $t_2, \dots, t_R$  be the time rounds  $t \geq 2$  such that  $E_t \neq E_{t-1}$ .

Using the second bound of Lemma 4 on  $t \in \mathcal{T}$  (since, for  $t \in \mathcal{T}$ ,  $\eta_t E_t \leq E_t/E_{t-1} = 1$ ) and the first bound of Lemma 4 on  $t \notin \mathcal{T}$ , which in this case reads  $\Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) \leq E_t$ , we get

$$\begin{aligned} \hat{X}_n - X_n^* &\geq -2 \max \left\{ E_n \ln N, (1/C) \sqrt{V_n \ln N} \right\} \\ &\quad - (e - 2) \sum_{t \in \mathcal{T}} \eta_t \text{Var } Z_t - \sum_{t \notin \mathcal{T}} E_t. \end{aligned} \quad (17)$$

We consider the  $r$ -th regime,  $r = 1, \dots, R$ , that is, the time steps  $s$  between  $t_r + 1$  and  $t_{r+1} - 1$  (with  $t_{R+1} = n$  by convention whenever  $t_R <$

$n$ ). For all these time steps  $s$ ,  $E_s = E_{t_r}$ . We use the same arguments that led to (14) and (15): denote by  $T_r$  the first time step  $s \geq t_r + 1$  when  $V_s > E_{t_r}^2/4$ . Then,

$$\sum_{s=t_r+1}^{t_{r+1}-1} \eta_t \text{Var } Z_t \leq \frac{E_{t_r}}{2} + \frac{C\sqrt{\ln N}}{\sqrt{2}-1} \left( \sqrt{V_{t_{r+1}-1}} - \sqrt{V_{T_r}} \right) .$$

Summing over  $r = 1, \dots, R$  and noting that a telescoping argument is given by  $V_{t_r} \leq V_{T_r}$ ,

$$\sum_{t \in T} \eta_t \text{Var } Z_t \leq \frac{C\sqrt{\ln N}}{\sqrt{2}-1} \sqrt{V_n} + \frac{1}{2} \sum_{r=1}^R E_{t_r} .$$

We deal with the last sum (also present in (17)) by noting that

$$\sum_{t \notin T} E_t = \sum_{r=1}^R E_{t_r} \leq \sum_{r=-\infty}^{\lceil \log_2 E \rceil} 2^r \leq 2^{1+\lceil \log_2 E \rceil} \leq 4E .$$

Putting things together,

$$\begin{aligned} \hat{X}_n - X_n^* &\geq -2 \max \left\{ E_n \ln N, (1/C) \sqrt{V_n \ln N} \right\} \\ &\quad - \frac{(e-2)C\sqrt{\ln N}}{\sqrt{2}-1} \sqrt{V_n} - 2e E . \end{aligned}$$

The proof is concluded, as the previous one, by noting that  $E_n \leq 2E$ .  $\square$

#### 4.3. RANDOMIZED PREDICTION AND ACTUAL REGRET

In this paper, the focus is on improved bounds for the expected regret. After choosing a probability distribution  $\mathbf{p}_t$  on the actions, the forecaster gets  $\hat{x}_t = x_{1,t}p_{1,t} + \dots + x_{N,t}p_{N,t}$  as a reward. In case randomized prediction is considered, after choosing  $\mathbf{p}_t$ , the forecaster draws an action  $I_t$  at random according to  $\mathbf{p}_t$  and gets the reward  $x_{I_t,t}$ , whose conditional expectation is  $\hat{x}_t$ . In this version of the game of prediction, the aim is now to minimize the (actual) regret, defined as the difference between  $x_{I_{1,1}} + \dots + x_{I_{n,n}}$  and  $X_n^*$ .

Bernstein's inequality for martingales (see, e.g., Freedman, 1975) shows however that the actual regret of any forecaster is bounded by the expected regret with probability  $1 - \delta$  up to deviations of the order of  $\sqrt{V_n \ln(n/\delta)} + M \ln(n/\delta)$ . These deviations are of the same order of magnitude as the bound of Theorem 6. Unless we are able to apply a sharper concentration result than Bernstein's inequality, no

further refinement of the above bounds is worthwhile. In particular, in view of the deviations from the expectations, as far as actual regret is concerned, we may prefer the results of Section 4 to those of Section 3. The next section, as well as Section 5, explain how bounds in terms of  $\sqrt{V_n}$  lead to many interesting bounds on the regret that do not depend on quantities related to the forecaster's rewards.

#### 4.4. BOUNDS ON THE FORECASTER'S CUMULATIVE VARIANCE

In this section we show a first way to deal with the dependency of the bound on  $V_n$ , the forecaster's cumulative variance. Section 5 will illustrate this further.

Recall that  $Z_t$  is the random variable which takes the value  $x_{i,t}$  with probability  $p_{i,t}$ , for  $i = 1, \dots, N$ . The main term of the bound stated in Theorem 6 contains  $V_n = \text{Var } Z_1 + \dots + \text{Var } Z_n$ . Note that  $V_n$  is therefore smaller than all quantities of the form

$$\sum_{t=1}^n \sum_{i=1}^N p_{i,t} (x_{i,t} - \mu_t)^2$$

where  $(\mu_t)_{t \geq 1}$  is any sequence of real numbers which may be chosen in *hindsight*, as it is not required for the definition of the forecaster. (The minimal value of the expression is obtained for  $\mu_t = \hat{x}_t$ .) This gives us a whole family of upper bounds, and we may choose for the analysis the most convenient sequence of  $\mu_t$ .

To provide a concrete example, recall the definition (1) of payoff effective range  $E_t$  and consider the choice  $\mu_t = \min_{j=1, \dots, N} x_{j,t} + E_t/2$ .

*Corollary 1.* The regret of the weighted majority forecaster with variable learning rate (16) satisfies

$$\hat{X}_n - X_n^* \geq -2 \sqrt{(\ln N) \sum_{t=1}^n E_t^2 - 4E \ln N - 6E}$$

where  $E$  is a bound on the payoff ranges,  $E = \max_{t=1, \dots, n} E_t$ .

The bound proposed by Corollary 1 shows that for an effective range of  $E$ , say if the payoffs all fall in  $[0, E]$ , the regret is lower bounded by a quantity equal to  $-2E\sqrt{n \ln N}$  (a closer look at the proof of Theorem 6 shows that this constant factor is less than 1.9, and could be made as close to  $2\sqrt{(e-2)} = \sqrt{2}\sqrt{2(e-2)}$  as desired). The best leading constant for such bounds is, to our knowledge,  $\sqrt{2}$  (see Cesa-Bianchi and Lugosi, 2006). This shows that the improved dependence in the bound does not come at a significant increase in the magnitude of the

leading coefficient. When the actual ranges are small, these bounds give a considerable advantage. Such a situation arises, for instance, in the setting of on-line portfolio selection, when we use linear upper bound on the regrets (see, e.g., the EG strategy by Helmbold et al., 1998). Moreover, we note that Corollary 1 improves on a result of Allenberg-Neeman and Neeman (2004), who show a regret bound, in terms of the cumulative effective range, whose main term is  $5.7\sqrt{2M(\ln N)\sum_{t=1}^n E_t}$ , for a given bound  $M$  over the payoffs.

