

Towards Optimal Algorithms for Prediction with Expert Advice

Nick Gravin*

Yuval Peres[†]

Balasubramanian Sivan[†]

Abstract

We study the classical problem of prediction with expert advice in the adversarial setting with a geometric stopping time. In 1965, Cover gave the optimal algorithm for the case of 2 experts. In this paper, we design the optimal algorithm, adversary and regret for the case of 3 experts. Further, we show that the optimal algorithm for 2 and 3 experts is a probability matching algorithm (analogous to Thompson sampling) against a particular randomized adversary. Remarkably, our proof shows that the probability matching algorithm is not only optimal against this particular randomized adversary, but also minimax optimal.

Our analysis develops upper and lower bounds simultaneously, analogous to the primal-dual method. Our analysis of the optimal adversary goes through delicate asymptotics of the random walk of a particle between multiple walls. We use the connection we develop to random walks to derive an improved algorithm and regret bound for the case of 4 experts, and, provide a general framework for designing the optimal algorithm and adversary for an arbitrary number of experts.

1 Introduction

Predicting future events based on past observations, a.k.a. prediction with expert advice, is a classic problem in learning. The experts framework was the first framework proposed for online learning and encompasses several applications as special cases. The underlying problem is an online optimization problem: a *player* has to make a decision at each time step, namely, decide which of the k experts' advice to follow. At every time t , an *adversary* sets gains for each expert: a gain of g_{it} for expert i at time t . Simultaneously, the player, seeing the gains from all previous steps except t , has to choose an action, i.e., decide on which expert to follow. If the player follows expert $j(t)$ at time t , he gains $g_{j(t),t}$. At the end of each step t , the gains associated with all experts are revealed to the player, and the player's choice is revealed to the

adversary. In the *finite horizon model*, this process is repeated for T steps, and the player's goal is to perform (achieve a cumulative gain) as close as possible to the best single action (best expert) in hindsight, i.e., to minimize his *regret* R_T :

$$R_T = \max_{1 \leq i \leq k} \sum_{t=1}^T g_{it} - \sum_{t=1}^T g_{j(t),t}.$$

Apart from assuming that the g_{it} 's are bounded in $[0, 1]$, we don't assume anything else about the gains. Just as natural as the finite horizon model is the model with a *geometric horizon*: the stopping time is a geometric random variable with expectation $\frac{1}{\delta}$. In other words, the process ends at any given step with probability δ , independently of the past. In this paper, we study both the finite horizon model and the geometric horizon model.

Questions and motivation. Given the breadth of applications and the significance of the experts problem in online learning, in this work we seek to understand and crisply characterize the structure of the optimal algorithm and the structure of the worst case input sequences. We ask:

1. What is the precisely optimal algorithm and regret values?
2. Does the optimal algorithm have a succinct and intuitive description (even for 2 experts)?
3. What are the hardest (adversarial) sequences of experts' gains and do they follow a succinct pattern?

Our motivation in exploring these questions include the following.

1. Half a century after Cover [9] described the optimal adversary for the case of 2 experts, we still do not have general insights about the structure of the optimal algorithm or the optimal adversarial sequences.
2. Several applications of the experts paradigm involve dealing with a small constant number of experts. What amount of gain can the optimal algorithm get over the multiplicative weights algorithm for a constant number of experts?

*Microsoft Research. One Memorial Drive, Cambridge, MA 02142. ngravin@gmail.com.

[†]Microsoft Research. One Microsoft Way, Redmond, WA 98052. peres@microsoft.com, balu2901@gmail.com.

3. The problem is theoretically clean and challenging: a priori it is not even clear if the algorithm and adversary are succinctly describable. It could well be that the optimal algorithm’s actions depend on various aspects of history in a manner that cannot be succinctly described.

Notation: We fix some notation before proceeding. We denote by G_{it} the cumulative gain of expert i after t steps. Namely, $G_{it} = \sum_{s=1}^t g_{is}$. We show that in the worst-case instance there is no benefit in using gains other than 0 and 1, so we restrict to $g_{it} \in \{0, 1\}$. The notion of optimality used for the experts framework is the minimax regret obtained against all possible adversarial sequences of experts’ predictions, *the adversary* for short. We study the optimal adversary that inflicts maximal regret (maxmin) against all possible algorithms.

Our contributions *Balanced adversary:* Our first general insight about the structure of the optimal adversary is that it is *balanced across all experts* at every time step, i.e., irrespective of the experts’ past gains, the adversary sets equal expected gain for each expert in this step. This insight is pervasive in this paper and greatly simplifies the problem in that it lets us describe the optimal (minimax/maxmin) regret without making any reference to the optimal algorithm. Indeed, *every* algorithm performs equally well against the maximin optimal adversary that equalizes the expected gains of all experts¹, and gets a gain of the average over all the k experts of the cumulative gain, namely, $\frac{1}{k}|\mathbf{G}|_1$. Given this, the adversary’s problem of maximizing regret can be reduced to maximizing the difference between maximum and average of the cumulative gains vector, i.e., $|\mathbf{G}|_\infty - \frac{1}{k}|\mathbf{G}|_1$.

Maximizing the number of collisions between the leading and second best expert: Our second insight about the adversary’s structure is that its objective, namely, $|\mathbf{G}|_\infty - \frac{1}{k}|\mathbf{G}|_1$ never changes in expectation, except when there is no unique expert with the largest cumulative gains. This is because, the maximin optimal adversary being balanced implies that all experts’ cumulative gains increase equally in expectation, including the cumulative gain of the expert who is currently leading. Thus the two quantities $|\mathbf{G}|_\infty$ and $\frac{1}{k}|\mathbf{G}|_1$ increase equally in

¹We clarify that this doesn’t mean all algorithms get the same minimax regret. It is only the maximin optimal adversary that is balanced and not every adversary. In other words, if the adversary is maximin optimal, all algorithms are equal. But if the algorithm is not minimax optimal, the optimal adversary for that algorithm is not maximin optimal and hence not necessarily balanced, and will inflict a larger regret on the said algorithm than a balanced adversary does.

expectation implying that a balanced adversary will find no way to increase its objective in expectation. The only time when this breaks is when the leading expert is not unique: this is because in this situation, the probability that the maximum cumulative gain increases is the probability of the cumulative gain of any one of the leading experts getting increased. Given this, the adversary’s problem essentially boils down to the crisp and challenging probability question of constructing gain sequences that maximizes the number of collisions between the leading and the second leading expert. Namely, construct a balanced random walk in \mathbb{Z}^k with the objective of maximizing the number of collisions between the largest and second largest coordinates.

Designing the random walk that maximizes the number of collisions: We use this insight about the adversary’s problem being a controlled random walk to construct such walks and hence succinct adversaries for the case of $k = 2$ and 3 experts (this also gives a simple alternative proof for Cover’s optimal adversary). While the case of $k = 2$ is special, the progress for $k = 3$ crucially relies on the above mentioned insights. While constructing the optimal such random walk for general k is still complicated, we believe that this reduction is powerful and gives a useful starting point for thinking about possible candidates that come close. For instance, the “comb adversary” described later in the introduction, has a simple but non-trivial structure that was inspired from the number-of-collisions characterization.

Probability matching algorithm: We establish a strong connection between the structure of the optimal algorithm and the optimal adversary. Namely, the optimal algorithm is a probability matching algorithm (analogous to the popular Thompson sampling procedure) that follows each expert with the probability that this expert finishes as the leader, when the sequence of gains is set by the optimal adversary .

We describe our results in detail below.

1. Two experts. The optimal adversary, designed by Cover, chooses one expert uniformly at random, sets a gain of 1 for that expert, and a gain of 0 for the other. We give a very simple characterization of the *unique* optimal algorithm in *both the finite horizon and the geometric horizon models*, namely, *follow each expert with the probability that he finishes as the leading expert² (the one with maximum cumulative gains), when gains are set*

²If there is a tie in the finite horizon model, we consider the leader to be the unique expert who did not have any expert ahead of him in the last two steps. In the geometric horizon model ties are just broken uniformly at random.

by Cover’s optimal adversary. Further, this probability of finishing as a leading expert has a simple analytical expression in the geometric horizon model, given in Algorithm 1 (see Theorem 4.1). The finite horizon has a simple expression too (see Theorem 3.1). The optimal algorithm achieves a regret of $\sqrt{\frac{T}{2\pi}}$ in the finite horizon model, and $\frac{1}{2} \frac{1}{\sqrt{2\delta}}$ in the geometric horizon model, respectively as $T \rightarrow \infty$ and $\delta \rightarrow 0$ (see Theorem 4.1 for the precisely optimal regret for every δ).

Algorithm 1 : Optimal Algorithm for Geometric Horizon Model with Two Experts

- 1: Initialize $\xi = \frac{1 - \sqrt{1 - (1 - \delta)^2}}{1 - \delta} \sim 1 - \sqrt{2\delta}$ as $\delta \rightarrow 0$
 - 2: Convention: Leading expert (larger cumulative gains) is numbered 1, and lagging expert is numbered 2
 - 3: **for** Each time step t till the game stops **do**
 - 4: Compute cumulative gains for both experts: $G_{1t} = \sum_{s=1}^t g_{1s}$, and $G_{2t} = \sum_{s=1}^t g_{2s}$
 - 5: Let $d = G_{1t} - G_{2t}$. Note that by definition $d \geq 0$
 - 6: Follow the leading expert with probability he will finish as leader, namely, $p_1(d) = 1 - \frac{1}{2}\xi^d$
 - 7: Follow the lagging expert with probability he will finish as leader, namely, $p_2(d) = \frac{1}{2}\xi^d$
 - 8: **end for**
-

2. Three experts. We derive the precisely optimal algorithm, adversary and regret values for three experts in the geometric horizon model (see Theorem 4.2). The optimal regret as $\delta \rightarrow 0$ is asymptotic to $\frac{2}{3} \frac{1}{\sqrt{2\delta}}$. The optimal adversary (as $\delta \rightarrow 0$)³ is as follows: it pairs up the middle and the lagging experts, and together this pair always disagrees with the leading expert. That is, the g_{it} ’s are of the form $(0, 1, 1)$ or $(1, 0, 0)$ where the ordering in the tuple captures the leading, middle and lagging experts (and do not refer to the identities of experts). The optimal algorithm is again a simple probability matching algorithm to the optimal adversary, that follows each expert with the probability that this expert finishes in the lead.

The case of 3 experts is significantly more complicated than 2 experts. In particular, while the optimal adversary for 2 experts was discovered back in 1965 ([9]), the optimal adversary for 3 experts was not known so far. The relative coordinate system we introduce, that numbers experts according to their cumulative gains, provides a convenient way to describe the optimal adversary.

3. Arbitrary number of experts. All our basic results continue to hold in both geometric and finite hori-

³The optimal adversary for all values of δ (asymptotic or not) is almost identical to this. We describe this in Section 4.

zon models. I.e., the optimal adversary (i) plays only $g_{it} \in \{0, 1\}$ (and not in $[0, 1]$), (ii) is balanced, (iii) at every step plays only one of the finitely many vertices of the convex polytope of balanced distributions. Prior to this work, given T and k , an algorithm for computing the precisely optimal adversary was not known. Result (iii) reduces the search space of the optimal adversary to finitely many balanced distributions, and thereby enables us to write a mundane dynamic program of size $O(T^k)$. Note that even after realizing that $g_{it} \in \{0, 1\}$ without loss of generality, the adversary has infinitely many balanced probability distributions available to choose from in every step, and thus, a priori it is not clear how to write a meaningful dynamic program.

Through this dynamic program, we found that the optimal adversary for $k \geq 4$ does not always have a simple description like for $k = 2, 3$. Further, we observe that unlike $k = 2, 3$, for 4 experts, the optimal adversary is already δ -dependent, and its actions at a given configuration of cumulative gains depend on the exact values of cumulative gains and not just their order. Nevertheless, inspired by results from this dynamic program, we conjecture that there is a simple adversary (“comb adversary”) which is asymptotically (in δ or T) optimal: split experts into two teams $\{1, 3, 5, \dots\}$ and $\{2, 4, 6, \dots\}$ and increment the gains of all experts in exactly one of these teams chosen uniformly at random. We analyze the comb adversary for $k = 4$ and show that as $\delta \rightarrow 0$, it inflicts a regret of $\frac{\pi}{4} \frac{1}{\sqrt{2\delta}}$ (see Theorem 5.2). We observed that for reasonably small values of δ , the optimal regret converges to our lower bound of $\frac{\pi}{4} \frac{1}{\sqrt{2\delta}}$ indicating that the comb adversary is indeed asymptotically optimal for $k = 4$.

