
Non-trivial two-armed partial-monitoring games are bandits

András Antos
 Computer and Automation
 Research Institute
 Hungarian Academy of Sciences

Gábor Bartók and Csaba Szepesvári
 Department of Computing Science
 University of Alberta
 Edmonton, Canada

Abstract

We consider online learning in partial-monitoring games against an oblivious adversary. We show that when the number of actions available to the learner is two and the game is nontrivial then it is reducible to a bandit-like game and thus the minimax regret is $\Theta(\sqrt{T})$.

1 Introduction

The *partial-monitoring games* we consider are defined as follows: Two players interact with each other in a sequential manner, *Learner* and *Nature*. In each time step *Learner* can choose one of N actions, while *Nature* can choose one of M actions. We use the notation $\underline{n} = \{1, \dots, n\}$ for any integer and denote the actions of both players by integers, starting from 1, so the action sets are \underline{N} and \underline{M} . At the beginning of the game both *Learner* and *Nature* are given a pair of $N \times M$ matrices, $\mathbf{G} = (\mathbf{L}, \mathbf{H})$, where \mathbf{L} is the *loss matrix* and \mathbf{H} is the *feedback matrix*. The elements ℓ_{ij} of \mathbf{L} are real numbers and in fact we shall assume that they belong to the $[0, 1]$ interval. The elements h_{ij} of \mathbf{H} could be chosen from any alphabet. However, for the sake of simplicity, and without loss of generality (w.l.o.g.), we may assume that the elements of \mathbf{H} are also real numbers. Now, still at the beginning of the game, *Nature* decides about the sequence of actions (J_1, J_2, \dots) to be played. These actions are kept private, *i.e.*, they are not revealed to *Learner*. *Nature's* actions will also be called *outcomes*.

The game is played in discrete time steps. At time step t ($t = 1, 2, \dots$), first *Learner* chooses an action I_t based on the information available to him up to time t . The choice of the action may be randomized. Upon announcing his action, *Learner* gets the feedback h_{I_t, J_t} and suffers the loss ℓ_{I_t, J_t} . The cycle is then repeated for time step $t + 1$. It is important to note that the loss suffered is not revealed to *Learner*.

The goal of *Learner* is to keep his cumulative loss

$$L_T = \sum_{t=1}^T \ell_{I_t, J_t}$$

small, where T denotes the time horizon. *Learner's* performance is evaluated by comparing his cumulative loss to the cumulative loss of the best fixed action from \underline{N} ,

$$L_T^* = \min_{i \in \underline{N}} \sum_{t=1}^T \ell_{i, J_t},$$

giving rise to the *cumulative expected regret* (or simply *regret*),

$$R_T(\mathcal{A}, \mathbf{G}) = \mathbb{E}[L_T - L_T^*],$$

where \mathcal{A} is the *strategy* *Learner* follows. Note that in the definition of L_T^* , the best fixed action is selected in hindsight. When the growth rate of regret is sublinear in T , *i.e.*, the average regret R_T/T converges to zero, in the long run, *Learner* can be said to perform as well as an oracle who can play this best action in hindsight.

The problem just described is of major importance in learning theory since it models a number of interesting scenarios including *apple tasting* [Bartók et al., 2010], a variant of *label efficient prediction*, and *dynamic pricing* [Cesa-Bianchi and Lugosi, 2006]. For further discussion and examples see Chapters 2–7 in the book by Cesa-Bianchi and Lugosi [2006].

