



**HAL**  
open science

# Robust approachability and regret minimization in games with partial monitoring

Shie Mannor, Vianney Perchet, Gilles Stoltz

► **To cite this version:**

Shie Mannor, Vianney Perchet, Gilles Stoltz. Robust approachability and regret minimization in games with partial monitoring. 2012. <hal-00595695v3>

**HAL Id: hal-00595695**

**<https://hal.science/hal-00595695v3>**

Preprint submitted on 15 Feb 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

# Robust approachability and regret minimization in games with partial monitoring

Shie Mannor

Israel Institute of Technology (Technion), Haifa, Israel  
email: [shie@ee.technion.ac.il](mailto:shie@ee.technion.ac.il) <http://webee.technion.ac.il/people/shie/>

Vianney Perchet

Université Paris-Diderot, Paris, France  
email: [vianney.perchet@normalesup.org](mailto:vianney.perchet@normalesup.org) <https://sites.google.com/site/vianneyperchet/home>

Gilles Stoltz

Ecole Normale Supérieure – CNRS – INRIA, Paris, France & HEC Paris – CNRS, Jouy-en-Josas, France  
email: [gilles.stoltz@ens.fr](mailto:gilles.stoltz@ens.fr) <http://www.math.ens.fr/~stoltz>

Approachability has become a standard tool in analyzing learning algorithms in the adversarial online learning setup. We develop a variant of approachability for games where there is ambiguity in the obtained reward that belongs to a set, rather than being a single vector. Using this variant we tackle the problem of approachability in games with partial monitoring and develop simple and efficient algorithms (i.e., with constant per-step complexity) for this setup. We finally consider external regret and internal regret in repeated games with partial monitoring and derive regret-minimizing strategies based on approachability theory.

*Key words:* Approachability; repeated games; partial monitoring; regret; online learning

*MSC2000 Subject Classification:* Primary: 91A20, 91A26, 62L12; Secondary: 68Q32

*OR/MS subject classification:* Primary: decision analysis: sequential; games/group decisions: noncooperative; secondary: computer science: artificial intelligence

---

**1. Introduction.** Blackwell’s approachability theory and its variants has become a standard and useful tool in analyzing online learning algorithms (Cesa-Bianchi and Lugosi [5]) and algorithms for learning in games (Hart and Mas-Colell [13, 14]). The first application of Blackwell’s approachability to learning in the online setup is due to Blackwell [3] himself. Numerous other contributions are summarized in the monograph by Cesa-Bianchi and Lugosi [5]. Blackwell’s approachability theory enjoys a clear geometric interpretation that allows it to be used in situations where online convex optimization or exponential weights do not seem to be easily applicable and, in some sense, to go beyond the minimization of the regret and/or to control quantities of a different flavor; e.g., in the article by Mannor et al. [20], to minimize the regret together with path constraints, and in the one by Mannor and Shimkin [18], to minimize the regret in games whose stage duration is not fixed. Recently, it has been shown by Abernethy et al. [1] that approachability and low regret learning are equivalent in the sense that efficient reductions exist from one to the other. Another recent paper by Rakhlin et al. [27] showed that approachability can be analyzed from the perspective of learnability using tools from learning theory.

In this paper we consider approachability and online learning with partial monitoring in games against Nature. In partial monitoring the decision maker does not know how much reward was obtained and only gets a (random) signal whose distribution depends on the action of the decision maker and the action of Nature. There are two extremes of this setup that are well studied. On the one extreme we have the case where the signal includes the reward itself (or a signal that can be used to unbiasedly estimate the reward), which is essentially the celebrated bandits setup. The other extreme is the case where the signal is not informative (i.e., it tells the decision maker nothing about the actual reward obtained); this setting then essentially consists of repeating the same situation over and over again, as no information is gained over time. We consider a setup encompassing these situations and more general ones, in which the signal is indicative of the actual reward, but is not necessarily a sufficient statistics thereof. The difficulty is that the decision maker cannot compute the actual reward he obtained nor the actions of Nature.

Regret minimization with partial monitoring has been studied in several papers in the learning theory community. Piccolboni and Schindelhauer [26], Mannor and Shimkin [17], Cesa-Bianchi et al. [6] study special cases where an accurate estimation of the rewards (or worst-case rewards) of the decision maker is possible thanks to some extra structure. A general policy with vanishing regret is presented by Lugosi et al. [16]. This policy is based on exponential weights and a specific estimation procedure

for the (worst-case) obtained rewards. In contrast, we provide approachability-based results for the problem of regret minimization. On route, we define a new type of approachability setup, which enables us to re-derive the extension of approachability to the partial monitoring vector-valued setting proposed by Perchet [23]. More importantly, we provide concrete algorithms for this approachability problem that are more efficient in the sense that, unlike previous works in the domain, their complexity is constant over all steps. Moreover, their rates of convergence are independent of the game at hand, as in the seminal paper by Blackwell [3] but for the first time in this general framework. For example, the recent purely theoretical (and fairly technical) study of approachability by Perchet and Quincampoix [25], which is based on somehow related arguments, does neither provide rates of convergence nor concrete algorithms for this matter.

**Outline.** The paper is organized as follows. In Section 2 we recall some basic facts from approachability theory in the standard vector-valued games setting where a decision maker is engaged in a repeated vector-valued game against an arbitrary opponent (or “Nature”). In Section 3 we propose a novel setup for approachability, termed “robust approachability,” where instead of obtaining a vector-valued reward, the decision maker obtains a set, that represents the ambiguity concerning his reward. We provide a simple characterization of approachable convex sets and an algorithm for the set-valued reward setup under the assumption that the set-valued reward functions are linear. In Section 4 we extend the robust approachability setup to problems where the set-valued reward functions are not linear, but rather concave in the mixed action of the decision maker and convex in the mixed action of Nature. In Section 5 we show how to apply the robust approachability framework to the repeated vector-valued games with partial monitoring. In Section 6 we consider a special type of games where the signaling structure possesses a special property, called bi-piecewise linearity, that can be exploited to derive efficient strategies. This type of games is rich enough as it encompasses several useful special cases. In Section 6.1 we provide a simple and constructive algorithm for these games. Previous results for approachability in this setup were either non-constructive (Rustichini [29]) or were highly inefficient as they relied on some sort of lifting to the space of probability measures on mixed actions (Perchet [23]) and typically required a grid that is progressively refined (leading to a step complexity that is exponential in the number  $T$  of past steps). In Section 6.2 we apply our results for both external-regret and internal-regret minimization in repeated games with partial monitoring. In both cases our proofs are simple, lead to algorithms with constant complexity at each step, and are accompanied with rates. Our results for external regret have rates similar to the ones obtained by Lugosi et al. [16], but our proof is direct and simpler. In Section 7 we mention the general signaling case and explain how it is possible to approach certain special sets such as polytopes efficiently and general convex sets although inefficiently.

**2. Some basic facts from approachability theory.** In this section we recall the most basic versions of Blackwell’s approachability theorem for vector-valued payoff functions.

We consider a vector-valued game between two players, a decision maker (first player) and Nature (second player), with respective finite action sets  $\mathcal{A}$  and  $\mathcal{B}$ , whose cardinalities are referred to as  $N_{\mathcal{A}}$  and  $N_{\mathcal{B}}$ . We denote by  $d$  the dimension of the reward vector and equip  $\mathbb{R}^d$  with the  $\ell^2$ -norm  $\|\cdot\|_2$ . The payoff function of the first player is given by a mapping  $m : \mathcal{A} \times \mathcal{B} \rightarrow \mathbb{R}^d$ , which is multi-linearly extended to  $\Delta(\mathcal{A}) \times \Delta(\mathcal{B})$ , the set of product-distributions over  $\mathcal{A} \times \mathcal{B}$ .

We consider two frameworks, depending on whether pure or mixed actions are taken.

**Pure actions taken and observed.** We denote by  $A_1, A_2, \dots$  and  $B_1, B_2, \dots$  the actions in  $\mathcal{A}$  and  $\mathcal{B}$  sequentially taken by each player; they are possibly given by randomized strategies, i.e., the actions  $A_t$  and  $B_t$  were obtained by random draws according to respective probability distributions denoted by  $\mathbf{x}_t \in \Delta(\mathcal{A})$  and  $\mathbf{y}_t \in \Delta(\mathcal{B})$ . For now, we assume that the first player has a full or bandit monitoring of the pure actions taken by the opponent player: at the end of round  $t$ , when receiving the payoff  $m(A_t, B_t)$ , either the pure action  $B_t$  (full monitoring) or only the indicated payoff (bandit monitoring) is revealed to him.

**DEFINITION 2.1** A set  $\mathcal{C} \subseteq \mathbb{R}^d$  is  $m$ -approachable with pure actions if there exists a strategy of the first

player such that, for all  $\varepsilon > 0$ , there exists an integer  $T_\varepsilon$  such that for all strategies of the second player,

$$\mathbb{P} \left\{ \forall T \geq T_\varepsilon, \inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T m(A_t, B_t) \right\|_2 \leq \varepsilon \right\} \geq 1 - \varepsilon.$$

In particular, the first player has a strategy that ensures that the average of his vector-valued payoffs converges almost surely to the set  $\mathcal{C}$  (uniformly with respect to the strategies of the second player).

The above convergence will be achieved in the course of this paper under two forms. Most often we will exhibit strategies such that, for all strategies of the second player, for all  $\delta > 0$ , with probability at least  $1 - \delta$ ,

$$\inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T m(A_t, B_t) \right\|_2 \leq \beta(T, \delta).$$

A union bound shows that such strategies  $m$ -approach  $\mathcal{C}$  as soon as there exists a positive sequence  $\varepsilon_T$  such that  $\sum \varepsilon_t$  is finite and  $\beta(T, \varepsilon_T) \rightarrow 0$ . Sometimes we will also deal with strategies directly ensuring that, for all strategies of the second player, for all  $\delta > 0$ , with probability at least  $1 - \delta$ ,

$$\sup_{\tau \geq T} \inf_{c \in \mathcal{C}} \left\| c - \frac{1}{\tau} \sum_{t=1}^{\tau} m(A_t, B_t) \right\|_2 \leq \beta(T, \delta).$$

Such strategies  $m$ -approach  $\mathcal{C}$  as soon as  $\beta(T, \delta) \rightarrow 0$  for all  $\delta > 0$ .

**Mixed actions taken and observed.** In this case, we denote by  $\mathbf{x}_1, \mathbf{x}_2, \dots$  and  $\mathbf{y}_1, \mathbf{y}_2, \dots$  the actions in  $\Delta(\mathcal{A})$  and  $\Delta(\mathcal{B})$  sequentially taken by each player. We also assume a full or bandit monitoring for the first player: at the end of round  $t$ , when receiving the payoff  $m(\mathbf{x}_t, \mathbf{y}_t)$ , either the mixed action  $\mathbf{y}_t$  (full monitoring) or the indicated payoff (bandit monitoring) is revealed to him.

**DEFINITION 2.2** *A set  $\mathcal{C} \subseteq \mathbb{R}^d$  is  $m$ -approachable with mixed actions if there exists a strategy of the first player such that, for all  $\varepsilon > 0$ , there exists an integer  $T_\varepsilon$  such that for all strategies of the second player,*

$$\mathbb{P} \left\{ \forall T \geq T_\varepsilon, \inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T m(\mathbf{x}_t, \mathbf{y}_t) \right\|_2 \leq \varepsilon \right\} \geq 1 - \varepsilon.$$

As indicated below, in this setting the first player may even have deterministic strategies such that, for all (deterministic or randomized) strategies of the second player,

$$\inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T m(\mathbf{x}_t, \mathbf{y}_t) \right\|_2 \leq \beta(T)$$

with probability 1, where  $\beta(T) \rightarrow 0$ .

**Necessary and sufficient condition for approachability.** For closed convex sets there is a simple characterization of approachability that is a direct consequence of the minimax theorem; the condition is the same for the two settings, whether pure or mixed actions are taken and observed.

**THEOREM 2.1 (THEOREM 3 OF BLACKWELL [2])** *A closed convex set  $\mathcal{C} \subseteq \mathbb{R}^d$  is approachable (with pure or mixed actions) if and only if*

$$\forall \mathbf{y} \in \Delta(\mathcal{B}), \exists \mathbf{x} \in \Delta(\mathcal{A}), \quad m(\mathbf{x}, \mathbf{y}) \in \mathcal{C}.$$

**An associated strategy (that is efficient depending on the geometry of  $\mathcal{C}$ ).** Blackwell suggested a simple strategy with a geometric flavor; it only requires a bandit monitoring.

Play an arbitrary  $\mathbf{x}_1$ . For  $t \geq 1$ , given the vector-valued quantities

$$\widehat{m}_t = \frac{1}{t} \sum_{s=1}^t m(A_s, B_s) \quad \text{or} \quad \widehat{m}_t = \frac{1}{t} \sum_{s=1}^t m(\mathbf{x}_s, \mathbf{y}_s),$$

depending on whether pure or mixed actions are taken and observed, compute the projection  $c_t$  (in  $\ell^2$ -norm) of  $\widehat{m}_t$  on  $\mathcal{C}$ . Find a mixed action  $\mathbf{x}_{t+1}$  that solves the minimax equation

$$\min_{\mathbf{x} \in \Delta(\mathcal{A})} \max_{\mathbf{y} \in \Delta(\mathcal{B})} \langle \widehat{m}_t - c_t, m(\mathbf{x}, \mathbf{y}) \rangle, \quad (1)$$

where  $\langle \cdot, \cdot \rangle$  is the Euclidian inner product in  $\mathbb{R}^d$ . In the case when pure actions are taken and observed, draw  $A_{t+1}$  at random according to  $\mathbf{x}_{t+1}$ .

The minimax problem used above to determine  $\mathbf{x}_{t+1}$  is easily seen to be a (scalar) zero-sum game and is therefore efficiently solvable using, e.g., linear programming: the associated complexity is polynomial in  $N_{\mathcal{A}}$  and  $N_{\mathcal{B}}$ . All in all, this strategy is efficient if the computations of the required projections onto  $\mathcal{C}$  in  $\ell^2$ -norm can be performed efficiently.

The strategy presented above enjoys the following rates of convergence for approachability.

**THEOREM 2.2** (THEOREM 3 OF BLACKWELL [2]; THEOREM II.4.3 OF MERTENS ET AL. [21]) *We denote by  $M$  a bound in norm over  $m$ , i.e.,*

$$\max_{(a,b) \in \mathcal{A} \times \mathcal{B}} \|m(a,b)\|_2 \leq M.$$

*With mixed actions taken and observed, the above strategy ensures that for all strategies of the second player, with probability 1,*

$$\inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T m(\mathbf{x}_t, \mathbf{y}_t) \right\|_2 \leq \frac{2M}{\sqrt{T}};$$

*while with pure actions taken and observed, for all  $\delta \in (0, 1)$  and for all strategies of the second player, with probability at least  $1 - \delta$ ,*

$$\sup_{\tau \geq T} \inf_{c \in \mathcal{C}} \left\| c - \frac{1}{\tau} \sum_{t=1}^{\tau} m(A_t, B_t) \right\|_2 \leq 2M \sqrt{\frac{2}{\delta T}}.$$

**An alternative strategy in the case where pure actions are taken and observed.** Convergence rates of a slightly different flavor (but still implying approachability) can be proved, in the full monitoring case, by modifying the above procedure as follows. For  $t \geq 1$ , consider instead the vector-valued quantity

$$\widehat{m}_t = \frac{1}{t} \sum_{s=1}^t m(\mathbf{x}_s, B_s),$$

compute its projection  $c_t$  (in  $\ell^2$ -norm) on  $\mathcal{C}$ , and solve the associated minimax problem (1).

This modified strategy enjoys the following rates of convergence for approachability when pure actions are taken and observed.

**THEOREM 2.3** (SECTION 7.7 AND EXERCISE 7.23 OF CESA-BIANCHI ET AL. [6]) *We denote by  $M$  a bound in norm over  $m$ , i.e.,*

$$\max_{(a,b) \in \mathcal{A} \times \mathcal{B}} \|m(a,b)\|_2 \leq M.$$

*With pure actions taken and observed, the above strategy ensures that for all strategies of the second player, with probability at least  $1 - \delta$ ,*

$$\inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T m(A_t, B_t) \right\|_2 \leq \frac{2M}{\sqrt{T}} \left(1 + 2\sqrt{\ln(2/\delta)}\right).$$

In the next section, we will rather resort to this slightly modified procedure as the form of the resulting bounds is closer to the one derived in the main section (Section 6) of this paper.

**3. Robust approachability for finite set-valued games.** In this section we extend the results from the previous section to set-valued payoff functions in the case of full monitoring. We denote by  $\mathcal{S}(\mathbb{R}^d)$  the set of all subsets of  $\mathbb{R}^d$  and consider a set-valued payoff function  $\overline{m} : \mathcal{A} \times \mathcal{B} \rightarrow \mathcal{S}(\mathbb{R}^d)$ .

**Pure actions taken and observed.** At each round  $t$ , the players choose simultaneously respective actions  $A_t \in \mathcal{A}$  and  $B_t \in \mathcal{B}$ , possibly at random according to mixed distributions  $\mathbf{x}_t$  and  $\mathbf{y}_t$ . Full monitoring takes place for the first player: he observes  $B_t$  at the end of round  $t$ . However, as a result, the first player gets the subset  $\overline{m}(A_t, B_t)$  as a payoff. This models the ambiguity or uncertainty associated with some true underlying payoff gained.

We extend  $\overline{m}$  multi-linearly to  $\Delta(\mathcal{A}) \times \Delta(\mathcal{B})$  and even to  $\Delta(\mathcal{A} \times \mathcal{B})$ , the set of joint probability distributions on  $\mathcal{A} \times \mathcal{B}$ , as follows. Let

$$\mu = (\mu_{a,b})_{(a,b) \in \mathcal{A} \times \mathcal{B}}$$

be such a joint probability distribution; then  $\overline{m}(\mu)$  is defined as a finite convex combination<sup>1</sup> of subsets of  $\mathbb{R}^d$ ,

$$\overline{m}(\mu) = \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{B}} \mu_{a,b} \overline{m}(a, b).$$

When  $\mu$  is the product-distribution of some  $\mathbf{x} \in \Delta(\mathcal{A})$  and  $\mathbf{y} \in \Delta(\mathcal{B})$ , we use the notation  $\overline{m}(\mu) = \overline{m}(\mathbf{x}, \mathbf{y})$ .

We denote by

$$\pi_T = \frac{1}{T} \sum_{t=1}^T \delta_{(A_t, B_t)}$$

the empirical distribution of the pairs  $(A_t, B_t)$  of actions taken during the first  $T$  rounds, and will be interested in the behavior of

$$\frac{1}{T} \sum_{t=1}^T \overline{m}(A_t, B_t),$$

which can also be rewritten here in a compact way as  $\overline{m}(\pi_T)$ , by linearity of the extension of  $\overline{m}$ .

The distance of this set  $\overline{m}(\pi_T)$  to the target set  $\mathcal{C}$  will be measured in a worst-case sense: we denote by

$$\varepsilon_T = \sup_{d \in \overline{m}(\pi_T)} \inf_{c \in \mathcal{C}} \|c - d\|_2$$

the smallest value such that  $\overline{m}(\pi_T)$  is included in an  $\varepsilon_T$ -neighborhood of  $\mathcal{C}$ . Robust approachability of a set  $\mathcal{C}$  with the set-valued payoff function  $\overline{m}$  then simply means that the sequence of  $\varepsilon_T$  tends almost-surely to 0, uniformly with respect to the strategies of the second player.