Finally, we note that using translations of payoffs for **prod**-type algorithms, as suggested by Section 5.1, may be worthwhile as well, see Corollary 4 below. However, unlike the approach presented here for the weighted majority based forecaster, there the payoffs have to be translated explicitly and on-line by the forecaster, and thus, each translation rule corresponds to a different forecaster.

#### 4.5. EXTENSION TO PROBLEMS WITH INCOMPLETE INFORMATION

An interesting issue is how the second-order bounds of this section extend to incomplete information problems. In the literature of this area, exponentially weighted averages of estimated cumulative payoffs play a key role (see, for instance, Auer et al., 2002 for the multiarmed bandit problem, Cesa-Bianchi, Lugosi, and Stoltz, 2005 for label-efficient prediction, and Piccolboni and Schindelhauer, 2001, Cesa-Bianchi, Lugosi, and Stoltz, 2004 for prediction under partial monitoring).

A careful analysis of the proofs therein shows that the order of magnitude of the bound on the regret is given by the root of the sum of the conditional variances of the estimates of the payoffs used for prediction,

$$\sqrt{(\ln N) \sum_{t=1}^n \mathbb{E}_t \left[ \sum_{i=1}^N p_{i,t} (\tilde{x}_{i,t})^2 - \left( \sum_{i=1}^N p_{i,t} \tilde{x}_{i,t} \right)^2 \right]}.$$

Here we denote by  $\tilde{x}_{i,t}$  the (unbiased) estimate available for  $x_{i,t}$  (whose form varies depending on the precise setup and the considered strategy), by  $\mathbf{p}_t = (p_{1,t}, \dots, p_{N,t})$  the probability distributions over the actions, and by  $\mathbb{E}_t$  the conditional expectation with respect to the information available up to round  $t$  (for instance, in multiarmed bandit problems, this information is the past payoffs). Note that the conditioning in  $\mathbb{E}_t$  determines the values of the payoffs  $\mathbf{x}_t = (x_{1,t}, \dots, x_{N,t})$  and of  $\mathbf{p}_t$ .

In setups with full monitoring, that is, for the setups considered in this paper, no estimation is needed,  $\tilde{x}_{i,t} = x_{i,t}$ , and the bound is exactly that of Theorem 6.

In multiarmed bandit problems (with payoffs in, say,  $[-M, M]$ ), the estimators are given by  $\tilde{x}_{i,t} = (x_{i,t}/p_{i,t})\mathbb{I}_{[I_t=i]}$  where  $I_t$  is the index of the chosen component of the payoff vector. Now,

$$\mathbb{E}_t [p_{i,t} \tilde{x}_{i,t}^2] = x_{i,t}^2 \leq M^2 . \quad (18)$$

Summing over  $i = 1, \dots, N$  and  $t = 1, \dots, n$  the bound  $M\sqrt{nN \ln N}$  of Auer et al. (2002) is recovered.

In label-efficient prediction problems,  $\tilde{x}_{i,t} = (x_{i,t}/\varepsilon)Z_t$ , where the  $Z_t$  are i.i.d. random variables distributed according to a Bernoulli distribution with parameter  $\varepsilon \sim m/n$ . Then,

$$\mathbb{E}_t [p_{i,t} \tilde{x}_{i,t}^2] = p_{i,t} \frac{x_{i,t}^2}{\varepsilon} \leq p_{i,t} \frac{M^2}{\varepsilon} .$$

Summing over  $i = 1, \dots, N$  and  $t = 1, \dots, n$  we recover the bound  $M\sqrt{(n/\varepsilon) \ln N} \sim Mn\sqrt{(\ln N)/m}$  of Cesa-Bianchi, Lugosi, and Stoltz (2005).

Finally, in games with partial monitoring, the quantity (18) is less than  $M^2 t^{-1/3} N^{2/3} (\ln N)^{-1/3}$ . Summing over  $i = 1, \dots, N$  and  $t = 1, \dots, n$  we recover the  $Mn^{2/3} N^{2/3} (\ln N)^{1/3}$  bound of Cesa-Bianchi, Lugosi, and Stoltz (2004).

In conclusion, the faster  $\sqrt{n}$  rate in bandit problems, as opposed to the  $n^{2/3}$  rate in problems of prediction under partial monitoring, is due to better statistical performances (i.e., smaller conditional variance) of the available estimators.

## 5. Using translations of the payoffs

We now consider the bounds derived from those of Sections 3 and 4 in the case when translations are performed on the payoffs (Section 5.1). We show that they lead to several improvements or extensions of earlier results (Section 5.2) and also relieve the forecaster from the need of any preliminary manipulation on the payoffs (Section 5.3).

### 5.1. ON-LINE TRANSLATIONS OF THE PAYOFFS

Note that any on-line forecasting strategy may be used by a meta-forecaster which, before applying the given strategy, may first translate the payoffs according to a prescribed rule that may depend on the past. More formally, the meta-forecaster runs the strategy with the payoffs  $r_{k,t} = x_{k,t} - \mu_t$ , where  $\mu_t$  is any quantity possibly based on the past payoffs  $x_{i,s}$ , for  $i = 1, \dots, N$  and  $s = 1, \dots, t$ .



The forecasting strategies of Section 4 (and the obtained bounds) are invariant by such translations. This is however not the case for the **prod**-type algorithms of Section 3. An interesting application is obtained in Section 5.2 by considering  $\mu_t = \hat{x}_t$  where we recall that  $\hat{x}_t = x_{1,t}p_{1,t} + \dots + x_{N,t}p_{N,t}$  is the forecaster's reward at time  $t$ . As the sums  $\mu_1 + \dots + \mu_n$  cancel out in the difference  $\hat{X}_n - X_{k,n}$ , we obtain the following corollary of Theorem 2. Note that the remainder term here is now expressed in terms of the effective ranges (1) of the payoffs.

*Corollary 2.* Given  $E > 0$ , for all  $n \geq 1$  and all sequences of payoffs with effective ranges  $E_t$  bounded by  $E$ , the cumulative reward of algorithm **prod-Q**( $E$ ) run using translated payoffs  $x_{k,t} - \hat{x}_t$  satisfies

$$\begin{aligned} \hat{X}_n \geq X_n^* &- 8\sqrt{(\ln N) \max_{s \leq n} R_s^*} \\ &- 2E \left( 1 + \log_4 n + 2(1 + \lfloor (\log_2 \ln N)/2 \rfloor) \ln N \right) . \end{aligned}$$

where the  $R_s^*$  are defined as follows. For  $1 \leq t \leq n$  and  $k = 1, \dots, N$ ,  $R_{k,t} = (x_{k,1} - \hat{x}_1)^2 + \dots + (x_{k,t} - \hat{x}_t)^2$  and  $R_t^* = R_{k_t^*,t}$ , where  $k_t^*$  is the index of the action achieving the best cumulative payoff at round  $t$  (ties are broken by choosing the action  $k$  with smallest associated  $R_{k,t}$ ).

