

Near Minimax Optimal Players for the Finite-Time 3-Expert Prediction Problem

Authored by:

Victor Gabillon
Peter L. Bartlett
Yasin Abbasi

Abstract

We study minimax strategies for the online prediction problem with expert advice. It has been conjectured that a simple adversary strategy, called COMB, is near optimal in this game for any number of experts. Our results and new insights make progress in this direction by showing that, up to a small additive term, COMB is minimax optimal in the finite-time three expert problem. In addition, we provide for this setting a new near minimax optimal COMB-based learner. Prior to this work, in this problem, learners obtaining the optimal multiplicative constant in their regret rate were known only when $K=2$ or $K \rightarrow \infty$. We characterize, when $K=3$, the regret of the game scaling as $\sqrt{8/(9\pi)T} \pm \log(T)^2$ which gives for the first time the optimal constant in the leading (\sqrt{T}) term of the regret.

1 Paper Body

This paper studies the online prediction problem with expert advice. This is a fundamental problem of machine learning that has been studied for decades, going back at least to the work of Hannan [12] (see [4] for a survey). As it studies prediction under adversarial data the designed algorithms are known to be robust and are commonly used as building blocks of more complicated machine learning algorithms with numerous applications. Thus, elucidating the yet unknown optimal strategies has the potential to significantly improve the performance of these higher level algorithms, in addition to providing insight into a classic prediction problem. The problem is a repeated two-player zero-sum game between an adversary and a learner. At each of the T rounds, the adversary decides the quality/gain of K experts' advice, while simultaneously the learner decides to follow the advice of one of the experts. The objective of the adversary is to maximize the regret of the learner, defined as the difference between the total gain of the learner and the total gain of the best fixed

expert. Open Problems and our Main Results. Previously this game has been solved asymptotically as both T and K tend to ∞ : asymptotically the upper bound on the performance of the state-of-the-art Multiplicative Weights Algorithm (MWA) for the learner matches the optimal multiplicative γ constant of the asymptotic minimax optimal regret rate $(T/2) \log K$ [3]. However, for finite K , this asymptotic quantity actually overestimates the finite-time value of the game. Moreover, Gravin et al. [10] proved a matching lower bound $(T/2) \log K$ on the regret of the classic version of MWA, additionally showing that the optimal learner does not belong to an extended MWA family. Already, Cover [5] proved that the value of the game is of order of $T/(2\gamma)$ when $K = 2$, meaning that the regret of a MWA learner is 47% larger than the optimal learner in this case. Therefore the question of optimality remains open for non-asymptotic K which are the typical cases in applications, and therefore progress in this direction is important. In studying a related setting with $K = 3$, where T is sampled from a geometric distribution with parameter γ , Gravin et al. [9] conjectured that, for any K , a simple adversary strategy, called the C OMB adversary, is asymptotically optimal ($T \rightarrow \infty$, or when $\gamma \rightarrow 0$), and also excessively competitive for finite-time fixed T . The C OMB strategy sorts the experts based on their cumulative 31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

gains and, with probability one half, assigns gain one to each expert in an odd position and gain zero to each expert in an even position. With probability one half, the zeros and ones are swapped. The simplicity and elegance of this strategy, combined with its almost optimal performance makes it very appealing and calls for a more extensive study of its properties. Our results and new insights make progress in this direction by showing that, for any fixed T and up to small additive terms, C OMB is minimax optimal in the finite-time three expert problem. Additionally and with similar guarantees, we provide for this setting a new near minimax optimal C OMB-based learner. For $K = 3$, the regret of a MWA learner is 39% larger than our new optimal learner.² In this paper we also characterize, when $K = 3$, the regret of the game as $8/(9\gamma)T \log(T)$ which γ gives for the first time the optimal constant in the leading (T) term of the regret. Note that the state-of-the-art non-asymptotic lower bound in [15] on the value of this problem is non informative as the lower bound for the case of $K = 3$ is a negative quantity. Related Works and Challenges. For the case of $K = 3$, Gravin et al. [9] proved the exact minimax optimality of a C OMB-related adversary in the geometrical setting, i.e. where T is not fixed in advance but rather sampled from a geometric distribution with parameter γ . However the connection between the geometrical setting and the original finite-time setting is not well understood, even asymptotically (possibly due to the large variance of geometric distributions with small γ). Addressing this issue, in Section 7 of [8], Gravin et al. formulate the ‘Finite vs Geometric Regret’ conjecture which states that the value of the game in the geometrical setting, V_γ , and the value of the game in the finite-time setting, V_T , verify $V_T = \gamma V_\gamma$. We resolve here the conjecture for $K = 3$. Analyzing the finite-time expert problem raises new challenges compared to the geometric setting. In