Given a game $\mathbf{G} = (\mathbf{L}, \mathbf{H})$, our goal is to find out the growth rate of the *minimax regret* associated with \mathbf{G} , and to design strategies that allow Learner to achieve this minimal growth rate. Let the worst-case regret of algorithm \mathcal{A} when used in \mathbf{G} for time horizon T be

$$\bar{R}_T(\mathcal{A}, \mathbf{G}) = \sup_{(J_1, \dots, J_T) \in \underline{M}^T} R_T(\mathcal{A}, \mathbf{G}),$$

where the supremum is taken over all outcome sequences. Formally, the *minimax regret* of game \mathbf{G} for time horizon T is defined by

$$R_T^*(\mathbf{G}) = \inf_{\mathcal{A}} \bar{R}_T(\mathcal{A}, \mathbf{G}) = \inf_{\mathcal{A}} \sup_{(J_1, \dots, J_T) \in \underline{M}^T} R_T(\mathcal{A}, \mathbf{G}),$$

where the infimum is taken over all strategies of Learner. Note that, since for constant outcome sequences $R_T(\mathcal{A}, \mathbf{G}) \geq 0$, also $\bar{R}_T(\mathcal{A}, \mathbf{G}) \geq 0$ and $R_T^*(\mathbf{G}) \geq 0$.

Definition 1 *A game is called trivial if the minimax regret is either 0 or scales linearly with the number of time steps.*

Lemma 1 *The following three statements are equivalent:*

- a) *The minimax regret is zero for each T .*
- b) *The minimax regret is zero for some T .*
- c) *There exists an action $i \in \underline{N}$ whose loss is not larger than the loss of any other action irrespectively of the choice of Nature's action.*

The proof is in the Appendix.

2 Previous work

The growth rate of the minimax regret is strongly influenced by the choice of \mathbf{L} and \mathbf{H} . Consider, for example, the case of so-called *full-information* games, where the feedback is sufficient for Learner to recover Nature's action in each round. In the simplest case, this is represented by $h_{ij} = j$. However, from the point of view of the information content of feedback, we get an equivalent situation when each row of \mathbf{H} is composed of pairwise distinct elements. The following result is known to hold for these games:

Theorem 2 *Consider a full-information game \mathbf{G} when Learner has N actions. Then there exists an algorithm \mathcal{A} such that for any time horizon T , $\bar{R}_T(\mathcal{A}, \mathbf{G}) \leq \sqrt{(T/2) \ln N}$.*

Algorithm \mathcal{A} in the theorem above can be the *Exponentially Weighted Average Forecaster* with appropriate tuning (see *e.g.*, Cesa-Bianchi and Lugosi [2006, Corollary 4.2]).

Another special case is when the only information that Learner receives is the loss of the action taken (*i.e.*, when $\mathbf{H} = \mathbf{L}$), which we call the *bandit case*, following Cesa-Bianchi and Lugosi [2006]. Then, the INF algorithm due to Audibert and Bubeck [2009] is known to achieve a constant multiple of the minimax regret:

Theorem 3 *Take a bandit game \mathbf{G} when Learner has N actions. Then there exists an algorithm \mathcal{A} such that $\bar{R}_T(\mathcal{A}, \mathbf{G}) \leq 15\sqrt{NT}$. Further, for any N there exists a game \mathbf{G} such that for any time horizon T , $R_T^*(\mathbf{G}) \geq 1/20\sqrt{NT}$.*

The lower bound on the minimax regret is due to Auer et al. [2002] (also, Cesa-Bianchi and Lugosi [2006, Theorem 6.11]), while the upper bound is due to Audibert and Bubeck [2009]. (The Exp3 algorithm due to Auer et al. [2002] achieves the same upper bound up to logarithmic factors.)

The following theorem, due to Antos et al. [2011], is a lower bound for any non-trivial game.