Remarks

1. In this work, we develop the optimal algorithm and adversary *simultaneously*, thereby, completely bridging the gap between upper and lower bounds for a small number of experts (our analysis obtains the optimal regret for $k = 2, 3$ experts for every value of δ , and not just asymptotically as $\delta \rightarrow 0$).
2. Although the optimal algorithm for $k = 2, 3$ experts does a probability matching with respect to a particular adversary it turns out that this algorithm is not only optimal against this adversary, but also minimax optimal against all possible adversaries. Our algorithms for $k = 2, 3$ experts (also our conjectured optimal algorithm for k experts) are simple and practical. One can implement our algorithms as follows: from any configuration of cumulative gains simulate the “comb adversary” till the end of the process and follow the expert who finishes in the lead. Simulating the comb adversary simply entails flipping a coin in every step

and incrementing by 1 the gains of the respective (odd or even numbered) team of experts.

3. The comb adversary that we introduce and analyze presents a simple to describe random process. But even for $k = 3, 4$ analyzing this process requires an understanding of non-trivial aspects of simple random walk (see Theorems 5.1 and 5.2). Developing a method to analyze this process for general k is a clean and challenging question on random walks.

Comparison with Multiplicative Weights Algorithm. It is known that for general bounded gains the widely used multiplicative weights algorithm (MWA), obtains a $\sqrt{\frac{T \ln k}{2}}$ regret and this regret is asymptotically optimal as both $\{T, k\} \rightarrow \infty$ (see Cesa-Bianchi et al. [8] and its generalization by Haussler et al. [14]). However, asymptotic analysis in k does not shed much light on the structure of the optimal algorithm and the hardest sequences of experts' gains: this is because the quantity $\sqrt{\frac{T \ln k}{2}}$ is insensitive if we employ 100 times as many experts and it is not optimal for a small (constant) number of experts. In this paper, we show that the optimal algorithm is not in the family of multiplicative weight algorithms. Namely, we show that the optimal algorithm cannot be expressed as a MWA or even as a convex combination of MWAs. We refer the reader to Appendix D for a detailed discussion and proof.

Related work. In this work, our goal is to identify the structure of the optimal algorithm and the adversary via a precise and efficient algorithmic description. There are recent works that *characterize* the optimal algorithm/adversary and regret as the supremum or infimum of some stochastic process, rather than give an efficient algorithmic description. There is also a significant body of recent work that either identify *approximately* optimal algorithms/adversaries, or, identify special cases where the optimal adversary can be precisely described. All of these relaxations allow for solving more general frameworks. But this body of work doesn't identify the precisely optimal algorithm/adversary/regret for the classical setting. In contrast, in our work we show that probability matching is precisely optimal, and we identify the precisely optimal adversary for the classical setting. We discuss all these lines of recent work after discussing some classics in this area.

As mentioned earlier, for the exact setting we consider in this work, the work of [9] is most closely related as it gives the optimal adversary and algorithm for the case of 2 experts.

Classic works: The book by Cesa-Bianchi and Lugosi [7] is an excellent source for both applications and references. The prediction with experts advice paradigm

was introduced by Littlestone and Warmuth [17] and Vovk [26]. The famous multiplicative weights update algorithm was introduced independently by these two works: as the weighted majority algorithm by Littlestone and Warmuth and as the aggregating algorithm by Vovk. The pioneering work of Cesa-Bianchi et al. [8] considered $\{0, 1\}$ outcome space for nature and showed that for the absolute loss function $\ell(x, y) = |x - y|$ (or $g(x, y) = 1 - |x - y|$), the asymptotically optimal regret is $\sqrt{\frac{T \ln k}{2}}$. This was later extended to $[0, 1]$ outcomes for nature by Haussler et al. [14]. The asymptotic optimality of $\sqrt{\frac{T \ln k}{2}}$ for arbitrary loss (gain) functions follows from the analysis of Cesa-Bianchi [6]. When it is known beforehand that the cumulative loss of the optimal expert is going to be small, the optimal regret can be considerably improved, and such results were obtained by Littlestone and Warmuth [17] and Freund and Schapire [11]. With certain assumptions on the loss function, the simplest possible algorithm of following the best expert already guarantees sub-linear regret Hannan [13]. Even when the loss functions are unbounded, if the loss functions are exponential concave, sub-linear regret can still be achieved Blum and Kalai [4].

Recent works: [18] consider a setting where the adversary is restricted to pick gain vectors from the basis vector space $\{e_1, \dots, e_k\}$. For this set of gain vectors, the only balanced adversary is to pick a random expert in every step. Since our analysis shows that the optimal adversary is balanced without loss of generality, it is immediate that a uniformly random adversary is optimal in this setting. [2] consider a different variant of experts problem where the game stops when cumulative loss of any expert exceeds given threshold. Here too there is a clear candidate for the optimal adversary: the same as in Luo and Schapire, namely, pick an expert uniformly at random at every step. They specify optimal algorithm in terms of the underlying random walk. The notable distinction of both [18, 2] from our setting is that their adversary is simple and static, i.e., it does not depend on the prior history. The random process to be analyzed in their setting is a standard random walk in \mathbb{Z}^k , while the random process in our setting even for $k = 3$ is non-trivial. [1] consider general convex games and compute the minimax regret exactly when the input space is a ball, and show that the algorithms of [27] and [15] are optimal w.r.t. minimax regret. [3] provide upper and lower bounds on the regret of an optimal strategy for several online learning problems without providing algorithms, by relating the optimal regret to the behavior of a certain stochastic process. [21] consider a continuous experts setting where the algorithm knows beforehand the

maximum number of mistakes of the best expert. [22] introduce the notion of sequential Rademacher complexity and use it to analyze the learnability of several problems in online learning w.r.t. minimax regret. [23] use the sequential Rademacher complexity introduced in [22] to analyze learnability w.r.t. general notions of regret (and not just minimax regret). Rakhlin et al. [24] use the notion of conditional sequential Rademacher complexity to find relaxations of problems like prediction with static experts that immediately lead to algorithms and associated regret guarantees. They show that the random playout strategy has a sound basis and propose a general method to design algorithms as a random playout. In our work, we show that random playout (probability matching) is not just a good strategy, but it is optimal, for the case of $k = 2, 3$ experts. Koolen [16] studies the regret w.r.t. every expert, rather than just the best expert in hindsight and considers tradeoffs in the pareto-frontier. [19] characterize the minimax optimal regret for online linear optimization games as the supremum over the expected value of a function of a martingale difference sequence, and similar characterizations for the minimax optimal algorithm and the adversary. [20] study online linear optimization in Hilbert spaces and characterize minimax optimal algorithms.

2 Preliminaries

Adversary. The adversary at each time t increases the gain of expert $i \in \{1, 2, \dots, k\}$ by a value $g_{it} \in [0, 1]$. Thus adversary decides on $\{g_{i \in [k]t}\}_{t=1}^T$. In particular, for each time t the adversary decides on the distribution \mathcal{D}_t to draw \mathbf{g}_t from. In general, the adversary could be adaptive: i.e., \mathcal{D}_t could depend, apart from the history of gains $\mathbf{g}_{[0,t-1]}$ till time $t - 1$, also on the player's past choices. But for the experts problem, it is known (Lemma 4.1 in [5]) that an *oblivious adversary*, whose distribution \mathcal{D}_t at time t is a function only of $\mathbf{g}_{[0,t-1]}$, is equally powerful⁴. Thus we focus on oblivious adversaries from now on. We denote the joint distribution for all $t \leq T$ as \mathcal{D} . We denote the cumulative gain till time t of expert i by $G_{it} = \sum_{s=1}^t g_{is}$. We denote the vector of cumulative gains at time t by $\mathbf{G}_t = (G_{1t}, \dots, G_{kt})$, and denote the entire history of cumulative gains by $\mathbf{G}_{[0,T]}$.

Player. Before making his decision at time t , the player observes all prior history, that is $\mathbf{g}_{[0,t-1]}$, but doesn't observe g_{it} . He decides on which expert to follow,

and, if the player follows expert i , he gains g_{it} at the end of step t . Specifically, the player decides on the distribution \mathcal{A}_t over experts $\{1, \dots, k\}$. In general, the player could be adaptive: i.e., his distribution \mathcal{A}_t could depend, apart from $\mathbf{g}_{[0,t-1]}$, on his own past choices. But an *oblivious player*, whose distribution \mathcal{A}_t at time t is a function only of $\mathbf{g}_{[0,t-1]}$, is equally powerful. Thus we focus on oblivious players from now on. We use $\mathcal{A}_t(\mathbf{g}_{[0,t-1]})$ to denote the gain of player at time t .

Regret. The stopping time T is known to both the algorithm and the adversary. If the adversary chooses $\mathbf{g}_{[0,T]}$ and the player plays \mathcal{A} , the regret is given by the expression:

$$(2.1) \quad R_T(\mathbf{g}_{[0,T]}, \mathcal{A}) = \max_{i \in [k]} G_{iT} - \sum_{t=1}^T \mathbf{E} \left[\mathcal{A}_t(\mathbf{g}_{[0,t-1]}) \right].$$

If the adversary uses a distribution \mathcal{D} , the regret is given by $R_T(\mathcal{D}, \mathcal{A}) = \mathbf{E}_{\mathbf{g}_{[0,T]} \sim \mathcal{D}} [R_T(\mathbf{g}_{[0,T]}, \mathcal{A})]$.

Minimax regret. The worst-case regret a player playing \mathcal{A} could experience is $\sup_{\mathcal{D}} R_T(\mathcal{D}, \mathcal{A})$. Hence a robust guarantee on the player's regret would be to optimize over \mathcal{A} for worst-case regret, namely, $\inf_{\mathcal{A}} \sup_{\mathcal{D}} R_T(\mathcal{D}, \mathcal{A})$. This is also referred to as the player's minimax regret as he tries to minimize his maximum regret.

Binary adversary. It turns out that an adversary that sets gains in $\{0, 1\}$ (that we call as a binary adversary) is as powerful as an adversary that sets gains in $[0, 1]$ (much like Theorem 10 in Luo and Schapire [18]). Formally, let $\mathcal{D}^{[0,1]}$ be an arbitrary adversary distribution with gains in $[0, 1]$ and let $\mathcal{D}^{\{0,1\}}$ be an arbitrary adversary distribution with gains in $\{0, 1\}$. Basically, we show that $\inf_{\mathcal{A}} \sup_{\mathcal{D}^{\{0,1\}}} R_T(\mathcal{D}^{\{0,1\}}, \mathcal{A}) = \inf_{\mathcal{A}} \sup_{\mathcal{D}^{[0,1]}} R_T(\mathcal{D}^{[0,1]}, \mathcal{A})$ (see Claim 4 in Appendix A). From now on, without loss of generality, we focus only on binary adversaries.

Minimax theorem. Our setting is naturally seen as a two player zero-sum game between the player and the adversary. The player and the adversary, though online in nature, can be described entirely upfront, i.e., by describing their (randomized) actions for every possible history. The set of deterministic strategies for the player is a (huge) finite set, and hence the set of player's randomized strategies is a (huge) simplex. Similarly, the set of adversary's randomized strategies is a (huge) simplex. The regret function is a bilinear function in the player's and adversary's strategies. Thus, the inf and sup can be replaced by min and max, and the famous minimax theorem due to von Neumann [25] applies, telling us that the minimax

⁴For the case of adaptive adversary, there is an alternative definition of regret known as policy regret [10], where this reduction does not apply. However in our setting, we don't use policy regret as it is too powerful and results in a linear regret. Also, for the bandits setting, where the player gets the feedback only about the gains of his chosen action (and not of every action) it is unknown whether adaptive adversaries are any more powerful than oblivious adversaries.

regret of the game is given by

$$(2.2) \quad \min_{\mathcal{A}} \max_{\mathcal{D}} [R_T(\mathcal{D}, \mathcal{A})] = \max_{\mathcal{D}} \min_{\mathcal{A}} [R_T(\mathcal{D}, \mathcal{A})].$$

We refer to the optimal algorithm that defines the LHS as the *minimax optimal algorithm* and similarly, the optimal adversary that defines the RHS as the *minimax optimal adversary*. The minimax optimal algorithm \mathcal{A}^* and the minimax optimal adversary \mathcal{D}^* form a Nash equilibrium: that is, they are mutual best responses.

Balanced adversary. We show that the minimax optimal adversary can, without loss of generality, be “balanced” (see Claim 5 in Appendix A). In other words, for every time t , irrespective of what the history $\mathbf{g}_{[0,t-1]}$ is, the minimax optimal adversary can pick \mathcal{D}_t such that

$\mathbf{E}_{\mathcal{D}_t(\mathbf{g}_{[0,t-1]})} [g_{it}]$ is the same for each expert i . I.e., the expected gains of all experts are equal at every step, irrespective of history.

Dependence on cumulative gains. The minimax optimal algorithm can also choose the distribution \mathcal{A}_t at time t , based only on \mathbf{G}_{t-1} instead of $\mathbf{g}_{[0,t-1]}$. Henceforth, we focus on such algorithms and adversaries, and denote the time t distributions by $\mathcal{A}_t(\mathbf{G}_{t-1})$ and $\mathcal{D}_t(\mathbf{G}_{t-1})$ respectively.