*Remark 1.* In one-sided games, for instance in gain games, the forecaster has always an incentive to translate the payoffs by the minimal payoff  $\mu_t$  obtained at each round  $t$ ,

$$\mu_t = \min_{j=1,\dots,N} x_{j,t} .$$

This is since for all  $j$  and  $t$ ,  $(x_{j,t} - \mu_t)^2 \leq x_{j,t}^2$  in a gain game. The issue is not so clear however for signed games, and it may be a delicate issue to determine beforehand if the payoffs should be translated, and if so, which translation rule should be used. See also Section 4.4, as well as Section 5.2.

## 5.2. IMPROVEMENTS FOR SMALL OR LARGE PAYOFFS

As recalled in Section 2.2, when all payoffs have the same sign Freund and Schapire (1997) first showed that Littlestone and Warmuth's weighted majority algorithm (1994) can be used to construct a forecasting strategy achieving a regret of order  $\sqrt{M|X_n^*| \ln N} + M \ln N$ , where  $N$  is the number of actions,  $M$  is a known upper bound on the magnitude of payoffs ( $|x_{i,t}| \leq M$  for all  $t$  and  $i$ ), and  $|X_n^*|$  is the absolute value of the cumulative payoff of the best action (i.e., the

largest cumulative payoff in a gain game or the smallest cumulative loss in a loss game), see also Auer, Cesa-Bianchi, and Gentile (2002).

This bound is good when  $|X_n^*|$  is small in the one-sided game; that is, when the best action has a small gain (in a gain game) or a small loss (in a loss game). However, one often expects the best expert to be effective (for instance, because we have many experts and at least one of them is accurate). An effective expert in a loss game suffers a small cumulative loss, but in a gain game, such an expert should get a large cumulative payoff  $X_n^*$ . To obtain a bound that is good when  $|X_n^*|$  is large one could apply the translation  $x'_{i,t} = x_{i,t} - M$  (from gains to losses) or the translation  $x'_{i,t} = x_{i,t} + M$  (from losses to gains). In both cases one would obtain a bound of the form  $\sqrt{M(Mn - |X_n^*|) \ln N}$ , which is now suited for effective experts in gain games and poor experts in loss games, but not for effective experts in loss games and poor experts in gain games. Since the original bound is not stable under the operation of conversion from one type of one-sided game into the other, the forecaster has to guess whether to play the original game or its translated version, depending on his beliefs on the quality of the experts and on the nature of the game (losses or gains).

In Corollary 4 we use the sharper bound of Corollary 2 to prove a (first-order) bound of the form

$$\sqrt{M \min\{|X_n^*|, Mn - |X_n^*|\} \ln N}.$$

This is indeed an improvement for small losses or large gains, though it requires knowledge of  $M$ . However, in Remark 2 we will indicate how to extend this result to the case when  $M$  is not known beforehand. Note that the (second-order) bound of Corollary 3 also yields the same result without any preliminary knowledge of  $M$ .

We thus recover an earlier result by Allenberg-Neeman and Neeman (2004). They proved, in a gain game, for a related algorithm, and with the previous knowledge of a bound  $M$  on the payoffs, a bound whose main term is  $11.4\sqrt{M \min\{\sqrt{X_n^*}, \sqrt{Mn - X_n^*}\}}$ . That algorithm was specifically designed to ensure a regret bound of this form, and is different from the algorithm whose performance we discussed before the statement of Corollary 1, whereas we obtain the improvements for small losses or large gains as corollaries of much more general bounds that have other consequences.

### 5.2.1. Analysis for exponentially weighted forecasters

The main drawback of  $V_n$ , used in Theorem 6, is that it is defined directly in terms of the forecaster's distributions  $\mathbf{p}_t$ . We now show how this dependence could be removed.

*Corollary 3.* Consider the weighted majority forecaster run with the time-varying learning rate (16). Then, for all sequences of payoffs in a one-sided game (i.e., payoffs are all nonpositive or all nonnegative),

$$\hat{X}_n \geq X_n^* - 4\sqrt{|X_n^*| \left(M - \frac{|X_n^*|}{n}\right) \ln N} - 39 M \max\{1, \ln N\}$$

where  $M = \max_{t=1,\dots,n} \max_{i=1,\dots,N} |x_{i,t}|$ .

*Proof.* We give the proof for a gain game. Since the payoffs are in  $[0, M]$ , we can write

$$\begin{aligned} V_n &\leq \sum_{t=1}^n \left( M \sum_{i=1}^N p_{i,t} x_{i,t} - \left( \sum_{i=1}^N p_{i,t} x_{i,t} \right)^2 \right) = \sum_{t=1}^n (M - \hat{x}_t) \hat{x}_t \\ &\leq n \left( \frac{M \hat{X}_n}{n} - \left( \frac{\hat{X}_n}{n} \right)^2 \right) = \hat{X}_n \left( M - \frac{\hat{X}_n}{n} \right) \end{aligned}$$

where we used the concavity of  $x \mapsto Mx - x^2$ . Assume that  $\hat{X}_n \leq X_n^*$  (otherwise the result is trivial). Then, Theorem 6 ensures that

$$\hat{X}_n - X_n^* \geq -4\sqrt{X_n^* \left(M - \frac{\hat{X}_n}{n}\right) \ln N} - \kappa$$

where  $\kappa = 4M \ln N + 6M$ . We solve for  $\hat{X}_n$  obtaining

$$\hat{X}_n - X_n^* \geq -4\sqrt{X_n^* \left(M - \frac{X_n^*}{n} + \frac{\kappa}{n}\right) \ln N} - \kappa - 16\frac{X_n^*}{n} \ln N.$$

Using the crude upper bound  $X_n^*/n \leq M$  and performing some simple algebra, we get the desired result.  $\square$

Similarly to the remark about constant factors in Section 4.4 the factor 4 in Corollary 3 can be made as close as desired to  $4\sqrt{e-2} = 2\sqrt{2}\sqrt{e-2}$ , which is not much larger than the best known leading constant for improvements for small losses,  $2\sqrt{2}$ , see Auer, Cesa-Bianchi, and Gentile (2002). But here, we have in addition an improvement for large losses, and deal with unknown ranges  $M$ . (Note, similarly to the discussion in Section 4.4, the presence of the same small factor  $\sqrt{2(e-2)} \approx 1.2$ .)

### 5.2.2. Analysis for prod-type forecasters

Quite surprisingly, a bound of the same form as the one shown in Corollary 3 can be derived from Corollary 2.