Theorem 4 *If \mathbf{G} is a non-trivial partial-monitoring game then there exists a constant $c > 0$ such that for any T , $R_T^*(\mathbf{G}) \geq c\sqrt{T}$.*

Now, consider the game $\mathbf{G} = (\mathbf{L}, \mathbf{H})$ with

$$\mathbf{L} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad \mathbf{H} = \begin{pmatrix} 1 & 2 \\ 1 & 1 \\ 1 & 1 \end{pmatrix}. \quad (1)$$

That is, the first action of Learner gives full information about the outcome, but it has a high cost, while the other two actions do not reveal any information. Further, the ordering of actions 2 and 3 by costs is reversed based on the choice of Nature. Then, the following holds [Cesa-Bianchi et al., 2006, Theorem 5.1]:

Theorem 5 *The above game has $R_T^*(\mathbf{G}) = \Omega(T^{2/3})$.*

This shows that the above game is intrinsically harder than a bandit problem. Further, the algorithm FEDEXP3 by Piccolboni and Schindelhauer [2001] is known to achieve this growth-rate [Cesa-Bianchi et al., 2006]:

Theorem 6 *Consider any partial-monitoring game $\mathbf{G} = (\mathbf{L}, \mathbf{H})$ such that $\mathbf{L} = \mathbf{KH}$ for some matrix \mathbf{K} . Then, there exist an algorithm \mathcal{A} such that $\bar{R}_T(\mathcal{A}, \mathbf{G}) \leq CT^{2/3}$, where C depends on N and $k^* \stackrel{\text{def}}{=} \max_{i,j} |k_{ij}|$.*

Thus, we see that the difficulty of a game depends on the structure of \mathbf{L} and \mathbf{H} . Recently, Bartók et al. [2010] classified almost all games by their difficulty when the number of actions available to Nature is limited to $M = 2$. In effect, they showed that the exponent in the dependence of the minimax regret on T in these games is one of $\{0, 1/2, 2/3, 1\}$.

In this short communication, we investigate the dual case when the number of actions available to Nature is not restricted, but the number of actions available to Learner is limited to $N = 2$.

3 Result

In this section we state and prove that, in essence, any non-trivial two-action game can be viewed as a bandit game.

We need some preparations. First, we will make use of the following concept:

Definition 2 *Take two games, $\mathbf{G} = (\mathbf{L}, \mathbf{H})$, $\mathbf{G}' = (\mathbf{L}', \mathbf{H}')$, where \mathbf{L} , \mathbf{L}' , \mathbf{H} , and \mathbf{H}' are $N \times M$ matrices. We say that \mathbf{G}' is simulation-and-regret-not-harder than \mathbf{G} (or easier for short, denoted by $\mathbf{G}' \leq \mathbf{G}$) when the following holds: Fix any algorithm \mathcal{A} . Then, one can find an algorithm \mathcal{A}' such that the behavior of \mathcal{A} on \mathbf{G} can be replicated by using \mathcal{A}' on \mathbf{G}' in the sense that for the same outcome sequences, the two algorithms will choose the same action sequences and the regret in the second case is at most the regret in the first case, that is, $R_T(\mathcal{A}', \mathbf{G}') \leq R_T(\mathcal{A}, \mathbf{G})$.*

We say that \mathbf{G} and \mathbf{G}' are simulation-and-regret-equivalent (or equivalent, $\mathbf{G}' \simeq \mathbf{G}$) when both $\mathbf{G}' \leq \mathbf{G}$ and $\mathbf{G} \leq \mathbf{G}'$.

Clearly, \leq is a preorder and \simeq is an equivalence relation on the set of $N \times M$ games, moreover, if $\mathbf{G}' \leq \mathbf{G}$ then $R_T^*(\mathbf{G}') \leq R_T^*(\mathbf{G})$, and if $\mathbf{G} \simeq \mathbf{G}'$ then their minimax regret is the same.