the geometric setting, at any time (round) t of the game, the expected number of remaining rounds before the end of the game is constant (does not depend on the current time t). This simplifies the problem to the point that, when $K = 3$, there exists an exactly minimax optimal adversary that ignores the time t and the parameter γ . As noted in [9], and noticeable from solving exactly small instances of the game with a computer, in the finite-time case, the exact optimal adversary seems to depend in a complex manner on time and state. It is therefore natural to compromise for a simpler adversary that is optimal up to a small additive error term. Actually, based on the observation of the restricted computer-based solutions, the additive error term of C OMB seems to vanish with larger T . Tightly controlling the errors made by C OMB is a new challenge with respect to [9], where the solution to the optimality equations led directly to the exact optimal adversary. The existence of such equations in the geometric setting crucially relies on the fact that the value-to-go of a given policy in a given state does not depend on the current time t (because geometric distributions are memoryless). To control the errors in the finite-time setting, our new approach solves the game by backward induction showing the approximate greediness of C OMB with respect to itself (read Section 2.1 for an overview of our new proof techniques and their organization). We use a novel exchangeability property, new connections to random walks and a close relation that we develop between C OMB and a TWIN-C OMB strategy. Additional connections with new related optimal strategies and random walks are used to compute the value of the game (Theorem 2). We discuss in Section 6 how our new techniques have more potential to extend to an arbitrary number of arms, than those of [9]. Additionally, we show how the approximate greediness of C OMB with respect to itself is key to proving that a learner based directly on the C OMB adversary is itself quasi-minimax-optimal. This is the first work to extend to the approximate case, approaches used to designed exactly optimal players in related works. In [2] a probability matching learner is proven optimal under the assumption that the adversary is limited to a fixed cumulative loss for the best expert. In [14] and [1], the optimal learner relies on estimating the value-to-go of the game through rollouts of the optimal adversary's plays. The results in these papers were limited to games where the optimal adversary was only playing canonical unit vector while our result holds for general gain vectors. Note also that a probability matching learner is optimal in [9].

Notation: Let $[a : b] = \{a, a + 1, \dots, b\}$ with $a, b \in \mathbb{N}$, $a \leq b$, and $[a] = [1 : a]$. For a vector $w \in \mathbb{R}^n$, $n \in \mathbb{N}$, $\|w\|_1 = \sum_{k \in [n]} |w_k|$. A vector indexed by both a time t and a specific element index k is $w_{t,k}$. An undiscounted Markov Decision Process (MDP) [13, 16] M is a 4-tuple (S, A, r, p) . S is the state space, A is the set of actions, $r : S \times A \rightarrow \mathbb{R}$ is the reward function, and the transition model $p(\cdot | s, a)$ gives the probability distribution over the next state when action a is taken in state s . A state is denoted by s or s_t if it is taken at time t . An action is denoted by a or a_t .

2

2

The Game

We consider a game, composed of T rounds, between two players, called a learner and an adversary. At each time/round t the learner chooses an index $I_t \in [K]$ from a distribution p_t on the K arms. Simultaneously, the adversary assigns a binary gain to each of the arms/experts, possibly at random from a distribution A_t , and we denote the vector of these gains by $g_t \in \{0, 1\}^K$. The adversary and the learner then observe I_t and g_t . For simplicity we use the notation $g[t] = (g_s)_{s=1, \dots, t}$. The value of one realization of such a game is the cumulative regret defined as $\sum_{t=1}^T \sum_{k=1}^K g_t(k) I_t(k) - \max_{k \in [K]} \sum_{t=1}^T g_t(k)$.

A state $s \in S = (\{0, 1\}^K)^T$ is a K -dimensional vector such that the k -th element is the cumulative sum of gains dealt by the adversary on arm k before the current time t . Here the state does not include I_t but is typically denoted for a specific time t as s_t and computed as $s_t(k) = \sum_{\tau=1}^t g_\tau(k)$. This definition is motivated by the fact that there exist minimax strategies for both players that rely solely on the state and time information as opposed to the complete history of plays, $g[t] \in \{0, 1\}^{K \times t}$. In state s , the set of leading experts, i.e., those with maximum cumulative gain, is $X(s) = \{k \in [K] : s_k = \max_{s'} s'_k\}$.

We use π_t to denote the (possibly non-stationary) strategy/policy used by the adversary, i.e., for any input state s and time t it outputs the gain distribution $\pi_t(s, t)$ played by the adversary at time t in state s . Similarly we use p_t to denote the strategy of the learner. As the state depends only on the adversary plays, we can sample a state s at time t from π_t . The expected regret of the game, $V_{\pi, p}$. Given an adversary π and a learner p , $V_{\pi, p} = \mathbb{E}[\sum_{t=1}^T \sum_{k=1}^K g_t(k) I_t(k) - \max_{k \in [K]} \sum_{t=1}^T g_t(k)]$. The learner tries to minimize the expected regret while the adversary π tries to maximize it. The value of the game is the minimax value $V_T = \min_{\pi} \max_{p} V_{\pi, p} = \max_{\pi} \min_{p} V_{\pi, p}$.

