
Minimizing Wide Range Regret with Time Selection Functions

Subhash Khot and Ashok Kumar Ponnuswami

New York University and Georgia Tech

khot@cs.nyu.edu and pashok@cc.gatech.edu

Abstract

We consider the problem of minimizing regret with respect to a given set \mathcal{S} of pairs of time selection functions and modifications rules. We give an online algorithm that has $O(\sqrt{T \log |\mathcal{S}|})$ regret with respect to \mathcal{S} when the algorithm is run for T time steps and there are N actions allowed. This improves the upper bound of $O(\sqrt{TN \log(|\mathcal{I}||\mathcal{F}|)})$ given by Blum and Mansour [BM07a] for the case when $\mathcal{S} = \mathcal{I} \times \mathcal{F}$ for a set \mathcal{I} of time selection functions and a set \mathcal{F} of modification rules. We do so by giving a simple reduction that uses an online algorithm for external regret as a black box.

1 Introduction

We consider the following online optimization problem. At the beginning of each day (a time step), we have to choose one of the N allowed actions. Instead of picking one action deterministically, we may come up with a distribution over the actions. At the end of the day, an adversary, with the knowledge of the distribution we picked, fixes a loss for each action. We give a concrete example from Cesa-Bianchi et al. [CBFH⁺97]. Suppose we want to predict the probability that it rains on a day based on the predictions of N weather forecasting websites. But we don't know which of these "experts" give good forecasts. We come up with some weights on the websites using an online algorithm and use the weighted prediction as our guess for the probability of raining. At the end of the day, based on whether or not it rained, everyone incurs a loss depending on how inaccurate their prediction was. Usually it is assumed that the loss for each action is picked from a fixed interval, like $[0, 1]$. For example, we could charge a person who predicts p as the probability of rain $1 - p$ if it rains and p if it does not. After T days, we compare the loss incurred by the online algorithm we used to the loss incurred if we had followed a simple strategy (like just picking the same action each day). Our goal is to minimize our *regret* for not following one of the simple strategies. One may also compare the algorithm's performance to the performance if the distribution over actions at each time step were modified using a certain set of rules. We consider the problem of designing algorithms with low regret with respect to a given set of strategies or modification rules.

The most basic regret studied is *external regret*, which is the difference between the loss incurred by the algorithm and the loss incurred by the best action in hindsight. Another kind of regret commonly studied is called *internal regret*. This was introduced by Foster and Vohra [FV98]. Here, we consider the set of modification rules where for each pair (a, b) of actions we have a rule of the kind: Every time the algorithm suggests picking a , pick b instead. The internal regret of the algorithm is the regret of not having applied one of these modification rules. Each rule here can be considered as a function $f_{a,b}$ that maps every action to itself, except action a which gets mapped to b . If we consider the set of modification rules corresponding to all functions mapping the set of actions into itself, we get the notion of *swap regret*. Finally, we can allow any subset of these mappings as the set of allowed modification rules which gives the notion of *wide range regret*. This was defined by Lehrer [Leh03]. Lehrer also associates *time selection* function with each rule that indicates whether a rule is "active" at a given time or not. A related model is that of "sleeping experts" or "specialists" defined in Freund et al. [FSSW97]. Here, at the beginning of time t , each specialist can decide whether or not the current situation is her area of speciality and make a prediction only if it does. In addition, Blum and Mansour [BM07a] consider the case where the experts can be "partially awake". One way to interpret the activeness function is that it measures degree of confidence that the corresponding rule will perform well at a given time. In this case, we weigh the loss incurred by the algorithm and the modified action with the time selection function to calculate the regret.

The first algorithm with external regret sublinear in T was developed by Hannan [Han57]. An algorithm whose external regret has only logarithmic dependence on N was given by Littlestone and Warmuth [LW94] and Cesa-Bianchi et al. [CBFH⁺97].