**DEFINITION 3.1** *A set  $\mathcal{C} \subseteq \mathbb{R}^d$  is  $\overline{m}$ -robust approachable with pure actions if there exists a strategy of the first player such that, for all  $\varepsilon > 0$ , there exists an integer  $T_\varepsilon$  such that for all strategies of the second player,*

$$\mathbb{P} \left\{ \forall T \geq T_\varepsilon, \sup_{d \in \overline{m}(\pi_T)} \inf_{c \in \mathcal{C}} \|c - d\|_2 \leq \varepsilon \right\} \geq 1 - \varepsilon.$$

**Mixed actions taken and observed.** At each round  $t$ , the players choose simultaneously respective mixed actions  $\mathbf{x}_t \in \Delta(\mathcal{A})$  and  $\mathbf{y}_t \in \Delta(\mathcal{B})$ . Full monitoring still takes place for the first player: he observes  $\mathbf{y}_t$  at the end of round  $t$ ; he however gets the subset  $\overline{m}(\mathbf{x}_t, \mathbf{y}_t)$  as a payoff (which, again, accounts for the uncertainty).

The product-distribution of two elements  $\mathbf{x} = (x_a)_{a \in \mathcal{A}} \in \Delta(\mathcal{A})$  and  $\mathbf{y} = (y_b)_{b \in \mathcal{B}} \in \Delta(\mathcal{B})$  will be denoted by  $\mathbf{x} \otimes \mathbf{y}$ ; it gives a probability mass of  $x_a y_b$  to each pair  $(a, b) \in \mathcal{A} \times \mathcal{B}$ . We consider the empirical joint distribution of mixed actions taken during the first  $T$  rounds,

$$\nu_T = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \otimes \mathbf{y}_t,$$

and will be interested in the behavior of

$$\frac{1}{T} \sum_{t=1}^T \overline{m}(\mathbf{x}_t, \mathbf{y}_t),$$

which can also be rewritten here in a compact way as  $\overline{m}(\nu_T)$ , by linearity of the extension of  $\overline{m}$ .

<sup>1</sup>For two sets  $S, T$  and  $\alpha \in [0, 1]$ , the convex combination  $\alpha S + (1 - \alpha)T$  is defined as  $\{\alpha s + (1 - \alpha)t, s \in S \text{ and } t \in T\}$ .

DEFINITION 3.2 *A set  $C \subseteq \mathbb{R}^d$  is  $\bar{m}$ -robust approachable with mixed actions if there exists a strategy of the first player such that, for all  $\varepsilon > 0$ , there exists an integer  $T_\varepsilon$  such that for all strategies of the second player,*

$$\mathbb{P} \left\{ \forall T \geq T_\varepsilon, \quad \sup_{d \in \bar{m}(\nu_T)} \inf_{c \in C} \|c - d\|_2 \leq \varepsilon \right\} \geq 1 - \varepsilon.$$

Actually, the bounds exhibited below in this setting will be of the form

$$\sup_{d \in \bar{m}(\nu_T)} \inf_{c \in C} \|c - d\|_2 \leq \beta(T)$$

with probability 1 and uniformly over all (deterministic or randomized) strategies of the second player, where  $\beta(T) \rightarrow 0$  and for deterministic strategies of the first player.

**A useful continuity lemma.** Before proceeding we provide a continuity lemma. It can be reformulated as indicating that for all joint distributions  $\mu$  and  $\nu$  over  $\mathcal{A} \times \mathcal{B}$ , the set  $\bar{m}(\mu)$  is contained in a  $M \|\mu - \nu\|_1$ -neighborhood of  $\bar{m}(\nu)$ , where  $M$  is a bound in  $\ell^2$ -norm on  $\bar{m}$ ; this is a fact that we will use repeatedly below.

LEMMA 3.1 *Let  $\mu$  and  $\nu$  be two probability distributions over  $\mathcal{A} \times \mathcal{B}$ . We assume that the set-valued function  $\bar{m}$  is bounded in norm by  $M$ , i.e., that there exists a real number  $M > 0$  such that*

$$\forall (a, b) \in \mathcal{A} \times \mathcal{B}, \quad \sup_{d \in \bar{m}(a, b)} \|d\|_2 \leq M.$$

Then

$$\sup_{d \in \bar{m}(\mu)} \inf_{c \in \bar{m}(\nu)} \|d - c\|_2 \leq M \|\mu - \nu\|_1 \leq M \sqrt{N_{\mathcal{A}} N_{\mathcal{B}}} \|\mu - \nu\|_2,$$

where the norms in the right-hand side are respectively the  $\ell^1$  and  $\ell^2$ -norms between probability distributions.

PROOF. Let  $d$  be an element of  $\bar{m}(\mu)$ ; it can be written as

$$d = \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{B}} \mu_{a, b} \theta_{a, b}$$

for some elements  $\theta_{a, b} \in \bar{m}(a, b)$ . We consider

$$c = \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{B}} \nu_{a, b} \theta_{a, b},$$

which is an element of  $\bar{m}(\nu)$ . Then by the triangle inequality,

$$\|d - c\|_2 = \left\| \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{B}} (\mu_{a, b} - \nu_{a, b}) \theta_{a, b} \right\|_2 \leq \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{B}} |\mu_{a, b} - \nu_{a, b}| \|\theta_{a, b}\|_2 \leq M \sum_{a \in \mathcal{A}} \sum_{b \in \mathcal{B}} |\mu_{a, b} - \nu_{a, b}|.$$

This entails the first claimed inequality. The second one follows from an application of the Cauchy-Schwarz inequality.  $\square$

COROLLARY 3.1 *When the set-valued function  $\bar{m}$  is bounded in norm, for all  $\mathbf{y} \in \Delta(\mathcal{B})$ , the mapping  $D_{\mathbf{y}} : \Delta(\mathcal{A}) \rightarrow \mathbb{R}$  defined by*

$$\forall \mathbf{x} \in \Delta(\mathcal{A}), \quad D_{\mathbf{y}}(\mathbf{x}) = \sup_{d \in \bar{m}(\mathbf{x}, \mathbf{y})} \inf_{c \in C} \|c - d\|_2$$

is continuous.

PROOF. We show that for all  $\mathbf{x}, \mathbf{x}' \in \Delta(\mathcal{A})$ , the condition  $\|\mathbf{x}' - \mathbf{x}\|_1 \leq \varepsilon$  implies that  $D_{\mathbf{y}}(\mathbf{x}) - D_{\mathbf{y}}(\mathbf{x}') \leq M\varepsilon$ , where  $M$  is the bound in norm over  $\bar{m}$ . Indeed, fix  $\delta > 0$  and let  $d_{\delta, \mathbf{x}} \in \bar{m}(\mathbf{x}, \mathbf{y})$  be such that

$$D_{\mathbf{y}}(\mathbf{x}) \leq \inf_{c \in C} \|c - d_{\delta, \mathbf{x}}\|_2 + \delta. \quad (2)$$

By Lemma 3.1 (with the choices  $\mu = \mathbf{x} \otimes \mathbf{y}$  and  $\nu = \mathbf{x}' \otimes \mathbf{y}$ ) there exists  $d_{\delta, \mathbf{x}'} \in \overline{m}(\mathbf{x}', \mathbf{y})$  such that  $\|d_{\delta, \mathbf{x}} - d_{\delta, \mathbf{x}'}\|_2 \leq M\varepsilon + \delta$ . The triangle inequality entails that

$$\inf_{c \in \mathcal{C}} \|c - d_{\delta, \mathbf{x}}\|_2 \leq \inf_{c \in \mathcal{C}} \|c - d_{\delta, \mathbf{x}'}\|_2 + M\varepsilon + \delta.$$

Substituting in (2), we get that

$$D_{\mathbf{y}}(\mathbf{x}) \leq M\varepsilon + 2\delta + \inf_{c \in \mathcal{C}} \|c - d_{\delta, \mathbf{x}'}\|_2 \leq M\varepsilon + 2\delta + D_{\mathbf{y}}(\mathbf{x}'),$$

which, letting  $\delta \rightarrow 0$ , proves our continuity claim.  $\square$

**Necessary and sufficient condition for robust approachability.** This condition reads as follows and will be referred to as (RAC), an acronym that stands for “robust approachability condition.”

**THEOREM 3.1** *Suppose that the set-valued function  $\overline{m}$  is bounded in norm by  $M$ . A closed convex set  $\mathcal{C} \subseteq \mathbb{R}^d$  is  $\overline{m}$ -approachable (with pure or mixed actions) if and only if the following robust approachability condition is satisfied,*

$$\forall \mathbf{y} \in \Delta(\mathcal{B}), \quad \exists \mathbf{x} \in \Delta(\mathcal{A}), \quad \overline{m}(\mathbf{x}, \mathbf{y}) \subseteq \mathcal{C}. \quad (\text{RAC})$$

**PROOF OF THE NECESSITY OF CONDITION (RAC).** If the condition does not hold, then there exists  $\mathbf{y}_0 \in \Delta(\mathcal{B})$  such that for every  $\mathbf{x} \in \mathcal{A}$ , the set  $\overline{m}(\mathbf{x}, \mathbf{y}_0)$  is not included in  $\mathcal{C}$ , i.e., it contains at least one point not in  $\mathcal{C}$ . We consider the mapping  $D_{\mathbf{y}_0}$  defined in the statement of Corollary 3.1. Since  $\mathcal{C}$  is closed, distances of given individual points to  $\mathcal{C}$  are achieved; therefore, by the choice of  $\mathbf{y}_0$ , we get that  $D_{\mathbf{y}_0}(\mathbf{x}) > 0$  for all  $\mathbf{x} \in \Delta(\mathcal{A})$ . Now, since  $D_{\mathbf{y}_0}$  is continuous on the compact set  $\Delta(\mathcal{A})$ , as asserted by the indicated corollary, it attains its minimum, whose value we denote by  $D_{\min} > 0$ .

Assume now that the second player chooses at each round  $\mathbf{y}_t = \mathbf{y}_0$  as his mixed action. In the case of mixed actions taken and observed, denoting

$$\overline{\mathbf{x}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t,$$

we get that  $\nu_t = \overline{\mathbf{x}}_T \otimes \mathbf{y}_0$ , and hence, for all strategies of the first player and for all  $T \geq 1$ ,

$$\sup_{d \in \overline{m}(\nu_T)} \inf_{c \in \mathcal{C}} \|c - d\|_2 = D_{\mathbf{y}_0}(\overline{\mathbf{x}}_T) \geq D_{\min} > 0,$$

which shows that  $\mathcal{C}$  is not approachable.

The case of pure actions taken and observed is treated similarly, with the sole addition of a concentration argument. By martingale convergence (e.g., repeated uses of the Hoeffding-Azuma inequality together with an application of the Borel-Cantelli lemma),  $\delta_T = \|\pi_T - \nu_T\|_1 \rightarrow 0$  almost surely as  $T \rightarrow \infty$ . By applying Lemma 3.1, we get

$$\sup_{d \in \overline{m}(\pi_T)} \inf_{c \in \mathcal{C}} \|c - d\|_2 \geq \sup_{d \in \overline{m}(\nu_T)} \inf_{c \in \mathcal{C}} \|c - d\|_2 - M\delta_T \geq D_{\min} - M\delta_T$$

and simply take the  $\liminf$  in the above inequalities to conclude the argument.  $\square$

That (RAC) is sufficient to get robust approachability is proved in a constructive way, by exhibiting suitable strategies. We identify probability distributions over  $\mathcal{A} \times \mathcal{B}$  with vectors in  $\mathbb{R}^{\mathcal{A} \times \mathcal{B}}$  and consider the vector-valued payoff function

$$m : (a, b) \in \mathcal{A} \times \mathcal{B} \mapsto \delta_{(a,b)} \in \mathbb{R}^{\mathcal{A} \times \mathcal{B}},$$

which we extend multi-linearly to  $\Delta(\mathcal{A}) \times \Delta(\mathcal{B})$ ; the target set will be

$$\tilde{\mathcal{C}} = \{\mu \in \Delta(\mathcal{A} \times \mathcal{B}) : \overline{m}(\mu) \subseteq \mathcal{C}\}. \quad (3)$$

Since  $\overline{m}$  is a linear function on  $\Delta(\mathcal{A} \times \mathcal{B})$  and  $\mathcal{C}$  is convex, the set  $\tilde{\mathcal{C}}$  is convex as well. In addition, since  $\mathcal{C}$  is closed,  $\tilde{\mathcal{C}}$  is also closed.

**LEMMA 3.2** *Condition (RAC) is equivalent to the  $m$ -approachability of  $\tilde{\mathcal{C}}$ .*

PROOF. This equivalence is immediate via Theorem 2.1. The latter indeed states that the  $m$ -approachability of  $\tilde{\mathcal{C}}$  is equivalent to the fact that for all  $\mathbf{y} \in \Delta(\mathcal{B})$ , there exists some  $\mathbf{x} \in \Delta(\mathcal{A})$  such that  $\mu = m(\mathbf{x}, \mathbf{y})$ , the product-distribution between  $\mathbf{x}$  and  $\mathbf{y}$ , belongs to  $\tilde{\mathcal{C}}$ , i.e., satisfies  $\overline{m}(\mu) = \overline{m}(\mathbf{x}, \mathbf{y}) \subseteq \mathcal{C}$ .  $\square$

The above definition of  $m$  entails the following rewriting,

$$\pi_T = \frac{1}{T} \sum_{t=1}^T m(A_t, B_t) \quad \text{and} \quad \nu_T = \frac{1}{T} \sum_{t=1}^T m(\mathbf{x}_t, \mathbf{y}_t).$$

Let  $P_{\tilde{\mathcal{C}}}$  denote the projection operator onto  $\tilde{\mathcal{C}}$ ; the quantities at hand in the definition of  $m$ -approachability of  $\tilde{\mathcal{C}}$  are given by

$$\varepsilon_T = \left\| \pi_T - P_{\tilde{\mathcal{C}}}(\pi_T) \right\|_2 = \inf_{\mu \in \tilde{\mathcal{C}}} \|\pi_T - \mu\|_2 \quad \text{and} \quad \varepsilon'_T = \left\| \nu_T - P_{\tilde{\mathcal{C}}}(\nu_T) \right\|_2 = \inf_{\mu \in \tilde{\mathcal{C}}} \|\nu_T - \mu\|_2.$$

We now relate the quantities of interest, i.e., the ones arising in the definition of  $\overline{m}$ -robust approachability of  $\mathcal{C}$ , to the former quantities.

LEMMA 3.3 *With pure actions taken and observed,*

$$\sup_{d \in \overline{m}(\pi_T)} \inf_{c \in \mathcal{C}} \|c - d\|_2 \leq M \sqrt{N_{\mathcal{A}} N_{\mathcal{B}}} \varepsilon_T.$$

*With mixed actions taken and observed,*

$$\sup_{d \in \overline{m}(\nu_T)} \inf_{c \in \mathcal{C}} \|c - d\|_2 \leq M \sqrt{N_{\mathcal{A}} N_{\mathcal{B}}} \varepsilon'_T.$$

PROOF. Lemma 3.1 entails that the sets  $\overline{m}(\pi_T)$  are included in  $M\sqrt{N_{\mathcal{A}}N_{\mathcal{B}}}\varepsilon_T$ -neighborhoods of  $\overline{m}(P_{\tilde{\mathcal{C}}}(\pi_T))$ . Since by definition of  $\tilde{\mathcal{C}}$ , one has  $\overline{m}(P_{\tilde{\mathcal{C}}}(\pi_T)) \subseteq \mathcal{C}$ , we get in particular that the sets  $\overline{m}(\pi_T)$  are included in  $M\sqrt{N_{\mathcal{A}}N_{\mathcal{B}}}\varepsilon_T$ -neighborhoods of  $\mathcal{C}$ , which is exactly what was stated. The argument can be repeated with the  $\nu_T$  to get the second bound in the statement of the lemma.  $\square$

PROOF OF THE SUFFICIENCY OF CONDITION (RAC). First, Lemma 3.2 shows that Condition (RAC) (via Theorems 2.2 or 2.3) ensures the existence of strategies  $m$ -approaching  $\tilde{\mathcal{C}}$ . Second, Lemma 3.3 indicates that these strategies also  $\overline{m}$ -robust approach  $\mathcal{C}$ . (It even translates the rates for the  $m$ -approachability of  $\tilde{\mathcal{C}}$  into rates for the  $\overline{m}$ -robust approachability of  $\mathcal{C}$ ; for instance, in the case of mixed actions taken and observed, the  $2/\sqrt{T}$  rate for the  $m$ -approachability of  $\tilde{\mathcal{C}}$  becomes a  $2M\sqrt{N_{\mathcal{A}}N_{\mathcal{B}}/T}$  rate for the  $\overline{m}$ -robust approachability of  $\mathcal{C}$ , a fact that we will use in the proof of Theorem 6.1.)  $\square$

**Two concluding remarks.** Note that, as explained around Equation (1), the considered strategies for  $m$ -approaching  $\tilde{\mathcal{C}}$ , or equivalently  $\overline{m}$ -robust approaching  $\mathcal{C}$ , are efficient as soon as projections in  $\ell^2$ -norm onto the set  $\tilde{\mathcal{C}}$  defined in (3) can be computed efficiently. The latter fact depends on the respective geometries of  $\overline{m}$  and  $\mathcal{C}$ . We will provide examples of favorable cases (see, e.g., Section 6.2.1 about minimization of external regret under partial monitoring).

A final remark is that the proposed strategies require full monitoring, as they rely on the observations of either the pair of played mixed actions  $m(\mathbf{x}_t, \mathbf{y}_t)$  or of played pure actions  $m(A_t, B_t)$ . They enjoy no obvious extension to a case where only a bandit monitoring of the played sets  $\overline{m}(\mathbf{x}_t, \mathbf{y}_t)$  or  $\overline{m}(A_t, B_t)$  would be available.

**4. Robust approachability for concave–convex set-valued games.** We consider in this section the same setting of mixed actions taken and observed as in the previous section, that is, we deal with set-valued payoff functions  $\overline{m} : \Delta(\mathcal{A}) \times \Delta(\mathcal{B}) \rightarrow \mathcal{S}(\mathbb{R}^d)$  under full monitoring. However, in the previous section  $\overline{m}$  was linear on  $\Delta(\mathcal{A}) \times \Delta(\mathcal{B})$ , an assumption that we now weaken while still having that (RAC) is the necessary and sufficient condition for robust approachability. The price to pay for this is the loss of the possible efficiency of the approachability strategies exhibited and the worsening of the convergence rates.

Formally, the functions  $\overline{m} : \Delta(\mathcal{A}) \times \Delta(\mathcal{B}) \rightarrow \mathcal{S}(\mathbb{R}^d)$  that we will consider will satisfy one or several of the following properties.

DEFINITION 4.1 A function  $\bar{m} : \Delta(\mathcal{A}) \times \Delta(\mathcal{B}) \rightarrow \mathcal{S}(\mathbb{R}^d)$  is uniformly continuous in its first argument if for all  $\varepsilon > 0$ , there exists  $\eta > 0$  such that for all  $\mathbf{x}, \mathbf{x}' \in \Delta(\mathcal{A})$  satisfying  $\|\mathbf{x} - \mathbf{x}'\|_1 \leq \eta$  and for all  $\mathbf{y} \in \Delta(\mathcal{B})$ , the set  $\bar{m}(\mathbf{x}', \mathbf{y})$  is included in an  $\varepsilon$ -neighborhood of  $\bar{m}(\mathbf{x}, \mathbf{y})$  in the Euclidian norm. Put differently,

$$\sup_{d \in \bar{m}(\mathbf{x}', \mathbf{y})} \inf_{c \in \bar{m}(\mathbf{x}, \mathbf{y})} \|d - c\|_2 \leq \varepsilon \quad \text{or} \quad \bar{m}(\mathbf{x}', \mathbf{y}) \subseteq \bar{m}(\mathbf{x}, \mathbf{y}) + \varepsilon \mathbf{B},$$

where  $\mathbf{B}$  is the unit Euclidian ball in  $\mathbb{R}^d$ .

Uniform continuity in the second argument is defined symmetrically.