CLAIM 1. *For any balanced adversary \mathcal{D} all algorithms will result in the same regret for the player. In particular, focusing on the algorithm that chooses an expert uniformly at random at every time t , the regret inflicted by \mathcal{D} is given by*

$$\begin{aligned} R_T(\mathcal{D}, \mathcal{A}) &= R_T(\mathcal{D}) \\ &= \mathbf{E}_{\mathbf{g}_{[0,T]} \sim \mathcal{D}} \left[\max_{i \in [k]} G_{iT} - \frac{\sum_{i \in [k]} G_{iT}}{k} \right]. \end{aligned}$$

Given that a minimax optimal adversary \mathcal{D} can always be balanced, the minimax optimal regret is given by

$$R_T(\mathcal{D}) = \mathbf{E}_{\mathbf{g}_{[0,T]} \sim \mathcal{D}} \left[\max_{i \in [k]} G_{iT} - \frac{\sum_{i \in [k]} G_{iT}}{k} \right].$$

Geometric horizon. We introduce the (almost identical) notation that we use for the geometric horizon setting in Section 4.

3 Finite horizon

Two experts: optimal adversary and regret. The optimal regret in the finite horizon setting for the case of $k = 2$ was derived by Cover [9], showing that as $T \rightarrow \infty$, the optimal regret approaches $\sqrt{\frac{T}{2\pi}}$. While Cover also gave an expression (involving a sum and binomial coefficients) for the algorithm’s probabilities, getting just the optimal adversary and the optimal regret value of

$\sqrt{\frac{T}{2\pi}}$ is simpler. We begin by rewriting the expression for the regret from Claim 1 for the case of two experts.

$$\begin{aligned} R_T(\mathcal{D}) &= \mathbf{E}_{\mathbf{g}_{[0,T]} \sim \mathcal{D}} \left[\max_{i=1,2} G_{iT} - \frac{G_{1T} + G_{2T}}{2} \right] \\ &= \mathbf{E}_{\mathbf{g}_{[0,T]} \sim \mathcal{D}} \left[\frac{|G_{1T} - G_{2T}|}{2} \right] \\ (3.3) \quad &= \frac{1}{2} \mathbf{E}_{\mathbf{g}_{[0,T]} \sim \mathcal{D}} \left[\left| \sum_{t=1}^T (g_{1t} - g_{2t}) \right| \right] \end{aligned}$$

The adversary’s optimization problem now is to construct these g_{it} ’s such that they maximize the RHS of equation (3.3), subject to being balanced. This problem is equivalent to the problem of designing a one-dimensional random walk, that respects the constraint that the probability of jumping one step left and one step right are the same, and maximizes the absolute distance from the origin. This equivalence is obtained by interpreting $g_{1t} - g_{2t}$ as the random-walk variable (which can take values of $-1, 1$ and 0), and the condition of being balanced translates to the constraint of jumping left and right with equal probability. We emphasize that the adversary has a control over this random walk, i.e., he can decide separately on the probabilities of jumping left or right at every time step t and every vector of gains \mathbf{G}_t with the restriction to be unbiased towards jumping left or right. Being balanced means that the only design choice left is the probability of staying still (not jumping left or right). To maximize the absolute distance from the origin, the latter probability has to be zero. Indeed, the adversary may as well postpone all his “staying still” turns until the deadline T . In such a case remaining still for the last few steps is not better in expectation than doing random walk. Thus, the optimal adversarial strategy in the 2 experts case is: at every step, choose an expert uniformly at random (with probability $1/2$) and set him to 1, and the other expert to 0.

Given the optimal adversarial strategy description above, the optimal regret in the finite horizon model with T steps is exactly half the expected distance travelled by a simple random walk in T steps, which approaches $\sqrt{\frac{T}{2\pi}}$ as $T \rightarrow \infty$. Thus, $R_T(\mathcal{D}) \rightarrow \sqrt{\frac{T}{2\pi}}$, when $T \rightarrow \infty$.

Two experts: optimal probability matching algorithm. It turns out that the optimal algorithm is precisely a probability matching algorithm, i.e, the algorithm picks each expert with the probability that the respective expert finishes in the lead (we break possible ties in favor of the unique expert who does not have any expert ahead of him in each of the last two steps).

We derive this from an explicit correspondence between simple random walk and the minimax regret value

of games with any given initial configuration of expert cumulative gains. The formal argument is given in Appendix C.2. The probability matching interpretation allows us to give the following simple and explicit description of the optimal algorithm for two experts in the finite horizon model.

THEOREM 3.1. *Let k be the number of remaining time steps and let X be a random variable with a binomial distribution $\text{Binom}(k, \frac{1}{2})$, when k is odd and $\text{Binom}(k - 1, \frac{1}{2})$, when k is even. The optimal algorithm computes the difference $d(\geq 0)$ of cumulative gains between the leading and lagging expert and chooses them with probabilities $p_1(d) = \mathbf{P}[X - \mathbf{E}[X] < d]$, and $p_2(d) = \mathbf{P}[X - \mathbf{E}[X] > d]$.*

k experts: optimal adversary and regret. The simplification afforded by the 2 experts case doesn't carry through for arbitrary k . In Appendix B we have a detailed technical description of the adversary's problem and useful observations about them, including the proofs of claims 2 and 3 below.

CLAIM 2. *For each time step t , the set of all possible distributions $\mathcal{D}_t(\mathbf{G}_{t-1})$ for a balanced adversary forms a convex polytope in 2^k -dimensional space.*

CLAIM 3. *There is a fixed finite set of distributions (over 2^k actions) such that at every time step t and every previous history, the minimax optimal adversary can always choose a distribution from this set.*

k experts: optimal algorithm. We refer the reader to appendix B for an expanded version of the discussion in this subsection. The main algorithmic question in the finite horizon case is if there is a simple description of the optimal algorithm. For the case of $k = 2$ experts, we show that the answer is yes: the optimal algorithm is a simple probability matching algorithm. I.e., the optimal algorithm follows expert i with the probability that, given the current cumulative gains of both the experts and the number of remaining steps, expert i will finish as the leading expert. The derivation of this optimal algorithm is related to how we derive the optimal algorithm for the geometric horizon case. So we do this in Appendix C.1 along with the derivation for geometric horizon.

4 Geometric horizon

Minimax theorem for the geometric horizon model. We use the same notation for the geometric horizon model and the finite horizon model except that we use R_δ for regret in the geometric model instead of R_T . Our setting in the geometric model is again a two player

zero-sum game between the player and the adversary, although the game is not finite now. But a slight generalization of von Neumann's minimax theorem guarantees that the minimax relation we need is still true. For any bilinear function $f(x, y)$ defined on $\mathcal{X} \times \mathcal{Y}$, where \mathcal{X} and \mathcal{Y} are convex and compact sets⁵, we have $\inf_{x \in \mathcal{X}} \sup_{y \in \mathcal{Y}} f(x, y) = \sup_{y \in \mathcal{Y}} \inf_{x \in \mathcal{X}} f(x, y)$. In our case, the space of strategies of the adversary and the algorithm can be easily shown to be convex compact sets. Thus, we have the expected minimax regret of the game given by:

$$(4.4) \quad \inf_{\mathcal{A}} \sup_{\mathcal{D}} [R_\delta(\mathcal{D}, \mathcal{A})] = \sup_{\mathcal{D}} \inf_{\mathcal{A}} [R_\delta(\mathcal{D}, \mathcal{A})].$$

Preliminary claims on the geometric horizon model. All the claims for the finite horizon model, namely, Claims 4, 5, 1, 2 and 3, also carry over to the geometric horizon model with appropriate modifications. We state the modified claims as Claim 6 in Appendix C.

We now derive the optimal adversary, regret and algorithm for the case of 2 experts in appendix C.1. We state our results here.

THEOREM 4.1. *In the geometric horizon model for 2 experts with parameter $\delta \in (0, 1)$:*

1. *The optimal adversary, at every time step, advances the leading expert alone with probability $\frac{1}{2}$ and lagging expert alone with probability $\frac{1}{2}$.*
2. *The optimal regret is $\frac{1-\delta}{2\sqrt{1-(1-\delta)^2}} \rightarrow \frac{1}{2\sqrt{2\delta}}$ as $\delta \rightarrow 0$.*
3. *The optimal algorithm, at every time step, computes the difference $d(\geq 0)$ of cumulative gains between the leading and lagging expert, and chooses them with probabilities $p_1(d) = 1 - \frac{1}{2}\xi^d$, and $p_2(d) = \frac{1}{2}\xi^d$. Here $\xi = \frac{1-\sqrt{1-(1-\delta)^2}}{1-\delta} \sim 1 - \sqrt{2\delta}$.*

We derive the optimal adversary, regret and algorithm for the case of 3 experts in appendix C.3. This derivation is significantly more involved for $k = 3$ when compared to $k = 2$. We state our results here.

THEOREM 4.2. *In the geometric horizon model for 3 experts with parameter $\delta \in (0, 1)$:*

⁵the space of pure strategies for the adversary consists of all infinite sequences of vector gains $\{v_1, v_2, \dots\}$, where each $v_t \in [0, 1]^k$. This space is compact, for example, in a normed space $\left\| \{v_t\}_{t=1}^\infty \right\|_2 = \sum_t \|v_t\|_2^2 / t^2$. Note that in such a normed space the regret is still a continuous function of the sequence $\{v_t\}_{t=1}^\infty$: this is because the geometric-infinite horizon results in a discount of $(1 - \delta)^t$ for round t 's utilities, and this decays much faster than $1/t^2$.

1. The optimal regret is $\frac{2}{3} \frac{1-\delta}{\sqrt{1-(1-\delta)^2}} \rightarrow \frac{2}{3} \frac{1}{\sqrt{2\delta}}$ as $\delta \rightarrow 0$.
2. The optimal algorithm, at every time step, computes the differences d_{ij} between the cumulative gains of experts (i denotes the expert with i th largest cumulative gains, and hence $d_{ij} \geq 0$ for all $i < j$). As a function of the d_{ij} 's the algorithm follows the leading expert with probability $p_1(\mathbf{d}) = 1 - \frac{\xi^{d_{12}}}{2} - \frac{\xi^{d_{13}+d_{23}}}{6}$, the second expert with probability $p_2(\mathbf{d}) = \frac{\xi^{d_{12}}}{2} - \frac{\xi^{d_{13}+d_{23}}}{6}$, and the lagging expert with probability $p_3(\mathbf{d}) = \frac{\xi^{d_{13}+d_{23}}}{3}$. Here $\xi = \frac{1-\sqrt{1-(1-\delta)^2}}{1-\delta} \sim \sqrt{2\delta}$.
3. The optimal adversary, at every time step, computes the differences d_{ij} 's, and follows the following strategy as a function of the d_{ij} 's. Here strategy $\{1\}\{2\}\{3\}$ means exclusively advancing expert 1 (leading expert) with probability $1/3$, expert 2 (middle expert) with probability $\frac{1}{3}$ and expert 3 (lagging expert) with probability $\frac{1}{3}$. Strategy $\{1\}\{23\}$ means advancing expert 1 alone with probability $\frac{1}{2}$ and experts 2 and 3 together with probability $\frac{1}{2}$.

$0 < d_{12} < d_{13} : \{1\}\{23\}$ (any mixture of $\{1\}\{23\}$ with $\{13\}\{2\}$ would also work).

$0 = d_{12} < d_{13} : \{1\}\{23\}$ (any mixture of $\{1\}\{23\}$ with $\{13\}\{2\}$ would also work).

$0 < d_{12} = d_{13} : \{1\}\{23\}$ (any mixture of $\{1\}\{23\}$ with $\{1\}\{2\}\{3\}$ would also work).

$0 = d_{12} = d_{13} : \{1\}\{2\}\{3\}$.

Interpretation as a probability matching algorithm. We show that the optimal algorithms for $k = 2, 3$ can be interpreted as following each expert with the probability he finishes as the leader (probability matching) when following an optimal adversary. We prove this respectively in Appendix C.1 and C.3.

5 Connections to random walk

We already saw for the case of two experts that the optimal strategy for the adversary has a direct connection to random walk. In this section we study larger number of experts, and show that this connection is deep and extends to nontrivial aspects of random walk. We state our results here and prove them (Theorems 5.1 and 5.2) in the full version [12].

Regret Lower Bounds for $k = 3, 4$ experts. While we already have shown in Section 4 that the optimal regret in the case of 3 experts is $\frac{2}{3} \frac{1}{\sqrt{2\delta}}$ as $\delta \rightarrow 0$, the adversary we used there was not the comb adversary. Here we derive

the same regret through the comb adversary. Next, we analyze the comb adversary for $k = 4$ experts and show that as $\delta \rightarrow 0$ it inflicts a regret that is asymptotic to $\frac{\pi}{4} \frac{1}{\sqrt{2\delta}}$.