Lemma 1 ([CBFH⁺97]) *There exists an online algorithm with external regret at most $O(\sqrt{T \log N})$ when the losses are picked from $[-1, +1]$. The running time is polynomial in T and N .*

The number of time steps T for which it will be run need not be provided as an input to the above algorithm. Stoltz and Lugosi [SL05] give a general method to convert any "weighted average predictor" algorithm for external regret to a low internal or swap regret algorithm. At a high level, they pretend there is an expert for each modification rule

who always suggests using that rule. At each time step, the expert is charged the loss that would be incurred if his modification rule were actually used. The weighted average predictor would give a distribution over the experts. The distribution over the actual actions is found by computing the fixed point of the expected modification rule picked from the distribution over the experts. This gives algorithms with $O(\sqrt{T \log N})$ internal regret and $O(\sqrt{TN \log N})$ swap regret. Our approach for wide range regret with time selection functions is based on the same idea. A drawback of a swap regret algorithm constructed this way is that it needs to maintain N^N weights. Blum and Mansour [BM07a] give an algorithm that has $O(\sqrt{TN \log N})$ swap regret and runs in time polynomial in N too. They also give an algorithm that has $O(\sqrt{TN \log(KM)})$ regret with respect to K modification rules and M time selection functions. Here, for each modification rule and time selection function, the regret of not having modified the algorithm's action by the rule with the losses weighed by the time selection function is considered. In this case, we can think of there being M people who are interested in following an algorithm's predictions. They have varying degrees of importance associated with each day (given by their corresponding time selection function) and want to minimize regret with respect to all the modification rules. The algorithm's goal is to minimize the maximum regret of a person. This is a bit different from the model considered in Lehrer [Leh03]. But with some effort, one can check that the result of Blum and Mansour [BM07a] can be generalized to the model of Lehrer [Leh03]. We refer the reader to [BM07b] for other bounds on the regret minimization and the relation of various kinds of regret to equilibriums in games.

The paper is organized as follows. In the next section, we define the model we work with formally and state our main result. We state the ideas we use from related results in Section 3. We prove our main result of an improved upper bound for wide range regret in Section 4. We conclude with a "first-order" upper bound in Section 5.

2 Our Model and Result

Let the set of actions be $[N] = \{1, 2, \dots, N\}$. Consider the following T round game between an online algorithm H and an adversary. At the beginning of time $t = 1, 2, \dots, T$, the algorithm picks a probability vector¹² $\mathbf{p}^t = (p_1^t, p_2^t, \dots, p_N^t)$. The adversary then picks the loss vector $\mathbf{l}^t = (l_1^t, l_2^t, \dots, l_N^t)$ for time t . The entries of \mathbf{l}^t are picked from a fixed interval. In this paper, we assume the losses are either picked from $[0, 1]$ or from $[-1, +1]$.

Define the regret of H with respect to action $a \in [N]$ to be

$$R_{H,a} = \sum_{t=1}^T \left(\sum_{b \in [N]} p_b^t l_b^t - l_a^t \right) = \sum_{t=1}^T \sum_{b \in [N]} p_b^t (l_b^t - l_a^t).$$

This can be interpreted as the difference between the expected loss of H and the loss of action a . Define the *external*

¹A probability vector is a vector in which the entries are non-negative and sum to 1.

²All vectors we consider are column vectors. We will use \top to denote the transpose.

regret of H to be

$$R_{H,ext} = \max_{a \in [N]} R_{H,a}.$$

We now define the model with time selection functions from Blum and Mansour [BM07a]. A time selection function is a function $I : \mathbb{N} \rightarrow [0, 1]$. Let \mathcal{I} be the set of time selection functions. At the beginning of time t , the adversary sets the values of $I(t)$ for each $I \in \mathcal{I}$. The algorithm then picks \mathbf{p}^t after which the adversary now picks \mathbf{l}^t as before. Given a modification rule $f : [N] \rightarrow [N]$, define \mathbf{M}_f to be the matrix with a 1 in column $f(i)$ of row i for all i and zeros everywhere else. Define the regret of H with respect to time selection function I and a modification rule f to be

$$\begin{aligned} R_{H,I,f} &= \sum_t I(t) \sum_{a \in [N]} p_a^t (l_a^t - l_{f(a)}^t) \\ &= \sum_t I(t) (\mathbf{p}^t \cdot \mathbf{l}^t - \mathbf{p}^{t\top} \mathbf{M}_f \mathbf{l}^t). \end{aligned}$$