In this work, we are interested in the search for optimal minimax strategies, which are adversary π^* , such that $V_T = \max_{\pi} V_{\pi, p^*}$, and learner strategies p^* such that $V_T = \min_{\pi} V_{\pi, p^*}$ and learner strategies p^* .

Summary of our Approach to Obtain the Near Greediness of C OMB

Most of our material is new. First, Section 3 recalls that Gravin et al. [9] have shown that the search for the optimal adversary π^* can be restricted to the finite family of balanced strategies (defined in the next section). When $K = 3$, the action space of a balanced adversary is limited to seven stochastic $\pi \in \{ \{1, 2\}, \{1, 3\}, \{2, 3\} \}$ (see Section 5.1 for their $C, V, 1$, actions (gain distributions), denoted by $B_3 = \{W \text{ description}\}$). The C OMB adversary repeats the gain distribution C at each time and in any state. In Section 4 we provide an explicit formulation of the problem as finding π^* inside an MDP with a specific reward function. Interestingly, we observe that another adversary, which we call T WIN π , has the same value as π^* (Section 5.1). C OMB and denote by π^* , which repeats the distribution W To control the errors

made by C OMB, the proof uses a novel and intriguing exchangeability property (Section 5.2). This exchangeability property holds thanks to the surprising role played by the T WIN ?, C OMB strategy. For any distributions $A? ? B? ? 3$ there exists a distribution $D? ?$, mixture of $C? ?$ and $W? ?$ and then $A? ?$ in terms of such that for almost all states, playing $A? ?$ and then $D? ?$ is the same as playing W the expected reward and the probabilities over the next states after these two steps. Using Bellman operators, this can be concisely written as: for any (value) function $f : S ?? R$, in (almost) any state s , we have that $[TA? ? [TD? ? f]](s) = [TW? ? [TA? ? f]](s)$. We solve the MDP with a backward induction in time from $t = T$. We show that playing $C? ?$ at time t is almost greedy with respect to playing $C? ?$ in later rounds $t? ? t$. The greedy error is defined as the difference of expected reward between always playing $C? ?$ and playing the best (greedy) ?rst action before playing C OMB. Bounding how these errors accumulate through the rounds relates the value of C OMB to the value of $C? ?$ (Lemma 16). To illustrate the main ideas, let us ?rst make two simplifying (but unrealistic) assumptions at time t : C OMB has been proven greedy w.r.t. itself in rounds $t? ? t$ and the exchangeability holds in all states. Then we would argue at time t that by the exchangeability property, instead of optimizing the greedy

$3 A? ? C? ? . . . C? ?$. Then action w.r.t. C OMB as $\max_{A? ?} B? ? 3 A? ? C? ? . . . C? ?$, we can study the optimizer of $\max_{A? ?} B? ? 3 W$ we use the induction property to conclude that $C? ?$ is the solution of the previous optimization problem. Unfortunately, the exchangeability property does not hold in one specific state denoted by $s? ?$. What saves us though is that we can directly compute the error of greedy?cation of any gain distribution with respect to C OMB in $s? ?$ and show that it diminishes exponentially fast as $T ? t$, the number of rounds remaining, increases (Lemma 7). This helps us to control how the errors accumulate during the induction. From one given state $st ?= s? ?$ at time t , ?rst, we use the exchangeability property once when trying to assess the ?quality? of an action $A? ?$ as a greedy action w.r.t. C OMB. This leads us to consider the quality of playing $A? ?$ in possibly several new states $\{st+1\}$ at time $t + 1$ reached following T WIN -C OMB in s . We use our exchangeability property repeatedly, starting from the state st until a subsequent state reaches $s? ?$, say at time $t? ?$, where we can substitute the exponentially decreasing greedy error computed at this time $t? ?$ in $s? ?$. Here the subsequent states are the states reached after having played T WIN -C OMB repetitively starting from the state st . If $s? ?$ is never reached we use the fact that C OMB is an optimal action everywhere else in the last round. The problem is then to determine at which time $t? ?$, starting from any state at time t and following a T WIN -C OMB strategy, we hit $s? ?$ for the ?rst time. This is translated into a classical gambler?s ruin problem, which concerns the hitting times of a simple random walk (Section 5.3). Similarly the value of the game is computed using the study of the expected number of equalizations of a simple random walk (Theorem 5.1).