*Corollary 4.* Given  $M > 0$ , for all  $n \geq 1$  and all sequences of payoffs bounded by  $M$ , i.e.,  $\max_{1 \leq i \leq N} \max_{1 \leq t \leq n} |x_{i,t}| \leq M$ , the cumulative reward of algorithm **prod-Q**( $2M$ ), run using translated payoffs  $x_{k,t} - \hat{x}_t$  in a one-sided game, is larger than

$$\begin{aligned} \hat{X}_n \geq X_n^* &- 8\sqrt{2M \min\{X_n^*, Mn - X_n^*\} \ln N} \\ &- 128 M \ln N - \kappa - 8\sqrt{2M(\ln N)\kappa} \end{aligned}$$

where

$$\begin{aligned} \kappa &= 4M \left( 1 + \log_4 n + 2(1 + \lfloor (\log_2 \ln N)/2 \rfloor) \ln N \right) \\ &= \Theta \left( M(\ln n) + M(\ln N)(\ln \ln N) \right). \end{aligned}$$

*Proof.* As in the proof of Corollary 3, it suffices to give the proof for a gain game. In fact, we apply below the bound of Corollary 2, which is invariant under the change  $\ell_{i,t} = M - x_{i,t}$  that converts bounded losses into bounded nonnegative payoffs.

The main term in the bound of Corollary 2, with the notations therein, involves

$$\max_{s \leq n} R_s^* \leq \min \left\{ M \left( X_n^* + \hat{X}_n \right), M \left( 2Mn - X_n^* - \hat{X}_n \right) \right\}. \quad (19)$$

Indeed, using that  $(a - b)^2 \leq a^2 + b^2$  for  $a, b \geq 0$ , we get on the one hand, for all  $1 \leq s \leq n$ ,

$$R_s^* \leq \sum_{t=1}^s x_{k_s^*, t}^2 + \hat{x}_s^2 \leq M \left( X_{k_s^*, s} + \hat{X}_s \right) \leq M \left( X_n^* + \hat{X}_n \right)$$

whereas on the other hand, the same techniques yield

$$\begin{aligned} R_s^* &= \sum_{t=1}^s \left( (M - x_{k_s^*, t}) - (M - \hat{x}_s^2) \right)^2 \\ &\leq M \left( (Ms - X_s^*) + (Ms - \hat{X}_s) \right). \end{aligned}$$

Now, we note that for all  $s$ ,  $X_{s+1}^* \leq X_s^* + M$ , and similarly,  $\hat{X}_{s+1} \leq \hat{X}_s + M$ . Thus we also have  $\max_{s \leq n} R_s^* \leq M(2Mn - X_n^* - \hat{X}_n)$ .

Corollary 2, combined with (19), yields

$$\hat{X}_n \geq \hat{X}_n^* - 8\sqrt{M(\ln N) \min \left\{ \left( X_n^* + \hat{X}_n \right), \left( 2Mn - X_n^* - \hat{X}_n \right) \right\}} - \kappa$$

where  $\kappa = 4M(1 + \log_4 n + 2(1 + \lfloor (\log_2 \ln N)/2 \rfloor) \ln N)$ . Without loss of generality, we may assume that  $\hat{X}_n \leq X_n^*$  and get

$$\hat{X}_n \geq \hat{X}_n^* - 8\sqrt{2M(\ln N) \min \left\{ X_n^*, (Mn - \hat{X}_n) \right\}} - \kappa.$$

Solving for  $\hat{X}_n$  and performing simple algebra in case the minimum is achieved by the term containing  $\hat{X}_n$  concludes the proof.  $\square$

*Remark 2.* The forecasting strategy of Theorem 4, when used by a meta-forecaster translating the payoffs by  $\hat{x}_t$ , achieves an improvement for small or large payoffs of the form

$$M\sqrt{\min \left\{ \max_{s \leq n} \frac{X_s^*}{M_s}, \max_{s \leq n} \frac{sM_s - X_s^*}{M_s} \right\}}$$

without previous knowledge of  $M$ .

### 5.2.3. The case of signed games

The proofs of Corollaries 3 and 4 reveal that the assumption of one-sidedness cannot be relaxed. However, we may also prove a version of the improvement for small losses or for large gains suited to signed games. Remember that, as explained in Section 2.3, a meta-forecaster may always convert a signed game into a one-sided game by performing a suitable translation on the payoffs, and then apply a strategy for one-sided games. Since Corollary 2 and Theorem 6 are stable under general translations, applying them to the payoffs  $x_{i,t}$  or to a translated version of them  $x'_{i,t}$  results in the same bounds. If the translated version  $x'_{i,t}$  correspond to a one-sided game, then the bounds of Corollaries 3 and 4 may be applied. Using  $x'_{i,t} = x_{i,t} - \min_{j=1,\dots,N} x_{j,t} \geq 0$  and  $x'_{i,t} = x_{i,t} - \max_{j=1,\dots,N} x_{j,t} \leq 0$  for the analysis, we may show, for instance, that for any signed game the forecaster of Theorem 6 ensures that the regret is bounded by a quantity whose main term is less than

$$\min \left\{ \sqrt{(\ln N) \max_{j=1,\dots,N} \left( \sum_{t=1}^n \left( x_{j,t} - \min_{i=1,\dots,N} x_{i,t} \right) \right)}, \right. \\ \left. \sqrt{(\ln N) \min_{j=1,\dots,N} \left( \sum_{t=1}^n \left( \max_{i=1,\dots,N} x_{i,t} - x_{j,t} \right) \right)} \right\}.$$

This bound is obtained without any previous knowledge of a bound  $M$  on the payoffs, and is sharper than both bounds (2) and (3). It may be interpreted as an improvement for small or large cumulative payoffs.

### 5.3. WHAT IS A “FUNDAMENTAL” BOUND?

Most of the known regret bounds are not stable under natural transformations of the payoffs, such as translations and rescalings.<sup>1</sup> If a regret bound is not stable, then a (meta-)prediction algorithm might be willing to manipulate the payoffs in order to achieve a better regret. However, in general it is hard to choose the payoff transformation that is best for a given and unknown sequence of payoffs. For this reason, we argue that regret bounds that are stable under payoff transformations are, in some sense, more fundamental than others. The bounds that we have derived in this paper are based on sums of squared payoffs. They are not only generally tighter than the previously known bounds, but also stable under different transformations, such as those described below (in what follows, we use  $x'_{i,t}$  to indicate a transformed payoff).

*Additive translations:*  $x'_{i,t} = x_{i,t} - \mu_t$ .

Note that the regret (of playing a fixed sequence  $\mathbf{p}_1, \mathbf{p}_2, \dots$ ) is not affected by this transformation. Hence, stable bounds should not change when payoffs are translated. As already explained in Section 5.2, translations can be used to turn a gain game into a loss game and vice versa.