Informally, we first weigh all the losses at time t by $I(t)$, the significance attached to time t . Then we look at the difference between the expected loss of H and the expected loss if the output of H were modified every time by applying f . That is, we measure the regret of not having played action $f(a)$ every time we played a . Given a set \mathcal{S} of pairs (I, f) , where I is a time selection function and f is a modification rule, the *wide range regret* of H with respect to \mathcal{S} is defined as

$$R_{H,\mathcal{S}} = \max_{(I,f) \in \mathcal{S}} R_{H,I,f}.$$

Let $\mathbb{1} : \mathbb{N} \rightarrow [0, 1]$ be the function that always outputs 1, i.e., $\mathbb{1}(t) = 1$. For simplicity of notation, we will use f to also denote the pair $(\mathbb{1}, f)$ when we are not concerned with time selection functions, in which case we assume that the adversary always sets $\mathbb{1}(t)$ to 1. It is easy to check that external regret is the same $R_{H,\mathcal{F}_{ext}}$ where $\mathcal{F}_{ext} = \{f_a\}_{a \in [N]}$ and $\forall b \in [N] : f_a(b) = a$. The *internal regret* of H is defined to be $R_{H,\mathcal{F}_{int}}$, where $\mathcal{F}_{int} = \{f_{a,b}\}_{a,b \in [N]}$ and $f_{a,b}(a) = b$ while $f_{a,b}(c) = c$ for $c \neq a$. The *swap regret* of H is defined to be $R_{H,\mathcal{F}_{swap}}$, where \mathcal{F}_{swap} is the set of all functions $f : [N] \rightarrow [N]$.

We prove the following theorem for minimizing wide range regret.

Theorem 2 *There exists an online algorithm H that for any given set \mathcal{S} satisfies*

- $R_{H,\mathcal{S}} = O(\sqrt{T \log |\mathcal{S}|})$ when the losses are picked from the $[0, 1]$.
- The running time of H is polynomial in T , N and $|\mathcal{S}|$.

Note that this matches (upto a constant) the results for external, internal and swap regret if we are not concerned with time selection functions. A drawback of our approach is that if the size of the set \mathcal{S} is large, the running time is high. For example, for swap regret with time selection functions, we may need time polynomial in T and N^N . But for this case, the result of Blum and Mansour already gives a more efficient algorithm with the same regret (upto a constant).

3 Previous Results

We use ideas from Stoltz and Lugosi [SL05] and Blum and Mansour [BM07a].

We first describe the approach of Stoltz and Lugosi [SL05] for internal regret. The idea is to simulate a low external regret algorithm for $N(N-1)$ imaginary experts. Start with any “weighted average predictor” H_{ext} with low external regret. There are $N(N-1)$ imaginary experts, one for each modification rule $f_{a,b}$. The expert corresponding to $f_{a,b}$ always suggests playing b instead of a . We will specify how the probability weights over the actual actions are calculated from the output of H_{ext} and how the losses are generated for the imaginary experts of H_{ext} .

At time t , suppose H_{ext} outputs probability $q_{a,b}^t$ for the expert corresponding to $f_{a,b}$. Then compute the probability vector $\mathbf{p}^t = (p_1^t, p_2^t, \dots, p_N^t)$ on the actual actions as a fixed point of

$$\mathbf{p}^t = \sum_{a,b \in [N]} q_{a,b}^t \mathbf{p}_{a \rightarrow b}^t,$$

where $\mathbf{p}_{a \rightarrow b}^t$ denotes the probability vector obtained from \mathbf{p}^t by changing the weight of action a to zero at putting it on action b . This can also be expressed as

$$\mathbf{p}^{t\top} = \sum_{a,b} q_{a,b}^t \mathbf{p}^{t\top} \mathbf{M}_{f_{a,b}} = \mathbf{p}^{t\top} \sum_{a,b} q_{a,b}^t \mathbf{M}_{f_{a,b}}.$$