We need a few simple lemmata on these relations of games:

Lemma 7 *The regret of a sequence of actions in a game does not change if the loss matrix is changed by subtracting the same real number from each coordinate of one of its columns (see e.g., Piccolboni and Schindelhauer [2001]). Therefore, letting $\mathbf{1} = (1, \dots, 1)^\top \in \mathbb{R}^N$, $\mathbf{v} \in \mathbb{R}^M$, and $\mathbf{G}' = (\mathbf{L} - \mathbf{1}\mathbf{v}^\top, \mathbf{H})$, we have that $\mathbf{G} \simeq \mathbf{G}'$.*

Lemma 8 *If $\mathbf{G} = (\mathbf{L}, \mathbf{H})$ and $\mathbf{G}' = (\mathbf{L}, \mathbf{H}')$ differ only in their feedback matrices and \mathbf{H}' can be obtained by $h'_{ij} = f_i(h_{ij})$ with the help of some mappings f_i ($i \in \underline{N}$) then $\mathbf{G} \leq \mathbf{G}'$. If each f_i is injective then $\mathbf{G} \simeq \mathbf{G}'$.*

In what follows, a transformation of some game into another game that takes either the first or the second form just defined shall be called an *admissible* transformation.

The following proposition shows that if a 2-armed partial-monitoring game is non-trivial then there is no loss in generality by assuming that $\mathbf{L} = \mathbf{KH}$ for some $\mathbf{K} \in \mathbb{R}^{2 \times 2}$. This statement for arbitrary N and most of the ideas for its proof could be extracted from the paper of Piccolboni and Schindelhauer [2001]. An exact detailed proof for $N = 2$ is included here for the sake of completeness.

Proposition 1 *Let $\mathbf{G}_0 = (\mathbf{L}_0, \mathbf{H}_0)$ be a non-trivial 2-armed partial-monitoring game. Then, there exist matrices $\mathbf{L}, \mathbf{H} \in \mathbb{R}^{2 \times M}$ such that $\mathbf{G}_0 \leq \mathbf{G} = (\mathbf{L}, \mathbf{H})$ and $\mathbf{L} = \mathbf{KH}$ for some $\mathbf{K} \in \mathbb{R}^{2 \times 2}$.*

Proof:[Proof of Proposition 1] First, we transform \mathbf{L}_0 to \mathbf{L} using Lemma 7 with \mathbf{v}^\top being its first row. Thus, the first row of \mathbf{L} becomes identically zero, and we get a non-trivial game $\mathbf{G}_1 = (\mathbf{L}, \mathbf{H}_0) \simeq \mathbf{G}_0$. Let ℓ denote the transpose of the second row of \mathbf{L} . In what follows we construct the matrix \mathbf{H} using an admissible transformation of \mathbf{H}_0 defined in Lemma 8.

We construct matrix \mathbf{A} in the following way. Assume that there are m_1 (m_2) distinct entries in the first (respectively, second) row of \mathbf{H}_0 , and transform \mathbf{H}_0 by two injective mappings (Lemma 8) such that the elements of its i^{th} row ($i \in \underline{2}$) are from $\underline{m_i}$. We define the matrices $\mathbf{A}_i \in \mathbb{R}^{m_i \times M}$ as

$$H_0 = \begin{pmatrix} 1 & 2 & 3 & 1 \\ 1 & 2 & 2 & 2 \end{pmatrix} \longrightarrow A = \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \end{pmatrix}$$

Figure 1: An example for the construction of matrix A used in the proof of Proposition 1. The first three rows of A are constructed from the first row of \mathbf{H}_0 which has three distinct elements, the remaining two rows are constructed from the second row of \mathbf{H}_0 . For more details, see the text.

follows: Let each row of A_i be the “indicator” row of the corresponding value of the i^{th} row of \mathbf{H}_0 , that is, $[A_i]_{jk} \stackrel{\text{def}}{=} \mathbb{I}_{\{[H_0]_{ik}=j\}}$. Define A by stacking these matrices on top of each other:

$$A = \begin{pmatrix} A_1 \\ A_2 \end{pmatrix}.$$

See Figure 1 for an example.

The following lemma, proven in the Appendix, is key to prove Proposition 1.