The invariance by general translations is the hardest to obtain, and this paper is the first one to show tight translation-invariant bounds that depend on the specific sequence of payoffs rather than just on its length (see Corollary 2, Theorem 6 and some of their corollaries, e.g., Corollary 1). It is also important to remark that, in a stable bound, not only the leading term, but also the smaller order terms, have to be stable under translations. This is why the smaller order terms of Corollary 2 and Theorem 6 involve bounds on the payoff ranges  $x_{i,t} - x_{j,t}$  rather than just on the payoffs  $x_{i,t}$ .

*Rescalings:*  $x'_{i,t} = \alpha x_{i,t}$ ,  $\alpha > 0$ .

As this transformation causes the regret to be multiplied by a factor of  $\alpha$ , stable bounds should only change by the same factor  $\alpha$ . Obtaining bounds that are stable under rescalings is not always easy when the payoff ranges are not known beforehand, or when we try to get bounds sharper than the basic zero-order bounds discussed in Section 2.1. For instance, the application of a doubling trick on the magnitude of the

---

<sup>1</sup> Here we do not distinguish between stable bounds and stable algorithms because all the stability properties we consider for the bounds are due to a corresponding stability of the prediction scheme they are derived from. When a stable algorithm does not achieve a stable bound, it suffices to optimize the bound in hindsight, thanks to the stability properties of the prediction scheme.

payoffs, or even the use of more sophisticated incremental techniques, may lead to small but undesirable  $M \ln(Mn)$  terms, which behave badly upon rescalings. This was the case with the remainder term  $M \ln(1 + |X_n^*|)$  in Theorem 2.1 by Auer, Cesa-Bianchi, and Gentile (2002) where they assume knowledge of the payoff range but seek sharper bounds.

Note also that forecasters with scaling-invariant bounds should require no previous knowledge on the payoff sequence (such as the payoff range) as this information is scale-sensitive. This is why, for instance, the bounds of Theorems 2 and 5 cannot be considered scaling-invariant. However, modifications of these forecasters that increase their adaptiveness lead to Theorems 4 and 6. There we could derive scaling-invariant bounds by using forecasters based on updates which are defined in terms of quantities that already have this type of invariance.

Whereas translation-invariant bounds that are also sharp are generally hard to obtain, we feel that any bound can be made stable with respect to rescalings via a reasonably accurate analysis.

Unstable bounds can lead the meta-forecaster to Cornelian dilemmas. Consider for the instance the bound (4) by Allenberg-Neeman and Neeman (2004). If we use a meta-forecaster that translates payoffs by a quantity  $\mu_t$  (possibly depending on past observations), then the bound takes the form

$$\sqrt{M(\ln N) \max_{t=1,\dots,n} \sum_{s=1}^t |x_{k_t^*,s} - \mu_s|} + M \ln N .$$

Note that the choice  $\mu_t = -M$  (or  $\mu_t = \min_{j=1,\dots,N} x_{j,t}$ ) yields the improvement for small payoffs (2) and the choice  $\mu_t = M$  (or  $\mu_t = \max_{j=1,\dots,N} x_{j,t}$ ) yields the improvement for large payoffs (3). In general, the above bound is tight if, for a large number of rounds, all payoffs  $x_{j,t}$  at a given round  $t$  are close to a common value, and we may guess this value to choose  $\mu_t$  accordingly. In Section 5.2.3, on the other hand, we show that Corollaries 3 and 4 propose bounds that need no preliminary choices of  $\mu_t$  and are better than both (2) and (3).

## 6. Discussion and open problems

We have analyzed forecasting algorithms that work indifferently in loss games, gain games, and signed games. In Corollary 2 and Theorem 6 we have shown, for these forecasters, sharp regret bounds that are stable under rescalings and general translations. These bounds lead to improvements for small or large payoffs in one-sided games (Corollaries 3

and 4) and do not assume any preliminary information about the payoff sequence.<sup>2</sup>

A practical advantage of the weighted majority forecaster is that its update rule is completely incremental and never needs to reset the weights. This in contrast to the forecaster **prod-MQ** of Theorem 4 that uses a nested doubling trick. On the other hand, the bound proposed in Theorem 6 is not in closed form, as it still explicitly depends through  $V_n$  on the forecaster's rewards  $\hat{x}_t$ . We therefore need to solve for the regrets as we did, for instance, in Sections 4.4 and 5.2. Finally, it was also noted in Section 4.4 that the weighted majority forecaster update is invariant under translations of the payoffs. This is not the case for the **prod**-type forecasters, which need to perform translations explicitly. Though in general it may be difficult to determine beforehand what a good translation could be, Corollaries 2 and 4, as well as Remark 1, indicate some general effective translation rules.

Several issues are left open:

- Design and analyze incremental updates for the **prod**-type forecasters of Section 3.
- Obtain second-order bounds with updates that are not multiplicative; for instance, updates based on the polynomial potentials (see Cesa-Bianchi and Lugosi, 2003). These updates could be used as basic ingredients to derive forecasters achieving optimal orders of magnitude on the regret when applied to problems such as nonstochastic multiarmed bandits, label-efficient prediction, and partial monitoring. Note that, to the best of our knowledge, in the literature about incomplete information problems only exponentially weighted averages have been able to achieve these optimal rates (see Section 4.5 and the references therein).
- Extend the analysis of **prod**-type algorithms to obtain an oracle inequality of the form

$$\hat{X}_n \geq \max_{k=1,\dots,N} \left( X_{k,n} - \gamma_1 \sqrt{Q_{k,n} \ln N} \right) - \gamma_2 M \ln N$$

where  $\gamma_1$  and  $\gamma_2$  are absolute constants. Inequalities of this form can be viewed as game-theoretic versions of the model selection bounds in statistical learning theory.

---

<sup>2</sup> Whereas the bound of Theorem 6 is already stated this way, we recall that it is easy to modify the forecaster used to prove Corollary 2 in order to dispense with the need of any preliminary knowledge of a bound  $E$  on the payoff ranges.



## References

- C. Allenberg-Neeman and B. Neeman. Full information game with gains and losses. Algorithmic Learning Theory, 15th International Conference, ALT 2004, Padova, Italy, October 2004, Proceedings, volume 3244 of Lecture Notes in Artificial Intelligence, pages 264–278. Springer, 2004.
- P. Auer, N. Cesa-Bianchi, Y. Freund, and R.E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32:48–77, 2002.
- 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.
- N. Cesa-Bianchi, Y. Freund, D.P. Helmbold, D. Haussler, R. Schapire, and M.K. Warmuth. How to use expert advice. *Journal of the ACM*, 3:427–485, 1997.
- N. Cesa-Bianchi and G. Lugosi. Potential-based algorithms in on-line prediction and game theory. *Machine Learning*, 51:239–261, 2003.
- N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- N. Cesa-Bianchi, G. Lugosi, and G. Stoltz. Minimizing regret with label efficient prediction. *IEEE Transactions on Information Theory*, 51:2152–2162, 2005.
- N. Cesa-Bianchi, G. Lugosi, and G. Stoltz. Regret minimization under partial monitoring. Submitted for journal publication, 2004.
- D. A. Freedman. On tail probabilities for martingales. *The Annals of Probability*, 3:100–118, 1975.
- 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*, 55(1):119–139, 1997.
- D. P. Helmbold, R. E. Schapire, Y. Singer, and M. K. Warmuth. On-line portfolio selection using multiplicative updates. *Mathematical Finance*, 8:325–344, 1998.
- N. Littlestone and M.K. Warmuth. The weighted majority algorithm. *Information and Computation*, 108:212–261, 1994.
- A. Piccolboni and C. Schindelhauer. Discrete prediction games with arbitrary feedback and loss. In *Proceedings of the 14th Annual Conference on Computational Learning Theory*, pages 208–223, 2001.
- V.G. Vovk. A game of prediction with expert advice. *Journal of Computer and System Sciences*, 56(2):153–73, 1998.