3

Solving for the Adversary Directly

In this section, we recall the results from [9] that, for arbitrary K , permit

us to directly search for the minimax optimal adversary in the restricted set of balanced adversaries while ignoring the learner. Definition 1. A gain distribution $A?$ is balanced if there exists a constant $cA?$, the mean gain of $A?$, such that $\forall k \in [K], cA? = E[g_k | A?]$. A balanced adversary uses exclusively balanced gain distributions. Lemma 1 (Claim 5 in [9]). There exists a minimax optimal balanced adversary. Use B to denote the set of all balanced strategies and $B?$ to denote the set of all balanced gain distributions. Interestingly, as demonstrated in [9], a balanced adversary σ incurs the same regret $T \sum_{p \in B} Vp, \sigma$ on every learner: If $\sigma \in B$, then $\sum_{p \in B} VT \sum_{p \in B} R : \sigma, p = VT$. (See Lemma 10) Therefore, given an adversary strategy σ , we can define the value-to-go $Vt?0(s)$ associated with σ from time t_0 in state s , $\forall t \geq 0$ (s) =

$$E[VsT+1 | \sigma, s]$$

$$sT+1$$

$$T \sum_{p \in B}$$

$$t=t_0$$

$$\sum_{p \in B} E[c?(st, t),$$

$$st$$

$$st+1 \sum_{p \in B} P(\cdot | st, \sigma(st, t), st_0 = s).$$

Another reduction comes from the fact that the set of balanced gain distributions can be seen as a convex combination of a finite set of balanced distributions [9, Claim 2 and 3]. We call this limited set the atomic gain distributions. Therefore the search for σ can be limited to this set. The set of convex combinations of the m distributions $A?_1, \dots, A?_m$ is denoted by $\Delta(A?_1, \dots, A?_m)$.

4

Reformulation as a Markovian Decision Problem

In this section we formulate, for arbitrary K , the maximization problem over balanced adversaries as an undiscounted MDP problem (S, A, r, p) . The state space S was defined in Section 2 and the action space is the set of atomic balanced distributions as discussed in Section 3. The transition model is defined by $p(\cdot | s, D?)$, which is a probability distribution over states given the current state s and a balanced distribution over gains $D?$. In this model, the transition dynamics are deterministic and entirely controlled by the adversary's action choices. However, the adversary is forced to choose stochastic actions (balanced gain distributions). The maximization problem can therefore also be thought of as designing a balanced random walk on states so as to maximize a sum of rewards (that are yet to be defined). First, we define $PA?$ the transition probability operator with respect to a gain distribution $A?$. Given function $f : S \rightarrow \mathbb{R}$, $PA?$ returns $[PA? f](s) = E[f(s') | s' \sim p(\cdot | s, A?)] = E[f(s + g)]$. $g \sim A?$

g is sampled in s according to $A?$. Given $A?$ in s , the per-step regret is denoted by $rA?(s)$ and defined as $rA?(s) = E[s' \sim p(\cdot | s, A?) - cA? \cdot s' | s, A?]$

4

Given an adversary σ strategy σ , starting in s at time t_0 , the cumulative per-step regret is $\sum_{t=t_0}^T Vt?0(s) = \sum_{t=t_0}^T E[r?(\sigma, t)(st) - st+1 \sum_{p \in B} P(\cdot | st, \sigma(st, t), st_0 = s)]$. The action-value function of σ at $(s, D?)$ and t is the expected sum of rewards received by starting from s , taking action $D?$, and then σ from t (st, D?)

$) = E [\sum_{t=0}^{\infty} \gamma^t (r_t - A_t)] = D \cdot \sum_{t=0}^{\infty} \gamma^t p(-st, A_t)$, $A_{t+1} = \gamma(st+1, t+1)$. following γ : $Q_t = \gamma V_{t+1} + \gamma \sum_{s'} p(s'|s, A_t) [r_t + \gamma V_{t+1}(s')]$. The Bellman operator of A_t , TA_t , is $[TA_t f](s) = rA_t(s) + \gamma \sum_{s'} p(s'|s, A_t) f(s')$. with $[T^*(s, t) V_{t+1}](s) = V_{t+1}(s)$. This per-step regret, $rA_t(s)$, depends on s and A_t and not on the time step t . Removing the time from the picture permits a simplified view of the problem that leads to a natural formulation of the exchangeability property that is independent of the time t . Crucially, this decomposition of the regret into per-step regrets is such that maximizing $\sum_{t=0}^{\infty} \gamma^t rA_t(s)$ over adversaries γ is equivalent, for all time t_0 and s , to maximizing over adversaries the original value of the game, the regret $\sum_{t=0}^{\infty} \gamma^t rA_t(s)$ (Lemma 2). Lemma 2. For any adversary strategy γ and any state s and time t_0 , $\sum_{t=t_0}^{\infty} \gamma^t rA_t(s) = \sum_{t=t_0}^{\infty} \gamma^t rA_t(s) + \sum_{t=t_0}^{\infty} \gamma^t rA_t(s)$.