Let the adversary return back \mathbf{l}^t as the loss vector at time t . The loss incurred at time t by each of the imaginary experts for $f_{a,b}$ is calculated as

$$l_{f_{a,b}}^t = \mathbf{l}^t \cdot \mathbf{p}_{i \rightarrow j}^t = \mathbf{p}^{t\top} \mathbf{M}_{f_{a,b}} \mathbf{l}^t.$$

This quantity can be thought of as the loss incurred if we followed the expert’s suggestion of playing b instead of a . Stoltz and Lugosi [SL05] showed that this achieves low internal regret. For an arbitrary set of modification rules \mathcal{F} , we have an expert for each modification rule $f \in \mathcal{F}$ and the probability and loss vectors are now calculated as

$$\mathbf{p}^t = \mathbf{p}^{t\top} \sum_{f \in \mathcal{F}} q_f^t \mathbf{M}_f$$

and

$$l_f^t = \mathbf{p}^{t\top} \mathbf{M}_f \mathbf{l}^t.$$

We now discuss the ideas we use from Blum and Mansour [BM07a]. We start with the case where $\mathcal{S} = \mathcal{I} \times \mathcal{F}_{ext}$ for some \mathcal{I} . In this case, there is an expert for each $(I, f_a) \in \mathcal{S}$. There is a weight $w_{I,a}^t$ associated with this expert at the end of time t where

$$w_{I,a}^t = \beta^{-\tilde{R}_{I,a}^t}$$

and

$$\tilde{R}_{I,a}^t = \sum_{t'=1}^t I(t') (\beta l_H^{t'} - l_a^{t'})$$

for some parameter $\beta \in (0, 1)$. Above, l_H^t is the actual loss incurred at time t . The quantity $\tilde{R}_{I,a}^t$ is called a “less-strict” external regret. The probability p_a^t associated with the action a at time t is then proportional to $\sum_{I \in \mathcal{I}} I(t) w_{I,a}^t$. By optimizing for the parameter β , Blum and Mansour [BM07a]

show that this achieves a low external regret with respect to all time selection functions.

To generalize this idea for wide range regret, where $\mathcal{S} = \mathcal{I} \times \mathcal{F}$, they introduce an expert for each $a \in [N]$, $I \in \mathcal{I}$ and $f \in \mathcal{F}$. There is a weight $w_{a,I,f}^t$ for each such expert. Note that this does not simplify to the reduction in the previous paragraph for the case when $\mathcal{F} = \mathcal{F}_{ext}$. Instead, in the next section we obtain a reduction where there are experts only for each $(I, f) \in \mathcal{S}$. Intuitively, this is where we remove the polynomial dependence of wide range regret on N and obtain a slightly simpler reduction.

4 A Reduction from Wide Range Regret to External Regret

We will prove Theorem 2 in this section. We first give an algorithm that when given a low external regret algorithm as a black box uses it to guarantee low wide range regret.

Theorem 3 *Given an algorithm H_{ext} with external regret $R(T, N)$ when the losses are from $[-1, +1]$, one can construct an algorithm H that when given losses from $[0, 1]$ satisfies:*

- $R_{H,\mathcal{S}} = R(T, |\mathcal{S}|)$
- The running time of H is polynomial in the running time of H_{ext} , T , N , and $|\mathcal{S}|$.

Idea: H will basically simulate an instance of H_{ext} with the elements of \mathcal{S} being the actions. Figure 1 shows the inputs and outputs of H and H_{ext} at time t . At time t , H_{ext} produces some $q_{I,f}^t$ for each $(I, f) \in \mathcal{S}$, where the $q_{I,f}^t$ form a probability distribution over \mathcal{S} . H will then use this to come up with a probability vector $\mathbf{p}^t = (p_1^t, p_2^t, \dots, p_N^t)$ on the actual actions. H will basically pick a random (I, f) with probability proportional to $I(t)q_{I,f}^t$. After this, it picks a vector \mathbf{p}^t over the actual actions such that \mathbf{p}^t is a fixed point of such a random f , i.e., modifying \mathbf{p}^t by f in expectation just yields \mathbf{p}^t . Intuitively, the loss passed to the black box H_{ext} for (I, f) is such that $q_{I,f}^t$ measures the regret with respect to time selection function I of not having modified the output of H using function f . Multiplying this by $I(t)$ takes care of the relevance of (I, f) at time t . Basically, the algorithm makes sure that if the regret with respect to (I, f) was large so far, then that regret doesn’t increase at the current step.