DEFINITION 4.2 A function  $\bar{m} : \Delta(\mathcal{A}) \times \Delta(\mathcal{B}) \rightarrow \mathcal{S}(\mathbb{R}^d)$  is concave in its first argument if for all  $\mathbf{x}, \mathbf{x}' \in \Delta(\mathcal{A})$ , all  $\mathbf{y} \in \Delta(\mathcal{B})$ , and all  $\alpha \in [0, 1]$ ,

$$\bar{m}(\alpha \mathbf{x} + (1 - \alpha) \mathbf{x}', \mathbf{y}) \subseteq \alpha \bar{m}(\mathbf{x}, \mathbf{y}) + (1 - \alpha) \bar{m}(\mathbf{x}', \mathbf{y}).$$

A function  $\bar{m} : \Delta(\mathcal{A}) \times \Delta(\mathcal{B}) \rightarrow \mathcal{S}(\mathbb{R}^d)$  is convex in its second argument if for all  $\mathbf{x} \in \Delta(\mathcal{A})$ , all  $\mathbf{y}, \mathbf{y}' \in \Delta(\mathcal{B})$ , and all  $\alpha \in [0, 1]$ ,

$$\alpha \bar{m}(\mathbf{x}, \mathbf{y}) + (1 - \alpha) \bar{m}(\mathbf{x}, \mathbf{y}') \subseteq \bar{m}(\mathbf{x}, \alpha \mathbf{y} + (1 - \alpha) \mathbf{y}').$$

An example of such a function  $\bar{m}$  is discussed in Lemma 5.1.

The following theorem indicates that (RAC) is the necessary and sufficient condition for the  $\bar{m}$ -robust approachability of a closed convex set  $\mathcal{C}$  with mixed actions when the payoff function  $\bar{m}$  satisfies all four properties stated above. (Boundedness of  $\bar{m}$  indeed follows from the continuity of  $\bar{m}$  in each variable.)

THEOREM 4.1 If  $\bar{m}$  is bounded, convex, and uniformly continuous in its second argument, then (RAC) entails that a closed convex set  $\mathcal{C}$  is  $\bar{m}$ -robust approachable with mixed actions.

On the contrary, if  $\bar{m}$  is concave and uniformly continuous in its first argument, then a closed convex set  $\mathcal{C}$  can be  $\bar{m}$ -robust approachable with mixed actions only if (RAC) is satisfied.

PROOF OF THE SECOND STATEMENT OF THEOREM 4.1. The proof of Corollary 3.1 extends to the case considered here and shows, thanks to the ad hoc consideration of the result stated in Lemma 3.1 as following from Definition 4.1, that for all  $\mathbf{y} \in \Delta(\mathcal{B})$ , the mapping  $D_{\mathbf{y}}$  is still continuous over  $\Delta(\mathcal{A})$ . We now proceed by contradiction and assume that (RAC) is not satisfied; the first part of the proof of the necessity of (RAC) in Theorem 3.1 also applies to the present case: there exists  $\mathbf{y}_0$  such that  $D_{\mathbf{y}_0} \geq D_{\min} > 0$  over  $\Delta(\mathcal{A})$ . It then suffices to note that whenever the second player resorts to  $\mathbf{y}_t = \mathbf{y}_0$  at all rounds  $t \geq 1$ , then for all strategies of the first player, the quantity of interest in robust approachability can be lower bounded as follows, thanks to the concavity in the first argument:

$$\begin{aligned} & \sup \left\{ \inf_{c \in \mathcal{C}} \|d - c\|_2 : d \in \frac{1}{T} \sum_{t=1}^T \bar{m}(\mathbf{x}_t, \mathbf{y}_0) \right\} \\ & \geq \sup \left\{ \inf_{c \in \mathcal{C}} \|d - c\|_2 : d \in \bar{m} \left( \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t, \mathbf{y}_0 \right) \right\} = D_{\mathbf{y}_0} \left( \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \right) \geq D_{\min} > 0. \end{aligned}$$

Therefore,  $\mathcal{C}$  is  $\bar{m}$ -robust approachable with mixed actions by no strategy of the first player.  $\square$

The proof of the first statement of Theorem 4.1 relies on the use of approximately calibrated strategies of the first player, as introduced and studied (among others) by Dawid [8], Foster and Vohra [9], Mannor and Stoltz [19]. Formally, given  $\eta > 0$ , an  $\eta$ -calibrated strategy of the first player considers some finite covering of  $\Delta(\mathcal{B})$  by  $N_\eta$  balls of radius  $\eta$  and abides by the following constraints. Denoting by  $\mathbf{y}^1, \dots, \mathbf{y}^{N_\eta}$  the centers of the balls in the covering (they form what will be referred to later on as an  $\eta$ -grid), such a strategy chooses only forecasts in  $\{\mathbf{y}^1, \dots, \mathbf{y}^{N_\eta}\}$ . We thus denote by  $L_t$  the index chosen in  $\{1, \dots, N_\eta\}$  at round  $t$  and by

$$N_T(\ell) = \sum_{t=1}^T \mathbb{I}_{\{L_t = \ell\}}$$

the total number of rounds within the first  $T$  ones when the element  $\ell$  of the grid was chosen. We denote by  $(\cdot)_+$  the function that gives the nonnegative part of a real number. The final condition to be satisfied

is that for all  $\delta > 0$ , there exists an integer  $T_\delta$  such that for all strategies of the second player, with probability at least  $1 - \delta$ , for all  $T \geq T_\delta$ ,

$$\sum_{\ell=1}^{N_\eta} \frac{N_T(\ell)}{T} \left( \left\| \mathbf{y}^\ell - \frac{1}{N_T(\ell)} \sum_{t=1}^T \mathbf{y}_t \mathbb{I}_{\{L_t=\ell\}} \right\|_1 - \eta \right)_+ \leq \delta. \quad (4)$$

This calibration criterion is slightly stronger than the classical  $\eta$ -calibration score usually considered in the literature, which consists of omitting nonnegative parts in the criterion above and ensuring that for all strategies of the second player, with probability at least  $1 - \delta$ , for all  $T \geq T_\delta$ ,

$$\sum_{\ell=1}^{N_\eta} \frac{N_T(\ell)}{T} \left\| \mathbf{y}^\ell - \frac{1}{N_T(\ell)} \sum_{t=1}^T \mathbf{y}_t \mathbb{I}_{\{L_t=\ell\}} \right\|_1 \leq \eta + \delta. \quad (5)$$

The existence of a calibrated strategy in the sense of (4) however follows from the same approachability-based construction studied in Mannor and Stoltz [19] to get (5) and is detailed in the appendix. In the sequel we will only use the following consequence of calibration: that for all strategies of the second player, with probability at least  $1 - \delta$ , for all  $T \geq T_\delta$ ,

$$\max_{\ell=1, \dots, N_\eta} \frac{N_T(\ell)}{T} \left( \left\| \mathbf{y}^\ell - \frac{1}{N_T(\ell)} \sum_{t=1}^T \mathbf{y}_t \mathbb{I}_{\{L_t=\ell\}} \right\|_1 - \eta \right)_+ \leq \delta. \quad (6)$$

**PROOF OF THE FIRST STATEMENT OF THEOREM 4.1.** The insight of this proof is similar to the one illustrated in Perchet [22]. We first note that it suffices to prove that for all  $\varepsilon > 0$ , the set  $\mathcal{C}_\varepsilon$  defined as the  $\varepsilon$ -neighborhood of  $\mathcal{C}$  is  $\bar{m}$ -robust approachable with mixed actions; this is so up to proceeding in regimes  $r = 1, 2, \dots$  each corresponding to a dyadic value  $\varepsilon_r = 2^{-r}$  and lasting for a number of rounds carefully chosen in terms of the length of the previous regimes.

Therefore, we fix  $\varepsilon > 0$  and associate with it a modulus of continuity  $\eta > 0$  given by the uniform continuity of  $\bar{m}$  in its second argument. We consider an  $\eta/2$ -calibrated strategy of the first player, which we will use as an auxiliary strategy. Since (RAC) is satisfied, we may associate with each element  $\mathbf{y}^\ell$  of the underlying  $\eta/2$ -grid a mixed action  $\mathbf{x}^\ell \in \Delta(\mathcal{A})$  such that  $\bar{m}(\mathbf{x}^\ell, \mathbf{y}^\ell) \subseteq \mathcal{C}$ . The main strategy of the first player then prescribes the use of  $\mathbf{x}_t = \mathbf{x}^{L_t}$  at each round  $t \geq 1$ . The intuition behind this definition is that if  $\mathbf{y}^{L_t}$  is forecast by the auxiliary strategy, then since the latter is calibrated, one should play as good as possible against  $\mathbf{y}^{L_t}$ ; in view of the aim at hand, which is approaching  $\mathcal{C}$ , such a good reply is given by  $\mathbf{x}^{L_t}$ .

To assess the constructed strategy, we group rounds according to the values  $\ell$  taken by the  $L_t$ ; to that end, we recall that  $N_T(\ell)$  denotes the number of rounds in which  $\mathbf{y}^\ell$  was forecast and  $\mathbf{x}^\ell$  was played. The average payoff up to round  $T$  is then rewritten as

$$\frac{1}{T} \sum_{t=1}^T \bar{m}(\mathbf{x}_t, \mathbf{y}_t) = \sum_{\ell=1}^{N_{\eta/2}} \frac{N_T(\ell)}{T} \left( \frac{1}{N_T(\ell)} \sum_{t=1}^T \bar{m}(\mathbf{x}^\ell, \mathbf{y}_t) \mathbb{I}_{\{L_t=\ell\}} \right).$$

We denote for all  $\ell$  such that  $N_T(\ell) > 0$  the average of their corresponding mixed actions  $\mathbf{y}_t$  by

$$\bar{\mathbf{y}}_T^\ell = \frac{1}{N_T(\ell)} \sum_{t=1}^T \mathbf{y}_t \mathbb{I}_{\{L_t=\ell\}}.$$

The convexity of  $\bar{m}$  in its second argument leads to the inclusion

$$\frac{1}{T} \sum_{t=1}^T \bar{m}(\mathbf{x}_t, \mathbf{y}_t) = \sum_{\ell=1}^{N_{\eta/2}} \frac{N_T(\ell)}{T} \left( \frac{1}{N_T(\ell)} \sum_{t=1}^T \bar{m}(\mathbf{x}^\ell, \mathbf{y}_t) \mathbb{I}_{\{L_t=\ell\}} \right) \subseteq \sum_{\ell=1}^{N_{\eta/2}} \frac{N_T(\ell)}{T} \bar{m}(\mathbf{x}^\ell, \bar{\mathbf{y}}_T^\ell).$$

To show that the above-defined strategy  $\bar{m}$ -robustly approaches  $\mathcal{C}_\varepsilon = \mathcal{C} + \varepsilon \mathbf{B}$ , it suffices to show that for all  $\delta > 0$ , there exists an integer  $T'_\delta$  such that for all strategies of the second player,

$$\mathbb{P} \left\{ \forall T \geq T'_\delta, \sum_{\ell=1}^{N_{\eta/2}} \frac{N_T(\ell)}{T} \bar{m}(\mathbf{x}^\ell, \bar{\mathbf{y}}_T^\ell) \subseteq \mathcal{C} + (\varepsilon + \delta) \mathbf{B} \right\} \geq 1 - \delta.$$

We denote by  $M$  a bound in  $\ell^2$ -norm on  $\bar{m}$ , i.e., for all  $\mathbf{x} \in \Delta(\mathcal{A})$  and  $\mathbf{y} \in \Delta(\mathcal{B})$ , the inclusion  $\bar{m}(\mathbf{x}, \mathbf{y}) \subseteq M\mathbf{B}$  holds. We let  $\delta' = \delta(\eta/2)/(M N_{\eta/2})$  and define  $T'_\delta$  as the time  $T_{\delta'}$  corresponding to (6). All statements that follow will be for all strategies of the second player and with probability at least  $1 - \delta' \geq 1 - \delta$ , for all  $T \geq T'_\delta$ , as required. For each index  $\ell$  of the grid, either  $\delta'T/N_T(\ell) \leq \eta/2$  or  $\delta'T/N_T(\ell) > \eta/2$ . In the first case, following (6),  $\|\mathbf{y}^\ell - \bar{\mathbf{y}}_T^\ell\| \leq \eta/2 + \delta'T/N_T(\ell) \leq \eta$ ; since  $\eta$  is the modulus of continuity for  $\varepsilon$ , we get that

$$\frac{N_T(\ell)}{T} \bar{m}(\mathbf{x}^\ell, \bar{\mathbf{y}}_T^\ell) \subseteq \frac{N_T(\ell)}{T} (\bar{m}(\mathbf{x}^\ell, \mathbf{y}^\ell) + \varepsilon\mathbf{B}) \subseteq \frac{N_T(\ell)}{T} (\mathcal{C} + \varepsilon\mathbf{B}),$$

where we used the definition of  $\mathbf{x}^\ell$  to get the second inclusion. In the second case, using the boundedness of  $\bar{m}$ , we simply write

$$\frac{N_T(\ell)}{T} \bar{m}(\mathbf{x}^\ell, \bar{\mathbf{y}}_T^\ell) \subseteq \frac{N_T(\ell)}{T} M\mathbf{B} \subseteq \frac{\delta'}{\eta/2} M\mathbf{B}.$$

Summing these bounds over  $\ell$  yields

$$\sum_{\ell=1}^{N_{\eta/2}} \frac{N_T(\ell)}{T} \bar{m}(\mathbf{x}^\ell, \bar{\mathbf{y}}_T^\ell) \subseteq \mathcal{C} + \varepsilon\mathbf{B} + \frac{N_{\eta/2}\delta'}{\eta/2} M\mathbf{B} = \mathcal{C} + (\varepsilon + \delta)\mathbf{B},$$

where we used the definition of  $\delta'$  in terms of  $\delta$ . This concludes the proof.  $\square$

**5. Approachability in games with partial monitoring: statement of the necessary and sufficient condition; links with robust approachability.** A repeated vector-valued game with partial monitoring is described as follows (see, e.g., Mertens et al. [21], Rustichini [29], and the references therein). The players have respective finite action sets  $\mathcal{I}$  and  $\mathcal{J}$ . We denote by  $r : \mathcal{I} \times \mathcal{J} \rightarrow \mathbb{R}^d$  the vector-valued payoff function of the first player and extend it multi-linearly to  $\Delta(\mathcal{I}) \times \Delta(\mathcal{J})$ . At each round, players simultaneously choose their actions  $I_t \in \mathcal{I}$  and  $J_t \in \mathcal{J}$ , possibly at random according to probability distributions denoted by  $\mathbf{p}_t \in \Delta(\mathcal{I})$  and  $\mathbf{q}_t \in \Delta(\mathcal{J})$ . At the end of a round, the first player does not observe  $J_t$  nor  $r(I_t, J_t)$  but only a signal. There is a finite set  $\mathcal{H}$  of possible signals; the feedback  $S_t$  that is given to the first player is drawn at random according to the distribution  $H(I_t, J_t)$ , where the mapping  $H : \mathcal{I} \times \mathcal{J} \rightarrow \Delta(\mathcal{H})$  is known by the first player.

EXAMPLE 5.1 *Examples of such partial monitoring games are provided by, e.g., Cesa-Bianchi et al. [6], among which we can cite the apple tasting problem, the label-efficient prediction constraint, and the multi-armed bandit settings.*

Some additional notation will be useful. We denote by  $R$  the norm of (the linear extension of)  $r$ ,

$$R = \max_{(i,j) \in \mathcal{I} \times \mathcal{J}} \|r(i, j)\|_2.$$

The cardinalities of the finite sets  $\mathcal{I}$ ,  $\mathcal{J}$ , and  $\mathcal{H}$  will be referred to as  $N_{\mathcal{I}}$ ,  $N_{\mathcal{J}}$ , and  $N_{\mathcal{H}}$ .

Definition 2.1 can be extended as follows in this setting; the only new ingredient is the signaling structure, the aim is unchanged.

DEFINITION 5.1 *Let  $\mathcal{C} \subseteq \mathbb{R}^d$  be some set;  $\mathcal{C}$  is  $r$ -approachable for the signaling structure  $H$  if there exists a strategy of the first player such that, for all  $\varepsilon > 0$ , there exists an integer  $T_\varepsilon$  such that for all strategies of the second player,*

$$\mathbb{P} \left\{ \forall T \geq T_\varepsilon, \inf_{\mathbf{c} \in \mathcal{C}} \left\| \mathbf{c} - \frac{1}{T} \sum_{t=1}^T r(I_t, J_t) \right\|_2 \leq \varepsilon \right\} \geq 1 - \varepsilon.$$

*That is, the first player has a strategy that ensures that the sequence of his average vector-valued payoffs converges to the set  $\mathcal{C}$  (uniformly with respect to the strategies of the second player), even if he only observes the random signals  $S_t$  as a feedback.*

**Our contributions.** A necessary and sufficient condition for  $r$ -approachability with the signaling structure  $H$  was stated and proved by Perchet [23]; we therefore need to indicate where our contribution lies. First, both proofs are constructive but our strategy can be efficient (as soon as some projection operator can be computed efficiently, e.g., in the cases of external and internal regret minimization described below) whereas the one of Perchet [23] relies on auxiliary strategies that are calibrated and that require a grid that is progressively refined (leading to a step complexity that is exponential in the number  $T$  of past steps); the latter construction is in essence the one used in Section 4. Second, we are able to exhibit convergence rates. Third, as far as elegance is concerned, our proof is short, compact, and more direct than the one of Perchet [23], which relied on several layers of notations (internal regret in games with partial monitoring, calibration of auxiliary strategies, etc.).

**5.1 Statement of the necessary and sufficient condition for approachability in games with partial monitoring.** To recall the mentioned approachability condition of Perchet [23] we need some additional notation: for all  $\mathbf{q} \in \Delta(\mathcal{J})$ , we denote by  $\tilde{H}(\mathbf{q})$  the element in  $\Delta(\mathcal{H})^{\mathcal{I}}$  defined as follows. For all  $i \in \mathcal{I}$ , its  $i$ -th component is given by the convex combination of probability distributions over  $\mathcal{H}$

$$\tilde{H}(\mathbf{q})_i = H(i, \mathbf{q}) = \sum_{j \in \mathcal{J}} q_j H(i, j).$$

Finally, we denote by  $\mathcal{F}$  the convex set of feasible vectors of probability distributions over  $\mathcal{H}$ :

$$\mathcal{F} = \left\{ \tilde{H}(\mathbf{q}) : \mathbf{q} \in \Delta(\mathcal{J}) \right\}.$$

A generic element of  $\mathcal{F}$  will be denoted by  $\sigma \in \mathcal{F}$  and we define the set-valued function  $\overline{m}$ , for all  $\mathbf{p} \in \Delta(\mathcal{I})$  and  $\sigma \in \mathcal{F}$ , by

$$\overline{m}(\mathbf{p}, \sigma) = \left\{ r(\mathbf{p}, \mathbf{q}') : \mathbf{q}' \in \Delta(\mathcal{J}) \text{ such that } \tilde{H}(\mathbf{q}') = \sigma \right\}.$$

The necessary and sufficient condition exhibited by Perchet [23] for the  $r$ -approachability of  $\mathcal{C}$  with the signaling structure  $H$  can now be recalled. In the sequel we will refer to this condition as Condition (APM), an acronym that stands for “approachability with partial monitoring.”

CONDITION 1 (REFERRED TO AS CONDITION (APM)) *The signaling structure  $H$ , the vector-payoff function  $r$ , and the set  $\mathcal{C}$  satisfy*

$$\forall \mathbf{q} \in \Delta(\mathcal{J}), \exists \mathbf{p} \in \Delta(\mathcal{I}), \forall \mathbf{q}' \in \Delta(\mathcal{J}), \quad \tilde{H}(\mathbf{q}) = \tilde{H}(\mathbf{q}') \Rightarrow r(\mathbf{p}, \mathbf{q}') \in \mathcal{C}.$$

*The condition can be equivalently reformulated as*

$$\forall \sigma \in \mathcal{F}, \exists \mathbf{p} \in \Delta(\mathcal{I}), \quad \overline{m}(\mathbf{p}, \sigma) \subseteq \mathcal{C}. \quad (\text{APM})$$

**This condition is necessary.** The subsequent sections show (in a constructive way) that Condition (APM) is sufficient for  $r$ -approachability of closed convex sets  $\mathcal{C}$  given the signaling structure  $H$ . That this condition is necessary was already proved in Section 3.1 of Perchet [23].