## Appendix

### Proof of Theorem 4

We use some additional notation for the proof:  $(r, s) - 1$  denotes the epoch right before  $(r, s)$ ; that is,  $(r, s - 1)$  when  $s > 0$ , and  $(r - 1, S_{r-1} - S_{r-2})$  when  $s = 0$ . For notational convenience,  $t_{(0,0)-1}$  is conventionally set to 0.

*Proof.* The proof combines the techniques from Theorems 2 and 3. As in the proof of Theorem 3, we denote by  $(R, S_R - S_{R-1})$  the index of the last epoch and let  $t_{(R, S_R - S_{R-1})} = n$ .

We assume  $R \geq 1$  and  $S_R \geq 1$ . Otherwise, if  $R = 0$ , this means that  $M_t = M^{(0)}$  for all  $t \leq n - 1$ , and the strategy, and thus the proposed bound, reduces to the one of Theorem 2. The case  $S_R = 0$  is dealt with at the end of the proof. In particular,  $S_R \geq 1$  implies that some epoch ended at time  $t$  when  $Q_t^* > 4^{S_R-1} M_t^2$ . This implies that  $q \geq 4^{S_R-1} (\geq 1)$ , which in turn implies  $2^{S_R} \leq 2\sqrt{q}$  and  $S_R \leq 1 + (\log_2 q)/2$ .

Denote  $M^{(R+1)} = M_n$ . Note that at time  $n$  we have either  $M_n \leq M^{(R)}$ , implying  $M_n = M^{(R+1)} = M^{(R)}$ , or we have  $M_n > M^{(R)}$ , implying  $M_n = M^{(R+1)} = 2M^{(R)}$ . In both cases,  $M^{(R)} \leq M^{(R+1)} \leq 2M$ . Furthermore,  $M^{(s)} \geq 2^{s-r} M^{(r)}$  for each  $0 \leq r \leq s \leq R$ , and thus (11) holds for  $s \leq R$  with  $M_{t_r}$  replaced by  $M^{(r)}$ .

Similar to the proof of Theorem 2, for each epoch  $(r, s)$ , let

$$X_k^{(r,s)} = \sum_{t=t_{(r,s)}-1+1}^{t_{(r,s)}-1} x_{k,t}, \quad Q_k^{(r,s)} = \sum_{t=t_{(r,s)}-1+1}^{t_{(r,s)}-1} x_{k,t}^2, \quad \hat{X}^{(r,s)} = \sum_{t=t_{(r,s)}-1+1}^{t_{(r,s)}-1} \hat{x}_t$$

where the sums are over all the time steps  $t$  in epoch  $(r, s)$  except the last one,  $t_{(r,s)}$ . We also denote  $k_{(r,s)} = k_{t_{(r,s)}-1}^*$  the index of the best overall expert up to time  $t_{(r,s)} - 1$  (one time step before the end of epoch  $(r, s)$ ).

We upper bound the cumulative payoff of the best action as

$$X_n^* \leq \sum_{r=0}^R \left( M^{(r+1)} + (S_r - S_{r-1}) M^{(r)} + \sum_{s=0}^{S_r - S_{r-1}} X_{k_{(r,s)}}^{(r,s)} \right) \quad (20)$$

by using the same argument by induction as in (10). More precisely, we write, for each  $(s, r)$ ,

$$\begin{aligned} X_{k_{(r,s)}, t_{(r,s)}-1} &= X_{k_{(r,s)}}^{(r,s)} + M_{t_{(r,s)}-1} + X_{k_{(r,s)}-1, t_{(r,s)}-1-1} \\ &\leq X_{k_{(r,s)}}^{(r,s)} + M_{t_{(r,s)}-1} + X_{k_{(r,s)}-1, t_{(r,s)}-1-1}. \end{aligned}$$

We note that  $M_{t_{(r,s)}-1} = M^{(r)}$  whenever  $0 \leq s < S_r - S_{r-1}$  and  $M_{t_{(r,s)}-1} = M^{(r+1)}$  otherwise. This and

$$X_n^* \leq X_{n-1}^* + M^{(R+1)} = X_{k_{(R,S_R)}, t_{(R,S_R)}-1} + M^{(R+1)}$$

show (20) by induction.

Let

$$\kappa = \sum_{r=0}^R \left( M^{(r+1)} + (S_r - S_{r-1}) M^{(r)} \right).$$

To show a bound on  $\kappa$  note that (11) implies

$$\sum_{r=0}^R M^{(r+1)} \leq 2M^{(R)} + M^{(R+1)} \leq 3M^{(R+1)} \leq 6M \quad (21)$$

and

$$\sum_{r=0}^R (S_r - S_{r-1}) M^{(r)} \leq 2MS_R \leq M(2 + \log_2 q) .$$

Thus,  $\kappa \leq (8 + \log_2 q)M$ .

Now, similarly to the above bound on  $X_n^*$ ,

$$\hat{X}_n \geq -\kappa + \sum_{r=0}^R \sum_{s=0}^{S_r - S_{r-1}} \hat{X}^{(r,s)}$$

so that the regret  $\hat{X}_n - X_n^*$  is larger than

$$\hat{X}_n - X_n^* \geq -2\kappa + \sum_{r=0}^R \sum_{s=0}^{S_r - S_{r-1}} \left( \hat{X}^{(r,s)} - X_{k(r,s)}^{(r,s)} \right) .$$

Now note that each time step  $t$  (but the last one) of epoch  $(r, s)$  satisfies  $M_t \leq M^{(r)}$  and  $\eta_{(r,s)} \leq 1/2M^{(r)}$ . Therefore, we can apply Lemma 2 to  $\hat{X}^{(r,s)} - X_{k(r,s)}^{(r,s)}$  for each epoch  $(r, s)$ . This gives

$$\hat{X}_n - X_n^* \geq -2\kappa - \sum_{r=0}^R \sum_{s=0}^{S_r - S_{r-1}} \left( \frac{\ln N}{\eta_{(r,s)}} + \eta_{(r,s)} Q_{k(r,s)}^{(r,s)} \right) .$$