THEOREM 5.1. *The regret inflicted by the adversary that advances experts 1 and 3 together with probability $\frac{1}{2}$, and, expert 2 with probability $\frac{1}{2}$, as $\delta \rightarrow 0$, is $\frac{2}{3} \frac{1}{\sqrt{2\delta}}$.*

THEOREM 5.2. *The regret inflicted by the adversary that advances experts 1 and 3 together with probability $\frac{1}{2}$, and, experts 2 and 4 together with probability $\frac{1}{2}$, as $\delta \rightarrow 0$, is $\frac{\pi}{4} \frac{1}{\sqrt{2\delta}}$.*

Main idea behind the analysis. We show a bijection between the random process defined by the comb adversary and the simple random walk of a particle between two walls. For $k = 3$, the two walls are “movable”, while for $k = 4$, one wall is “fixed” and the other is movable. I.e., when the particle coincides with the wall and tries to penetrate it in the next step, a movable wall moves one step in the direction of particle’s movement while the particle doesn’t move, but a fixed wall doesn’t move and the particle bounces one step back. The comb adversary’s regret maps to half of the expected number of visits of the particle to one of the movable walls for $k = 3$, and the fixed wall for $k = 4$. Computing the expected number of visits leads to interesting asymptotic analysis.

References

- [1] Jacob Abernethy, Peter L. Bartlett, Alexander Rakhlin, and Ambuj Tewari. Optimal strategies and minimax lower bounds for online convex games. In *21st Annual Conference on Learning Theory - COLT 2008, Helsinki, Finland, July 9-12, 2008*, pages 415–424, 2008.
- [2] Jacob Abernethy, Manfred K. Warmuth, and Joel Yellin. When random play is optimal against an adversary. In *COLT*, pages 437–446, 2008.
- [3] Jacob Abernethy, Alekh Agarwal, Peter L. Bartlett, and Alexander Rakhlin. A stochastic view of optimal regret through minimax duality. In *COLT 2009 - The 22nd Conference on Learning Theory, Montreal, Quebec, Canada, June 18-21, 2009*, 2009.
- [4] Avrim Blum and Adam Kalai. Universal portfolios with and without transaction costs. *Machine Learning*, 35(3):193–205, June 1999. ISSN 0885-6125.
- [5] Jeremy I. Bulow and Jonathan Levin. Matching and price competition. *American Economic Review*, 96: 652–668, 2006.

- [6] Nicolò Cesa-Bianchi. Analysis of two gradient-based algorithms for on-line regression. In *Proceedings of the Tenth Annual Conference on Computational Learning Theory, COLT '97*, pages 163–170, New York, NY, USA, 1997. ACM.
- [7] Nicolò Cesa-Bianchi and Gabor Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, New York, NY, USA, 2006. ISBN 0521841089.
- [8] 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, May 1997. ISSN 0004-5411.
- [9] Thomas M. Cover. Behavior of sequential predictors of binary sequences. In *Proceedings of the 4th Prague Conference on Information Theory, Statistical Decision Functions, Random Processes*, pages 263–272, 1965.
- [10] Ofer Dekel, Ambuj Tewari, and Raman Arora. Online bandit learning against an adaptive adversary: from regret to policy regret. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*, 2012.
- [11] Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.*, 55(1): 119–139, August 1997. ISSN 0022-0000.
- [12] Nick Gravin, Yuval Peres, and Balasubramanian Sivan. Towards optimal algorithms for prediction with expert advice. *CoRR*, abs/1409.3040, 2014. URL <http://arxiv.org/abs/1409.3040>.
- [13] James Hannan. Approximation to bayes risk in repeated play. *Contributions to the Theory of Games*, 3:97–139, 1957.
- [14] David Haussler, Jyrki Kivinen, and Manfred K. Warmuth. Tight worst-case loss bounds for predicting with expert advice. In *EuroCOLT*, pages 69–83, 1995.
- [15] Elad Hazan, Adam Kalai, Satyen Kale, and Amit Agarwal. Logarithmic regret algorithms for online convex optimization. In *Learning Theory, 19th Annual Conference on Learning Theory, COLT 2006, Pittsburgh, PA, USA, June 22-25, 2006, Proceedings*, pages 499–513, 2006.
- [16] Wouter M. Koolen. The pareto regret frontier. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.*, pages 863–871, 2013.
- [17] Nick Littlestone and Manfred K. Warmuth. The weighted majority algorithm. *Information and Computation*, 108(2):212–261, February 1994. ISSN 0890-5401.
- [18] Haipeng Luo and Robert E. Schapire. Towards minimax online learning with unknown time horizon. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, pages 226–234, 2014.
- [19] H. Brendan McMahan and Jacob Abernethy. Minimax optimal algorithms for unconstrained linear optimization. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.*, pages 2724–2732, 2013.
- [20] H. Brendan McMahan and Francesco Orabona. Unconstrained online linear learning in hilbert spaces: Minimax algorithms and normal approximations. In *Proceedings of The 27th Conference on Learning Theory, COLT 2014, Barcelona, Spain, June 13-15, 2014*, pages 1020–1039, 2014.
- [21] Indraneel Mukherjee and Robert E. Schapire. Learning with continuous experts using drifting games. *Theor. Comput. Sci.*, 411(29-30):2670–2683, 2010.
- [22] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning: Random averages, combinatorial parameters, and learnability. In *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010. Proceedings of a meeting held 6-9 December 2010, Vancouver, British Columbia, Canada.*, pages 1984–1992, 2010.
- [23] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning: Beyond regret. In *COLT 2011 - The 24th Annual Conference on Learning Theory, June 9-11, 2011, Budapest, Hungary*, pages 559–594, 2011.
- [24] Alexander Rakhlin, Ohad Shamir, and Karthik Sridharan. Relax and randomize : From value to algorithms. In *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural*

Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States., pages 2150–2158, 2012.

- [25] John von Neumann. Zur theorie der gesellschaftsspiele. *Math Annalen*, 100:295–320, 1928.
- [26] Volodimir G. Vovk. Aggregating strategies. In *Proceedings of the Third Annual Workshop on Computational Learning Theory*, COLT '90, pages 371–386, 1990. ISBN 1-55860-146-5.
- [27] Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Machine Learning, Proceedings of the Twentieth International Conference (ICML 2003), August 21-24, 2003, Washington, DC, USA*, pages 928–936, 2003.

A Proofs from Section 2

CLAIM 4. (BINARY ADVERSARY) *The minimax regret defined by the class of binary adversaries is exactly the same as that defined by general adversaries:*

$$\inf_{\mathcal{A}} \sup_{\mathcal{D}^{\{0,1\}}} R_T(\mathcal{D}^{\{0,1\}}, \mathcal{A}) = \inf_{\mathcal{A}} \sup_{\mathcal{D}^{[0,1]}} R_T(\mathcal{D}^{[0,1]}, \mathcal{A}).$$

Proof. Given that the class of general adversaries is larger than the class of binary adversaries, it immediately follows that $\inf_{\mathcal{A}} \sup_{\mathcal{D}^{\{0,1\}}} R_T(\mathcal{D}^{\{0,1\}}, \mathcal{A}) \leq \inf_{\mathcal{A}} \sup_{\mathcal{D}^{[0,1]}} R_T(\mathcal{D}^{[0,1]}, \mathcal{A})$. It is therefore enough to show that $\inf_{\mathcal{A}} \sup_{\mathcal{D}^{\{0,1\}}} R_T(\mathcal{D}^{\{0,1\}}, \mathcal{A}) \geq \inf_{\mathcal{A}} \sup_{\mathcal{D}^{[0,1]}} R_T(\mathcal{D}^{[0,1]}, \mathcal{A})$. This can be seen as follows: consider the minimax optimal algorithm \mathcal{A}^* for the class of binary adversaries. When faced with a $[0, 1]$ adversary, \mathcal{A}^* , in every round, “discretizes” the gains set by the adversary by independently rounding them to 0 or 1 so that the expectation is equal to the gain g_{it} set by the adversary: i.e., a gain of g_{it} is set to 1 with probability g_{it} and 0 with the remaining probability. From the algorithm \mathcal{A}^* ’s point of view, whether the adversary originally used a distribution with gains in $[0, 1]$ that \mathcal{A}^* discretized to $\{0, 1\}$, or the adversary already set gains in $\{0, 1\}$ with the *same* distribution doesn’t make a difference. Both result in exactly the same expected gains for the algorithm. However, using the discretized version could possibly help the adversary. We see this as follows.

For some step t and history $\mathbf{g}_{[0,t-1]}$, let the adversary set expert’s i gain to be $g_{it} \notin \{0, 1\}$ with non zero probability. Consider the following step-by-step discretization by the adversary. It changes random variable g_{it} (only

for expert i and time t and history $\mathbf{g}_{[0,t-1]}$) to be $\{0, 1\}$ while preserving expectations. While performing this discretization the adversary does not change the distribution in future steps, i.e., it chooses future distributions as if the discretization was not performed. We now show that the expected gain of the best expert can only increase. For each fixed value g_{it} let us denote by ξ a random variable that takes value 1 with probability g_{it} and 0 with probability $1 - g_{it}$. Let us fix all choices of the adversary other than ξ . Then our substitution of constant g_{it} by a random variable ξ can only increase the gain of the best expert $\max_{i \in [k]} G_{iT}$. Indeed, this follows from the inequality

$$\max(\mathbf{E}[\xi] + c_1, c_2) \leq \mathbf{E}[\max(\xi + c_1, c_2)],$$

where c_1 and c_2 are two constants determined by a fixed set of adversary’s random choices. Hence, our modification may only increase the total expected regret, proving the theorem.

CLAIM 5. (BALANCED ADVERSARY) *For each time t and for every possible history $\mathbf{g}_{[0,t-1]}$, the minimax optimal adversary can pick $\mathcal{D}_t(\mathbf{g}_{[0,t-1]})$, such that*

$$\mathbf{E}_{\mathcal{D}_t(\mathbf{g}_{[0,t-1]})}[g_{it}] \text{ is the same for each expert } i.$$

Proof. Given an adversary that is not balanced, we modify it so that algorithm cannot improve, but the expected gain of the best expert $\max_{i \in [k]} G_{iT}$ may only increase. For the minimax optimal adversary \mathcal{D} , one best response algorithm is to choose an expert

$$i^* \in \operatorname{argmax}_{i \in [k]} \mathbf{E}_{\mathcal{D}_t(\mathbf{g}_{[0,t-1]})}[g_{it}].$$

The adversary can modify distribution $\mathcal{D}_t(\mathbf{g}_{[0,t-1]})$ so that for all experts $[k] \setminus \{i^*\}$

$$\mathbf{E}_{\mathcal{D}_t(\mathbf{g}_{[0,t-1]})}[g_{it}] = \mathbf{E}_{\mathcal{D}_t(\mathbf{g}_{[0,t-1]})}[g_{i^*t}],$$

by switching some of the gains from 0 to 1 for $i \in [k] \setminus \{i^*\}$. While making such transformation the adversary does not change \mathcal{D} in the future time steps after t , i.e., the adversary continues as if there was no transformation at time t . The adversary also reveals to the algorithm the value of g_{it} as it was drawn in the original $\mathcal{D}_t(\mathbf{g}_{[0,t-1]})$.

The best response algorithm described above cannot improve its gain at time t , as the expected gain of the best expert does not change. We also note that the algorithm cannot improve in time before t , nor it can improve for the time steps after t , as the knowledge of the algorithm about prior history and the adversary distribution do not change for these times.

On the other hand, the expected gain of the best expert $\max_{i \in [k]} G_{iT}$ could only improve for every such modification of $\mathcal{D}_t(\mathbf{g}_{[0,t-1]})$.

B Proofs and Results from Section 3

B.1 k experts, finite horizon: optimal adversary and regret As mentioned before, the simplification afforded by the 2 experts case doesn't carry through for arbitrary k . Here is the design problem faced by the optimal adversary: for every time step t , given the gains \mathbf{G}_{t-1} at time $t-1$, the adversary has to compute the distribution $\mathcal{D}_t(\mathbf{G}_{t-1})$ at time t so as to maximize the expression for regret given by

$$R_T(\mathcal{D}) = \mathbf{E}_{\mathbf{g}_{[0,T]} \sim \mathcal{D}} \left[\max_{i \in [k]} G_{iT} - \frac{\sum_{i \in [k]} G_{iT}}{k} \right].$$

Note that given any vector of gains \mathbf{G}_{t-1} after $t-1$ time steps, the adversary's distribution \mathcal{D}_t at time t is over 2^k actions corresponding to "setting gain to 0" or "setting gain to 1" for each expert with the restriction that the expected gain of each expert is the same. This design problem of the adversary can be thought of as the design of a controlled random walk on \mathbb{Z}^k so that the advance in each dimension in expectation is the same at every step, with the objective of maximizing the regret expression above.