**5.2 Links with robust approachability.** As will become clear in the proof of Theorem 6.1, the key in our problem will be to ensure the robust approachability of  $\mathcal{C}$  with the following non-linear set-valued payoff function, that is however concave–convex in the sense of Definition 4.2.

LEMMA 5.1 *The function*

$$(\mathbf{p}, \mathbf{q}) \in \Delta(\mathcal{I}) \times \Delta(\mathcal{J}) \longmapsto \overline{m}(\mathbf{p}, H(\mathbf{q})).$$

*is concave in its first argument and convex in its second argument.*

Unfortunately, efficient strategies for robust approachability were only proposed in the linear case, not in the concave–convex case. But we illustrate in the next example (and provide a general theory in the next section) how working in lifted spaces can lead to linearity and hence to efficiency.

EXAMPLE 5.2 *We consider a game in which the second player (the column player) can force the first player (the row player) to play a game of matching pennies in the dark by choosing actions  $L$  or  $M$ ; in the matrix below, the real numbers denote the payoff while  $\clubsuit$  and  $\heartsuit$  denote the two possible signals. The respective sets of actions are  $\mathcal{I} = \{T, B\}$  and  $\mathcal{J} = \{L, M, R\}$ .*

|     | $L$    | $M$    | $R$   |
|-----|--------|--------|-------|
| $T$ | 1 / ♣  | -1 / ♣ | 2 / ♥ |
| $B$ | -1 / ♣ | 1 / ♣  | 3 / ♥ |

In this example we only study the mapping  $\mathbf{p} \mapsto \overline{m}(\mathbf{p}, \clubsuit)$  and show that it is piecewise linear on  $\Delta(\mathcal{I})$ , thus, is induced by a linear mapping defined on a lifted space.

We introduce a set  $\mathcal{A} = \{\mathbf{p}_T, \mathbf{p}_B, \mathbf{p}_{1/2}\}$  of possibly mixed actions extending the set  $\mathcal{I} = \{T, B\}$  of pure actions; the set  $\mathcal{A}$  is composed of

$$\mathbf{p}_T = \delta_T, \quad \mathbf{p}_B = \delta_B, \quad \text{and} \quad \mathbf{p}_{1/2} = \frac{1}{2}\delta_T + \frac{1}{2}\delta_B.$$

Each mixed action in  $\Delta(\mathcal{I})$  can be uniquely written as  $\mathbf{p}_\lambda = \lambda\delta_B + (1-\lambda)\delta_T$  for some  $\lambda \in [0, 1]$ . Now, for  $\lambda \geq 1/2$ , first,

$$\mathbf{p}_\lambda = (2\lambda - 1)\delta_B + (1 - (2\lambda - 1))\mathbf{p}_{1/2};$$

second, by definition of  $\overline{m}$ ,

$$\overline{m}(\mathbf{p}_\lambda, \clubsuit) = [1 - 2\lambda, 2\lambda - 1];$$

since in particular  $\overline{m}(\mathbf{p}_{1/2}, \clubsuit) = \{0\}$  and  $\overline{m}(\delta_B, \clubsuit) = [-1, 1]$ , we have the convex decomposition

$$\overline{m}(\mathbf{p}_\lambda, \clubsuit) = (2\lambda - 1)\overline{m}(\delta_B, \clubsuit) + (1 - (2\lambda - 1))\overline{m}(\mathbf{p}_{1/2}, \clubsuit),$$

which can be restated as

$$\overline{m}(\mathbf{p}_\lambda, \clubsuit) = \overline{m}\left((2\lambda - 1)\delta_B + (1 - (2\lambda - 1))\mathbf{p}_{1/2}, \clubsuit\right) = (2\lambda - 1)\overline{m}(\delta_B, \clubsuit) + (1 - (2\lambda - 1))\overline{m}(\mathbf{p}_{1/2}, \clubsuit).$$

That is,  $\overline{m}(\cdot, \clubsuit)$  is linear on the subset of  $\Delta(\mathcal{I})$  corresponding to mixed actions  $\mathbf{p}_\lambda$  with  $\lambda \geq 1/2$ .

A similar property holds the subset of distributions with  $\lambda \leq 1/2$ , so that we have proved that  $\overline{m}(\cdot, \clubsuit)$  is piecewise linear on  $\Delta(\mathcal{I})$ .

The linearity on a lifted space comes from the following observation:  $\overline{m}$  is induced by the linear extension to  $\Delta(\mathcal{A})$  of the restriction of  $\overline{m}$  to  $\mathcal{A}$  (see Definition 6.1 for a more formal statement).

**6. Application of robust approachability to games with partial monitoring: for a particular class of games encompassing regret minimization.** In this section we consider the case where the signaling structure has some special properties described below (linked to linearity properties on lifted spaces) and that can be exploited to get efficient strategies. The case of general signaling structures is then considered in Section 7 but the particular class of games considered here is already rich enough to encompass the minimization of external and internal regret.

**6.1 Approachability in bi-piecewise linear games.** To define bi-piecewise linearity of a game, we start from a technical lemma that shows that  $\overline{m}(\mathbf{p}, \sigma)$  can be written as a *finite* convex combination of sets of the form  $\overline{m}(\mathbf{p}, b)$ , where  $b$  belongs to some finite set  $\mathcal{B} \subseteq \mathcal{F}$  that depends on the game. Under the additional assumption of piecewise linearity of the thus-defined mappings  $\overline{m}(\cdot, b)$ , we then describe a (possibly) efficient strategy for approachability followed by convergence rate guarantees.

**6.1.1 Bi-piecewise linearity of a game – A preliminary technical result.**

**LEMMA 6.1** *For any game with partial monitoring, there exists a finite set  $\mathcal{B} \subset \mathcal{F}$  and a piecewise-linear (injective) mapping  $\Phi : \mathcal{F} \rightarrow \Delta(\mathcal{B})$  such that*

$$\forall \sigma \in \mathcal{F}, \quad \forall \mathbf{p} \in \Delta(\mathcal{I}), \quad \overline{m}(\mathbf{p}, \sigma) = \sum_{b \in \mathcal{B}} \Phi_b(\sigma) \overline{m}(\mathbf{p}, b),$$

where we denoted the convex weight vector  $\Phi(\sigma) \in \Delta(\mathcal{B})$  by  $(\Phi_b(\sigma))_{b \in \mathcal{B}}$ .

**PROOF.** Since  $\tilde{H}$  is linear on the polytope  $\Delta(\mathcal{J})$ , Proposition 2.4 in Rambau and Ziegler [28] implies that its inverse application  $\tilde{H}^{-1}$  is a piecewise linear mapping of  $\mathcal{F}$  into the subsets of  $\Delta(\mathcal{J})$ . This means

that there exists a finite decomposition of  $\mathcal{F}$  into polytopes  $\{P_1, \dots, P_K\}$  each on which  $\tilde{H}^{-1}$  is linear. Up to a triangulation (see, e.g., Chapter 14 in [12]), we can assume that each  $P_k$  is a simplex. Denote by  $\mathcal{B}_k \subseteq \mathcal{F}$  the set of vertices of  $P_k$ ; then, the finite subset stated in the lemma is

$$\mathcal{B} = \bigcup_{k=1}^K \mathcal{B}_k,$$

the set of all vertices of all the simplices.

Fix any  $\sigma \in \mathcal{F}$ . It belongs to some simplex  $P_k$ , so that there exists a convex decomposition  $\sigma = \sum_{b \in \mathcal{B}_k} \lambda_b b$ ; this decomposition is unique within the simplex  $P_k$ . If  $\sigma$  belongs to two different simplices, then it actually belongs to their common face and the two possible decompositions coincide (some coefficients  $\lambda_b$  in the above decomposition are null). All in all, with each  $\sigma \in \mathcal{F}$ , we can associate a unique decomposition in  $\mathcal{B}$ ,

$$\sigma = \sum_{b \in \mathcal{B}} \Phi_b(\sigma) b,$$

where the coefficients  $(\Phi_b(\sigma))_{b \in \mathcal{B}}$  form a convex weight vector over  $\mathcal{B}$ , i.e., belong to  $\Delta(\mathcal{B})$ ; in addition,  $\Phi_b(\sigma) > 0$  only if  $b \in \mathcal{B}_k$ , where  $k$  is such that  $\sigma \in P_k$ .

Since  $\tilde{H}^{-1}$  is linear on each simplex  $P_1, \dots, P_K$ , we therefore get

$$\tilde{H}^{-1}(\sigma) = \sum_{b \in \mathcal{B}} \Phi_b(\sigma) \tilde{H}^{-1}(b).$$

Finally, the result is a consequence of the fact that

$$\overline{m}(\mathbf{p}, \sigma) = r\left(\mathbf{p}, \tilde{H}^{-1}(\sigma)\right) = r\left(\mathbf{p}, \sum_{b \in \mathcal{B}} \Phi_b(\sigma) \tilde{H}^{-1}(b)\right),$$

which implies, by linearity of  $r$ , that

$$\overline{m}(\mathbf{p}, \sigma) = \sum_{b \in \mathcal{B}} \Phi_b(\sigma) r\left(\mathbf{p}, \tilde{H}^{-1}(b)\right) = \sum_{b \in \mathcal{B}} \Phi_b(\sigma) \overline{m}(\mathbf{p}, b),$$

which concludes the proof.  $\square$

**REMARK 6.1** *The proof shows that  $\Phi$  is piecewise linear on a finite decomposition of  $\mathcal{F}$ ; it is therefore Lipschitz on  $\mathcal{F}$ . We denote by  $\kappa_\Phi$  its Lipschitz constant with respect to the  $\ell^2$ -norms.*

The main contribution of this subsection (Definition 6.1) relies on the following additional assumption.

**ASSUMPTION 6.1** *A game is bi-piecewise linear if  $\overline{m}(\cdot, b)$  is piecewise linear on  $\Delta(\mathcal{I})$  for every  $b \in \mathcal{B}$ .*

Assumption 6.1 means that for all  $b \in \mathcal{B}$  there exists a decomposition of  $\Delta(\mathcal{I})$  into polytopes each on which  $\overline{m}(\cdot, b)$  is linear. Since  $\mathcal{B}$  is finite, there exists a finite number of such decompositions, and thus there exists a decomposition to polytopes that refines all of them. (The latter is generated by the intersection of all considered polytopes as  $b$  varies.) By construction, every  $\overline{m}(\cdot, b)$  is linear on any of the polytopes of this common decomposition. We denote by  $\mathcal{A} \subset \Delta(\mathcal{I})$  the finite subset of all their vertices: a construction similar to the one used in the proof of Lemma 6.1 (provided below) then leads to a piecewise linear (injective) mapping  $\Theta : \Delta(\mathcal{I}) \rightarrow \Delta(\mathcal{A})$ , where  $\Theta(\mathbf{p})$  is the decomposition of  $\mathbf{p}$  on the vertices of the polytope(s) of the decomposition to which it belongs, satisfying

$$\forall b \in \mathcal{B}, \quad \forall \mathbf{p} \in \Delta(\mathcal{I}), \quad \overline{m}(\mathbf{p}, b) = \sum_{a \in \mathcal{A}} \Theta_a(\mathbf{p}) \overline{m}(a, b),$$

where we denoted the convex weight vector  $\Theta(\mathbf{p}) \in \Delta(\mathcal{A})$  by  $(\Theta_a(\mathbf{p}))_{a \in \mathcal{A}}$ . This, Lemma 6.1, and Assumption 6.1 show that on a lifted space,  $\overline{m}$  coincides with a bi-linear mapping  $\overline{\overline{m}}$ , as is made formal in the next definition.

**DEFINITION 6.1** *We denote by  $\overline{\overline{m}}$  the linear extension to  $\Delta(\mathcal{A} \times \mathcal{B})$  of the restriction of  $\overline{m}$  to  $\mathcal{A} \times \mathcal{B}$ , so that for all  $\mathbf{p} \in \Delta(\mathcal{I})$  and  $\sigma \in \mathcal{F}$ ,*

$$\overline{m}(\mathbf{p}, \sigma) = \overline{\overline{m}}(\Theta(\mathbf{p}), \Phi(\sigma)).$$

---

Approaching Strategy in Games with Partial Monitoring

---

*Parameters:* an integer block length  $L \geq 1$ , an exploration parameter  $\gamma \in [0, 1]$ , a strategy  $\Psi$  for  $\overline{\overline{m}}$ -robust approachability of  $\mathcal{C}$

*Notation:*  $\mathbf{u} \in \Delta(\mathcal{I})$  is the uniform distribution over  $\mathcal{I}$ ,  $P_{\mathcal{F}}$  denotes the projection operator in  $\ell^2$ -norm of  $\mathbb{R}^{\mathcal{H} \times \mathcal{I}}$  onto  $\mathcal{F}$

*Initialization:* compute the finite set  $\mathcal{B}$  and the mapping  $\Phi : \mathcal{F} \rightarrow \Delta(\mathcal{B})$  of Lemma 6.1, compute the finite set  $\mathcal{A}$  and the mapping  $\Theta : \Delta(\mathcal{I}) \rightarrow \Delta(\mathcal{A})$  defined based on Assumption 6.1, pick an arbitrary  $\boldsymbol{\theta}_1 \in \Delta(\mathcal{A})$

For all blocks  $n = 1, 2, \dots$ ,

- (i) define  $\mathbf{x}_n = \sum_{a \in \mathcal{A}} \theta_{n,a} a$  and  $\mathbf{p}_n = (1 - \gamma) \mathbf{x}_n + \gamma \mathbf{u}$ ;
- (ii) for rounds  $t = (n - 1)L + 1, \dots, nL$ ,
  - 2.1 draw an action  $I_t \in \mathcal{I}$  at random according to  $\mathbf{p}_n$ ;
  - 2.2 get the signal  $S_t$ ;
- (iii) form the estimated vector of probability distributions over signals,

$$\tilde{\sigma}_n = \left( \frac{1}{L} \sum_{t=(n-1)L+1}^{nL} \frac{\mathbb{I}_{\{S_t=s\}} \mathbb{I}_{\{I_t=i\}}}{p_{I_t,n}} \right)_{(i,s) \in \mathcal{I} \times \mathcal{H}} ;$$

- (iv) compute the projection  $\hat{\sigma}_n = P_{\mathcal{F}}(\tilde{\sigma}_n)$ ;
  - (v) choose  $\boldsymbol{\theta}_{n+1} = \Psi(\boldsymbol{\theta}_1, \Phi(\hat{\sigma}_1), \dots, \boldsymbol{\theta}_n, \Phi(\hat{\sigma}_n))$ .
- 

Figure 1: The proposed strategy, which plays in blocks.

**6.1.2 Construction of a strategy to approach  $\mathcal{C}$ .** The approaching strategy for the original problem is based on a strategy  $\Psi$  for  $\overline{\overline{m}}$ -approachability of  $\mathcal{C}$ , provided by Theorem 3.1; we therefore first need to prove the existence of such a  $\Psi$ .

LEMMA 6.2 *Under Condition (APM), the closed convex set  $\mathcal{C}$  is  $\overline{\overline{m}}$ -robust approachable.*

PROOF. We show that Condition (RAC) in Theorem 3.1 is satisfied, that is, that for all  $\mathbf{y} \in \Delta(\mathcal{B})$ , there exists some  $\mathbf{x} \in \Delta(\mathcal{A})$  such that  $\overline{\overline{m}}(\mathbf{x}, \mathbf{y}) \subseteq \mathcal{C}$ . With such a given  $\mathbf{y} \in \Delta(\mathcal{B})$ , we associate<sup>2</sup> the feasible vector of signals  $\sigma = \sum_{b \in \mathcal{B}} y_b b \in \mathcal{F}$  and let  $\mathbf{p}$  be given by Condition (APM), so that  $\overline{\overline{m}}(\mathbf{p}, \sigma) \subseteq \mathcal{C}$ . By linearity of  $\overline{\overline{m}}$  (for the first equality), by convexity of  $\overline{\overline{m}}$  in its second argument (for the first inclusion), by Lemma 6.1 (for the second and fourth equalities), by construction of  $\mathcal{A}$  (for the third equality),

$$\begin{aligned} \overline{\overline{m}}(\Theta(\mathbf{p}), \mathbf{y}) &= \sum_{a \in \mathcal{A}} \Theta_a(\mathbf{p}) \sum_{b \in \mathcal{B}} y_b \overline{\overline{m}}(a, b) \subseteq \sum_{a \in \mathcal{A}} \Theta_a(\mathbf{p}) \overline{\overline{m}}(a, \sigma) = \sum_{a \in \mathcal{A}} \Theta_a(\mathbf{p}) \sum_{b \in \mathcal{B}} \Phi_b(\sigma) \overline{\overline{m}}(a, b) \\ &= \sum_{b \in \mathcal{B}} \Phi_b(\sigma) \overline{\overline{m}}(\mathbf{p}, b) = \overline{\overline{m}}(\mathbf{p}, \sigma) \subseteq \mathcal{C}, \end{aligned}$$

which concludes the proof. □

We consider the strategy described in Figure 1. It forces exploration at a  $\gamma$  rate, as is usual in situations with partial monitoring. One of its key ingredient, that conditionally unbiased estimators are available, is extracted from Section 6 in the article by Lugosi et al. [16]: in block  $n$  we consider sums of elements of the form

$$\hat{H}_t = \left( \frac{\mathbb{I}_{\{S_t=s\}} \mathbb{I}_{\{I_t=i\}}}{p_{I_t,n}} \right)_{(i,s) \in \mathcal{I} \times \mathcal{H}} \in \mathbb{R}^{\mathcal{H} \times \mathcal{I}};$$

averaging over the respective random draws of  $I_t$  and  $S_t$  according to  $\mathbf{p}_n$  and  $H(I_t, J_t)$ , i.e., taking the conditional expectation  $\mathbb{E}_t$  with respect to  $\mathbf{p}_n$  and  $J_t$ , we get

$$\mathbb{E}_t[\hat{H}_t] = \tilde{H}(\delta_{J_t}). \tag{7}$$

---

<sup>2</sup>Note however that we do not necessarily have that  $\Phi(\sigma)$  and  $\mathbf{y}$  are equal, as  $\Phi$  is not a one-to-one mapping (it is injective but not surjective).

Indeed, the conditional expectation of the component  $i$  of  $\widehat{H}_t$  equals

$$\mathbb{E}_t \left[ \left( \frac{\mathbb{I}_{\{S_t=s\}} \mathbb{I}_{\{I_t=i\}}}{p_{I_t,n}} \right)_{s \in \mathcal{H}} \right] = \mathbb{E}_t \left[ \frac{H(I_t, J_t) \mathbb{I}_{\{I_t=i\}}}{p_{I_t,n}} \right] = \frac{H(i, J_t)}{p_{i,n}} \mathbb{E}_t [\mathbb{I}_{\{I_t=i\}}] = H(i, J_t),$$

where we first took the expectation over the random draw of  $S_t$  (conditionally to  $\mathbf{p}_n$ ,  $J_t$ , and  $I_t$ ) and then over the one of  $I_t$ . Consequently, concentration-of-the-measure arguments can show that for  $L$  large enough,

$$\tilde{\sigma}_n = \frac{1}{L} \sum_{t=(n-1)L+1}^{nL} \widehat{H}_t \quad \text{is close to} \quad \widetilde{H}(\widehat{\mathbf{q}}_n), \quad \text{where} \quad \widehat{\mathbf{q}}_n = \frac{1}{L} \sum_{t=(n-1)L+1}^{nL} \delta_{J_t}.$$

Actually, since  $\mathcal{F} \subseteq \Delta(\mathcal{H})^{\mathcal{I}}$ , we have a natural embedding of  $\mathcal{F}$  into  $\mathbb{R}^{\mathcal{H} \times \mathcal{I}}$  and we can define  $P_{\mathcal{F}}$ , the convex projection operator onto  $\mathcal{F}$  (in  $\ell^2$ -norm). Instead of using directly  $\tilde{\sigma}_n$ , we consider in our strategy  $\widehat{\sigma}_n = P_{\mathcal{F}}(\tilde{\sigma}_n)$ , which is even closer to  $\widetilde{H}(\widehat{\mathbf{q}}_n)$ .