By definition of the algorithm, for all epochs  $(r, s)$ ,

$$Q_{k(r,s)}^{(r,s)} \leq Q_{k(r,s), t(r,s)-1} = Q_{t(r,s)-1}^* \leq 4^{S_{r-1}+s} (M^{(r)})^2$$

and

$$\eta_{(r,s)} \leq \sqrt{\ln N} / \left( 2^{S_{r-1}+s} M^{(r)} \right) .$$

Therefore,

$$\begin{aligned} & \sum_{r=0}^R \sum_{s=0}^{S_r - S_{r-1}} \eta_{(r,s)} Q_{k(r,s)}^{(r,s)} \\ & \leq \sum_{r=0}^R \sum_{s=0}^{S_r - S_{r-1}} 2^{S_{r-1}+s} M^{(r)} \sqrt{\ln N} \\ & \leq \sum_{r=0}^R \sum_{s=1}^{S_r - S_{r-1}} 2^{S_{r-1}+s} (2M) \sqrt{\ln N} + \sum_{r=0}^R 2^{S_{r-1}} M^{(r)} \sqrt{\ln N} \end{aligned}$$

$$\begin{aligned}
&\leq (2M) \sum_{s=1}^{S_R} 2^s \sqrt{\ln N} + 2^{S_R} \sum_{r=0}^R M^{(r)} \sqrt{\ln N} \\
&\leq (2M) 2^{S_R+1} \sqrt{\ln N} + 2^{S_R} (4M) \sqrt{\ln N} \\
&\quad \text{(using (11) and } M^{(R)} \leq 2M) \\
&\leq (16M) \sqrt{q \ln N}
\end{aligned} \tag{22}$$

since  $q \geq 4^{S_R-1}$  implies  $2^{S_R} \leq 2\sqrt{q}$ .

We now turn our attention to the remaining sum

$$\sum_{r=0}^R \sum_{s=0}^{S_r-S_{r-1}} \frac{\ln N}{\eta_{(r,s)}}.$$

By definition of the algorithm,

$$\eta_{(r,s)} = \begin{cases} 1/(2M^{(r)}) & \text{if } S_{r-1} + s \leq \lceil (\log_2 \ln N)/2 \rceil \\ \sqrt{\ln N} / (2^{S_{r-1}+s} M^{(r)}) & \text{otherwise.} \end{cases}$$

We denote by  $(r^*, s^*)$  the last couple  $(r, s)$  for which  $\eta_{r,s} = 1/(2M^{(r)})$ . With obvious notation, a crude overapproximation leads to

$$\begin{aligned}
&\sum_{r=0}^R \sum_{s=0}^{S_r-S_{r-1}} \frac{\ln N}{\eta_{(r,s)}} \\
&\leq \sum_{(r,s) \leq (r^*, s^*)} 2M^{(r)} \ln N + \sum_{r=0}^R \sum_{s=0}^{S_r-S_{r-1}} 2^{S_{r-1}+s} M^{(r)} \sqrt{\ln N}.
\end{aligned}$$

We already have the upper bound  $(16M)\sqrt{q \ln N}$  for the second sum. For the first one, we write

$$\begin{aligned}
&\sum_{(r,s) \leq (r^*, s^*)} 2M^{(r)} \ln N \\
&= \sum_{r=0}^{r^*} 2M^{(r)} \ln N + \sum_{r=0}^{r^*-1} (S_r - S_{r-1}) (2M^{(r)}) \ln N \\
&\quad + s^* (2M^{(r^*)}) \ln N \\
&\leq \sum_{r=0}^R 2M^{(r)} \ln N + (S_{r^*-1} + s^*) (4M) \ln N \\
&\leq 2M(\ln N) (3 + 2\lceil (\log_2 \ln N)/2 \rceil)
\end{aligned}$$

where we used (21). The proof is concluded in the case  $S_R \geq 1$  by putting things together and performing some overapproximation.

When  $S_R = 0$ ,  $q = 1$ ,  $\kappa$  is simply less than  $6M$ , (22) is less than  $8M\sqrt{\ln N}$ , so that the bound holds as well in this case.  $\square$

### Proof of Lemma 3

We first note that Jensen's inequality implies that  $\Phi$  is nonnegative.

The proof below is a simple modification of an argument first proposed in Auer, Cesa-Bianchi, and Gentile (2002). Note that we consider real-valued (non necessarily nonnegative) payoffs in what follows. For  $t = 1, \dots, n$ , we rewrite  $p_{i,t} = w_{i,t}/W_t$ , where  $w_{i,t} = e^{\eta_t X_{i,t-1}}$  and  $W_t = \sum_{j=1}^N w_{j,t}$  (the payoffs  $X_{i,0}$  are understood to equal 0, and thus,  $\eta_1$  may be any positive number satisfying  $\eta_1 \geq \eta_2$ ). Use  $w'_{i,t} = e^{\eta_{t-1} X_{i,t-1}}$  to denote the weight  $w_{i,t}$  where the parameter  $\eta_t$  is replaced by  $\eta_{t-1}$ . The associated normalization factor will be denoted by  $W'_t = \sum_{j=1}^N w'_{j,t}$ . Finally, we use  $j_t^*$  to denote the expert with the largest cumulative payoff after the first  $t$  rounds (ties are broken by choosing the expert with smallest index). That is,  $X_{j_t^*,t} = \max_{i \leq N} X_{i,t}$ . We also make use of the following technical lemma.

*Lemma 5.* (Auer, Cesa-Bianchi, and Gentile, 2002) For all  $N \geq 2$ , for all  $\beta \geq \alpha \geq 0$ , and for all  $d_1, \dots, d_N \geq 0$  such that  $\sum_{i=1}^N e^{-\alpha d_i} \geq 1$ ,

$$\ln \frac{\sum_{i=1}^N e^{-\alpha d_i}}{\sum_{j=1}^N e^{-\beta d_j}} \leq \frac{\beta - \alpha}{\alpha} \ln N .$$