CLAIM 2. *For each time step t , the set of all possible distributions $\mathcal{D}_t(\mathbf{G}_{t-1})$ for a balanced adversary forms a convex polytope in 2^k -dimensional space.*

Proof. First, note that if two distributions are feasible, a convex combination of them is also feasible. Thus the set of feasible distributions is convex. Second, the feasibility conditions can all be described with linear equalities/inequalities. Finally, the set of feasible distributions is bounded. Thus the set of feasible distributions is a convex polytope.

CLAIM 3. *There is a fixed finite set of distributions (over 2^k actions) such that at every time step t and every previous history, the minimax optimal adversary can always choose a distribution from this set.*

Proof. Given that the set of possible distributions is a convex polytope, at every t , it is a weakly dominant strategy for the adversary to choose from one among the vertices of this polytope. This is because the expression for regret (which the adversary maximizes) is linear in distributions, i.e., a convex combination of two distributions will yield a regret which is the convex combination of the corresponding regrets. Furthermore, this convex polytope of possible distributions remains the same, independent of t and previous history.

REMARK 1. *Note that this polytope of possible distributions has exponentially many vertices. This is easy to see:*

for every subset S of $\{1, \dots, k\}$, treat experts in S as a group and those in \bar{S} as a group. With probability half, set the gains of experts in S to be 1 and those in \bar{S} to be 0, and with the remaining probability do the opposite. Each such distribution is a vertex, and there are exponentially many of them.

REMARK 2. *For concreteness, for the case of $k=3$ and $k=4$, we list all the vertices of the distribution polytopes. While describing a distribution, we shall list only actions in its support, as it turns out the respective probabilities can be reconstructed from the balanced condition for any extremal distribution in our convex polytope. While describing an action, we list the set of experts whom we advance. For instance, the list $\{1\}, \{23\}$ reads as "advance expert 1 with probability 0.5; advance experts 2 and 3 (but not 1) with remaining probability". Similarly, the list $\{234\}\{12\}\{13\}\{14\}$ reads as "with probability $2/5$ advance experts 2, 3, and 4; with probability $1/5$ advance experts 1 and 2; with probability $1/5$ advance experts 1 and 3; with probability $1/5$ advance experts 1 and 4." For $k=3$ and $k=4$ the lists are (excluding the trivial distribution $\{\}$ that advances no experts at all, the distributions $\{123\}$ for $k=3$ and $\{1234\}$ for $k=4$ that advance all the experts together):*

$k = 4$	
$\{123\}\{4\}$	$\{1\}\{2\}\{34\}$
$\{124\}\{3\}$	$\{1\}\{3\}\{24\}$
$\{134\}\{2\}$	$\{1\}\{4\}\{23\}$
$\{234\}\{1\}$	$\{2\}\{3\}\{14\}$
$\{12\}\{34\}$	$\{2\}\{4\}\{13\}$
$\{13\}\{24\}$	$\{3\}\{4\}\{12\}$
$\{14\}\{23\}$	$\{12\}\{134\}\{234\}$
$\{1\}\{23\}\{24\}\{34\}$	$\{13\}\{124\}\{234\}$
$\{2\}\{13\}\{14\}\{34\}$	$\{14\}\{123\}\{234\}$
$\{3\}\{12\}\{14\}\{24\}$	$\{23\}\{124\}\{134\}$
$\{4\}\{12\}\{13\}\{23\}$	$\{24\}\{123\}\{134\}$
$\{123\}\{124\}\{134\}\{234\}$	$\{34\}\{123\}\{124\}$
$\{1\}\{2\}\{3\}\{4\}$	$\{123\}\{14\}\{24\}\{34\}$
	$\{124\}\{13\}\{23\}\{34\}$
	$\{134\}\{12\}\{23\}\{24\}$
	$\{234\}\{12\}\{13\}\{14\}$

$k = 3$
$\{1\}\{23\}$
$\{2\}\{13\}$
$\{3\}\{12\}$
$\{1\}\{2\}\{3\}$
$\{12\}\{13\}\{23\}$

Encouraged by a very simple optimal adversary for $k=2$, one may think that similar behavior extends to

3 or more experts. Unfortunately, this is not the case. The optimal adversary will be time dependent for $k = 3$. For instance, if only one step remains before deadline the optimal adversary would do the following:

- if $G_{1T-1} = G_{2T-1} = G_{3T-1}$, then $\{1\}\{2\}\{3\}$;
- if $G_{1T-1} = G_{2T-1} > G_{3T-1}$, then $\{1\}\{23\}$ or $\{13\}\{2\}$;
- if $G_{1T-1} > G_{2T-1}$, then any balanced strategy.

B.2 k experts, finite horizon: optimal algorithm We note that given a finite time horizon T and finite list of balanced distributions for the adversary, one can write a dynamic program for the maximal value of the regret at any time period $t \leq T$ and initial vector of gains $\mathbf{G}_t \in [T]^k$. We can solve this program by using backward induction over time and furthermore given the regret function at every time step $t \in [T]$ and vector of gains $\mathbf{G}_t \in [T]^k$ we can compute the best strategy for the algorithm. The running time of such an algorithm would be $O(T^k)$. This approach gives us the answer for a small number of experts and reasonably small time horizon T . On the other hand, it becomes impractical as T and especially k get larger, and furthermore, it does not tell us much about intrinsic structure of the optimal algorithm and the optimal adversary.

Even for $k = 2$ the optimal algorithm depends on the time remaining before the deadline T . For example, if the leading expert is ahead of the lagging expert by more than the number of remaining time steps, then the optimal algorithm should always choose the leading expert; on the other hand, if the difference between leading and lagging experts is smaller than the time remaining, then there should be non zero chance of selecting the lagging expert.

The probability matching algorithm. Given this, the main question in the finite horizon case is if there is a simple description of the optimal algorithm. For the case of $k = 2$ experts, we show that the answer is yes: the optimal algorithm is a simple probability matching algorithm. I.e., the optimal algorithm follows expert i with the probability that, given the current cumulative gains of both the experts and the number of remaining steps, expert i will finish as the leading expert.

The derivation of this optimal algorithm is related to how we derive the optimal algorithm for the geometric horizon case. So we do this towards the end of Section C.2.

C Geometric Horizon Model

CLAIM 6. *Observations on the geometric horizon model. The following statements are true:*

1. *The minimax regret defined by the class of binary adversaries is exactly the same as that defined by general adversaries: $\inf_{\mathcal{A}} \sup_{\mathcal{D}^{\{0,1\}}} R_{\delta}(\mathcal{D}^{\{0,1\}}, \mathcal{A}) = \inf_{\mathcal{A}} \sup_{\mathcal{D}^{[0,1]}} R_{\delta}(\mathcal{D}^{[0,1]}, \mathcal{A})$.*
2. *For each time step t and for every possible history $\mathbf{g}_{[0,t-1]}$, the minimax optimal adversary can pick $\mathcal{D}_t(\mathbf{g}_{[0,t-1]})$, such that $\mathbf{E}_{\mathcal{D}_t(\mathbf{g}_{[0,t-1]})}[g_{it}]$ is the same for each expert i .*
3. *A balanced adversary \mathcal{D} inflicts the same regret on every algorithm \mathcal{A} . Since the minimax optimal adversary can always be balanced, the minimax optimal regret is given by:*

$$\begin{aligned} R_{\delta}(\mathcal{D}, \mathcal{A}) &= R_{\delta}(\mathcal{D}) \\ &= \sum_{t=0}^{\infty} \delta \cdot (1 - \delta)^t \\ &= \mathbf{E}_{\mathbf{g}_{[0,t]} \sim \mathcal{D}} \left[\left(\max_{i \in [k]} G_{it} - \frac{\sum_{i \in [k]} G_{it}}{k} \right) \right] \end{aligned}$$

4. *For each time t , the set of all possible distributions $\mathcal{D}_t(\mathbf{G}_{t-1})$ for the adversary forms a convex polytope with exponentially many (in k) vertices.*
5. *There is a fixed finite set of distributions (over 2^k actions) such that at every time t and every previous history, the minimax optimal adversary can always choose a distribution from this set.*

REMARK 3. *The third point in Claim 6 above says that the minimax optimal adversary makes all algorithms achieve the same regret. In particular, if the precise realization of the stopping time random variable was leaked to the algorithm, the minimax optimal regret is not influenced in any way. On the other hand, if the adversary knew the realization of the stopping time information, it could potentially increase the minimax optimal regret. This proves that the algorithm does not benefit from knowing the precise realization of the stopping time information, whereas the adversary could potentially benefit from it.*

C.1 Two experts: optimal algorithm, adversary and regret

Notational convention. At each time step t we always enumerate experts in the decreasing order of their cumulative gains \mathbf{G}_t , i.e., experts 1 and 2 don't refer to identities of experts but to the leading expert and trailing expert respectively. Observe that the strategy of the

optimal adversary at any moment t should not change if cumulative gains \mathbf{G}_t of all experts are changed by the same amount for every expert. Thereby, at every time step t we shall always adjust the total gains \mathbf{G}_t of our experts, so that the leading expert 1 has zero cumulative gain $\widehat{G}_{1t} = 0$. We denote the adjusted gain of the lagging expert $\widehat{G}_{2t} = G_{2t} - G_{1t}$ by x (note that $x \leq 0$).

We denote by $f(x)$ the optimal regret the adversary can obtain for an initial configuration of $\mathbf{G} = (0, x)$, i.e., leading expert has 0 gain and lagging expert has x gain (again, recall that $x \leq 0$). The useful thing about this notation is that if we start at $(0, x)$ for any x , and the game immediately ends at that round, we get a regret of 0: the max expert gain is 0, and the algorithm didn't get any chance to get any gain because the game ended right away. So $0 - 0 = 0$ is the regret.

System of Equations. We are now ready to write our system of equations connecting these $f(x)$'s. Our discussion of Cover's result in Section 3 showed that the minimax optimal adversary in the finite horizon model was independent of the horizon T , and advanced expert 1 or 2 mutually exclusively with probability $\frac{1}{2}$ each. The independence from time horizon T in the finite horizon model immediately means that this adversary is also minimax optimal for the geometric horizon model: it doesn't care when the game ends. This adversary advances the leading expert with probability $\frac{1}{2}$ and lagging expert with probability $\frac{1}{2}$. Thus, starting from the $(0, x)$ configuration, we go to the $(1, x)$ configuration with probability $\frac{1}{2}$ (corresponds to adversary advancing the leading expert), and go to the $(0, x + 1)$ configuration with probability $\frac{1}{2}$ (corresponds to adversary advancing the lagging expert). In the meanwhile, the algorithm would have gained 1 with probability $\frac{1}{2}$ regardless of which expert was advanced. This can be transcribed to an equation right away except that the $(1, x)$ configuration is not in our standard format: our standard format normalizes the largest gain to 0. To perform such a normalization here, notice that the paths of the optimal adversary starting at $(1, x)$ and $(0, x - 1)$ are identical except that the "max expert gain" is precisely one larger when starting from $(1, x)$ than when starting from $(0, x - 1)$. We take this into account in our equations. Summarising this as an equation we get,

$$\begin{aligned} f(x) &= \delta \cdot 0 + (1 - \delta) \\ &\cdot \left[\frac{1}{2} (f(x - 1) + 1) + \frac{1}{2} f(x + 1) - \frac{1}{2} \right] \\ &= (1 - \delta) \cdot \left[\frac{f(x - 1) + f(x + 1)}{2} \right]. \end{aligned}$$

When x is 0 we have to take special care because $(0, x + 1)$ is just $(0, 1)$. First we rewrite gains in the

descending order to obtain the $(1, 0)$ configuration. But this is not in standard format: so we go to the $(0, -1)$ format and add a 1 to the regret in this process. Thus the difference of $(0, 0)$ from $(0, x)$ is that normalization has to be done for both choices of adversary, as against for just one choice. We get,

$$\begin{aligned} f(0) &= \delta \cdot 0 + (1 - \delta) \\ &\cdot \left[\frac{1}{2} (f(0 - 1) + 1) + \frac{1}{2} (f(0 + 1) + 1) - \frac{1}{2} \right] \\ &= (1 - \delta) \cdot \left[f(-1) + \frac{1}{2} \right]. \end{aligned}$$

Combining these two equations, we get the following system:

$$(C.1) \quad f(x) = (1 - \delta) \cdot \frac{f(x - 1) + f(x + 1)}{2}$$

$$(C.2) \quad f(0) = (1 - \delta) \cdot \left(f(-1) + \frac{1}{2} \right)$$

Optimal regret. Thus we need to solve this recurrence relation for $f(x)$. The characteristic polynomial of this recurrence is $x^2 - \frac{2}{1-\delta}x + 1 = 0$, which has two real roots $\xi_1 > 1 > \xi_2$, and $\xi_1 \cdot \xi_2 = 1$, given by $\frac{1 \pm \sqrt{1 - (1-\delta)^2}}{1-\delta}$. The solution to our recurrence relation is then of the form $f(x) = c_1 \cdot \xi_1^x + c_2 \cdot \xi_2^x$. As the regret cannot grow faster than a linear function and cannot be negative, it follows that c_2 must be 0. Combining $f(x) = c_1 \xi_1^x$ with equation (C.2), we get $c_1 \cdot \xi_1^0 = (1 - \delta) \cdot \left(c_1 \cdot \xi_1^{-1} + \frac{1}{2} \right)$. This gives us that $c_1 = \frac{1}{\xi_1 - \xi_2}$. The optimal regret is simply the regret starting at $(0, 0)$, which is given by $f(0)$. Thus the optimal regret is $f(0) = c_1 \xi_1^0 = c_1 = \frac{1}{\xi_1 - \xi_2} = \frac{1 - \delta}{2\sqrt{1 - (1-\delta)^2}}$. Thus, as $\delta \rightarrow 0$, the optimal regret $f(0) \rightarrow \frac{1}{2} \frac{1}{\sqrt{2\delta}}$.