The proof of Lemma 2 is in Section 8. In the following, our focus will be on maximizing $\sum_{t=0}^{\infty} \gamma^t rA_t(s)$ in any state s . We now show some basic properties of the per-step regret that holds for an arbitrary number of experts K and discuss their implications. The proofs are in Section 9. γ for all s, t , we have $0 \leq rA_t(s) \leq 1$. Furthermore if $\sum_{s'} p(s'|s, A_t) = 1$, $rA_t(s) = 0$. Lemma 3. Let $A_t \in \mathcal{A}$, Lemma 3 shows that a state s in which the reward is not zero contains at least two equal leading experts, $\sum_{s'} p(s'|s, A_t) = 1$. Therefore the goal of maximizing the reward can be rephrased into finding a policy that visits the states with $\sum_{s'} p(s'|s, A_t) = 1$ as often as possible, while still taking into account that the per-step reward increases with $\sum_{s'} p(s'|s, A_t)$. The set of states with $\sum_{s'} p(s'|s, A_t) = 1$ is called the "reward wall". Lemma 4. In any state s , with $\sum_{s'} p(s'|s, A_t) = 2$, for any balanced gain distribution D_t such that with probability one exactly one of the leading expert receives a gain of 1, $rD_t(s) = \max_{A_t \in \mathcal{A}} rA_t(s)$.

5

The Case of $K = 3$

5.1

Notations in the 3-Experts Case, the C OMB and the T WIN -C OMB Adversaries

First we define the state space in the 3-expert case. The experts are sorted with respect to their cumulative gains and are named in decreasing order, the leading expert, the middle expert and the lagging expert. As mentioned in [9], in our search for the minimax optimal adversary, it is sufficient for any K to describe our state only using d_{ij} that denote the difference between the cumulative gains of consecutive sorted experts i and $j = i + 1$. Here, i denotes the expert with i th largest cumulative gains, and hence $d_{ij} \geq 0$ for all $i \leq j$. Therefore one notation for a state, that will be used throughout this section, is $s = (x, y) = (d_{12}, d_{23})$. We distinguish four types of states C_1, C_2, C_3, C_4 as detailed below in Figure 1. In the same figure, in the center, the states are represented on a 2d-grid. C_4 contains only the state denoted $s^* = (0, 0)$. Reward Wall

$s \in C_1, d_{12} > 0, d_{23} > 0$ $s \in C_2, d_{12} = 0, d_{23} > 0$ $s \in C_3, d_{12} > 0, d_{23} = 0$ $s \in C_4, d_{12} = 0, d_{23} = 0$

1
 1
 1
 2
 1
 1
 1
 1
 2
 1
 1
 1
 4
 3
 3
 3
 d12
 Atomic A? {1}{23} {2}{13} {3}{12} {1}{2}{3} {12}{13}{23}
 Symbol ? W ? C ? V 1? 2?
 cA? 1/2 1/2 1/2 1/3 2/3

Figure 1: 4 types of states (left), their location on the 2d grid of states (center) and 5 atomic A? (right) Concerning the action space, the gain distributions use brackets. The group of arms in the same bracket receive gains together and each group receive gains with equal probability. For instance, {1}{2}{3} exclusively deals a gain to expert 1 (leading expert) with probability 1/3, expert 2 (middle expert) with probability 1/3, and expert 3 (lagging expert) with probability 1/3, whereas {1}{23} means dealing a gain to expert 1 alone with probability 1/2 and experts 2 and 3 together with probability 1/2. As discussed in Section 3, we are searching for a π^* using mixtures of atomic balanced distributions. π^* , π^* , π^* , π^* , π^* , π^* , π^* When $K = 3$ there are seven atomic distributions, denoted by B^* , B^* , B^* , B^* , B^* , B^* , B^* and described in Figure 1 (right). Moreover, in Figure 2, we report in detail in a table (left) and 5

s
 rC? (s)
 C1 C2 C3 C4
 0 1/2 0 1/2
 Distribution of next state $s^* = p(s^* | s, C^*)$ with $s = (x, y)$ $P(s^* = (x+1, y+1)) = P(s^* = (x+1, y^*1)) = .5$ $P(s^* = (x+1, y)) = P(s^* = (x+1, y^*1)) = .5$ $P(s^* = (x, y+1)) = P(s^* = (x^*1, y+1)) = .5$ $P(s^* = (x, y+1)) = P(s^* = (x+1, y)) = .5$
 d23
 .5
 d12
 1
 .5
 0
 2