More precisely, the following result can be extracted from the proof of Theorem 6.1 in Lugosi et al. [16]. The proof is provided in Appendix B.

LEMMA 6.3 *With probability  $1 - \delta$ ,*

$$\left\| \widehat{\sigma}_n - \widetilde{H}(\widehat{\mathbf{q}}_n) \right\|_2 \leq \sqrt{N_{\mathcal{I}} N_{\mathcal{H}}} \left( \sqrt{\frac{2N_{\mathcal{I}}}{\gamma L} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta}} + \frac{1}{3} \frac{N_{\mathcal{I}}}{\gamma L} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta} \right).$$

**6.1.3 A performance guarantee for the strategy of Figure 1.** For the sake of simplicity, we provide first a performance bound for fixed parameters  $\gamma$  and  $L$  tuned as functions of  $T$ . Adaptation to  $T \rightarrow \infty$  is then described in the next section; note that it cannot be performed by simply proceeding in regimes, as the approachability guarantees offered by the second part of the theorem are only at time round  $T$ . (This is so because the considered strategy depends on  $T$  via the parameters  $\gamma$  and  $L$ .)

THEOREM 6.1 *Consider a closed convex set  $\mathcal{C}$  and a game  $(r, H)$  for which Condition (APM) is satisfied and that is bi-piecewise linear in the sense of Assumption 6.1. Then, for all  $T \geq 1$ , the strategy of Figure 1, run with parameters  $\gamma \in [0, 1]$  and  $L \geq 1$  and fed with a strategy  $\Psi$  for  $\overline{m}$ -approachability of  $\mathcal{C}$  (provided by Lemma 6.2) is such that, with probability at least  $1 - \delta$ ,*

$$\inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T r(I_t, J_t) \right\|_2 \leq \frac{2L}{T} R + 4R \sqrt{\frac{\ln((2T)/(L\delta))}{T}} + 2\gamma R + \frac{2R}{\sqrt{T/L-1}} \sqrt{N_{\mathcal{A}} N_{\mathcal{B}}} \\ + R\kappa_{\Phi} \sqrt{N_{\mathcal{I}} N_{\mathcal{H}} N_{\mathcal{A}}} \left( \sqrt{\frac{2N_{\mathcal{I}}}{\gamma L} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}} T}{L\delta}} + \frac{1}{3} \frac{N_{\mathcal{I}}}{\gamma L} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}} T}{L\delta} \right).$$

*In particular, for all  $T \geq 1$ , the choices of  $L = \lceil T^{3/5} \rceil$  and  $\gamma = T^{-1/5}$  imply that with probability at least  $1 - \delta$ ,*

$$\inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T r(I_t, J_t) \right\|_2 \leq \square \left( T^{-1/5} \sqrt{\ln \frac{T}{\delta}} + T^{-2/5} \ln \frac{T}{\delta} \right)$$

*for some constant  $\square$  depending only on  $\mathcal{C}$  and on the game  $(r, H)$  at hand.*

The efficiency of the strategy of Figure 1 depends on whether it can be fed with an efficient approachability strategy  $\Psi$ , which in turn depends on the respective geometries of  $\overline{m}$  and  $\mathcal{C}$ , as was indicated before the statement of Theorem 3.1. (Note that the projection onto  $\mathcal{F}$  can be performed in polynomial time, as the latter closed convex set is defined by finitely many linear constraints, and that the computation of  $\mathcal{A}$ ,  $\mathcal{B}$ , and  $\overline{m}$  can be performed beforehand.) In any case, the per-round complexity is constant (though possibly large).

PROOF. We write  $T$  as  $T = NL + k$  where  $N$  is an integer and  $0 \leq k \leq L - 1$  and will show successively that (possibly with overwhelming probability only) the following statements hold.

$$\frac{1}{T} \sum_{t=1}^T r(I_t, J_t) \quad \text{is close to} \quad \frac{1}{NL} \sum_{t=1}^{NL} r(I_t, J_t); \quad (8)$$

$$\frac{1}{NL} \sum_{t=1}^{NL} r(I_t, J_t) \quad \text{is close to} \quad \frac{1}{N} \sum_{n=1}^N r(\mathbf{p}_n, \hat{\mathbf{q}}_n); \quad (9)$$

$$\frac{1}{N} \sum_{n=1}^N r(\mathbf{p}_n, \hat{\mathbf{q}}_n) \quad \text{is close to} \quad \frac{1}{N} \sum_{n=1}^N r(\mathbf{x}_n, \hat{\mathbf{q}}_n); \quad (10)$$

$$\frac{1}{N} \sum_{n=1}^N r(\mathbf{x}_n, \hat{\mathbf{q}}_n) = \frac{1}{N} \sum_{n=1}^N \sum_{a \in \mathcal{A}} \theta_{n,a} r(a, \hat{\mathbf{q}}_n) \quad \text{belongs to the set} \quad \frac{1}{N} \sum_{n=1}^N \sum_{a \in \mathcal{A}} \theta_{n,a} \bar{m}(a, \tilde{H}(\hat{\mathbf{q}}_n));$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{a \in \mathcal{A}} \theta_{n,a} \bar{m}(a, \tilde{H}(\hat{\mathbf{q}}_n)) \quad \text{is equal to the set} \quad \frac{1}{N} \sum_{n=1}^N \bar{m}(\boldsymbol{\theta}_n, \Phi(\tilde{H}(\hat{\mathbf{q}}_n)));$$

$$\frac{1}{N} \sum_{n=1}^N \bar{m}(\boldsymbol{\theta}_n, \Phi(\tilde{H}(\hat{\mathbf{q}}_n))) \quad \text{is close to the set} \quad \frac{1}{N} \sum_{n=1}^N \bar{m}(\boldsymbol{\theta}_n, \Phi(\hat{\sigma}_n)); \quad (11)$$

$$\frac{1}{N} \sum_{n=1}^N \bar{m}(\boldsymbol{\theta}_n, \Phi(\hat{\sigma}_n)) \quad \text{is close to the set} \quad \mathcal{C}; \quad (12)$$

where we recall that the notation  $\hat{\mathbf{q}}_n$  was defined above and is referring to the empirical distribution of the  $J_t$  in the  $n$ -th block. Actually, we will show below the numbered statements only. The first unnumbered statement is immediate by the definition of  $\mathbf{x}_n$ , the linearity of  $r$ , and the very definition of  $\bar{m}$ ; while the second one follows from Definition 6.1:

$$\frac{1}{N} \sum_{n=1}^N \sum_{a \in \mathcal{A}} \theta_{n,a} \bar{m}(a, \tilde{H}(\hat{\mathbf{q}}_n)) = \frac{1}{N} \sum_{n=1}^N \sum_{(a,b) \in \mathcal{A} \times \mathcal{B}} \theta_{n,a} \Phi_b(\tilde{H}(\hat{\mathbf{q}}_n)) \bar{m}(a,b) = \frac{1}{N} \sum_{n=1}^N \bar{m}(\boldsymbol{\theta}_n, \Phi(\tilde{H}(\hat{\mathbf{q}}_n))).$$

**Step 1: Assertion (8).** A direct calculation decomposing the sum over  $T$  elements into a sum over the  $NL$  first elements and the  $k$  remaining ones shows that

$$\left\| \frac{1}{T} \sum_{t=1}^T r(I_t, J_t) - \frac{1}{NL} \sum_{t=1}^{NL} r(I_t, J_t) \right\|_2 \leq R \left( \frac{k}{T} + \left( \frac{1}{NL} - \frac{1}{T} \right) NL \right) = \frac{2k}{T} R \leq \frac{2L}{T} R.$$

**Step 2: Assertion (9).** We note that by defining  $\mathbb{E}_t$  the conditional expectation with respect to  $(I_1, S_1, J_1), \dots, (I_{t-1}, S_{t-1}, J_{t-1})$  and  $J_t$ , which fixes the values of the distribution  $\mathbf{p}'_t$  of  $I_t$  and the value of  $J_t$ , we have

$$\mathbb{E}_t[r(I_t, J_t)] = r(\mathbf{p}'_t, J_t).$$

We note that by definition of the forecaster,  $\mathbf{p}'_t = \mathbf{p}_n$  if  $t$  belongs to the  $n$ -th block. By a version of the Hoeffding-Azuma inequality for sums of Hilbert space-valued martingale differences stated as<sup>3</sup> Lemma 3.2 in Chen and White [7], we therefore get that with probability at least  $1 - \delta$ ,

$$\left\| \frac{1}{NL} \sum_{t=1}^{NL} r(I_t, J_t) - \frac{1}{N} \sum_{n=1}^N r(\mathbf{p}_n, \hat{\mathbf{q}}_n) \right\|_2 \leq 4R \sqrt{\frac{\ln(2/\delta)}{T}}.$$

**Step 3: Assertion (10).** Since by definition  $\mathbf{p}_n = (1 - \gamma) \mathbf{x}_n + \gamma \mathbf{u}$ , we get

$$\left\| \frac{1}{N} \sum_{n=1}^N r(\mathbf{p}_n, \hat{\mathbf{q}}_n) - \frac{1}{N} \sum_{n=1}^N r(\mathbf{x}_n, \hat{\mathbf{q}}_n) \right\|_2 \leq 2\gamma R.$$

**Step 4: Assertion (11).** We fix a given block  $n$ . Lemma 6.3 indicates that with probability  $1 - \delta$ ,

$$\left\| \hat{\sigma}_n - \tilde{H}(\hat{\mathbf{q}}_n) \right\|_2 \leq \sqrt{N_{\mathcal{I}} N_{\mathcal{H}}} \left( \sqrt{\frac{2N_{\mathcal{I}}}{\gamma L} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta}} + \frac{1}{3} \frac{N_{\mathcal{I}}}{\gamma L} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta} \right). \quad (13)$$

Since  $\Phi$  is Lipschitz (see Remark 6.1), with a Lipschitz constant in  $\ell^2$ -norms denoted by  $\kappa_{\Phi}$ , we get that with probability  $1 - \delta$ ,

$$\left\| \Phi(\hat{\sigma}_n) - \Phi(\tilde{H}(\hat{\mathbf{q}}_n)) \right\|_2 \leq \kappa_{\Phi} \sqrt{N_{\mathcal{I}} N_{\mathcal{H}}} \left( \sqrt{\frac{2N_{\mathcal{I}}}{\gamma L} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta}} + \frac{1}{3} \frac{N_{\mathcal{I}}}{\gamma L} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta} \right).$$

<sup>3</sup>Together with the fact that  $\sqrt{u} e^{-u} \leq e^{-u/2}$  for all  $u \geq 0$ .

By a union bound, the above bound holds for all blocks  $n = 1, \dots, N$  with probability at least  $1 - N\delta$ . Finally, an application of Lemma 3.1 shows that

$$\frac{1}{N} \sum_{n=1}^N \overline{m} \left( \boldsymbol{\theta}_n, \Phi \left( \tilde{H}(\hat{\boldsymbol{q}}_n) \right) \right) \quad \text{is in a } \varepsilon_T\text{-neighborhood (in } \ell^2\text{-norm) of } \frac{1}{N} \sum_{n=1}^N \overline{m} \left( \boldsymbol{\theta}_n, \Phi(\hat{\boldsymbol{\sigma}}_n) \right),$$

where

$$\varepsilon_T = R\sqrt{N_{\mathcal{B}}} \times \kappa_{\Phi} \sqrt{N_{\mathcal{I}}N_{\mathcal{H}}} \left( \sqrt{\frac{2N_{\mathcal{I}}}{\gamma L} \ln \frac{2N_{\mathcal{I}}N_{\mathcal{H}}}{\delta}} + \frac{1}{3} \frac{N_{\mathcal{I}}}{\gamma L} \ln \frac{2N_{\mathcal{I}}N_{\mathcal{H}}}{\delta} \right).$$

**Step 5: Assertion (12).** Since  $\mathcal{C}$  is  $\overline{m}$ -robust approachable and by definition of the choices of the  $\boldsymbol{\theta}_n$  in Figure 1, we get by (the proof of the sufficiency part of) Theorem 3.1, with probability 1,

$$\inf_{c \in \mathcal{C}} \left\| c - \frac{1}{N} \sum_{n=1}^N \overline{m} \left( \boldsymbol{\theta}_n, \Phi(\hat{\boldsymbol{\sigma}}_n) \right) \right\|_2 \leq \frac{2R}{\sqrt{N}} \sqrt{N_{\mathcal{A}}N_{\mathcal{B}}} \leq \frac{2R}{\sqrt{T/L - 1}} \sqrt{N_{\mathcal{A}}N_{\mathcal{B}}},$$

since  $T/L \leq N + k/L \leq N + 1$ .

**Conclusion of the proof.** The proof is concluded by putting the pieces together, thanks to a triangle inequality and by considering  $L\delta/T \leq \delta/(N + 1)$  instead of  $\delta$ .  $\square$

#### 6.1.4 Uniform guarantees over time for a time-adaptive version of the strategy of Figure 1.

We present here a variant of the strategy of Figure 1 for which the lengths  $L_n$  of blocks  $n$  and the exploration rates  $\gamma_n$  are no longer constant. To do so, we need the following generalization of Theorem 2.2 to polynomial averages; this result is of independent interest. We only state the result for mixed actions taken and observed, but the generalization for pure actions follows easily.

Consider the setting of Theorem 2.2. The studied strategy relies on a parameter  $\alpha \geq 0$ . It plays an arbitrary  $\boldsymbol{x}_1$ . For  $t \geq 1$ , it forms at stage  $t + 1$  the vector-valued polynomial average

$$\hat{m}_t^\alpha = \frac{1}{T_t^\alpha} \sum_{s=1}^t s^\alpha m(\boldsymbol{x}_s, \boldsymbol{y}_s) \quad \text{where} \quad T_t^\alpha = \sum_{s=1}^t s^\alpha,$$

computes its projection  $c_t^\alpha$  onto  $\mathcal{C}$ , and resorts to a mixed action  $\boldsymbol{x}_{t+1}$  solving the minimax equation

$$\min_{\boldsymbol{x} \in \Delta(\mathcal{A})} \max_{\boldsymbol{y} \in \Delta(\mathcal{B})} \langle \hat{m}_t^\alpha - c_t^\alpha, m(\boldsymbol{x}, \boldsymbol{y}) \rangle.$$

**THEOREM 6.2** *We denote by  $M$  a bound in norm over  $m$ , i.e.,*

$$\max_{(a,b) \in \mathcal{A} \times \mathcal{B}} \|m(a,b)\|_2 \leq M.$$

*For all  $\alpha \geq 0$ , when  $\mathcal{C}$  is an approachable closed convex set, the above strategy ensures that for all strategies of the second player, with probability 1, for all  $T \geq 1$ ,*

$$\inf_{c \in \mathcal{C}} \left\| c - \frac{1}{\sum_{t=1}^T t^\alpha} \sum_{t=1}^T t^\alpha m(\boldsymbol{x}_t, \boldsymbol{y}_t) \right\|_2 \leq 2M \frac{\sqrt{\sum_{t=1}^T t^{2\alpha}}}{\sum_{t=1}^T t^\alpha} \leq \frac{2MK_\alpha}{\sqrt{T}}, \quad (14)$$

where  $K_\alpha$  is a constant depending only  $\alpha$ .

It is interesting to note that the convergence rate are independent of  $\alpha$  and are the same as standard approachability ( $1/\sqrt{T}$ ).

**PROOF.** The proof is a slight modification of the one of Theorem 2.2. We denote by  $d_t^\alpha$  the squared distance of  $\hat{m}_t^\alpha$  to  $\mathcal{C}$ ,

$$d_t^\alpha = \inf_{c \in \mathcal{C}} \|c - \hat{m}_t^\alpha\|^2 = \|c_t^\alpha - \hat{m}_t^\alpha\|^2$$

and use the shortcut notation  $m_t = m(\mathbf{x}_t, \mathbf{y}_t)$  for all  $t \geq 1$ . Then,

$$\begin{aligned} d_{t+1}^\alpha &\leq \|\widehat{m}_{t+1}^\alpha - c_t^\alpha\|^2 = \left\| \widehat{m}_t^\alpha - c_t^\alpha + \frac{(t+1)^\alpha}{T_{t+1}^\alpha} (m_{t+1} - \widehat{m}_t^\alpha) \right\|^2 \\ &\leq \|\widehat{m}_t^\alpha - c_t^\alpha\|^2 + \frac{2(t+1)^\alpha}{T_{t+1}^\alpha} \langle \widehat{m}_t^\alpha - c_t^\alpha, m_{t+1} - m_t^\alpha \rangle + \left( \frac{(t+1)^\alpha}{T_{t+1}^\alpha} \right)^2 \|m_{t+1} - \widehat{m}_t^\alpha\|^2 \\ &\leq d_t^\alpha + \frac{2(t+1)^\alpha}{T_{t+1}^\alpha} \underbrace{\langle \widehat{m}_t^\alpha - c_t^\alpha, m_{t+1} - c_t^\alpha \rangle}_{\leq 0} + \frac{2(t+1)^\alpha}{T_{t+1}^\alpha} \langle \widehat{m}_t^\alpha - c_t^\alpha, c_t^\alpha - m_t^\alpha \rangle + \left( \frac{(t+1)^\alpha}{T_{t+1}^\alpha} \right)^2 4M^2 \\ &\leq d_t^\alpha \left( 1 - \frac{2(t+1)^\alpha}{T_{t+1}^\alpha} \right) + \left( \frac{(t+1)^\alpha}{T_{t+1}^\alpha} \right)^2 4M^2, \end{aligned}$$

where we used in the third inequality the same convex projection inequality as in the proof of Theorem 2.2.

The first inequality in (14) then follows by induction: the bound  $2M$  for  $t = 1$  is by boundedness of  $m$ . If the stated bound holds for  $d_t^\alpha$ , then

$$d_{t+1}^\alpha \leq \left( 2M \frac{\sqrt{\sum_{s=1}^t s^{2\alpha}}}{\sum_{s=1}^t s^\alpha} \right)^2 \left( 1 - \frac{2(t+1)^\alpha}{T_{t+1}^\alpha} \right) + \left( \frac{(t+1)^\alpha}{T_{t+1}^\alpha} \right)^2 4M^2 \leq 4M^2 \frac{\sum_{s=1}^{t+1} s^{2\alpha}}{(T_{t+1}^\alpha)^2},$$

as desired, since

$$\frac{1}{(T_t^\alpha)^2} \left( 1 - \frac{2(t+1)^\alpha}{T_{t+1}^\alpha} \right) = \frac{1}{T_{t+1}^\alpha (T_t^\alpha)^2} (T_t^\alpha - (t+1)^\alpha) \leq \frac{1}{T_{t+1}^\alpha (T_t^\alpha)^2} \frac{(T_t^\alpha)^2 - (t+1)^{2\alpha}}{T_t^\alpha + (t+1)^\alpha} \leq \frac{1}{(T_{t+1}^\alpha)^2}.$$

The second inequality in (14) can be proved as follows. First, for all  $\alpha \geq 0$ , by comparing sums and integrals, we get that for all  $t \geq 1$ ,

$$\frac{t^{\alpha+1}}{\alpha+1} = \int_0^t s^\alpha ds \leq \sum_{s=1}^t s^\alpha \leq \int_1^{t+1} s^\alpha ds \leq \frac{(t+1)^{\alpha+1}}{\alpha+1} \leq \frac{(2t)^{\alpha+1}}{\alpha+1}.$$

Therefore,

$$\frac{\sqrt{\sum_{s=1}^t s^{2\alpha}}}{\sum_{s=1}^t s^\alpha} \leq \frac{\alpha+1}{\sqrt{2\alpha+1}} \frac{\sqrt{(2t)^{\alpha+1}}}{t^{\alpha+1}} = K_\alpha \frac{1}{\sqrt{t}}$$

for

$$K_\alpha = \frac{\alpha+1}{\sqrt{2\alpha+1}} \sqrt{2^{\alpha+1}}.$$

This concludes the proof.  $\square$

The extension to polynomially weighted averages can also be obtained in the context of robust approachability as the key to Theorem 3.1 is Lemma 3.3, which indicates that to get robust approachability, it suffices to approach, in the usual sense,  $\widetilde{\mathcal{C}}$ ; both can thus be performed with polynomially weighted averages.