Optimal algorithm. Note that because of minimax principle, we were able to compute the precise regret without even knowing anything about the algorithm. We now proceed to compute the optimal algorithm for $k = 2$. This will reveal how even without knowing the optimal adversary a priori, we can simultaneously discover both the optimal adversary, optimal regret and the optimal algorithm (a useful exercise to the significantly more complicated case of $k = 3$).

Given configuration $(0, x)$ (with $x \leq 0$ as usual), the optimal algorithm assigns probabilities $p_1(x)$ and $p_2(x) = 1 - p_1(x)$ respectively for choosing leading and lagging experts. We drop the arguments for probabilities when it is clear from context. The adversary has four choices, namely advancing expert 1 alone, or expert 2 alone, or both experts, or none of the experts. For $x < 0$

this corresponds to decreasing x by 1, increasing x by 1, not changing x for the last two choices. When $x = 0$, we have to take care of the fact that advancing 1 alone and 2 alone are similar, in that both of them need the normalizing $+1$. Putting what we just described into equations, we get (note the extreme RHS corner gives the adversary's actions corresponding to each expression, and this is common for both $x < 0$ and $x = 0$):

$$(C.3) \quad f(x) = (1-\delta) \cdot \max \begin{cases} f(x-1) + 1 - p_1 & //\{1\} \\ f(x+1) - p_2 & //\{2\} \\ f(x) + 1 - p_1 - p_2 & //\{12\} \\ f(x) & //\{\} \end{cases}$$

$$f(0) = (1-\delta) \cdot \max \begin{cases} f(0-1) + 1 - p_1 & //\{1\} \\ f(0-1) + 1 - p_2 & //\{2\} \\ f(0) + 1 - p_1 - p_2 & //\{12\} \\ f(0) & //\{\} \end{cases}$$

We realize that $p_1(0) = p_2(0) = \frac{1}{2}$ by symmetry (the optimal algorithm is indifferent when the expert gains are the same). By removing strictly suboptimal actions $\{\}, \{12\}$ of the adversary, we obtain

$$f(0) = (1-\delta) \left(f(-1) + \frac{1}{2} \right).$$

Similarly, the first part of expression (C.3) for $x < 0$ boils down to

$$(C.4) \quad f(x) = (1-\delta) \cdot \max \left(f(x-1) + 1 - p_1, f(x+1) - p_2 \right).$$

We further simplify equation (C.4). Notice that the optimal algorithm in minimax equilibrium must make the adversary indifferent between any two actions the adversary is randomizing over. In this case it means that for each $x < 0$ the probabilities $p_1(x)$ and $p_2(x)$ must be chosen by the optimal algorithm in such a way that $f(x) = (1-\delta)(f(x-1) + 1 - p_1(x)) = (1-\delta)(f(x+1) - p_2(x))$. Note that adding these two equations and dividing by 2, we get equation (C.1). Thus we can solve for optimal regret. Additionally, solving for $p_1(x)$ and $p_2(x)$ we obtain

$$(C.5) \quad p_1(x) = 1 - \frac{1}{2} \xi_1^x; \quad p_2(x) = \frac{1}{2} \xi_1^x$$

This proves the optimality of the algorithm 1 for $k = 2$. We summarise our results for $k = 2$ in the following theorem. For convenience we replace the negative number x by positive $d = -x$, and also replace ξ_1 by $\xi = \xi_2 = \frac{1}{\xi_1} \sim 1 - \sqrt{2\delta}$.

THEOREM 4.1. *In the geometric horizon model for 2 experts with parameter $\delta \in (0, 1)$:*

1. *The optimal adversary, at every time step, advances the leading expert alone with probability $\frac{1}{2}$ and lagging expert alone with probability $\frac{1}{2}$.*
2. *The optimal regret is $\frac{1-\delta}{2\sqrt{1-(1-\delta)^2}} \rightarrow \frac{1}{2} \frac{1}{\sqrt{2\delta}}$ as $\delta \rightarrow 0$.*
3. *The optimal algorithm, at every time step, computes the difference $d(\geq 0)$ of cumulative gains between the leading and lagging expert, and chooses them with probabilities $p_1(d) = 1 - \frac{1}{2}\xi^d$, and $p_2(d) = \frac{1}{2}\xi^d$. Here $\xi = \frac{1-\sqrt{1-(1-\delta)^2}}{1-\delta} \sim 1 - \sqrt{2\delta}$.*

C.2 Two experts: interpretation as a probability matching algorithm

Geometric horizon model. The quantity ξ turns out to be precisely equal to the probability that a simple random walk that starts at 1 will reach 0 before the geometric process gets killed. To see this, just note that it is the root of the equation which captures the probability of the above event $\xi = (1-\delta) \cdot 0 + \delta \cdot \frac{1}{2} \cdot (1 + \xi^2)$ (the root that is smaller than 1), which is $\xi = \frac{1-\sqrt{1-(1-\delta)^2}}{1-\delta}$. Now, note that the minimax optimal adversary advances one of the experts uniformly at random and doesn't advance the other. This means that the gap between the cumulative gains of the leading and the lagging experts evolves as a random walk, and the probability that given a separation of d , the lagging expert will match the leading expert is precisely ξ^d . Once they match, each expert has an equal probability $\frac{1}{2}$ of being the leading expert⁶. This means, the probability that the currently lagging expert will finish as the leading expert is precisely $\frac{1}{2}\xi^d$, and the probability that the currently leading expert will finish as the leading expert is $1 - \frac{1}{2}\xi^d$.

Finite horizon model. We now show that for the finite horizon case too, the optimal algorithm is precisely a probability matching algorithm, i.e., the algorithm picks each expert with the probability that the respective expert finishes in the lead (we break possible ties in favor of the unique expert who doesn't have any expert ahead of him in each of the last two steps). We set up equations very similar to (C.3) and (C.4) in the finite horizon model

⁶If the experts are tied, the leader is chosen uniformly at random

except that f now will be a function of both x and the number of time steps left $\ell = T - t$ until the deadline. Thus

$$(C.6) \quad \begin{aligned} f(x, 0) &= 0, & \text{if } x \leq 0 \\ f(x, \ell) &= \frac{f(x+1, \ell-1) + f(x-1, \ell-1)}{2}, & \text{if } \ell > 0, \text{ and } x < 0 \\ f(0, \ell) &= f(-1, \ell-1) + \frac{1}{2}, & \text{if } \ell > 0, x = 0. \end{aligned}$$

We consider a simple random walk $SRW(x, \ell)$ that starts from position x and does ℓ steps (we also use $SRW(x, \ell)$ to denote the location of this walk after ℓ steps). It turns out that $g(x, \ell) = \frac{\mathbf{E}[|SRW(x, \ell)| - |x|]}{2}$ satisfies exactly the same set of equations (C.6) as $f(x, \ell)$ does. Thus $f(x, \ell) = g(x, \ell)$. Analogously to the geometric model, we can also derive that $p_2(x, \ell) = f(x+1, \ell-1) - f(x, \ell)$ for $x < 0$ and $p_2(0, \ell) = p_1(0, \ell) = \frac{1}{2}$. We immediately get the desired probability matching result for $x = 0$. To get the same for $x < 0$, we do a natural coupling of random walks $SRW(x+1, \ell-1)$ and $SRW(x, \ell)$ in the expression $p_2(x, \ell) = g(x+1, \ell-1) - g(x, \ell)$. When $SRW(x+1, \ell-1)$ arrives at y in this coupling, $SRW(x, \ell)$ does one more iteration from the location $y-1$. The expression $g(x+1, \ell-1) - g(x, \ell)$, given that $SRW(x+1, \ell-1)$ arrives at y can be written as:

$$\begin{aligned} & \frac{|y| - |x+1|}{2} - \frac{\frac{1}{2}|y-2| + \frac{1}{2}|y| - |x|}{2} \\ &= \frac{\frac{1}{2}(|y| - |y-2|) + 1}{2} \\ &= \begin{cases} 0 & \text{if } y \leq 0 \\ 1/2 & \text{if } y = 1 \\ 1 & \text{if } y > 1 \end{cases} \end{aligned}$$

The first line in the RHS of the above expression corresponds to the situations where $SRW(x, \ell)$ arrives at $y-1 < 0$ at step $\ell-1$, i.e., the second expert does not reach the leader till step $\ell-1$ (and therefore the first expert is the unique one who didn't lag in steps $\ell-1$ and ℓ); the second line corresponds to the situations where the second expert reaches the leader at step $\ell-1$ (and therefore overtakes him with probability $1/2$ in the last step); the third line represents situations when the second expert is the unique leader after $\ell-1$ steps (and therefore the second expert is the unique one that didn't lag in steps $\ell-1$ and ℓ). This yields the desired probability matching result.

Uniqueness of the optimal algorithm. In the geometric horizon model, we explicitly solve the infinite system of equations and realize that they have a unique solution proving the uniqueness of the optimal algorithm. In the finite horizon model, although we don't explicitly solve the system of equations, the discussion in the previous paragraph shows that the probabilities chosen by the optimal algorithm are unique, and hence the optimal algorithm is unique.

C.3 Three experts, geometric horizon: optimal algorithm, adversary and regret We derive the optimal adversary, algorithm and regret here. We restate Theorem 4.2 for ease of reading.

THEOREM 4.2. *In the geometric horizon model for 3 experts with parameter $\delta \in (0, 1)$:*

1. *The optimal regret is $\frac{2}{3} \frac{1-\delta}{\sqrt{1-(1-\delta)^2}} \rightarrow \frac{2}{3} \frac{1}{\sqrt{2\delta}}$ as $\delta \rightarrow 0$.*
2. *The optimal algorithm, at every time step, computes the differences d_{ij} between the cumulative gains of experts (i denotes the expert with i th largest cumulative gains, and hence $d_{ij} \geq 0$ for all $i < j$). As a function of the d_{ij} 's the algorithm follows the leading expert with probability $p_1(\mathbf{d}) = 1 - \frac{\xi^{d_{12}}}{2} - \frac{\xi^{d_{13}+d_{23}}}{6}$, the second expert with probability $p_2(\mathbf{d}) = \frac{\xi^{d_{12}}}{2} - \frac{\xi^{d_{13}+d_{23}}}{6}$, and the lagging expert with probability $p_3(\mathbf{d}) = \frac{\xi^{d_{13}+d_{23}}}{3}$. Here $\xi = \frac{1-\sqrt{1-(1-\delta)^2}}{1-\delta} \sim 1 - \sqrt{2\delta}$.*
3. *The optimal adversary, at every time step, computes the differences d_{ij} 's, and follows the following strategies below as a function of the d_{ij} 's. Strategy $\{1\}\{2\}\{3\}$ means: exclusively advancing with probability $1/3$ expert 1 (leading expert), expert 2 (middle expert), and expert 3 (lagging expert). Strategy $\{1\}\{23\}$ means: advancing with probability $\frac{1}{2}$ expert 1 alone, or advancing experts 2 and 3 together.*

$$0 < d_{12} < d_{13} : \{1\}\{23\}, \text{ or } \{13\}\{2\}.$$

$$0 = d_{12} < d_{13} : \{1\}\{23\}, \text{ or } \{13\}\{2\}.$$

$$0 < d_{12} = d_{13} : \{1\}\{23\}, \text{ or } \{1\}\{2\}\{3\}.$$

$$0 = d_{12} = d_{13} : \{1\}\{2\}\{3\}.$$

Notational convention. At each time period t we always enumerate experts in the decreasing order of their cumulative gains \mathbf{G}_t . We observe that the strategy of the adversary at any moment t should not change if cumulative gains \mathbf{G}_t of all experts are changed by the same amount for every expert. Thereby, at every time

step t we shall always adjust the total gains \mathbf{G}_t of our experts, so that the leading expert 1 has zero cumulative gain $G_{1t} = 0$. We denote the adjusted gain $G_{i+1t} - G_{1t}$ by $x_i(t)$ for each $i \in [2]$; we denote by $\mathbf{x}(t) = (x_1(t), x_2(t))$ the vector of adjusted gains. Note that both $x_1(t)$ and $x_2(t)$ are negative.