Proof: We first specify how H computes \mathbf{p}^t and \mathbf{l}^t at time t . To compute \mathbf{p}^t , get \mathbf{q}^t from H_{ext} . If $\sum_{(I,f) \in \mathcal{S}} I(t)q_{I,f}^t = 0$, then output any probability vector \mathbf{p} . Otherwise define \mathbf{p}^t to be any vector satisfying

$$\mathbf{p}^{t\top} = \mathbf{p}^{t\top} \left(\frac{\sum_{(I,f) \in \mathcal{S}} I(t)q_{I,f}^t \mathbf{M}_f}{\sum_{(I,f) \in \mathcal{S}} I(t)q_{I,f}^t} \right). \quad (1)$$

This is well defined since $\sum_{(I,f) \in \mathcal{S}} I(t)q_{I,f}^t \neq 0$. Such a vector \mathbf{p}^t exists since every row of

$$\frac{\sum_{(I,f) \in \mathcal{S}} I(t)q_{I,f}^t \mathbf{M}_f}{\sum_{(I,f) \in \mathcal{S}} I(t)q_{I,f}^t} \quad (2)$$

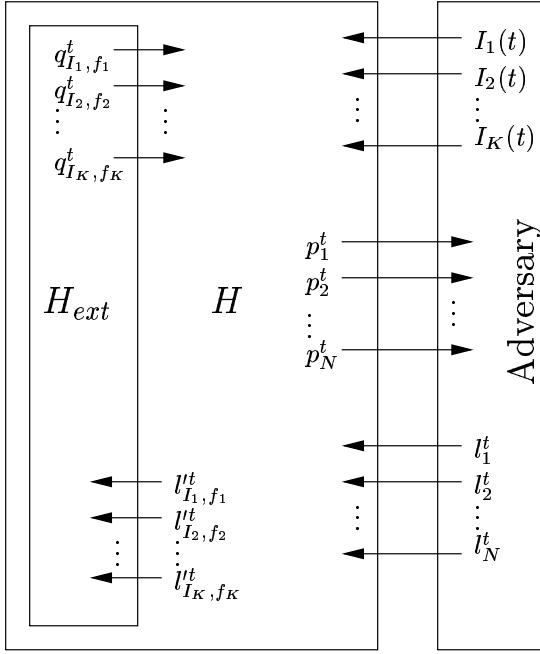


Figure 1: The reduction from wide range to external regret.

is a probability vector because \mathbf{M}_f has exactly one 1 in each row. That is, (2) defines the transition matrix of a Markov chain. When H gets back loss vector \mathbf{l}^t , it computes

$$l_{I,f}^t = I(t) \sum_{a \in N} p_a^t (l_{f(a)}^t - l_a^t) = I(t) \mathbf{p}^{t\top} (\mathbf{M}_f - \mathbf{I}) \mathbf{l}^t$$

where \mathbf{I} is the identity matrix. This yields

$$\sum_t l_{I,f}^t = \sum_t I(t) \mathbf{p}^{t\top} (\mathbf{M}_f - \mathbf{I}) \mathbf{l}^t = -R_{H,I,f}. \quad (3)$$

That is, $l_{I,f}^t$ is exactly the decrease at time t of the regret with respect to (I, f) . It is easy to check that $l_{I,f}^t \in [-1, +1]$.

From the low external regret guarantee of H_{ext} , for all $(I, f) \in \mathcal{S}$:

$$\sum_t \sum_{(J,g) \in \mathcal{S}} q_{J,g}^t l_{J,g}^t \leq \sum_t l_{I,f}^t + R(T, |\mathcal{S}|). \quad (4)$$

We will next show that

$$\sum_{(J,g) \in \mathcal{S}} q_{J,g}^t l_{J,g}^t = 0. \quad (5)$$

Together with (3) and (4), this will show that for all $(I, f) \in \mathcal{S}$,

$$0 \leq -R_{H,I,f} + R(T, |\mathcal{S}|),$$

or $R_{H,I,f} \leq R(T, |\mathcal{S}|)$ which proves the theorem.