Consider now the variant of the strategy of Figure 1 for which the length of the  $n$ -th block, denoted by  $L_n$ , is equal to  $n^\alpha$ , the exploration rate on this block comes at a rate  $\gamma_n = n^{-\alpha/3}$  and  $\Psi$  is an  $\overline{m}$ -robust approachability strategy of  $\mathcal{C}$  with respect to polynomially weighted averages with parameter  $\alpha = 3/2$ . We call it a time-adaptive version of this strategy; note that it does not depend anymore on any time horizon  $T$ , hence guarantees can be obtained for all  $T$ .

**THEOREM 6.3** *The time-adaptive version of the strategy described in Figure 1 (with  $L_n = n^\alpha$  and  $\gamma_n = n^{-\alpha/3}$  for  $\alpha = 3/2$ ) ensures that, for all  $T \geq 1$ , with probability at least  $1 - \delta$ ,*

$$\inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T r(I_t, J_t) \right\|_2 \leq \square \left( T^{-1/5} \sqrt{\ln \frac{T}{\delta}} + T^{-2/5} \ln \frac{T}{\delta} \right)$$

for some constant  $\square$  depending only on  $\mathcal{C}$  and the game  $(r, H)$  at hand.

PROOF. The proof follows closely the one of Theorem 6.1. We choose  $N$  so as to write  $T = T_N^\alpha + k$  where  $0 \leq k \leq L_{N+1} - 1$ . We adapt step 1 as follows,

$$\left\| \frac{1}{T} \sum_{t=1}^T r(I_t, J_t) - \frac{1}{T_N^\alpha} \sum_{t=1}^{T_N^\alpha} r(I_t, J_t) \right\|_2 \leq R \left( \frac{k}{T} + \left( \frac{1}{T_N^\alpha} - \frac{1}{T} \right) T_N^\alpha \right) = \frac{2k}{T} R \leq \frac{2L_{N+1}}{T} R.$$

Second, as in step 2, we resort again to the Hoeffding-Azuma inequality for sums of Hilbert space-valued martingale differences; with probability at least  $1 - \delta$ ,

$$\left\| \frac{1}{T_N^\alpha} \sum_{t=1}^{T_N^\alpha} r(I_t, J_t) - \frac{1}{T_N^\alpha} \sum_{n=1}^N n^\alpha r(\mathbf{p}_n, \hat{\mathbf{q}}_n) \right\|_2 \leq 4R \sqrt{\frac{\ln(2/\delta)}{T_N^\alpha}} \leq 4R \sqrt{\frac{\ln(2/\delta)}{T}}.$$

In view of the choice  $\gamma_n = n^{-\alpha/3}$ , step 3 translates here to

$$\left\| \frac{1}{T_N^\alpha} \sum_{n=1}^N n^\alpha r(\mathbf{p}_n, \hat{\mathbf{q}}_n) - \frac{1}{T_N^\alpha} \sum_{n=1}^N n^\alpha r(\mathbf{x}_n, \hat{\mathbf{q}}_n) \right\|_2 \leq 2R \frac{\sum_{n=1}^N n^\alpha \gamma_n}{T_N^\alpha} = 2R \frac{\sum_{n=1}^N n^{2\alpha/3}}{T_N^\alpha} = 2R \frac{T_N^{(2\alpha/3)}}{T_N^\alpha}.$$

The same argument as the one at the beginning of the proof of Theorem 6.1 shows that

$$\frac{1}{T_N^\alpha} \sum_{n=1}^N n^\alpha r(\mathbf{x}_n, \hat{\mathbf{q}}_n) \in \frac{1}{T_N^\alpha} \sum_{n=1}^N n^\alpha \bar{m} \left( \boldsymbol{\theta}_n, \Phi(\tilde{H}(\hat{\mathbf{q}}_n)) \right).$$

Step 4 starts also by an application of Lemma 6.3 together with the Lipschitzness of  $\Phi$  to get that for all regimes  $n = 1, \dots, N$ , with probability at least  $1 - \delta$ ,

$$\left\| \Phi(\hat{\sigma}_n) - \Phi(\tilde{H}(\hat{\mathbf{q}}_n)) \right\|_2 \leq \kappa_\Phi \sqrt{N_{\mathcal{I}} N_{\mathcal{H}}} \left( \sqrt{\frac{2N_{\mathcal{I}}}{\gamma_n L_n} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta}} + \frac{1}{3} \frac{N_{\mathcal{I}}}{\gamma_n L_n} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta} \right).$$

By a union bound, the above bound holds for all regimes  $n = 1, \dots, N$  with probability at least  $1 - N\delta$ . Then, an application of Lemma 3.1 shows that

$$\frac{1}{T_N^\alpha} \sum_{n=1}^N n^\alpha \bar{m} \left( \boldsymbol{\theta}_n, \Phi(\tilde{H}(\hat{\mathbf{q}}_n)) \right) \quad \text{is in a } \varepsilon_N\text{-neighborhood of} \quad \frac{1}{T_N^\alpha} \sum_{n=1}^N n^\alpha \bar{m} \left( \boldsymbol{\theta}_n, \Phi(\hat{\sigma}_n) \right),$$

where, substituting the values of  $L_n = n^\alpha$  and  $\gamma_n = n^{-\alpha/3}$ ,

$$\begin{aligned} \varepsilon_N &= R \sqrt{N_{\mathcal{B}}} \times \kappa_\Phi \sqrt{N_{\mathcal{I}} N_{\mathcal{H}}} \frac{1}{T_N^\alpha} \sum_{n=1}^N n^\alpha \left( \sqrt{\frac{2N_{\mathcal{I}}}{\gamma_n L_n} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta}} + \frac{1}{3} \frac{N_{\mathcal{I}}}{\gamma_n L_n} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta} \right) \\ &= R \sqrt{N_{\mathcal{B}}} \times \kappa_\Phi \sqrt{N_{\mathcal{I}} N_{\mathcal{H}}} \left( \frac{T_N^{(2\alpha/3)}}{T_N^\alpha} \sqrt{2N_{\mathcal{I}} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta}} + \frac{T_N^{(\alpha/3)}}{T_N^\alpha} \frac{N_{\mathcal{I}}}{3} \ln \frac{2N_{\mathcal{I}} N_{\mathcal{H}}}{\delta} \right). \end{aligned}$$

It then suffices, as in step 5 of the original proof, to write the convergence rates for robust approachability guaranteed by the strategy  $\Psi$ . By combining the result of Lemma 3.3 with Theorem 6.2 and Lemma 3.1, we get

$$\inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T_n} \sum_{n=1}^N n^\alpha \bar{m} \left( \boldsymbol{\theta}_n, \Phi(\hat{\sigma}_n) \right) \right\|_2 \leq \frac{2R K_\alpha}{\sqrt{N}} \sqrt{N_{\mathcal{A}} N_{\mathcal{B}}}.$$

Putting all things together and applying a union bound, we obtain that with probability at least  $1 - \delta$ ,

$$\begin{aligned} \inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T r(I_t, J_t) \right\|_2 \\ = O \left( \frac{(N+1)^\alpha}{T} + \sqrt{\frac{\ln(N/\delta)}{T_n^\alpha}} + \frac{T_N^{(2\alpha/3)}}{T_N^\alpha} + \frac{T_N^{(2\alpha/3)}}{T_N^\alpha} \sqrt{\ln \frac{N}{\delta}} + \frac{T_N^{(\alpha/3)}}{T_N^\alpha} \ln \frac{N}{\delta} + \frac{1}{\sqrt{N}} \right). \end{aligned}$$

Since (as proved at the end of Theorem 6.2)  $T_N^\beta \sim N^{\beta+1}/(\beta+1)$  for all  $\beta \geq 0$ , we get that

$$N \sim ((\alpha+1)T)^{1/(\alpha+1)} \quad \text{and} \quad T_N^\beta \sim \frac{N^{\beta+1}}{\beta+1} \sim \kappa_{\alpha,\beta} T^{(\beta+1)/(\alpha+1)},$$

where  $\kappa_{\alpha,\beta}$  is a constant that only depends on  $\alpha$  and  $\beta$ . Choosing  $\alpha = 3/2$  and substituting these equivalences ensures the result.  $\square$

**6.2 Application to regret minimization.** In this section we analyze external and internal regret minimization in repeated games with partial monitoring from the approachability perspective. We show how to—in particular—efficiently minimize regret in both setups using the results developed for vector-valued games with partial monitoring; to do so, we indicate why the assumption of bi-piecewise linearity (Assumption 6.1) is satisfied.

**6.2.1 External regret.** We consider in this section the framework and aim introduced by Rustichini [29] and studied, sometimes in special cases, by Piccolboni and Schindelhauer [26], Mannor and Shimkin [17], Cesa-Bianchi et al. [6], Lugosi et al. [16]. We show that our general strategy can be used for regret minimization.

Scalar payoffs are obtained (but not observed) by the first player, i.e.,  $d = 1$ : the payoff function  $r$  is a mapping  $\mathcal{I} \times \mathcal{J} \rightarrow \mathbb{R}$ ; we still denote by  $R$  a bound on  $|r|$ . We define in this section

$$\hat{\mathbf{q}}_T = \frac{1}{T} \sum_{t=1}^T \delta_{J_t}$$

as the empirical distribution of the actions taken by the second player during the first  $T$  rounds. (This is in contrast with the notation  $\hat{\mathbf{q}}_T$  used in the previous section to denote such an empirical distribution, but only taken within regime  $n$ .)

The external regret of the first player at round  $T$  equals by definition

$$R_T^{\text{ext}} = \max_{\mathbf{p} \in \Delta(\mathcal{I})} \rho(\mathbf{p}, \tilde{H}(\hat{\mathbf{q}}_T)) - \frac{1}{T} \sum_{t=1}^T r(I_t, J_t),$$

where  $\rho : \Delta(\mathcal{I}) \times \mathcal{F}$  is defined as follows: for all  $\mathbf{p} \in \Delta(\mathcal{I})$  and  $\sigma \in \mathcal{F}$ ,

$$\rho(\mathbf{p}, \sigma) = \min \left\{ r(\mathbf{p}, \mathbf{q}) : \mathbf{q} \text{ such that } \tilde{H}(\mathbf{q}) = \sigma \right\}.$$

The function  $\rho$  is continuous in its first argument and therefore the supremum in the defining expression of  $R_T^{\text{ext}}$  is a maximum.

We recall briefly why, intuitively, this is the natural notion of external regret to consider in this case. Indeed, the first term in the definition of  $R_T^{\text{ext}}$  is (close to) the worst-case average payoff obtained by the first player when playing consistently a mixed action  $\mathbf{p}$  against a sequence of mixed actions inducing on average the same laws on the signals as the sequence of actions actually played.

The following result is an easy consequence of Theorem 6.3, as is explained below; it corresponds to the main result of Lugosi et al. [16], with the same convergence rate but with a different strategy. (However, Section 2.3 of Perchet [24] exhibited an efficient strategy achieving a convergence rate of order  $T^{-1/3}$ , which is optimal; a question that remains open is thus whether the rates exhibited in Theorem 6.3 could be improved.)

**COROLLARY 6.1** *The first player has a strategy such that for all  $T$  and all strategies of the second player, with probability at least  $1 - \delta$ ,*

$$R_T^{\text{ext}} \leq \square \left( T^{-1/5} \sqrt{\ln \frac{T}{\delta}} + T^{-2/5} \ln \frac{T}{\delta} \right)$$

for some constant  $\square$  depending only on the game  $(r, H)$  at hand.

The proof below is an extension to the setting of partial monitoring of the original proof and strategy of Blackwell [3] for the case of external regret under full monitoring: in the latter case the vector-payoff function  $\underline{r}$  and the set  $\mathcal{C}$  considered in our proof are equal to the ones considered by Blackwell.

**PROOF.** We embed  $\mathcal{F}$  into  $\mathbb{R}^{\mathcal{I} \times \mathcal{H}}$  so that in this proof we will be working in the vector space  $\mathbb{R}^d = \mathbb{R} \times \mathbb{R}^{\mathcal{I} \times \mathcal{H}}$ . We consider the closed convex set  $\mathcal{C}$  and the vector-valued payoff function  $\underline{r}$  respectively defined by

$$\mathcal{C} = \left\{ (z, \sigma) \in \mathbb{R} \times \mathcal{F} : z \geq \max_{\mathbf{p} \in \Delta(\mathcal{I})} \rho(\mathbf{p}, \sigma) \right\} \quad \text{and} \quad \underline{r}(i, j) = \begin{bmatrix} r(i, j) \\ \tilde{H}(\delta_j) \end{bmatrix},$$

for all  $(i, j) \in \mathcal{I} \times \mathcal{J}$ .

We first show that Condition (APM) is satisfied for the considered convex set  $\mathcal{C}$  and game  $(\underline{r}, H)$ . To do so, by continuity of  $\rho$  in its first argument, we associate with each  $\mathbf{q} \in \Delta(\mathcal{J})$  an element  $\phi(\mathbf{q}) \in \Delta(\mathcal{I})$  such that

$$\phi(\mathbf{q}) \in \operatorname{argmax}_{\mathbf{p} \in \Delta(\mathcal{I})} \rho(\mathbf{p}, \tilde{H}(\mathbf{q})).$$

Then, given any  $\mathbf{q} \in \Delta(\mathcal{J})$ , we note that for all  $\mathbf{q}'$  satisfying  $\tilde{H}(\mathbf{q}') = \tilde{H}(\mathbf{q})$ , we have by definition of  $\rho$ ,

$$r(\phi(\mathbf{q}), \mathbf{q}') \geq \rho(\phi(\mathbf{q}), \tilde{H}(\mathbf{q}')) = \max_{\mathbf{p} \in \Delta(\mathcal{I})} \rho(\mathbf{p}, \tilde{H}(\mathbf{q}')),$$

which shows that  $\underline{r}(\phi(\mathbf{q}), \mathbf{q}') \in \mathcal{C}$ . The required condition is thus satisfied.

We then show that Assumption 6.1 is satisfied. To do so, we will actually prove the stronger property that the mappings  $\bar{m}(\cdot, \sigma)$  are piecewise linear for all  $\sigma \in \mathcal{F}$ ; we fix such a  $\sigma$  in the sequel. Only the first coordinate  $r$  of  $\underline{r}$  depends on  $\mathbf{p}$ , so the desired property is true if and only if the mapping  $\bar{m}_1(\cdot, \sigma)$  defined by

$$\mathbf{p} \in \Delta(\mathcal{I}) \longmapsto \bar{m}_1(\mathbf{p}, \sigma) = \left\{ r(\mathbf{p}, \mathbf{q}) : \mathbf{q} \in \Delta(\mathcal{J}) \text{ such that } \tilde{H}(\mathbf{q}) = \sigma \right\}$$

is piecewise linear. Since  $\tilde{H}$  is linear, the set

$$\left\{ \mathbf{q} \in \Delta(\mathcal{J}) \text{ such that } \tilde{H}(\mathbf{q}) = \sigma \right\}$$

is a polytope, thus, the convex hull of some finite set  $\{\mathbf{q}_{\sigma,1}, \dots, \mathbf{q}_{\sigma,M}\} \subset \Delta(\mathcal{J})$ . Therefore, for every  $\mathbf{p} \in \Delta(\mathcal{I})$ , by linearity of  $r$  (and by the fact that it takes one-dimensional values),

$$\bar{m}_1(\mathbf{p}, \sigma) = \operatorname{co} \left\{ r(\mathbf{p}, \mathbf{q}_{\sigma,1}), \dots, r(\mathbf{p}, \mathbf{q}_{\sigma,M}) \right\} = \left[ \min_{k \in \{1, \dots, M\}} r(\mathbf{p}, \mathbf{q}_{\sigma,k}), \max_{k' \in \{1, \dots, M\}} r(\mathbf{p}, \mathbf{q}_{\sigma,k'}) \right], \quad (15)$$

where  $\operatorname{co}$  stands for the convex hull. Since all applications  $r(\cdot, \mathbf{q}_{\sigma,k})$  are linear, their minimum and their maximum are piecewise linear functions, thus  $\bar{m}_1(\cdot, \sigma)$  is also piecewise linear. Assumption 6.1 is thus satisfied, as claimed.

Theorem 6.1 can therefore be applied to exhibit the convergence rates; we simply need to relate the quantity of interest here to the one considered therein. To that end we use the fact that the mapping

$$\sigma \in \mathcal{F} \longmapsto \max_{\mathbf{p} \in \Delta(\mathcal{I})} \rho(\mathbf{p}, \sigma)$$

is Lipschitz, with Lipschitz constant in  $\ell^2$ -norm denoted by  $L_\rho$ ; the proof of this fact is detailed below.

Now, the regret is non positive as soon as  $\sum_{t=1}^T \underline{r}(I_t, J_t)/T$  belongs to  $\mathcal{C}$ ; we therefore only need to consider the case when this average is not in  $\mathcal{C}$ . In the latter case, we denote by  $(\tilde{r}_T, \tilde{\sigma}_T)$  its projection in  $\ell^2$ -norm onto  $\mathcal{C}$ . We have first that the defining inequality of  $\mathcal{C}$  is an equality on its border, so that

$$\tilde{r}_T = \max_{\mathbf{p} \in \Delta(\mathcal{I})} \rho(\mathbf{p}, \tilde{\sigma}_T);$$

and second, that

$$\begin{aligned} R_T^{\text{ext}} &= \max_{\mathbf{p} \in \Delta(\mathcal{I})} \rho(\mathbf{p}, \tilde{H}(\hat{\mathbf{q}}_T)) - \frac{1}{T} \sum_{t=1}^T r(I_t, J_t) \\ &\leq \left| \max_{\mathbf{p} \in \Delta(\mathcal{I})} \rho(\mathbf{p}, \tilde{H}(\hat{\mathbf{q}}_T)) - \max_{\mathbf{p} \in \Delta(\mathcal{I})} \rho(\mathbf{p}, \tilde{\sigma}_T) \right| + \left| \tilde{r}_T - \frac{1}{T} \sum_{t=1}^T r(I_t, J_t) \right| \\ &\leq L_\rho \left\| \tilde{\sigma}_T - \tilde{H}(\hat{\mathbf{q}}_T) \right\|_2 + \left| \tilde{r}_T - \frac{1}{T} \sum_{t=1}^T r(I_t, J_t) \right| \\ &\leq \sqrt{2} \max\{L_\rho, 1\} \left\| \begin{bmatrix} \tilde{r}_T \\ \tilde{\sigma}_T \end{bmatrix} - \frac{1}{T} \sum_{t=1}^T \underline{r}(I_t, J_t) \right\|_2 \\ &= \sqrt{2} \max\{L_\rho, 1\} \inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T \underline{r}(I_t, J_t) \right\|_2. \end{aligned}$$

The claimed rates are now seen to follow from the ones indicated in Theorem 6.3.