We denote by $f(\mathbf{x})$ the optimal regret the adversary can obtain for an initial configuration $\mathbf{G} = (0, x_1, x_2)$ (where again $x_1, x_2 \leq 0$). Much like the case of $k = 2$,

the advantage of this convention is that if we start from configuration $(0, x_1, x_2)$ and stop immediately, the “max-expert-gain - algorithm’s gain” is just 0.

Algorithm assigns probabilities $p_1(\mathbf{x})$, $p_2(\mathbf{x})$, and $p_3(\mathbf{x}) = 1 - p_1(\mathbf{x}) - p_2(\mathbf{x})$ respectively to the leading, middle and lagging experts. Similarly to the case $k = 2$ the adversary now has eight choices and the regret satisfies the following expression for each $\mathbf{x} : 0 > x_1 > x_2$.

$$(C.7) \quad f(x_1, x_2) = \delta \cdot 0 + (1 - \delta) \cdot \max \begin{cases} f(x_1 - 1, x_2 - 1) + 1 - p_1 & // \{1\} \\ f(x_1 + 1, x_2 + 1) - p_2 - p_3 & // \{23\} \\ f(x_1 - 1, x_2) + 1 - p_1 - p_3 & // \{13\} \\ f(x_1 + 1, x_2) - p_2 & // \{2\} \\ f(x_1, x_2 - 1) + 1 - p_1 - p_2 & // \{12\} \\ f(x_1, x_2 + 1) - p_3 & // \{3\} \\ f(x_1, x_2) + 1 - p_1 - p_2 - p_3 & // \{123\} \\ f(x_1, x_2) & // \{\} \end{cases}$$

We note that we can omit the lines $\{123\}$ and $\{\}$ in the RHS of the expression above. For the boundary points $\mathbf{x} : 0 = x_1 > x_2$ and $\mathbf{x} : 0 > x_1 = x_2$ we need to take into account in (C.7) the possibility that the order $0 \geq x_1 \geq x_2$ might change:

$$(C.8) \quad \frac{f(0, x_2)}{1 - \delta} = \max \begin{cases} f(-1, x_2 - 1) + 1 - p_1 \\ f(-1, x_2) + 1 - p_2 - p_3 \\ f(-1, x_2) + 1 - p_1 - p_3 \\ f(-1, x_2 - 1) + 1 - p_2 \\ f(0, x_2 - 1) + 1 - p_1 - p_2 \\ f(0, x_2 + 1) - p_3 \end{cases} \quad \frac{f(x, x)}{1 - \delta} = \max \begin{cases} f(x - 1, x - 1) + 1 - p_1 & // \{1\} \\ f(x + 1, x + 1) - p_2 - p_3 & // \{23\} \\ f(x, x - 1) + 1 - p_1 - p_3 & // \{13\} \\ f(x + 1, x) - p_2 & // \{2\} \\ f(x, x - 1) + 1 - p_1 - p_2 & // \{12\} \\ f(x + 1, x) - p_3 & // \{3\} \end{cases}$$

For $\mathbf{x} : 0 = x_1 = x_2$ we have

$$(C.9) \quad f(0, 0) = (1 - \delta) \max \begin{cases} f(-1, -1) + 1 - p_1 & // \{1\} \\ f(0, -1) + 1 - p_2 - p_3 & // \{23\} \\ f(0, -1) + 1 - p_1 - p_3 & // \{13\} \\ f(-1, -1) + 1 - p_2 & // \{2\} \\ f(0, -1) + 1 - p_1 - p_2 & // \{12\} \\ f(-1, -1) + 1 - p_3 & // \{3\} \end{cases}$$

Our approach here will be a guess and verify approach. While there are several strategies possible for the adversary, we discovered that the optimal strategy for the adversary is to play $\{1\}$, $\{23\}$ or $\{13\}$, $\{2\}$ for most of the \mathbf{x} (at least for those $\mathbf{x} : 0 > x_1 > x_2$). We will now compute the consequences of this being the optimal adversary and finally verify if our guess was true. So playing $\{1\}$, $\{23\}$ or $\{13\}$, $\{2\}$ for most of the time means that for $\mathbf{x} : 0 > x_1 > x_2$ we have

$$(C.10) \quad \begin{aligned} f(x_1, x_2) &= \frac{1-\delta}{2} \left[f(x_1-1, x_2) + f(x_1+1, x_2) \right] \\ f(x_1, x_2) &= \frac{1-\delta}{2} \left[f(x_1-1, x_2-1) + \right. \\ &\quad \left. f(x_1+1, x_2+1) \right]. \end{aligned}$$

One can write generating function for $f(x_1, x_2)$:

$$G(u, v) = \sum_{x_1, x_2} f(x_1, x_2) u^{x_1} v^{x_2}$$

We can write two functional relations on $G(u, v)$ from expression (C.10) and further derive a parametric expression for $f(x_1, x_2)$:

$$f(x_1, x_2) = c_1 \cdot \xi_1^{x_1} + c_2 \cdot \xi_1^{2x_2-x_1} + c_3 \cdot \xi_1^{-x_1} + c_4 \cdot \xi_1^{x_1-2x_2},$$

where c_1, c_2, c_3, c_4 are unknown parameters and $\xi_1 > 1 > \xi_2$ are the roots of the characteristic polynomial $x^2 - \frac{2}{1-\delta}x + 1 = 0$. It turns out that, as the regret cannot grow faster than a linear function and cannot be negative, it follows that c_4 and also c_3 must be 0.

From the algorithm's point of view, the probabilities $p_1(\mathbf{x})$, $p_2(\mathbf{x})$, and $p_3(\mathbf{x})$ must be chosen in such a way that adversary will be indifferent between playing $\{1\}$, $\{23\}$, $\{13\}$, and $\{2\}$ for $\mathbf{x} : 0 > x_1 > x_2$. From this condition we derive that

$$\begin{aligned} 1 - p_1 &= \left(\frac{\xi_1 - \xi_2}{2} \right) \left(c_1 \cdot \xi_1^{x_1} + c_2 \cdot \xi_1^{2x_2-x_1} \right) \\ p_2 &= \left(\frac{\xi_1 - \xi_2}{2} \right) \left(c_1 \cdot \xi_1^{x_1} - c_2 \cdot \xi_1^{2x_2-x_1} \right) \\ p_3 &= \left(\frac{\xi_1 - \xi_2}{2} \right) \left(2c_2 \cdot \xi_1^{2x_2-x_1} \right) \end{aligned}$$

We also assume that the above formula for $\mathbf{p}(\mathbf{x})$ extends to the points of the form $\mathbf{x} : 0 = x_1 > x_2$ and $\mathbf{x} : 0 > x_1 = x_2$ and $\mathbf{x} : 0 = x_1 = x_2$. We equate

$p_1(\mathbf{x})$ and $p_2(\mathbf{x})$ for $\mathbf{x} : 0 = x_1 > x_2$, as now leading and middle experts are identical from the adversary's point of view. Similarly, we equate $p_2(\mathbf{x})$ and $p_3(\mathbf{x})$ for $\mathbf{x} : 0 > x_1 = x_2$; and equate $p_1(\mathbf{x})$, $p_2(\mathbf{x})$ and $p_3(\mathbf{x})$ for $\mathbf{x} : 0 = x_1 = x_2$. From these equations we deduce that

$$c_1 = \frac{1}{\xi_1 - \xi_2}; \quad c_2 = \frac{1}{3(\xi_1 - \xi_2)}.$$

This results in the following expression for the regret $f(\mathbf{x})$:

$$(C.11) \quad f(\mathbf{x}) = \frac{\xi_1^{x_1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2-x_1}}{3(\xi_1 - \xi_2)},$$

which gives us regret of $\frac{4}{3(\xi_1 - \xi_2)}$ at $\mathbf{x} = (0, 0)$. As $\delta \rightarrow 0$ the regret

$$R_\delta = \frac{4}{3(\xi_1 - \xi_2)} = \frac{2(1-\delta)}{3\sqrt{\delta} \cdot (2-\delta)} \xrightarrow{\delta \rightarrow 0} \frac{2}{3\sqrt{2\delta}}.$$

THEOREM C.1. *Equation (C.11) gives the precise value of the regret for every normalized $\mathbf{x} : 0 \geq x_1 \geq x_2$. Moreover, the optimal algorithm chooses leading, middle and lagging experts respectively with the following probabilities $p_1(\mathbf{x})$, $p_2(\mathbf{x})$, and $p_3(\mathbf{x})$:*

$$(C.12) \quad \begin{aligned} 1 - p_1(\mathbf{x}) &= \frac{\xi_1^{x_1}}{2} + \frac{\xi_1^{2x_2-x_1}}{6} \\ p_2(\mathbf{x}) &= \frac{\xi_1^{x_1}}{2} - \frac{\xi_1^{2x_2-x_1}}{6} \\ p_3(\mathbf{x}) &= \frac{\xi_1^{2x_2-x_1}}{3} \end{aligned}$$

Proof. To prove this theorem we shall first verify that the function $f(\cdot)$ given by (C.11) together with the probabilities (C.12) satisfies combined system of equations (C.7), (C.8), (C.9) for every $\mathbf{x} : 0 \geq x_1 \geq x_2$.

Then the expression (C.11) immediately gives us an upper bound on the regret function $f(\mathbf{x})$. Indeed, if we fix strategy of the algorithm to be as in (C.12), then $f(\mathbf{x})$ would be an upper bound on the regret that the best response adversary (with respect to this fixed algorithm) could get.

Finally, to show matching lower bound we will consider the best response strategy of the adversary in (C.11), i.e., those lines in RHS of (C.7), (C.8), (C.9) which are equal to LHS. We will make sure that among these strategies the adversary can always compose a mixed strategy which is balanced, i.e. the one that makes algorithm completely indifferent between all experts. Assume that we have restricted our adversary to these mixed strategies.

Then any algorithm will be the best response algorithm, in particular the algorithm defined by (C.12). Hence, this particular restricted strategy of the adversary provides a lower bound given by (C.11) on the regret function $f(\mathbf{x})$.

We begin by verifying (C.7) for the interior points

$\mathbf{x} : 0 > x_1 > x_2$.

LEMMA C.1. Equation (C.7) holds true for the interior points $\mathbf{x} : 0 > x_1 > x_2$.

$$\frac{\xi_1^{x_1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} = (1 - \delta) \cdot \max \left\{ \begin{array}{ll} \frac{\xi_1^{x_1 - 1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1 - 1}}{3(\xi_1 - \xi_2)} + \frac{\xi_1^{x_1}}{2} + \frac{\xi_1^{2x_2 - x_1}}{6} & // \{1\} \\ \frac{\xi_1^{x_1 + 1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1 + 1}}{3(\xi_1 - \xi_2)} - \frac{\xi_1^{x_1}}{2} - \frac{\xi_1^{2x_2 - x_1}}{6} & // \{23\} \\ \frac{\xi_1^{x_1 - 1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1 + 1}}{3(\xi_1 - \xi_2)} + \frac{\xi_1^{x_1}}{2} - \frac{\xi_1^{2x_2 - x_1}}{6} & // \{13\} \\ \frac{\xi_1^{x_1 + 1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1 - 1}}{3(\xi_1 - \xi_2)} - \frac{\xi_1^{x_1}}{2} + \frac{\xi_1^{2x_2 - x_1}}{6} & // \{2\} \\ \frac{\xi_1^{x_1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1 - 2}}{3(\xi_1 - \xi_2)} + \frac{\xi_1^{2x_2 - x_1}}{3} & // \{12\} \\ \frac{\xi_1^{x_1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1 + 2}}{3(\xi_1 - \xi_2)} - \frac{\xi_1^{2x_2 - x_1}}{3} & // \{3\} \\ \frac{\xi_1^{x_1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} & // \{123\} \\ \frac{\xi_1^{x_1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} & // \{\} \end{array} \right.$$

Proof. Clearly, the lines $\{\}$ and $\{123\}$ in the RHS are smaller than the LHS. Since we have chosen $f(\mathbf{x})$ according to (C.10), it immediately follows that the average of the lines $\{1\}$ and $\{23\}$ in RHS as well as average of the lines $\{13\}$ and $\{2\}$ in RHS are equal to the the LHS. We further notice that p_1, p_2, p_3 were chosen so that the RHS expressions in lines $\{1\}$ and $\{23\}$ are equal as well as are equal expressions in lines $\{13\}$ and $\{2\}$. This makes the

expressions in lines 1-4 in the RHS to be equal to the LHS.

We are only left to verify that LHS is greater than or equal to the expressions in the lines 5 and 6 in the RHS. We recall that $\xi_1 > 1 > \xi_2$ are the roots of the polynomial $x^2 - \frac{2}{1-\delta}x + 1$, so that $\xi_1 \cdot \xi_2 = 1$ and $\xi_1 + \xi_2 = \frac{2}{1-\delta}$.