We now proceed to prove (5).

$$\begin{aligned} \sum_{(J,g) \in \mathcal{S}} q_{J,g}^t l_{J,g}^t &= \sum_{(J,g)} q_{J,g}^t J(t) \mathbf{p}^{t\top} (\mathbf{M}_g - \mathbf{I}) \mathbf{l}^t \\ &= \sum_{(J,g)} q_{J,g}^t J(t) \mathbf{p}^{t\top} \mathbf{M}_g \mathbf{l}^t - \sum_{(J,g)} q_{J,g}^t J(t) \mathbf{p}^{t\top} \mathbf{l}^t \\ &= \mathbf{p}^{t\top} \left(\sum_{(J,g)} J(t) q_{J,g}^t \mathbf{M}_g \right) \mathbf{l}^t - \left(\sum_{(J,g)} q_{J,g}^t J(t) \right) (\mathbf{p}^{t\top} \mathbf{l}^t). \end{aligned}$$

(Case 1:) Suppose $\sum_{(J,g)} J(t) q_{J,g}^t \neq 0$. In this case we can use (1) to get

$$\begin{aligned} \sum_{(J,g) \in \mathcal{S}} q_{J,g}^t l_{J,g}^t &= \left(\sum_{(J,g)} q_{J,g}^t J(t) \right) (\mathbf{p}^{t\top} \mathbf{l}^t) \\ &\quad - \left(\sum_{(J,g)} q_{J,g}^t J(t) \right) (\mathbf{p}^{t\top} \mathbf{l}^t) \\ &= 0. \end{aligned}$$

(Case 2:) Assume $\sum_{(J,g)} J(t) q_{J,g}^t = 0$. Then $J(t) q_{J,g}^t = 0$ for all pairs (J, g) since $J(t)$ and $q_{J,g}^t$ are all non-negative, which implies

$$\sum_{(J,g) \in \mathcal{S}} q_{J,g}^t l_{J,g}^t = 0. \quad \blacksquare$$

It can be seen easily that Theorem 3 and Lemma 1 imply Theorem 2.

5 A First-Order Bound for Wide Range Regret

If we are only concerned with regret bounds as a function of T and N (called “zero-order” bounds in Cesa-Bianchi et al. [CBMS05]), Theorem 2 matches (up to a constant) the known upper bounds for external, internal and swap regret. One can also try to obtain “first-order” bounds, bounds that depend on the sum of payoffs of actions instead of the time. For example, Blum and Mansour [BM07a] show a $O(\sqrt{L_{min} \log(NM)} + \log(NM))$ upper bound for minimizing external regret with respect to a set \mathcal{I} of M time selection functions, where $L_{min} = \max_I \min_a L_{I,a}$ and $L_{I,a} = \sum_t I(t) l_a^t$. For the case when there is at least one “real” expert that does well most of the time, such a bound will be much tighter than a zero-order bound. One can hope to use external regret algorithms with good first-order bounds like the following to come up with good first-order bounds for wide range regret.

Lemma 4 (Cesa-Bianchi et al. [CBFH⁺97]) *There exists an algorithm with running time polynomial in T and N and external regret $O(\sqrt{L_{min} \log N} + \log N)$ when the losses are picked from $[0, 1]$.*

We need an algorithm that can handle losses from the interval $[-1, +1]$ in Theorem 3. One way to use the algorithm from Lemma 4 is to map the losses $l_{I,f}^t$ to the interval $[0, 1]$ by a linear transformation. But this also changes the loss of best action and makes the first order bound obtained very weak. Another alternative is to tinker with the quantity that $l_{I,f}^t$ signifies. If we are concerned only with modification rules (and not time selection functions), we can redefine l_f^t as

$$l_f^t = \sum_{a \in N} p_a^t l_{f(a)}^t.$$

But for technical reasons, this can’t be done if we are also working with time selection functions. Note that the only term in (4) that depends on I and f is $l_{I,f}^t$, and hence it must

capture all the terms that depend on *either* I or f in the definition of $R_{H,I,f}$. So we give a method based on the approach of Blum and Mansour [BM07a]. The main idea is to define a *reduced regret* for each pair (I, f) .