*Proof (of Lemma 5).* We begin by writing

$$\begin{aligned} \ln \frac{\sum_{i=1}^N e^{-\alpha d_i}}{\sum_{j=1}^N e^{-\beta d_j}} &= \ln \frac{\sum_{i=1}^N e^{-\alpha d_i}}{\sum_{j=1}^N e^{(\alpha-\beta)d_j} e^{-\alpha d_j}} \\ &= -\ln \mathbb{E} \left[ e^{(\alpha-\beta)D} \right] \\ &\leq (\beta - \alpha) \mathbb{E} [D] \end{aligned}$$

where we applied Jensen inequality to the random variable  $D$  taking value  $d_i$  with probability  $e^{-\alpha d_i} / \sum_{j=1}^N e^{-\alpha d_j}$  for each  $j = 1, \dots, N$ . Since  $D$  takes at most  $N$  distinct values, its entropy  $H(D)$  is at most  $\ln N$ . Therefore

$$\begin{aligned} \ln N \geq H(D) &= \frac{\sum_{i=1}^N e^{-\alpha d_i}}{\sum_{j=1}^N e^{-\beta d_j}} \left( \alpha d_i + \ln \sum_{j=1}^N e^{-\beta d_j} \right) \\ &= \alpha \mathbb{E} [D] + \ln \sum_{j=1}^N e^{-\beta d_j} \geq \alpha \mathbb{E} [D] \end{aligned}$$

where the last inequality holds since  $\sum_{i=1}^N e^{-\alpha d_i} \geq 1$ . Hence  $\mathbb{E}[D] \leq (\ln N)/\alpha$ . As  $\beta > \alpha$  by hypothesis, we can plug the bound on  $\mathbb{E}[D]$  in the upper bound above and conclude the proof.  $\square$

*Proof of Lemma 3.* As it is usual in the analysis of the exponentially weighted average predictor, we study the evolution of  $\ln(W_{t+1}/W_t)$ . However, here we need to couple this term with  $\ln(w_{j_{t-1}^*,t}/w_{j_t^*,t+1})$  including in both terms the time-varying parameters  $\eta_t, \eta_{t+1}$ . Tracking the currently best expert  $j_t^*$  is used to lower bound the weight  $\ln(w_{j_t^*,t+1}/W_{t+1})$ . In fact, the weight of the overall best expert (after  $n$  rounds) could get arbitrarily small during the prediction process. We thus obtain the following

$$\begin{aligned} & \frac{1}{\eta_t} \ln \frac{w_{j_{t-1}^*,t}}{W_t} - \frac{1}{\eta_{t+1}} \ln \frac{w_{j_t^*,t+1}}{W_{t+1}} \\ &= \left( \frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \ln \frac{W_{t+1}}{w_{j_t^*,t+1}} + \frac{1}{\eta_t} \ln \frac{w'_{j_t^*,t+1}/W'_{t+1}}{w_{j_t^*,t+1}/W_{t+1}} + \frac{1}{\eta_t} \ln \frac{w_{j_{t-1}^*,t}/W_t}{w'_{j_t^*,t+1}/W'_{t+1}} \\ &= (A) + (B) + (C) . \end{aligned}$$

We now bound separately the three terms on the right-hand side. The term (A) is easily bounded by using  $\eta_{t+1} \leq \eta_t$  and using the fact that  $j_t^*$  is the index of the expert with largest payoff after the first  $t$  rounds. Therefore,  $w_{j_t^*,t+1}/W_{t+1}$  must be at least  $1/N$ . Thus we have

$$(A) = \left( \frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \ln \frac{W_{t+1}}{w_{j_t^*,t+1}} \leq \left( \frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \ln N .$$

We proceed to bounding the term (B) as follows

$$\begin{aligned} (B) &= \frac{1}{\eta_t} \ln \frac{w'_{j_t^*,t+1}/W'_{t+1}}{w_{j_t^*,t+1}/W_{t+1}} = \frac{1}{\eta_t} \ln \frac{\sum_{i=1}^N e^{-\eta_{t+1}(X_{j_t^*,t} - X_{i,t})}}{\sum_{j=1}^N e^{-\eta_t(X_{j_t^*,t} - X_{j,t})}} \\ &\leq \frac{\eta_t - \eta_{t+1}}{\eta_t \eta_{t+1}} \ln N = \left( \frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \ln N \end{aligned}$$

where the inequality is proven by applying Lemma 5 with  $d_i = X_{j_t^*,t} - X_{i,t}$ . Note that  $d_i \geq 0$  since  $j_t^*$  is the index of the expert with largest payoff after the first  $t$  rounds and  $\sum_{i=1}^N e^{-\eta_{t+1}d_i} \geq 1$  as for  $i = j_t^*$  we have  $d_i = 0$ .

The term (C) is first split as follows,

$$(C) = \frac{1}{\eta_t} \ln \frac{w_{j_{t-1}^*,t}/W_t}{w'_{j_t^*,t+1}/W'_{t+1}} = \frac{1}{\eta_t} \ln \frac{w_{j_{t-1}^*,t}}{w'_{j_t^*,t+1}} + \frac{1}{\eta_t} \ln \frac{W'_{t+1}}{W_t} .$$

We bound separately each one of the two terms on the right-hand side. For the first one, we have

$$\frac{1}{\eta_t} \ln \frac{w_{j_{t-1}^*, t}^*}{w_{j_t^*, t+1}^*} = \frac{1}{\eta_t} \ln \frac{e^{\eta_t X_{j_{t-1}^*, t-1}^*}}{e^{\eta_t X_{j_t^*, t}^*}} = X_{j_{t-1}^*, t-1}^* - X_{j_t^*, t}^* .$$

The second term is handled by using the very definition of  $\Phi$ ,

$$\begin{aligned} \frac{1}{\eta_t} \ln \frac{W'_{t+1}}{W_t} &= \frac{1}{\eta_t} \ln \frac{\sum_{i=1}^N w_{i,t} e^{\eta_t x_{i,t}}}{W_t} = \frac{1}{\eta_t} \ln \sum_{i=1}^N p_{i,t} e^{\eta_t x_{i,t}} \\ &= \sum_{i=1}^N p_{i,t} x_{i,t} + \Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) . \end{aligned}$$

Finally, we plug back in the main equation the bounds on the first two terms (A) and (B), and the bounds on the two parts of the term (C). After rearranging we obtain

$$\begin{aligned} 0 \leq & \left( X_{j_{t-1}^*, t-1}^* - X_{j_t^*, t}^* \right) + \sum_{i=1}^N p_{i,t} x_{i,t} + \Phi(\mathbf{p}_t, \eta_t, \mathbf{x}_t) \\ & - \frac{1}{\eta_{t+1}} \ln \frac{w_{j_t^*, t+1}^*}{W_{t+1}} + \frac{1}{\eta_t} \ln \frac{w_{j_{t-1}^*, t}^*}{W_t} \\ & + 2 \left( \frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \ln N . \end{aligned}$$

We apply the above inequalities to each  $t = 1, \dots, n$  and sum up using

$$\begin{aligned} \sum_{t=1}^n X_{j_{t-1}^*, t-1}^* - X_{j_t^*, t}^* &= - \max_{j=1, \dots, N} X_{j,n} \\ \text{and } \sum_{t=1}^n \left( -\frac{1}{\eta_{t+1}} \ln \frac{w_{j_t^*, t+1}^*}{W_{t+1}} + \frac{1}{\eta_t} \ln \frac{w_{j_{t-1}^*, t}^*}{W_t} \right) &\leq -\frac{1}{\eta_1} \ln \frac{w_{j_0^*, 1}^*}{W_1} = \frac{\ln N}{\eta_1} \end{aligned}$$

to conclude the proof.  $\square$