For the line $\{12\}$ in RHS we need to verify the following.

$$\frac{1}{1 - \delta} \left(\frac{\xi_1^{x_1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} \right) \geq \frac{\xi_1^{x_1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1 - 2}}{3(\xi_1 - \xi_2)} + \frac{\xi_1^{2x_2 - x_1}}{3}$$

Equivalently, we need to show

$$\begin{aligned} \left(\frac{\delta}{1 - \delta} \right) \frac{\xi_1^{x_1}}{\xi_1 - \xi_2} &\geq \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} \left(\xi_2^2 - \frac{1}{1 - \delta} + \xi_1 - \xi_2 \right) \\ &= \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} \left(\xi_2^2 - \frac{1}{1 - \delta} + \frac{2}{1 - \delta} - 2\xi_2 \right) \\ &= \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} \left(\frac{2}{1 - \delta} \xi_2 - 1 - 2\xi_2 + \frac{1}{1 - \delta} \right) \\ &= \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} \left(\frac{\delta}{1 - \delta} \right) (2\xi_2 + 1). \end{aligned}$$

The last inequality holds true as

$$\xi_1^{2x_1 - 2x_2} \geq 1 \geq \frac{2\xi_2 + 1}{3}.$$

Similarly for the line $\{3\}$ in RHS we need to verify that

$$\frac{1}{1 - \delta} \left(\frac{\xi_1^{x_1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} \right) \geq \frac{\xi_1^{x_1}}{\xi_1 - \xi_2} + \frac{\xi_1^{2x_2 - x_1 + 2}}{3(\xi_1 - \xi_2)} - \frac{\xi_1^{2x_2 - x_1}}{3}$$

After some transformation we need to show that

$$\begin{aligned}
\left(\frac{\delta}{1-\delta}\right) \frac{\xi_1^{x_1}}{\xi_1 - \xi_2} &\geq \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} \left(\xi_1^2 - \frac{1}{1-\delta} - \xi_1 + \xi_2\right) \\
&= \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} \left(\xi_1^2 - \frac{1}{1-\delta} + \frac{2}{1-\delta} - 2\xi_1\right) \\
&= \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} \left(\frac{2}{1-\delta}\xi_1 - 1 - 2\xi_1 + \frac{1}{1-\delta}\right) \\
&= \frac{\xi_1^{2x_2 - x_1}}{3(\xi_1 - \xi_2)} \left(\frac{\delta}{1-\delta}\right) (2\xi_1 + 1).
\end{aligned}$$

We further compare LHS with RHS of the last inequality. We need to prove that

$$\xi_1^{2x_1 - 2x_2} \geq \frac{2\xi_1 + 1}{3}.$$

Since $x_1 > x_2$, we observe that $\xi_1^{2x_1 - 2x_2} \geq \xi_1^2$. The desired inequality is true, as $\xi_1^2 \geq \frac{2\xi_1 + 1}{3}$.

We next consider a few cases for boundary points when there are ties between leading, middle and legging experts. However, if there is no change in the order of experts after adversary's action, most of our derivations in Lemma C.1 applies to the boundary cases as well.

LEMMA C.2. *Equation (C.8) holds true for the boundary points $\mathbf{x} : 0 = x_1 > x_2$.*

Proof. We note that the expressions in the lines $\{1\}$, $\{13\}$, $\{12\}$, and $\{3\}$ of (C.8) are the same as in Lemma C.1 for $x_1 = 0$, since the order of leading, middle, and legging experts does not change for any of these choices of the adversary.

We observe that $p_1 = 1 - \frac{1}{2} - \frac{\xi_1^{2x_2}}{6} = \frac{1}{2} - \frac{\xi_1^{2x_2}}{6} = p_2$ for boundary points $\mathbf{x} : 0 = x_1 > x_2$. Furthermore, as first two experts are the same, leading and middle expert are equivalent from the perspective of the adversary. It implies that lines $\{1\}$ and $\{2\}$ as well as lines $\{13\}$ and $\{23\}$ in the RHS of (C.8) are identical. We conclude the proof by observing that

1. LHS is equal to the line $\{1\}$ in RHS (same argument as in Lemma C.1), which is equal to the expression in the line $\{2\}$ of RHS.
2. LHS is equal to the expression in the line $\{13\}$ in RHS (same argument as in Lemma C.1), which is the same as the line $\{23\}$ in RHS.
3. LHS is at least the expressions in the lines $\{12\}$ and $\{3\}$ (same argument as in Lemma C.1). Indeed, we

only used the fact that $x_1 - x_2 > 0$ and analytically all the rest derivations remain the same as in Lemma C.1.

LEMMA C.3. *Equation (C.8) holds true for the boundary points $\mathbf{x} : 0 > x = x_1 = x_2$.*

Proof. We note that the expressions in the lines $\{1\}$, $\{2\}$, $\{12\}$, and $\{23\}$ of (C.8) are the same as in Lemma C.1 for $x = x_1 = x_2$, since the order of leading, middle, and legging experts does not change for any of these choices of the adversary.

We observe that $p_2 = \frac{\xi_1^x}{2} - \frac{\xi_1^x}{6} = \frac{\xi_1^x}{3} = p_3$ for boundary points $\mathbf{x} : 0 > x = x_1 = x_2$. Furthermore, as last two experts are the same, middle and legging experts are equivalent from the perspective of the adversary. It implies that lines $\{2\}$ and $\{3\}$ as well as lines $\{13\}$ and $\{12\}$ in the RHS of (C.8) are identical. We conclude the proof by observing that

1. LHS is equal to the expression in the line $\{2\}$ in RHS (same argument as in Lemma C.1), which is equal to the expression in the line $\{3\}$ of RHS.
2. LHS is at least the expression in the line $\{12\}$ in RHS (same argument as in Lemma C.1), which is the same as the line $\{13\}$ in RHS. Indeed, for the line $\{12\}$ we don't need x_1 to be strictly greater than x_2 and our derivations as in Lemma C.1 do not change.
3. LHS is equal to the expressions in the lines $\{12\}$ and $\{3\}$ (same argument as in Lemma C.1).

LEMMA C.4. *Equation (C.9) holds true for the boundary point $\mathbf{x} : 0 = x_1 = x_2$.*

Proof. We observe that $p_1 = 1 - \frac{\xi_1^0}{2} - \frac{\xi_1^0}{6} = \frac{\xi_1^0}{2} - \frac{\xi_1^0}{6} = p_2 = \frac{\xi_1^0}{3} = p_3$. Therefore, the lines $\{1\}$, $\{2\}$, and $\{3\}$ in RHS are identical, similarly are identical the lines

{12}, {23}, and {13}. We also notice that expression in the lines {1} and {12} in (C.9) are special cases of the corresponding expressions in (C.7) for $x_1 = x_2 = 0$. Expression in the line {12} in RHS is not greater than LHS, because our derivations from Lemma C.1 analytically remain the same and for the expression in line {12} we only need $x_1 \geq x_2$.

We conclude the proof by observing that

1. LHS is equal to the expression in the line {1} in RHS (same argument as in Lemma C.1), which is equal to the expressions in the lines {2} and {3} in RHS.
2. LHS is at least the expression in the line {12} in RHS (same argument as in Lemma C.1), which is the same as the lines {13} and {23} in RHS.

We summarize below the best choices for the adversary (lines in RHS of (C.7),(C.8),(C.9) which are equal to LHS).

$$\mathbf{x} : 0 > x_1 > x_2 \quad \{1\}, \{23\}, \{13\}, \{2\}.$$

$$\mathbf{x} : 0 = x_1 > x_2 \quad \{1\}, \{2\}, \{13\}, \{23\}.$$

$$\mathbf{x} : 0 > x_1 = x_2 \quad \{1\}, \{23\}, \{2\}, \{3\}.$$

$$\mathbf{x} : 0 = x_1 = x_2 \quad \{1\}, \{2\}, \{3\}.$$

The corresponding mixed balanced strategies of the adversary are:

$$\mathbf{x} : 0 > x_1 > x_2 \quad \{1\}\{23\}, \text{ or } \{13\}\{2\}.$$

$$\mathbf{x} : 0 = x_1 > x_2 \quad \{1\}\{23\}, \text{ or } \{13\}\{2\}.$$

$$\mathbf{x} : 0 > x_1 = x_2 \quad \{1\}\{23\}, \text{ or } \{1\}\{2\}\{3\}.$$

$$\mathbf{x} : 0 = x_1 = x_2 \quad \{1\}\{2\}\{3\}.$$

This concludes the proof of Theorem C.1 and hence Theorem 4.2.

D Comparison with multiplicative weights algorithm

In this section, we show that the optimal algorithm is not in the family of multiplicative weight algorithms.

Multiplicative weights algorithm (MWA). Given cumulative gains $G_{1t-1}, \dots, G_{kt-1}$ for the k experts after $t-1$ steps, MWA computes the exponentials of these cumulative gains and follows expert i with probability proportional to these exponentials. Formally, MWA at time t follows expert i with probability $\frac{\exp(\eta G_{it-1})}{\sum_j \exp(\eta G_{jt-1})}$, where η is a parameter that can be tuned. For the special case of 2 experts, this description can be simplified: let $d(t-1) = G_{1t-1} - G_{2t-1}$ where we use 1 and 2 denote the leading and lagging experts respectively. Then, MWA follows the leading expert with probability $\frac{e^{\eta d(t-1)}}{e^{\eta d(t-1)} + 1}$ and the lagging expert with probability $\frac{1}{e^{\eta d(t-1)} + 1}$.

D.1 Optimal algorithm is not in the MWA family

Even for $k = 2$ experts, the optimal algorithm is not in the MWA family. From the tuple representation of MWA and OPT, namely,

$$\text{MWA: } \left(\frac{e^{\eta d}}{e^{\eta d} + 1}, \frac{1}{e^{\eta d} + 1} \right), \text{ and OPT: } \left(1 - \frac{1}{2}\xi^d, \frac{1}{2}\xi^d \right),$$

it is clear that the optimal algorithm cannot be expressed as a multiplicative weights algorithm. We now show that even a convex combination of MWAs cannot express it.

FACT D.1. *No convex combination of multiplicative weight algorithms can express the optimal algorithm.*

Proof. We show that even if at every step, the parameter η was allowed to be drawn from a measure μ , MWA cannot express the optimal algorithm, i.e., for no measure μ can we have that for all integer $d \geq 0$, $\int_{-\infty}^{\infty} \frac{d\mu(\eta)}{e^{\eta d} + 1} = \frac{1}{2}\xi^d$, or equivalently, $\int_{-\infty}^{\infty} \frac{d\mu(\eta)}{(\xi e^{\eta})^d + \xi^d} = \frac{1}{2}$. Let η_0 be such that $\xi e^{\eta_0} = 1$ (note that $\xi < 1$). The measure on $\{\eta : \eta < \eta_0\}$ should be 0 for otherwise the denominator in the integral goes to 0 as $d \rightarrow \infty$, which will make the integral go to ∞ where as the RHS is just $\frac{1}{2}$. Likewise, any measure on $\eta : \eta > \eta_0$ doesn't contribute to the integral as $d \rightarrow \infty$ since the integral will anyway be 0 in the region (η_0, ∞) . Thus, for the integral to be $\frac{1}{2}$ as $d \rightarrow \infty$, we need $\mu(\eta_0) = \frac{1}{2}$. We now expand twice the LHS, namely, the integral $\int_{\eta_0}^{\infty} \frac{2d\mu(\eta)}{(\xi e^{\eta})^d + \xi^d}$ by splitting it into two terms: the first term is the integral at η_0 where the measure is $\mu(\eta_0) = \frac{1}{2}$ and $\xi e^{\eta_0} = 1$, and the second term is the integral in the region (η_0, ∞) . So we have $\int_{\eta_0}^{\infty} \frac{2d\mu(\eta)}{(\xi e^{\eta})^d + \xi^d} = \frac{1}{\xi^d + 1} + \int_{(\eta_0, \infty)} \frac{2d\mu(\eta)}{(\xi e^{\eta})^d + \xi^d}$. If this were to be equal to twice the RHS, namely 1, we need $\frac{\xi^d}{\xi^d + 1} = \int_{(\eta_0, \infty)} \frac{2d\mu(\eta)}{(\xi e^{\eta})^d + \xi^d}$. Or equivalently, we need $\frac{1}{\xi^d + 1} = \int_{(\eta_0, \infty)} \frac{2d\mu(\eta)}{(\xi^2 e^{\eta})^d + \xi^{2d}}$. Now, as $d \rightarrow \infty$, the LHS approaches 1. For the RHS to approach 1 as $d \rightarrow \infty$, we need that $\mu(\eta') = \frac{1}{2}$, where η' is such that $\xi^2 \eta' = 1$. This completes the proof.