0
3
.5
1/2
.5 .5
.5
.5
4
.5 1/2

Figure 2: The per-step regret and transition probabilities of the gain distribution C ? an illustration (right) on the 2-D state grid?the properties of the C OMB gain distribution C ? . The remaining atomic distributions are similarly reported in the appendix in Figures 5 to 8. In the case of three experts, the C OMB distribution is simply playing $\{2\}\{13\}$ in any state. We use γ to denote the strategy that plays $\{1\}\{23\}$ in any state and refer to it as the T WIN -C OMB strategy. W The C OMB and T WIN -C OMB strategies (as opposed to the distributions) repeat their respective gain distributions in any state and any time. They are respectively denoted γ_C , γ_W . The Lemma 5 shows that the C OMB strategy γ_C , the T WIN -C OMB strategy γ_W and therefore any mixture of both, have the same expected cumulative per-step regret. The proof is reported to Section 11. Lemma 5. For all states s at time t , we have $V^{\gamma_C}(s) = V^{\gamma_W}(s)$. 5.2

The Exchangeability Property

γ) such that for any $s \neq s'$, and for any $f : S \rightarrow R$, Lemma 6. Let $A \in B^3$, there exists $D \in \gamma(C, W [TA \rightarrow TD \rightarrow f])(s) = [TW \rightarrow TA \rightarrow f](s)$.

γ , $A = \{\}$ or $A = \{123\}$, use $D = W$. If $A = C$, use Lemma 11 and 12. Proof. If $A = W$) with $s \in C_3$ then $s \in C_3 \cap C_4$. Case 1. $A = V$: V is equal to C in $C_3 \cap C_4$ and if $s \in p(\cdot, s, W)$ So when $s \in C_3$ we reuse the case $A = C$ above. When $s \in C_1 \cap C_2$, we consider two cases. γ which is $\{1\}\{23\}$. If $s \in p(\cdot, s, V)$ with $s \in C_2$ then Case 1.1. $s = (0, 1)$: We choose $D = W$ $s \in C_2$. Similarly, if $s \in p(\cdot, s, V)$ with $s \in C_1$ then $s \in C_1 \cap C_3$. Moreover D modifies similarly the coordinates (d_{12}, d_{23}) of $s \in C_1$ and $s \in C_3$. Therefore the effect in terms of transition probability and reward of D is the same whether it is done before or after the actions chosen by V . If $s \in p(\cdot, s, D)$ with $s \in C_1 \cap C_2$ then $s \in C_1 \cap C_2$. Moreover V modifies similarly the coordinates (d_{12}, d_{23}) of $s \in C_1$ and $s \in C_2$. Therefore the effect in terms of the transition probability of V is the same whether it is done before or after the action D . In terms of reward, notice that in the states $s \in C_1 \cap C_2$, V has 0 per-step regret and using V does not make s leave or enter the reward wall. γ . One can check from the tables in Figures 7 and 8 that Case 1.2 $s = (0, 1)$: We can chose $D = W$ exchangeability holds. Additionally we provide an illustration of the exchangeability equality in the 2d-grid in Figure 1. The starting state $s = (0, 1)$, is graphically represented by γ . We show on the grid the effect of the gain distribution V (in dashed red) followed (left picture) or preceded (right picture) by the gain distribution D (in plain blue). The illustration shows that $V \rightarrow D$ and $D \rightarrow V$ lead to the

1
7
1
5
3
.5
3
10
8
d12
6
4
.5

Consider a random walk that starts from state $s_0 = s$ and is generated by the TWIN-COMB strategy, $s_{t+1} = p(\cdot | s_t, W_t)$. Define the random variable $T_s = \min\{t \in \mathbb{N} \setminus \{0\} : s_t = s^*\}$. This random variable is the number of steps of the random walk before hitting s^* for the first time. Then, let $P^*(s, t)$ be the probability that s^* is reached after t steps: $P^*(s, t) = P(T_s = t)$. (top) & G random walks (bottom) Lemma 8 controls the COMB greedy error in s_t in relation to $P^*(s, t)$. Lemma 9 derives a state-independent upper-bound for $P^*(s, t)$. Lemma 8. For any time $t \in [T]$ and state $s, s^* \in \mathcal{S}$, $P^*(s, t) \leq P^*(s^*, t)$.

Proof. If $s = s^*$, this is a direct application of Lemma 7 as $P^*(s^*, t) = 0$ for $t \neq 0$. When $s \neq s^*$, the following proof is by induction.