Lemma 9 *If $\ell \notin \text{Im } A^\top$ then \mathbf{G}_1 is trivial.*

Using the assumption that \mathbf{G}_1 is non-trivial, we have from Lemma 9 that $\ell \in \text{Im } A^\top$ must hold. That is, ℓ can be written as a linear combination of the rows of A :

$$\ell = \sum_{i=1}^m \lambda_i \mathbf{a}_i,$$

where $m = m_1 + m_2$ and the vectors \mathbf{a}_i^\top are the rows of A . Let

$$\mathbf{h}_1 = \sum_{i=1}^{m_1} \lambda_i \mathbf{a}_i \quad \text{and} \quad \mathbf{h}_2 = \sum_{i=m_1+1}^m \lambda_i \mathbf{a}_i.$$

Finally, let

$$\mathbf{H} = \begin{pmatrix} \mathbf{h}_1^\top \\ \mathbf{h}_2^\top \end{pmatrix}$$

and $\mathbf{G} = (\mathbf{L}, \mathbf{H})$. Now if the k^{th} and k'^{th} entries of the first row of \mathbf{H}_0 are identical then $[\mathbf{a}_i]_k = [\mathbf{a}_i]_{k'}$ for $1 \leq i \leq m_1$, hence also $[\mathbf{h}_1]_k = [\mathbf{h}_1]_{k'}$. The same holds for the second row of \mathbf{H}_0 and \mathbf{h}_2 . Thus, \mathbf{H} can be obtained by appropriate mappings from \mathbf{H}_0 , and Lemma 8 implies $\mathbf{G}_1 \leq \mathbf{G}$.

On the other hand, setting

$$\mathbf{K} = \begin{pmatrix} 0 & 0 \\ 1 & 1 \end{pmatrix}, \tag{2}$$

we have that $\mathbf{L} = \mathbf{KH}$. ■

The following Proposition is more than what we need, but it is interesting in itself:

Proposition 2 *Let $\mathbf{G} = (\mathbf{L}, \mathbf{H})$ be a 2-armed partial-monitoring game such that $\mathbf{L} = \mathbf{KH}$ for some $\mathbf{K} \in \mathbb{R}^{2 \times 2}$. Then, there exist a $2 \times M$ bandit game \mathbf{G}' such that $\mathbf{G} \leq \mathbf{G}'$. If \mathbf{K} is given by (2) then $\mathbf{G} \simeq \mathbf{G}'$.*

Proof: We will construct a bandit game $\mathbf{G}' = (\mathbf{L}', \mathbf{H}') \geq \mathbf{G}$ that satisfies $\mathbf{L}' = \mathbf{H}'$. Let $\mathbf{K} = [k_{ij}]_{2 \times 2}$ and

$$\mathbf{D} = \text{diag}(k_{11} - k_{21}, k_{22} - k_{12})$$

be a 2×2 diagonal matrix, and define the feedback matrix of \mathbf{G}' by $\mathbf{H}' = \mathbf{DH}$. Then, both rows of \mathbf{H}' are scalar multiples of the corresponding rows of \mathbf{H} . Hence, by these mappings and Lemma 8, $\mathbf{G} \leq (\mathbf{L}, \mathbf{H}')$. If \mathbf{K} is given by (2) then $\mathbf{D} = \text{diag}(-1, 1)$, thus both mappings are injective and $\mathbf{G} \simeq (\mathbf{L}, \mathbf{H}')$. On the other hand, $\mathbf{K} - \mathbf{D} = \mathbf{1}\mathbf{k}^\top$ where $\mathbf{k}^\top = (k_{21}, k_{12})$. Consider the loss matrix

$$\mathbf{L}' \stackrel{\text{def}}{=} \mathbf{L} - \mathbf{1}(\mathbf{k}^\top \mathbf{H}).$$