It only remains to prove the indicated Lipschitzness. (All Lipschitzness statements that follow will be with respect to the  $\ell^2$ -norms.) We have by Definition 6.1 that for all  $\mathbf{p} \in \Delta(\mathcal{I})$  and  $\sigma \in \mathcal{F}$ ,

$$\rho(\mathbf{p}, \sigma) = \min \bar{m}_1(\mathbf{p}, \Phi(\sigma)),$$

where the linear  $\bar{m}_1$  is indifferently either relative to  $\bar{m}_1$  or is the projection onto the first component of the function  $\bar{m}$  relative to  $\bar{m}$ . By Remark 6.1 the mapping  $\sigma \in \mathcal{F} \mapsto \Phi(\sigma)$  is  $\kappa_\Phi$ -Lipschitz; this entails, by Lemma 3.1, that for all  $\mathbf{p} \in \Delta(\mathcal{I})$ , the mapping  $\sigma \in \mathcal{F} \mapsto \rho(\mathbf{p}, \sigma)$  is  $R\sqrt{N_B} \kappa_\Phi$ -Lipschitz. In particular, since the latter Lipschitz constant is independent of  $\mathbf{p}$ , the mapping

$$\sigma \in \mathcal{F} \mapsto \max_{\mathbf{p} \in \Delta(\mathcal{I})} \rho(\mathbf{p}, \sigma)$$

is  $R\sqrt{N_B} \kappa_\Phi$ -Lipschitz as well, which concludes the proof.  $\square$

A similar argument to the one in Perchet [24] shows that the convex set  $\mathcal{C}$  is defined by a finite number of piecewise linear equations, it is therefore a polyhedron so the projection onto it, and as well the computation of the strategy, can be done efficiently. We sketch the argument below, and refer the reader to Perchet [24] for details. Equation (15) indicates a priori that for each  $\sigma \in \mathcal{F}$ , there exist a finite number  $M_\sigma$  (depending on  $\sigma$ ) of mixed actions  $\mathbf{q}_{\sigma,1}, \dots, \mathbf{q}_{\sigma,M_\sigma}$  such that for all  $\mathbf{p} \in \Delta(\mathcal{I})$ , we have  $\rho(\mathbf{p}, \sigma) = \min\{r(\mathbf{p}, \mathbf{q}_{\sigma,1}), \dots, r(\mathbf{p}, \mathbf{q}_{\sigma,M_\sigma})\}$ . But by an argument stated in Perchet [24],

$$\sigma \mapsto \left\{ \mathbf{q} \in \Delta(\mathcal{J}) \text{ such that } \tilde{H}(\mathbf{q}) = \sigma \right\}$$

evolves in a piecewise linear way and thus there exist a finite number  $M$  of piecewise linear functions  $\sigma \mapsto \mathbf{q}'_{\sigma,k}$ , with  $k = 1, \dots, M$ , such that, for all  $\sigma \in \mathcal{F}$ ,

$$\left\{ \mathbf{q}_{\sigma,1}, \dots, \mathbf{q}_{\sigma,M_\sigma} \right\} = \left\{ \mathbf{q}'_{\sigma,1}, \dots, \mathbf{q}'_{\sigma,M} \right\}.$$

(There can be some redundancies between the  $\mathbf{q}'_{\sigma,k}$ .) Because of this, we have that for all  $\mathbf{p} \in \Delta(\mathcal{I})$  and  $\sigma \in \mathcal{F}$ ,

$$\rho(\mathbf{p}, \sigma) = \min\{r(\mathbf{p}, \mathbf{q}'_{\sigma,1}), \dots, r(\mathbf{p}, \mathbf{q}'_{\sigma,M})\}.$$

Each function  $\sigma \mapsto \mathbf{q}'_{\sigma,k}$  being piecewise linear, one can construct a finite set  $\{\mathbf{p}_1, \dots, \mathbf{p}_K\} \subset \Delta(\mathcal{I})$  such that, for any  $\sigma \in \mathcal{F}$ , the mapping  $\mathbf{p} \mapsto \rho(\mathbf{p}, \sigma)$  is maximized at one of these  $\mathbf{p}_k$ . The convex set  $\mathcal{C}$  is therefore defined by a finite number of piecewise linear equations, it is therefore a polyhedron; therefore the projection onto it, hence the computation of the proposed strategy, can be done efficiently.

**6.2.2 Internal / swap regret.** Foster and Vohra [10] defined internal regret with full monitoring as follows. A player has no internal regret if, for every action  $i \in \mathcal{I}$ , he has no external regret on the stages when this specific action  $i$  was played. In other words,  $i$  is the best response to the empirical distribution of action of the other player on these stages.

With partial monitoring, the first player evaluates his payoffs in a pessimistic way through the function  $\rho$  defined above. This function is not linear over  $\Delta(\mathcal{I})$  in general (it is concave), so that the best responses are not necessarily pure actions  $i \in \mathcal{I}$  but mixed actions, i.e., elements of  $\Delta(\mathcal{I})$ . Following Lehrer and Solan [15] one therefore can partition the stages not depending on the pure actions actually played but on the mixed actions  $\mathbf{p}_t \in \Delta(\mathcal{I})$  used to draw them. To this end, it is convenient to assume that the strategies of the first player need to pick these mixed actions in a finite (but possibly thin) grid of  $\Delta(\mathcal{I})$ , which we denote by  $\{\mathbf{p}_g, g \in \mathcal{G}\}$ , where  $\mathcal{G}$  is a finite set. At each round  $t$ , the first player picks an index  $G_t \in \mathcal{G}$  and uses the distribution  $\mathbf{p}_{G_t}$  to draw his action  $I_t$ . Up to a standard concentration-of-the-measure argument, we will measure the payoff at round  $t$  with  $r(\mathbf{p}_{G_t}, J_t)$  rather than with  $r(I_t, J_t)$ .

For each  $g \in \mathcal{G}$ , we denote by  $N_T(g)$  the number of stages in  $\{1, \dots, T\}$  for which we had  $G_t = g$  and, whenever  $N_T(g) > 0$ ,

$$\hat{\mathbf{q}}_{T,g} = \frac{1}{N_T(g)} \sum_{t:G_t=g} \delta_{J_t}.$$

We define  $\hat{\mathbf{q}}_{T,g}$  is an arbitrary way when  $N_T(g) = 0$ . The internal regret of the first player at round  $T$  is measured as

$$R_T^{\text{int}} = \max_{g,g' \in \mathcal{G}} \frac{N_T(g)}{T} \left( \rho(\mathbf{p}_{g'}, \tilde{H}(\hat{\mathbf{q}}_{T,g})) - r(\mathbf{p}_g, \hat{\mathbf{q}}_{T,g}) \right).$$

Actually, our proof technique rather leads to the minimization of some swap regret (see Blum and Mansour [4] for the definition of swap regret in full monitoring):

$$R_T^{\text{swap}} = \sum_{g \in \mathcal{G}} \frac{N_T(g)}{T} \left( \max_{g' \in \mathcal{G}} \rho(\mathbf{p}_{g'}, \tilde{H}(\hat{\mathbf{q}}_{T,g})) - r(\mathbf{p}_g, \hat{\mathbf{q}}_{T,g}) \right)_+.$$

Again, the following bound on the swap regret easily follows from Theorem 6.1; the latter constructs a simple and direct strategy to control the swap regret, thus also the internal regret. It therefore improves the results of Lehrer and Solan [15] and Perchet [22], two articles that presented more involved and less efficient strategies to do so (strategies based on auxiliary strategies using grids that need to be refined over time and whose complexities is exponential in the size of these grids; ideas all in all similar to what is done in calibration, see the references provided in Section 4). Moreover, we provide convergence rates.

**COROLLARY 6.2** *The first player has an explicit strategy such that for all  $T$  and all strategies of the second player, with probability at least  $1 - \delta$ ,*

$$R_T^{\text{swap}} \leq \square \left( T^{-1/5} \sqrt{\ln \frac{T}{\delta}} + T^{-2/5} \ln \frac{T}{\delta} \right)$$

for some constant  $\square$  depending only on the game  $(r, H)$  at hand and on the size of the finite grid  $\mathcal{G}$ .

**PROOF.** The proof of this corollary is based on ideas similar to the ones used in the proof of Corollary 6.1;  $\mathcal{G}$  will play the role of the action set of the first player. The proof proceeds in four steps. In the first step, we construct an approachability setup and show that Condition (APM) applies. In the second step, we show that Assumption 6.1 is satisfied. In the third step we analyze the convergence rates of the swap regret. In the fourth and final step, we show that the set we are approaching possess some smoothness properties by providing a uniform Lipschitz bound on certain functions.

**Step 1:** We denote by

$$\mathcal{F}_{\text{cone}} = \{ \lambda \sigma, \sigma \in \mathcal{F}, \lambda \in \mathbb{R}_+ \}$$

the cone generated by  $\mathcal{F}$  and extend linearly  $\rho : \Delta(\mathcal{I}) \times \mathcal{F} \rightarrow \mathbb{R}$  into a mapping  $\rho : \Delta(\mathcal{I}) \times \mathcal{F}_{\text{cone}} \rightarrow \mathbb{R}$  as follows: for all  $\mathbf{p} \in \Delta(\mathcal{I})$ , for all  $\lambda \geq 0$  with  $\lambda \neq 1$ , and all  $\sigma \in \mathcal{F}$ ,

$$\rho(\mathbf{p}, \lambda \sigma) = \begin{cases} 0 & \text{if } \lambda = 0, \\ \lambda \rho(\mathbf{p}, \sigma) & \text{if } \lambda > 0. \end{cases}$$

In the sequel, we embed  $\mathcal{F}_{\text{cone}}$  into  $\mathbb{R}^{\mathcal{I} \times \mathcal{H}}$ .

The closed convex set  $\mathcal{C}$  and the vector-valued payoff function  $\underline{r}$  are then respectively defined by

$$\mathcal{C} = \left\{ (z_g, \mathbf{v}_g)_{g \in \mathcal{G}} \in (\mathbb{R} \times \mathcal{F}_{\text{cone}})^{\mathcal{G}} : \forall g \in \mathcal{G}, z_g \geq \max_{g' \in \mathcal{G}} \rho(\mathbf{p}_{g'}, \mathbf{v}_g) \right\}$$

and, for all  $(g, j) \in \mathcal{G} \times \mathcal{J}$ ,

$$\underline{r}(g, j) = \begin{bmatrix} r(\mathbf{p}_g, j) \mathbb{I}_{\{g'=g\}} \\ \tilde{H}(\delta_j) \mathbb{I}_{\{g'=g\}} \end{bmatrix}_{g' \in \mathcal{G}}.$$

To show that  $\mathcal{C}$  is  $\underline{r}$ -approachable, we associate with each  $\mathbf{q} \in \Delta(\mathcal{J})$  an element  $g^*(\mathbf{q}) \in \mathcal{G}$  such that

$$g^*(\mathbf{q}) \in \operatorname{argmax}_{g \in \mathcal{G}} \rho(\mathbf{p}_g, \tilde{H}(\mathbf{q})).$$

Then, given any  $\mathbf{q} \in \Delta(\mathcal{J})$ , we note that for all  $\mathbf{q}'$  satisfying  $\tilde{H}(\mathbf{q}') = \tilde{H}(\mathbf{q})$ , the components of the vector  $\underline{r}(g^*(\mathbf{q}), \mathbf{q}')$  are all null but the ones corresponding to  $g^*(\mathbf{q})$ , for which we have

$$r(\mathbf{p}_{g^*(\mathbf{q})}, \mathbf{q}') \geq \rho(\mathbf{p}_{g^*(\mathbf{q})}, \tilde{H}(\mathbf{q}')) = \rho(\mathbf{p}_{g^*(\mathbf{q})}, \tilde{H}(\mathbf{q})) = \max_{g' \in \mathcal{G}} \rho(\mathbf{p}_{g'}, \tilde{H}(\mathbf{q})) = \max_{g' \in \mathcal{G}} \rho(\mathbf{p}_{g'}, \tilde{H}(\mathbf{q}')),$$

where the first inequality is by definition of  $\rho$ . Therefore,  $\underline{r}(g^*(\mathbf{q}), \mathbf{q}') \in \mathcal{C}$ . Condition (APM) in Lemma 6.2 and Theorem 6.1 is thus satisfied, so that we have approachability.

**Step 2:** We then show that Assumption 6.1 is satisfied. It suffices to show that for all  $\sigma \in \mathcal{F}$ , the application

$$\pi = (\pi_g)_{g \in \mathcal{G}} \in \Delta(\mathcal{G}) \longmapsto \bar{m}_1(\pi, \sigma) = \left\{ (\pi_g r(\mathbf{p}_g, \mathbf{q}))_{g \in \mathcal{G}} : \mathbf{q} \in \Delta(\mathcal{J}) \text{ such that } \tilde{H}(\mathbf{q}) = \sigma \right\}$$

is piecewise linear (as the other components in the definition of  $\bar{m}$  are linear in  $\pi$ ). This is the case since for each  $g$ , the application

$$\pi \in \Delta(\mathcal{G}) \longmapsto \left\{ \pi_g r(\mathbf{p}_g, \mathbf{q}) : \mathbf{q} \in \Delta(\mathcal{J}) \text{ such that } \tilde{H}(\mathbf{q}) = \sigma \right\}$$

is seen to be piecewise linear, by using the same one-dimensional argument as in the proof of Corollary 6.1.

**Step 3:** We now exhibit the convergence rates. In view of the form of the defining set of constraints for  $\mathcal{C}$ , the coordinates of the elements in  $\mathcal{C}$  can be grouped according to each  $g \in \mathcal{G}$  and projections onto  $\mathcal{C}$  can therefore be done separately for each such group. The group  $g$  of coordinates of  $\sum_{t=1}^T \mathbf{r}(G_t, J_t)/T$  is formed by

$$\frac{N_T(g)}{T} r(\mathbf{p}_g, \hat{\mathbf{q}}_{T,g}) \quad \text{and} \quad \frac{N_T(g)}{T} \tilde{H}(\hat{\mathbf{q}}_{T,g});$$

when

$$\frac{N_T(g)}{T} r(\mathbf{p}_g, \hat{\mathbf{q}}_{T,g}) \geq \max_{g' \in \mathcal{G}} \rho\left(\mathbf{p}_{g'}, \frac{N_T(g)}{T} \tilde{H}(\hat{\mathbf{q}}_{T,g})\right),$$

we denote these quantities by  $\tilde{r}_{T,g}$  and  $\tilde{\mathbf{v}}_{T,g}$ . Otherwise, we project this pair on the set

$$\mathcal{C}_g = \left\{ (z_g, \mathbf{v}_g) \in \mathbb{R} \times \mathcal{F}_{\text{cone}} : z_g \geq \max_{g' \in \mathcal{G}} \rho(\mathbf{p}_{g'}, \mathbf{v}_g) \right\}$$

and denote by  $\tilde{r}_{T,g}$  and  $\tilde{\mathbf{v}}_{T,g}$  the coordinates of the projection; they satisfy the defining inequality of  $\mathcal{C}_g$  with equality,

$$\tilde{r}_{T,g} = \max_{g' \in \mathcal{G}} \rho(\mathbf{p}_{g'}, \tilde{\mathbf{v}}_{T,g}).$$

By distinguishing for each  $g$  according to which of the two cases above arose (for the first inequality), we may decompose and upper bound the swap regret as follows,

$$\begin{aligned} R_T^{\text{swap}} &= \sum_{g \in \mathcal{G}} \frac{N_T(g)}{T} \left( \max_{g' \in \mathcal{G}} \rho(\mathbf{p}_{g'}, \tilde{H}(\hat{\mathbf{q}}_{T,g})) - r(\mathbf{p}_g, \hat{\mathbf{q}}_{T,g}) \right)_+ \\ &= \sum_{g \in \mathcal{G}} \left( \max_{g' \in \mathcal{G}} \rho\left(\mathbf{p}_{g'}, \frac{N_T(g)}{T} \tilde{H}(\hat{\mathbf{q}}_{T,g})\right) - \frac{N_T(g)}{T} r(\mathbf{p}_g, \hat{\mathbf{q}}_{T,g}) \right)_+ \\ &\leq \sum_{g \in \mathcal{G}} \left| \max_{g' \in \mathcal{G}} \rho\left(\mathbf{p}_{g'}, \frac{N_T(g)}{T} \tilde{H}(\hat{\mathbf{q}}_{T,g})\right) - \max_{g' \in \mathcal{G}} \rho(\mathbf{p}_{g'}, \tilde{\mathbf{v}}_{T,g}) \right| + \sum_{g \in \mathcal{G}} \left| \tilde{r}_{T,g} - \frac{N_T(g)}{T} r(\mathbf{p}_g, \hat{\mathbf{q}}_{T,g}) \right| \\ &\leq \sum_{g \in \mathcal{G}} \bar{L}_\rho \left\| \frac{N_T(g)}{T} \tilde{H}(\hat{\mathbf{q}}_{T,g}) - \tilde{\mathbf{v}}_{T,g} \right\|_2 + \sum_{g \in \mathcal{G}} \left| \tilde{r}_{T,g} - \frac{N_T(g)}{T} r(\mathbf{p}_g, \hat{\mathbf{q}}_{T,g}) \right|, \end{aligned}$$

where we used a fact proved below, that the application

$$\mathbf{v} \in \mathcal{F}_{\text{cone}} \longmapsto \max_{g' \in \mathcal{G}} \rho(\mathbf{p}_{g'}, \mathbf{v})$$

is  $\bar{L}_\rho$ -Lipschitz. In the last inequality we had a sum of  $\ell^2$ -norms, which can be bounded by a single  $\ell^2$ -norm,

$$\begin{aligned} R_T^{\text{swap}} &\leq \max\{\bar{L}_\rho, 1\} \sqrt{2N_{\mathcal{G}}} \left\| \begin{bmatrix} \tilde{r}_{T,g} \\ \tilde{\mathbf{v}}_{T,g} \end{bmatrix}_{g \in \mathcal{G}} - \frac{1}{T} \sum_{t=1}^T \mathbf{r}(I_t, J_t) \right\|_2 \\ &\leq \max\{\bar{L}_\rho, 1\} \sqrt{2N_{\mathcal{G}}} \inf_{c \in \mathcal{C}} \left\| c - \frac{1}{T} \sum_{t=1}^T \mathbf{r}(I_t, J_t) \right\|_2, \end{aligned}$$

where we denoted by  $N_{\mathcal{G}}$  the cardinality of  $\mathcal{G}$ . Resorting to the convergence rate stated in Theorem 6.3 concludes the proof, up to the claimed Lipschitzness, which we now prove. (All Lipschitzness statements that follow will be with respect to the  $\ell^2$ -norms.)