Initialization: Let $t = T$. At the last round only the last per-step regret matters (for all states $s, s^* \in \mathcal{S}$, $r_D(s) = r_D(s^*)$). As $s \neq s^*$, s is such that $r_D(s) > r_D(s^*)$ then $r_D(s) = \max_{s' \in \mathcal{S}} r_A(s')$ because of QALB Lemma 4 and Lemma 3. Therefore the statement holds. Induction: Let $t \leq T$. We assume the statement is true at time $t + 1$. We distinguish two cases. For all gain distributions $D \in \mathcal{B}$, (b) $r_D(s, A) = \max_{s' \in \mathcal{S}} r_D(s', A)$ (a) $r_D(s, D) = [TD + TE + V_{t+2}](s) = [TW + [TD + V_{t+2}]](s) = [TW + Q_{t+1} + T_{t+1}](s) = [TW + Q_{t+1} + T_{t+1}](s) = [TW + \max_{s' \in \mathcal{S}} Q_{t+1}(s', A)](s) = [TW + \max_{s' \in \mathcal{S}} Q_{t+1}(s', A)](s) = [TW + \max_{s' \in \mathcal{S}} Q_{t+1}(s', A)](s) = [TW + \max_{s' \in \mathcal{S}} Q_{t+1}(s', A)](s) = [TW + \max_{s' \in \mathcal{S}} Q_{t+1}(s', A)](s)$

- (d) $r_D(s, A) = \max_{s' \in \mathcal{S}} [TW + Q_{t+1}(s', A)](s)$
- (b) $r_D(s, A) = \max_{s' \in \mathcal{S}} Q_{t+1}(s', A)$
- (e) $r_D(s, A) = \max_{s' \in \mathcal{S}} Q_{t+1}(s', A)$
- T
- t1
- $r_D(s, A) = [PW + P_{t+1}(s, A)](s)$
- $r_D(s, A) = [PW + P_{t+1}(s, A)](s)$
- $r_D(s, A) = [PW + P_{t+1}(s, A)](s)$

7

π_t and this step holds because of Lemma 5, (b) holds where in (a) E_t is any distribution in $\Pi(C, W)$ because of the exchangeability property of Lemma 6, (c) is true by induction and monotonicity of Bellman operator, in (d) the max operators change from being specific to any next state s' at time $t + 1$ to being just one max operator that has to choose a single optimal gain distribution in state s at time t , (e) holds by definition as for any t_2 , (here the last equality holds because $s' = s$) $[P W P'(\cdot, t_2)](s) = E s' p(\cdot, s, W) [P'(s', t_2)] = E s' p(\cdot, s, W) [P'(T, s' = t_2)] = P(s, t_2 + 1)$. Lemma 9. For $t \geq 0$ and any s ,

$$P(s, t) = \sum_{s'} P(s', t) \pi_t(s')$$

?

$$P(s, t) = \sum_{s'} P(s', t) \pi_t(s')$$

Proof. Using the connection between the TWIN-C OMB strategy and a simple random walk in Proposition 1, a formula can be found for $P(s, t)$ from the classical "Gambler's ruin" problem, where one wants to know the probability that the gambler reaches ruin (here state s') at any time t given an initial capital in dollars (here is as defined in Proposition 1). The gambler has an equal probability to win or lose one dollar at each round and has no upper bound on his capital during the game. Using [7] (Chapter XIV, Equation 4.14) or [18] we have $P(s, t) = \sum_{s'} P(s', t) \pi_t(s')$, where the 2 binomial coefficient is 0 if t and s' are not of the same parity. The technical Lemma 14 completes the proof. We now state our main result, connecting the value of the C OMB adversary to the value of the game. Theorem 1. Let $K = 3$, the regret of C OMB strategies against any learner p , $\pi_t \in C$, satisfies $\sum_{t=1}^T \min_{p \in C} V_p(s) - V_t(s) \leq 12 \log(T + 1) + 2$.

We also characterize the minimax regret of the game. Theorem 2. Let $K = 3$, for even T , we have that $\sum_{t=1}^T \min_{p \in C} V_p(s) - V_t(s) \leq T + 2 \sqrt{T} + 1 \leq 12 \log_2(T + 1) + 3 \sqrt{T} + 2$.

with

?

$$\sum_{t=1}^T \min_{p \in C} V_p(s) - V_t(s) \leq T + 2 \sqrt{T} + 1 \leq 12 \log_2(T + 1) + 3 \sqrt{T} + 2$$

In Figure 4 we introduce a C OMB-based learner that is denoted by p^C . Here a state is represented by a vector of 3 integers. The three arms/experts are ordered as (1) (2) (3), breaking ties arbitrarily. We connect the value of the C OMB-based learner to the value of the game. Theorem 3. Let $K = 3$, the regret of C OMB-based p^C learner against any adversary π , $\max_{\pi \in \Pi} \sum_{t=1}^T V_{\pi}(s) - V_t(s) \leq VT + 36 \log_2(T + 1)$.