Theorem 5 *There exists an online algorithm that for any \mathcal{S} satisfies:*

- The wide range regret with respect to \mathcal{S} is at most $O(\sqrt{L_{min} \log |\mathcal{S}|} + \log |\mathcal{S}|)$, where

$$L_{min} = \max_I \min_{(I,f) \in \mathcal{S}} \sum_t I(t) \mathbf{p}^{t\top} \mathbf{M}_f \mathbf{l}^t.$$

- The running time is polynomial in T , N , and $|\mathcal{S}|$.

Proof: Define the loss of H with respect to I till time t as

$$L_{H,I}^t = \sum_{t'=1}^t I(t') \mathbf{p}^{t'} \cdot \mathbf{l}^{t'},$$

and the loss of H with respect to (I, f) till time t as

$$L_{H,I,f}^t = \sum_{t'=1}^t I(t') \mathbf{p}^{t'\top} \mathbf{M}_f \mathbf{l}^{t'}.$$

We assume that at any time t , not all $I(t)$ are zero. This is without loss of generality since in this case, the losses defined above don't change at time t . For some $\beta \in (0, 1)$ to be fixed later, we basically run an exponentially weighted predictor with a weight for each pair (I, f) . The weight of (I, f) at the end of time t is $w_{I,f}^t = \beta^{-\tilde{R}_{H,I,f}^t}$, where

$$\tilde{R}_{H,I,f}^t = \beta L_{H,I}^t - L_{H,I,f}^t.$$

That is, $\tilde{R}_{H,I,f}^t$ is a regret of H with respect to (I, f) where the incurred loss is reduced by a factor β . We define $q_{I,f}^t = w_{I,f}^{t-1} / W^{t-1}$, where $W^t = \sum_{(I,f) \in \mathcal{S}} w_{I,f}^t$ is the sum of the weights.

At time t , the algorithm does the following. It computes $q_{I,f}^t$ as above. The probability vector \mathbf{p}^t over the actual actions is picked as in (2). This is well defined since $w_{I,f}^t$ (and hence $q_{I,f}^t$) are all non-zero and at least one of the $I(t)$ is also non-zero (by assumption). Then the algorithm updates all the losses and weights when it gets back \mathbf{l}^t from the adversary. We first show that the sum of the weights can not increase at any time.

Claim 6

$$\forall t: \sum_{(I,f) \in \mathcal{S}} w_{I,f}^t \leq \sum_{(I,f) \in \mathcal{S}} w_{I,f}^{t-1}$$

Proof: We will use the fact that for any $\beta \in (0, 1)$ and $x \in [0, 1]$, $\beta^x \leq 1 - (1 - \beta)x$ and $\beta^{-x} \leq 1 + (1 - \beta)x/\beta$. This

gives

$$\begin{aligned} \sum_{(I,f) \in \mathcal{S}} w_{I,f}^t &= \sum_{(I,f)} w_{I,f}^{t-1} \beta^{I(t) (\mathbf{p}^{t\top} \mathbf{M}_f \mathbf{l}^t - \beta \mathbf{p}^t \cdot \mathbf{l}^t)} \\ &\leq \sum_{(I,f)} \left[w_{I,f}^{t-1} \left(1 - (1 - \beta) I(t) \mathbf{p}^{t\top} \mathbf{M}_f \mathbf{l}^t \right) \right. \\ &\quad \left. \times \left(1 + (1 - \beta) I(t) \mathbf{p}^t \cdot \mathbf{l}^t \right) \right] \\ &\leq \sum_{(I,f)} w_{I,f}^{t-1} - \left[(1 - \beta) W^{t-1} \sum_{(I,f)} q_{I,f}^t I(t) \mathbf{p}^{t\top} \mathbf{M}_f \mathbf{l}^t \right] \\ &\quad + \left[(1 - \beta) W^{t-1} \sum_{(I,f)} q_{I,f}^t I(t) \mathbf{p}^t \cdot \mathbf{l}^t \right] \\ &= \sum_{(I,f)} w_{I,f}^{t-1} - \left[(1 - \beta) W^{t-1} \mathbf{p}^{t\top} \left(\sum_{(I,f)} q_{I,f}^t I(t) \mathbf{M}_f \right) \mathbf{l}^t \right] \\ &\quad + \left[(1 - \beta) W^{t-1} \left(\sum_{(I,f)} q_{I,f}^t I(t) \right) (\mathbf{p}^t \cdot \mathbf{l}^t) \right] \\ &= \sum_{(I,f)} w_{I,f}^{t-1}. \end{aligned}$$