By Lemma 7, $\mathbf{G}' = (\mathbf{L}', \mathbf{H}') \simeq (\mathbf{L}, \mathbf{H}')$. Moreover,

$$\mathbf{L}' = \mathbf{L} - (\mathbf{1}\mathbf{k}^\top)\mathbf{H} = \mathbf{L} - (\mathbf{K} - \mathbf{D})\mathbf{H} = \mathbf{DH} = \mathbf{H}'.$$
■

Remark 1 It is worth to consider why the above proof works only for $N = 2$. We used that from any 2×2 matrix \mathbf{K} we can subtract a diagonal matrix resulting in a matrix with identical rows. For $N \geq 3$, this obviously does not hold (there is not enough “degrees of freedom”). Indeed, for $N \geq 3$, we have regret rates between $\Theta(\sqrt{T})$ and $\Theta(T)$, for example, Theorem 5 and 6 show that the game in (1) has minimax regret rate $\Theta(T^{2/3})$.

Now, we are ready to prove our main result.

Theorem 10 Each non-trivial 2-armed partial-monitoring game is easier than an appropriate $2 \times M$ bandit game. Consequently, its minimax regret is $\Theta(\sqrt{T})$, where T is the number of time steps.

Proof: According to Proposition 1 and 2, if \mathbf{G}_0 is non-trivial then we can construct first $\mathbf{G} = (\mathbf{L}, \mathbf{H})$ such that $\mathbf{L} = \mathbf{K}\mathbf{H}$ and $\mathbf{G}_0 \leq \mathbf{G}$, then a $2 \times M$ bandit game \mathbf{G}' such that $\mathbf{G} \leq \mathbf{G}'$. Thus $\mathbf{G}_0 \leq \mathbf{G}'$, that implies $R_T^*(\mathbf{G}_0) \leq R_T^*(\mathbf{G}') = O(\sqrt{T})$ by Theorem 3, finishing the proof. \blacksquare

Appendix

Proof:[Proof of Lemma 1] a) \rightarrow b) is obvious

b) \rightarrow c) For any \mathcal{A} ,

$$\begin{aligned} \sup_{(J_1, \dots, J_T) \in \underline{M}^T} R_T(\mathcal{A}, \mathbf{G}) &\geq \sup_{j \in \underline{M}, J_1 = \dots = J_T = j} \mathbb{E}[L_T - L_T^*] \\ &= \sup_{j \in \underline{M}} \mathbb{E} \left[\sum_{t=1}^T \ell_{I_t, j} - T \min_{i \in \underline{N}} \ell_{ij} \right] \\ &\geq \sup_{j \in \underline{M}} \left(\mathbb{E}[\ell_{I_1, j}] - \min_{i \in \underline{N}} \ell_{ij} \right) \stackrel{\text{def}}{=} f(\mathcal{A}). \end{aligned}$$

b) leads to

$$0 = R_T^*(\mathbf{G}) = \inf_{\mathcal{A}} \sup_{(J_1, \dots, J_T) \in \underline{M}^T} R_T(\mathcal{A}, \mathbf{G}) \geq \inf_{\mathcal{A}} f(\mathcal{A}).$$

Observe that $f(\mathcal{A})$ depends on \mathcal{A} through only the distribution of I_1 on \underline{N} denoted by $q = q(\mathcal{A})$ now, that is, $f(\mathcal{A}) = f'(q)$. This dependence is continuous on the compact domain of q , hence the infimum can be replaced by minimum. Thus $\min_q f'(q) \leq 0$, that is, there is a q that for all $j \in \underline{M}$, $\mathbb{E}[\ell_{I_1, j}] = \min_{i \in \underline{N}} \ell_{ij}$. This implies that the support of q contains only actions whose loss is not larger than the loss of any other action irrespectively of the choice of Nature’s action.

c) \rightarrow a) The algorithm that always plays i has zero regret for all outcome sequences and T . \blacksquare

Proof:[Proof of Lemma 9] $\ell \notin \text{Im } A^\top$ implies $\langle \ell \rangle \not\subseteq \text{Im } A^\top$, that is equivalent to $\ell^\perp \not\subseteq \text{Ker } A$, which can be seen by taking the orthogonal complement of both sides and using $(\text{Ker } A)^\perp = \text{Im } A^\top$. The latter implies that there exists v such that $v \in \text{Ker } A$ but $\ell^\top v \neq 0$. We may assume w.l.o.g. that $\ell^\top v > 0$ (otherwise take $-v$). Note that, since the first m_1 rows of A add up to $\mathbf{1}$ and $v \in \text{Ker } A$, the coordinates of v sum to zero.