Figure 4: A C OMB learner, p^C Similarly to [2] and [14], this strategy can be efficiently computed using rollouts/simulations from the C OMB adversary in order to estimate the value $V_t(s)$ of C in s at time t .

6

Discussion and Future Work

The main objective is to generalize our new proof techniques to higher dimensions. In our case, the MDP formulation and all the results in Section 4 already holds for general K . Interestingly, Lemma 3 and 4 show that the C OMB distribution is the balanced distribution with highest per-step regret in all the states s such that $\|X(s)\|_2 \leq 2$, for arbitrary K . Then assuming an ideal exchangeability property that gives $\max_{A \in \mathcal{A}} \sum_{t=1}^T \langle A, C^t \rangle = \max_{A \in \mathcal{A}} \sum_{t=1}^T \langle A, C^t \rangle$, a distribution would be greedy w.r.t the C OMB strategy at an early round of the game if it maximizes the per-step regret at the last round of the game. The C OMB policy specifically tends to visit almost exclusively states $\|X(s)\|_2 \leq 2$, states where C OMB itself is the maximizer of the per-step regret (Lemma 3). This would give that C OMB is greedy w.r.t. itself and therefore optimal. To obtain this result for larger K , we will need to extend the exchangeability property to higher K and therefore understand how the C OMB and TWIN-C OMB families extend to higher dimensions. One could also borrow ideas from the link with pde approaches made in [6]. 8

Acknowledgements We gratefully acknowledge the support of the NSF through grant IIS-1619362 and of the Australian Research Council through an Australian Laureate Fellowship (FL110100281) and through the Australian Research Council Centre of Excellence for Mathematical and Statistical Frontiers (ACEMS). We would like to thank Nate Eldredge for pointing us to the results in [18]!

2 References

- [1] Jacob Abernethy and Manfred K. Warmuth. Repeated games against budgeted adversaries. In *Advances in Neural Information Processing Systems (NIPS)*, pages 179, 2010.
- [2] Jacob Abernethy, Manfred K. Warmuth, and Joel Yellin. Optimal strategies from random walks. In *21st Annual Conference on Learning Theory (COLT)*, pages 437-446, 2008.
- [3] Nicolò Cesa-Bianchi, Yoav Freund, David Haussler, David P. Helmbold, Robert E. Schapire, and Manfred K. Warmuth. How to use expert advice. *Journal of the ACM (JACM)*, 44(3):427-485, 1997.
- [4] Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- [5] Thomas M. Cover. Behavior of sequential predictors of binary sequences. In *4th Prague Conference on Information Theory, Statistical Decision Functions, Random Processes*, pages 263-272, 1965.
- [6] Nadeja Drenska. A pde approach to mixed strategies prediction with expert advice. <http://www.gtcenter.org/Downloads/Conf/Drenska2708.pdf>. (Extended abstract).
- [7] William Feller. *An Introduction to Probability Theory and its Applications*, volume 2. John Wiley & Sons, 2008.
- [8] Nick Gravin, Yuval Peres, and Balasubramanian Sivan. Towards optimal algorithms for prediction with expert advice. In *arXiv preprint arXiv:1603.04981*, 2014.
- [9] Nick Gravin, Yuval Peres, and Balasubramanian Sivan. Towards optimal algorithms for prediction with expert advice. In *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 528-547, 2016.
- [10] Nick Gravin, Yuval Peres, and Balasubramanian Sivan. Tight Lower Bounds for Multiplicative Weights Algorithmic Families. In *44th International*

Colloquium on Automata, Languages, and Programming (ICALP), volume 80, pages 48:1?48:14, 2017. [11] Charles Miller Grinstead and James Laurie Snell. Introduction to probability. American Mathematical Soc., 2012. [12] James Hannan. Approximation to bayes risk in repeated play. Contributions to the Theory of Games, 3:97?139, 1957. [13] Ronald A. Howard. Dynamic Programming and Markov Processes. The MIT Press, Cambridge, MA, 1960. [14] Haipeng Luo and Robert E. Schapire. Towards minimax online learning with unknown time horizon. In Proceedings of The 31st International Conference on Machine Learning (ICML), pages 226?234, 2014. [15] Francesco Orabona and David P. L. Optimal non-asymptotic lower bound on the minimax regret of learning with expert advice. arXiv preprint arXiv:1511.02176, 2015. [16] Martin L. Puterman. Markov Decision Processes. Wiley, New York, 1994. [17] Pantelimon Stanica. Good lower and upper bounds on binomial coefficients. Journal of Inequalities in Pure and Applied Mathematics, 2(3):30, 2001. [18] Remco van der Hofstad and Michael Keane. An elementary proof of the hitting time theorem. The American Mathematical Monthly, 115(8):753?756, 2008. 9