Above, the second inequality follows from the definition of $q_{I,f}^t$ and the last equality follows from (2). ■

We now get back to the proof of the theorem. The claim implies that for all $(I, f) \in \mathcal{S}$,

$$\beta^{-(\beta L_{H,I}^T - L_{H,I,f}^T)} = \beta^{-\tilde{R}_{H,I,f}^T} = w_{I,f}^T \leq \sum_{(J,g) \in \mathcal{S}} w_{J,g}^0 = |\mathcal{S}|$$

which gives

$$(\beta L_{H,I}^T - L_{H,I,f}^T) \log(1/\beta) \leq \log |\mathcal{S}|$$

or

$$L_{H,I} \leq \frac{L_{H,I,f} + \frac{\log |\mathcal{S}|}{\log(1/\beta)}}{\beta}.$$

Since for a given I , the statement is true for all f such that $(I, f) \in \mathcal{S}$, we can rewrite it as:

$$L_{H,I} \leq \frac{L_{H,I,min} + \frac{\log |\mathcal{S}|}{\log(1/\beta)}}{\beta}$$

where

$$L_{H,I,min} = \min_{f:(I,f) \in \mathcal{S}} L_{H,I,f}.$$

Setting β so that

$$\beta^{-1} = 1 + \min \left\{ \sqrt{\frac{\log |\mathcal{S}|}{L_{min}}}, \frac{1}{2} \right\}$$

gives the theorem. ■

Acknowledgments

The authors would like to thank Yishay Mansour for comments on an early draft of the paper.

References

- [BM07a] Avrim Blum and Yishay Mansour. From external to internal regret. *J. Mach. Learn. Res.*, 8:1307–1324, 2007.
- [BM07b] Avrim Blum and Yishay Mansour. Learning, regret minimization and equilibria. In Noam Nisan, Tim Roughgarden, Eva Tardos, and Vijay Vazirani, editors, *Algorithmic Game Theory*, chapter 4. Cambridge University Press, 2007.
- [CBFH⁺97] Nicolò Cesa-Bianchi, Yoav Freund, David Haussler, David P. Helmbold, Robert E. Schapire, and Manfred K. Warmuth. How to use expert advice. *J. ACM*, 44(3):427–485, 1997.
- [CBMS05] Nicolò Cesa-Bianchi, Yishay Mansour, and Gilles Stoltz. Improved second-order bounds for prediction with expert advice. In *COLT*, pages 217–232, 2005.
- [FSSW97] Yoav Freund, Robert E. Schapire, Yoram Singer, and Manfred K. Warmuth. Using and combining predictors that specialize. In *STOC '97: Proceedings of the twenty-ninth annual ACM symposium on Theory of computing*, pages 334–343, 1997.
- [FV98] Dean Foster and Rakesh V. Vohra. Asymptotic calibration. *Biometrika*, 85:379–390, 1998.
- [Han57] J. Hannan. Approximation to bayes risk in repeated plays. In M.Dresher, A. Tucker, and P. Wolfe, editors, *Contributions to the Theory of Games*, volume 3, pages 97–139. Princeton University Press, 1957.
- [Leh03] Ehud Lehrer. A wide range no-regret theorem. *Games and Economic Behavior*, 42(1):101–115, 2003.
- [LW94] Nick Littlestone and Manfred K. Warmuth. The weighted majority algorithm. *Inf. Comput.*, 108(2):212–261, 1994.
- [SL05] Gilles Stoltz and Gábor Lugosi. Internal regret in on-line portfolio selection. *Mach. Learn.*, 59(1-2):125–159, 2005.