Let $\Delta^M \subseteq \mathbb{R}^M$ denote the M -dimensional probability simplex. If $p \in \Delta^M$ is a distribution over Nature’s actions \underline{M} , then it is easy to see that the first m_1 coordinates of Ap give the probability distribution of observing the different values of the first row of \mathbf{H}_0 while Learner chooses action 1 assuming Nature chooses her actions from p . The same applies to the last m_2 coordinates of Ap and action 2. It follows that if $Ap_1 = Ap_2$ for two distributions then no algorithm can distinguish them. We find such p_1, p_2 and apply this idea as follows:

If for all $p \in \Delta^M$, $\ell^\top p \geq 0$ (or $\ell^\top p \leq 0$), then \mathbf{G}_1 has zero minimax regret and thus it is trivial. Otherwise, there exist p_+ and p_- with $\ell^\top p_+ > 0$ and $\ell^\top p_- < 0$. Now either there exists $p_0 \in \text{Int}(\Delta^M)$ such that $\ell^\top p_0 = 0$, or we can assume w.l.o.g. that one of p_+ and p_- is in $\text{Int}(\Delta^M)$, in which case there must be again a $p_0 \in \text{Int}(\Delta^M)$ on the segment $\overline{p_+ p_-}$ such that $\ell^\top p_0 = 0$ by the continuity of $\ell^\top p$ in p . In other words, we have a distribution p_0 over \underline{M} such that p_0 is not on the boundary of the probability simplex and the expected loss of the two actions are equal.

Now let $p_1 = p_0 + \varepsilon v$ and $p_2 = p_0 - \varepsilon v$ for some $\varepsilon > 0$. If ε is small enough then both p_1 and p_2 are on the probability simplex Δ^M . Since $Av = 0$ we have that $Ap_1 = Ap_2$.

Given a $p \in \Delta^M$, we use randomization such that J_1, \dots, J_T is replaced by a vector $\tilde{J}_1, \dots, \tilde{J}_T \in \underline{M}^T$ of i.i.d. random variables distributed according to p , independent of the randomization in the algorithm. Let \mathcal{A} be an arbitrary strategy of Learner. For $k \in \underline{2}$, given that the outcome distribution

is p_k , let $\mathbb{P}_k[\cdot]$ be the probability of an event and $\mathbb{E}_k[\cdot]$ be the expectation of a random variable. Then the worst case regret of \mathcal{A} is

$$\begin{aligned}
\sup_{(J_1, \dots, J_T) \in \underline{M}^T} R_T(\mathcal{A}, \mathbf{G}_1) &\geq \mathbb{E}_k[R_T(\mathcal{A}, \mathbf{G}_1)] \\
&= \mathbb{E}_k \left[\sum_{t=1}^T \ell_{I_t, \tilde{J}_t} - \min_{i \in \underline{2}} \sum_{t=1}^T \ell_{i, \tilde{J}_t} \right] \\
&= \mathbb{E}_k \left[\sum_{t=1}^T \mathbb{I}_{\{I_t=2\}} \ell_{\tilde{J}_t} - \min \left(\sum_{t=1}^T \ell_{\tilde{J}_t}, 0 \right) \right] \\
&\quad (\ell_{1j} = 0, \ell_{2j} = \ell_j) \\
&\geq \sum_{t=1}^T \mathbb{E}_k [\mathbb{I}_{\{I_t=2\}}] \mathbb{E}_k \ell_{\tilde{J}_t} - \min \left(\sum_{t=1}^T \mathbb{E}_k \ell_{\tilde{J}_t}, 0 \right) \\
&\quad (\text{by the independence of } I_t \text{ and } \tilde{J}_t, \text{ and Jensen's inequality for min}) \\
&= \sum_{t=1}^T \mathbb{P}_k[I_t = 2] \ell^\top p_k + \left(- \sum_{t=1}^T \ell^\top p_k \right)^+ \\
&= \ell^\top p_k \mu_{T2} + T(-\ell^\top p_k)^+,
\end{aligned}$$

where

$$\mu_{T2} = \mu_{T2}(\mathcal{A}) \stackrel{\text{def}}{=} \sum_{t=1}^T \mathbb{P}_k[I_t = 2] \in [0, T]$$

is the expected number of times \mathcal{A} chooses action 2 under p_k up to time T . Observe that $Ap_1 = Ap_2$ means that for both actions, the feedback distribution is the same under outcome distributions p_1 and p_2 , implying (by induction) that for each $t \geq 1$, $\mathbb{P}_1[I_t = 2] = \mathbb{P}_2[I_t = 2]$. This leads to $\mu_{T1} = \mu_{T2} \stackrel{\text{def}}{=} \mu_T = \mu_T(\mathcal{A})$. Moreover, using $\ell^\top p_0 = 0$ and $\ell^\top v > 0$,

$$\ell^\top p_k \mu_T + T(-\ell^\top p_k)^+ = \begin{cases} \varepsilon \ell^\top v \mu_T & \text{if } k = 1, \\ \varepsilon \ell^\top v (T - \mu_T) & \text{if } k = 2. \end{cases}$$

Thus we have

$$\begin{aligned}
R_T^*(\mathbf{G}_1) &= \inf_{\mathcal{A}} \sup_{(J_1, \dots, J_T) \in \underline{M}^T} R_T(\mathcal{A}, \mathbf{G}_1) \geq \inf_{\mathcal{A}} \max_{k \in \underline{2}} (\ell^\top p_k \mu_T + T(-\ell^\top p_k)^+) \\
&= \varepsilon \ell^\top v \inf_{\mathcal{A}} \max(\mu_T, T - \mu_T) \geq \varepsilon \ell^\top v T/2,
\end{aligned}$$

that is, \mathbf{G}_1 is trivial. ■

References

- A. Antos, G. Bartók, D. Pál, and Cs. Szepesvári. Toward a classification of finite partial-monitoring games, 2011. <http://arxiv.org/abs/1102.2041>.
- J.-Y. Audibert and S. Bubeck. Minimax policies for adversarial and stochastic bandits. In S. Dasgupta and A. Klivans, editors, *22nd Annual Conference on Learning Theory*, 2009.
- P. Auer, N. Cesa-Bianchi, Y. Freund, and R.E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32:48–77, 2002.
- G. Bartók, D. Pál, and Cs. Szepesvári. Toward a classification of finite partial-monitoring games. In M. Hutter and T. Zeugmann, editors, *21st International Conference on Algorithmic Learning Theory*, Lecture Notes in Computer Science. Springer, 2010.
- N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, New York, NY, USA, 2006.
- N. Cesa-Bianchi, G. Lugosi, and G. Stoltz. Regret minimization under partial monitoring. *Mathematics of Operations Research*, 31:562–580, 2006.

A. Piccolboni and C. Schindelhauer. Discrete prediction games with arbitrary feedback and loss. In D.P. Helmbold and B. Williamson, editors, *14th Annual Conference on Computational Learning Theory and 5th European Conference on Computational Learning Theory*, volume 2111 of *Lecture Notes in Computer Science*, pages 208–223. Springer, 2001. ISBN 3-540-42343-5.