**Step 4:** To do so, it suffices to show that for all fixed elements  $\mathbf{p} \in \Delta(\mathcal{I})$ , the functions  $\mathbf{v} \in \mathcal{F}_{\text{cone}} \mapsto \rho(\mathbf{p}, \mathbf{v})$  are Lipschitz, with a Lipschitz constant  $\bar{L}_\rho$  that is independent of  $\mathbf{p}$ . Note that we already proved at the end of the proof of Corollary 6.1 that  $\sigma \in \mathcal{F} \mapsto \rho(\mathbf{p}, \sigma)$  is Lipschitz, with a Lipschitz constant  $L_\rho$  independent of  $\mathbf{p}$ . Consider now two elements  $\mathbf{v}, \mathbf{v}' \in \mathcal{F}_{\text{cone}}$ , which we write as  $\mathbf{v} = \lambda\sigma$  and  $\mathbf{v}' = \lambda'\sigma'$ , with  $\sigma, \sigma' \in \mathcal{F}$  and  $\lambda, \lambda' \in \mathbb{R}_+$ . Using triangle inequalities, the Lipschitzness of  $\rho$  on  $\mathcal{F}$ , and the fact that  $r$  thus  $\rho$  are bounded by  $R$ ,

$$\begin{aligned} |\rho(\mathbf{p}, \lambda\sigma) - \rho(\mathbf{p}, \lambda'\sigma')| &\leq |\lambda(\rho(\mathbf{p}, \sigma) - \rho(\mathbf{p}, \sigma'))| + |(\lambda - \lambda')\rho(\mathbf{p}, \sigma')| \\ &\leq \lambda L_\rho \|\sigma - \sigma'\|_2 + R|\lambda - \lambda'| \\ &\leq L_\rho \|\lambda\sigma - \lambda'\sigma' + (\lambda' - \lambda)\sigma'\|_2 + R|\lambda - \lambda'| \\ &\leq L_\rho \|\lambda\sigma - \lambda'\sigma'\|_2 + (R + L_\rho N_{\mathcal{I}}) |\lambda - \lambda'|, \end{aligned}$$

where we used also for the last inequality that since  $\sigma$  is a vector of  $N_{\mathcal{I}}$  probability distributions over the signals,  $\|\sigma\|_2 \leq \|\sigma\|_1 = N_{\mathcal{I}}$ . To conclude the argument, we simply need to show that  $|\lambda - \lambda'|$  can be bounded by  $\|\lambda\sigma - \lambda'\sigma'\|_2$  up to some universal constant, which we do now. We resort again to the fact that  $\|\sigma\|_1 = \|\sigma'\|_1 = N_{\mathcal{I}}$  and can thus write, thanks to a triangle inequality and assuming with no loss of generality that  $\lambda' < \lambda$ , that

$$|\lambda - \lambda'| = \frac{1}{N_{\mathcal{I}}} (\lambda \|\sigma\|_1 - \lambda' \|\sigma'\|_1) \leq \frac{1}{N_{\mathcal{I}}} \|\lambda\sigma - \lambda'\sigma'\|_1 \leq \frac{\sqrt{N_{\mathcal{H}} N_{\mathcal{I}}}}{N_{\mathcal{I}}} \|\lambda\sigma - \lambda'\sigma'\|_2,$$

where we used the Cauchy-Schwarz inequality for the final step. One can thus take, for instance,

$$\bar{L}_\rho = L_\rho + (R + L_\rho N_{\mathcal{I}}) \sqrt{\frac{N_{\mathcal{H}}}{N_{\mathcal{I}}}}.$$

This concludes the proof.  $\square$

**7. Approachability in the case of general games with partial monitoring.** Unfortunately, as is illustrated in the following example, there exist games with partial monitoring that are not bi-piecewise linear.

**EXAMPLE 7.1** *The following game (with the same action and signal sets as in Example 5.2) is not bi-piecewise linear.*

|     | $L$                        | $M$                        | $R$                         |
|-----|----------------------------|----------------------------|-----------------------------|
| $T$ | $(1, 0, 0, 0) / \clubsuit$ | $(0, 0, 1, 0) / \clubsuit$ | $(2, 0, 4, 0) / \heartsuit$ |
| $B$ | $(0, 1, 0, 0) / \clubsuit$ | $(0, 0, 0, 1) / \clubsuit$ | $(0, 3, 0, 5) / \heartsuit$ |

**PROOF.** We denote mixed actions of the first player by  $(p, 1 - p)$ , where  $p \in [0, 1]$  denotes the probability of playing  $T$  and  $1 - p$  is the probability of playing  $B$ . It is immediate that  $\bar{m}((p, 1 - p), \clubsuit)$  can be identified with the set of all product distributions on  $2 \times 2$  elements with first marginal distribution  $(p, 1 - p)$ . The proof of Lemma 6.1 shows that the set  $\mathcal{B}$  associated with any game always contains the Dirac masses on each signal; that is,  $\delta_{\clubsuit} \in \mathcal{B}$ . But for  $p \neq p'$  and  $\lambda \in (0, 1)$ , denoting  $\bar{p} = \lambda p + (1 - \lambda)p'$ , one necessarily has that

$$\bar{m}((\bar{p}, 1 - \bar{p}), \clubsuit) \subsetneq \lambda \bar{m}((p, 1 - p), \clubsuit) + (1 - \lambda) \bar{m}((p', 1 - p'), \clubsuit);$$

the inclusion  $\subseteq$  holds by concavity of  $\bar{m}$  in its first argument (Lemma 5.1) but this inclusion is always strict here since the left-hand side is formed by product distributions while the right-hand side also contains distributions with correlations. Hence, bi-piecewise linearity cannot hold for this game.  $\square$

However, we will show that if Condition (APM) holds there exist strategies with a constant per-round complexity to approach polytopes even when the game is not bi-piecewise linear. That is, by considering simpler closed convex sets  $\mathcal{C}$ , no assumption is needed on the pair  $(r, H)$ .

We will conclude this section by indicating that thanks to a doubling trick, Condition (APM) is still sufficient for approachability in the most general case when no assumption is made neither on  $(r, H)$  nor on  $\mathcal{C}$ —at the cost, however, of inefficiency.

**7.1 Approachability of the negative orthant in the case of general games.** For the sake of simplicity, we start with the case of the negative orthant  $\mathbb{R}_-^d$ . Our argument will be based on Lemma 6.1; we use in the sequel the objects and notation introduced therein. We denote by  $r = (r_k)_{1 \leq k \leq d}$  the components of the  $d$ -dimensional payoff function  $r$  and introduce, for all  $k \in \{1, \dots, d\}$ , the set-valued mapping  $\tilde{m}_k$  defined by

$$\tilde{m}_k : (\mathbf{p}, b) \in \Delta(\mathcal{I}) \times \mathcal{B} \mapsto \tilde{m}_k(\mathbf{p}, b) = \left\{ r_k(\mathbf{p}, \mathbf{q}) : \mathbf{q} \in \Delta(\mathcal{J}) \text{ such that } \tilde{H}(\mathbf{q}) = b \right\}.$$

The mapping  $\tilde{m}$  is then defined as the Cartesian product of the  $\tilde{m}_k$ ; formally, for all  $\mathbf{p} \in \Delta(\mathcal{I})$  and  $b \in \mathcal{B}$ ,

$$\tilde{m}(\mathbf{p}, b) = \left\{ (z_1, \dots, z_d) : \forall k \in \{1, \dots, d\}, z_k \in \tilde{m}_k(\mathbf{p}, b) \right\}.$$

We then linearly extend this mapping into a set-valued mapping  $\tilde{m}$  defined on  $\Delta(\mathcal{I}) \times \Delta(\mathcal{B})$  and finally consider the set-valued mapping  $\check{m}$  defined on  $\Delta(\mathcal{I}) \times \mathcal{F}$  by

$$\forall \sigma \in \mathcal{F}, \quad \forall \mathbf{p} \in \Delta(\mathcal{I}), \quad \check{m}(\mathbf{p}, \sigma) = \tilde{m}(\mathbf{p}, \Phi(\sigma)) = \sum_{b \in \mathcal{B}} \Phi_b(\sigma) \tilde{m}(\mathbf{p}, b),$$

where  $\Phi$  refers to the mapping defined in Lemma 6.1 (based on  $\overline{m}$ ). The lemma below indicates why  $\check{m}$  is an excellent substitute to  $\overline{m}$  in the case of the approachability of the orthant  $\mathbb{R}_-^d$ .

LEMMA 7.1 *The set-valued mappings  $\check{m}$  and  $\overline{m}$  satisfy that for all  $p \in \Delta(\mathcal{I})$  and  $\sigma \in \mathcal{F}$ ,*

- (i) *the inclusion  $\overline{m}(\mathbf{p}, \sigma) \subseteq \check{m}(\mathbf{p}, \sigma)$  holds;*
- (ii) *if  $\overline{m}(\mathbf{p}, \sigma) \subseteq \mathbb{R}_-^d$ , then one also has  $\check{m}(\mathbf{p}, \sigma) \subseteq \mathbb{R}_-^d$ .*

The interpretations of these two properties are that: 1.  $\check{m}$ -robust approaching a set  $\mathcal{C}$  is more difficult than  $\overline{m}$ -robust approaching it; and 2. that if Condition (APM) holds for  $\overline{m}$  and  $\mathbb{R}_-^d$ , it also holds for  $\check{m}$  and  $\mathbb{R}_-^d$ .

PROOF. For property 1., note that by the component-wise construction of  $\tilde{m}$ ,

$$\forall b \in \mathcal{B}, \quad \forall \mathbf{p} \in \Delta(\mathcal{I}), \quad \overline{m}(\mathbf{p}, b) \subseteq \tilde{m}(\mathbf{p}, b);$$

Lemma 6.1, the linear extension of  $\tilde{m}$ , and the definition of  $\check{m}$  then show that

$$\forall \sigma \in \mathcal{F}, \quad \forall \mathbf{p} \in \Delta(\mathcal{I}), \quad \overline{m}(\mathbf{p}, \sigma) = \sum_{b \in \mathcal{B}} \Phi_b(\sigma) \overline{m}(\mathbf{p}, b) \subseteq \tilde{m}(\mathbf{p}, \Phi(\sigma)) = \check{m}(\mathbf{p}, \sigma).$$

As for property 2., it suffices to work component-wise. Note that (by Lemma 6.1 again) the stated assumption exactly means that  $\sum_{b \in \mathcal{B}} \Phi_b(\sigma) \overline{m}(\mathbf{p}, b) \subseteq \mathbb{R}_-^d$ . In particular, rewriting the non-positivity constraint for each of the  $d$  components of the payoff vectors, we get

$$\sum_{b \in \mathcal{B}} \Phi_b(\sigma) \tilde{m}_k(\mathbf{p}, b) \subseteq \mathbb{R}_-,$$

for all  $k \in \{1, \dots, d\}$ ; thus, in particular,  $\sum_{b \in \mathcal{B}} \Phi_b(\sigma) \tilde{m}(\mathbf{p}, b) = \check{m}(\mathbf{p}, \sigma) \subseteq \mathbb{R}_-^d$ .  $\square$

We can then extend the result of the previous section without the bi-pieceswise linearity assumption.

THEOREM 7.1 *If Condition (APM) is satisfied for  $\overline{m}$  and  $\mathbb{R}_-^d$ , then there exists a strategy for  $(r, H)$ -approaching  $\mathbb{R}_-^d$  at a rate of the order of  $T^{-1/5}$ , with a constant per-round complexity.*

PROOF. The assumption of the theorem and Property 2. of Lemma 7.1 imply that Condition (APM) holds for  $\mathbb{R}_-^d$  and  $\check{m}$ ; furthermore, the latter corresponds to a bi-pieceswise linear game as can be seen by noting, similarly to what was done in the section devoted to regret minimization (Section 6.2), that each  $\tilde{m}_k$ , being based on the scalar payoff function  $r_k$ , is a pieceswise linear function. Thus,  $\check{m}$  is also a pieceswise linear function.

Therefore, the steps between Equations (10)–(12) of the proof of Theorem 6.1 (or the corresponding statements in the proof of Theorem 6.3) can be adapted by replacing  $\overline{m}$  and  $\overline{\overline{m}}$  by, respectively,  $\tilde{m}$ ,  $\check{m}$ , and its extension corresponding to Definition 6.1. The result follows.  $\square$

**7.2 Approachability of polytopes in the case of general games.** If that the target set  $\mathcal{C}$  is a polytope, then  $\mathcal{C}$  can be written as the intersection of a finite number of half-planes, i.e., there exists a finite family  $\{(e_k, f_k) \in \mathbb{R}^d \times \mathbb{R}, k \in \mathcal{K}\}$  such that

$$\mathcal{C} = \{z \in \mathbb{R}^d : \langle z, e_k \rangle \leq f_k, \forall k \in \mathcal{K}\}.$$

Given the original (not necessarily bi-piecewise linear) game  $(r, H)$ , we introduce another game  $(r_{\mathcal{C}}, H)$ , whose payoff function  $r_{\mathcal{C}} : \mathcal{I} \times \mathcal{J} \rightarrow \mathbb{R}^{\mathcal{K}}$  is defined as

$$\forall i \in \mathcal{I}, \quad \forall j \in \mathcal{J}, \quad r_{\mathcal{C}}(i, j) = \left[ \langle r(i, j), e_k \rangle - f_k \right]_{k \in \mathcal{K}}.$$

The following lemma follows by rewriting the above.

**LEMMA 7.2** *Given a polytope  $\mathcal{C}$ , the  $(r, H)$ -approachability of  $\mathcal{C}$  and the  $(r_{\mathcal{C}}, H)$ -approachability of  $\mathbb{R}_-^d$  are equivalent in the sense that every strategy for one problem translates to a strategy for the other problem. In addition, Condition (APM) holds for  $(r, H)$  and  $\mathcal{C}$  if and only if it holds for  $(r_{\mathcal{C}}, H)$  and  $\mathbb{R}_-^d$ .*

Via the lemma above, Theorem 7.1 indicates that Condition (APM) for  $(r, H)$  and  $\mathcal{C}$  is a sufficient condition for the  $(r, H)$ -approachability of  $\mathcal{C}$  and provides a strategy to do so. (The per-round complexity of this strategy depends in particular at least linearly on the cardinality of  $\mathcal{K}$ .)

**7.3 Approachability of general convex sets in the case of general games.** A general closed convex set can always be approximated arbitrarily well by a polytope (where the number of vertices of the latter however increases as the quality of the approximation does). Therefore, via playing in regimes, Condition (APM) is also seen to be sufficient to  $(r, H)$ -approach any general closed convex set  $\mathcal{C}$ . However, the computational complexity of the resulting strategy is much larger: the per-round complexity increases over time (as the numbers of vertices of the approximating polytopes do).

### Appendix A. An auxiliary result of calibration.

We prove here (4) for a given  $\eta > 0$  and do so by following the methodology of Mannor and Stoltz [19]. (Note that this result is of independent interest.)

We actually assume that the covering  $\mathbf{y}^1, \dots, \mathbf{y}^{N_\eta}$  is slightly finer than what was required around (4) and that it forms an  $\eta/N_{\mathcal{B}}$ -grid of  $\Delta(\mathcal{B})$ , i.e., that for all  $\mathbf{y} \in \Delta(\mathcal{B})$ , there exists  $\ell \in \{1, \dots, N_\eta\}$  such that  $\|\mathbf{y} - \mathbf{y}^\ell\|_1 \leq \eta/N_{\mathcal{B}}$ .

We recall that elements  $\mathbf{y} \in \mathcal{B}$  are denoted by  $\mathbf{y} = (y_b)_{b \in \mathcal{B}}$  and we identify  $\Delta(\mathcal{B})$  with a subset of  $\mathbb{R}^{N_{\mathcal{B}}}$ . In particular,  $\mathbb{I}_b$ , the Dirac mass on a given  $b \in \mathcal{B}$ , is a binary vector whose only non-null component is the one indexed by  $b$ . Finally, we denote by

$$\underline{\mathbf{0}} = (0, \dots, 0) \quad \text{and} \quad \underline{\mathbf{1}} = (1, \dots, 1)$$

the elements of  $\mathbb{R}^{\mathcal{B}}$  respectively formed by zeros and ones only.

We consider a vector-valued payoff function  $C : \{1, \dots, N_\eta\} \times \mathcal{B} \rightarrow \mathbb{R}^{2N_\eta N_{\mathcal{B}}}$  defined as follows; for all  $\ell \in \{1, \dots, N_\eta\}$  and for all  $b \in \mathcal{B}$ ,

$$C(\ell, b) = \left( \underline{\mathbf{0}}, \dots, \underline{\mathbf{0}}, \mathbf{y}^\ell - \mathbb{I}_b - \frac{\eta}{N_{\mathcal{B}}} \underline{\mathbf{1}}, \mathbb{I}_b - \mathbf{y}^\ell - \frac{\eta}{N_{\mathcal{B}}} \underline{\mathbf{1}}, \underline{\mathbf{0}}, \dots, \underline{\mathbf{0}} \right),$$

which is a vector of  $2N_\eta$  elements of  $\mathbb{R}^{\mathcal{B}}$  composed by  $2(N_\eta - 1)$  occurrences of the zero element  $\underline{\mathbf{0}} \in \mathbb{R}^{\mathcal{B}}$  and two non-zero elements, located in the positions indexed by  $2\ell - 1$  and  $2\ell$ .

We now show that the closed convex set  $(\mathbb{R}_-)^{2N_\eta N_{\mathcal{B}}}$  is  $C$ -approachable; to do so, we resort to the characterization stated in Theorem 2.1. To each  $\mathbf{y} \in \Delta(\mathcal{B})$  we will associate a pure action  $\ell_{\mathbf{y}}$  in  $\{1, \dots, N_\eta\}$  so that  $C(\ell_{\mathbf{y}}, \mathbf{y}) \in (\mathbb{R}_-)^{2N_\eta N_{\mathcal{B}}}$ ; note that to satisfy the necessary and sufficient condition, it is not necessary here to resort to mixed actions of the first player. The index  $\ell_{\mathbf{y}}$  is any index  $\ell$  such that  $\|\mathbf{y} - \mathbf{y}^\ell\|_1 \leq \eta/N_{\mathcal{B}}$ ; such an index always exists as noted at the beginning of this proof. Indeed, one then has in particular that for each component  $b \in \mathcal{B}$ ,

$$|y_b^{\ell_{\mathbf{y}}} - y_b| \leq \|\mathbf{y}^{\ell_{\mathbf{y}}} - \mathbf{y}\|_1 \leq \eta/N_{\mathcal{B}}.$$

A straightforward adaptation of the proof of Theorem 2.2 then yields a strategy such that for all  $\delta \in (0, 1)$  and for all strategies of the second player, with probability at least  $1 - \delta$ ,

$$\sup_{\tau \geq T} \inf_{c \in (\mathbb{R}_-)^{2N_\eta N_B}} \left\| c - \frac{1}{\tau} \sum_{t=1}^{\tau} C(L_t, \mathbf{y}_t) \right\|_2 \leq 2M \sqrt{\frac{2}{\delta T}}, \quad (16)$$

where  $M$  is a bound in Euclidian norm over  $C$ , e.g.,  $M = 4 + 2\eta$ . The quantities of interest can be rewritten as

$$\frac{1}{\tau} \sum_{t=1}^{\tau} C(L_t, \mathbf{y}_t) = \left( \frac{N_\tau(\ell)}{\tau} (\mathbf{y}^\ell - \bar{\mathbf{y}}_\tau^\ell) - \frac{N_\tau(\ell)}{\tau} \frac{\eta}{N_B} \mathbf{1}, \frac{N_\tau(\ell)}{\tau} (\bar{\mathbf{y}}_\tau^\ell - \mathbf{y}^\ell) - \frac{N_\tau(\ell)}{\tau} \frac{\eta}{N_B} \mathbf{1} \right)_{\ell \in \{1, \dots, N_\eta\}},$$

where we recall that we denoted for all  $\ell$  such that  $N_\tau(\ell) > 0$  the average of their corresponding mixed actions  $\mathbf{y}_t$  by

$$\bar{\mathbf{y}}_\tau^\ell = \frac{1}{N_\tau(\ell)} \sum_{t=1}^{\tau} \mathbf{y}_t \mathbb{I}_{\{L_t = \ell\}}.$$

The projection in  $\ell^2$ -norm of quantity of interest onto  $(\mathbb{R}_-)^{2N_\eta N_B}$  is formed by its non-positive components, so that its square distance to  $(\mathbb{R}_-)^{2N_\eta N_B}$  equals

$$\inf_{c \in (\mathbb{R}_-)^{2N_\eta N_B}} \left\| c - \frac{1}{\tau} \sum_{t=1}^{\tau} C(L_t, \mathbf{y}_t) \right\|_2^2 = \sum_{\ell=1}^{N_\eta} \left( \frac{N_\tau(\ell)}{\tau} \right)^2 \sum_{b \in \mathcal{B}} \underbrace{\left( \left( y_b^\ell - \bar{y}_{\tau,b}^\ell - \frac{\eta}{N_B} \right)_+^2 + \left( \bar{y}_{\tau,b}^\ell - y_b^\ell - \frac{\eta}{N_B} \right)_+^2 \right)}_{= (\bar{y}_{\tau,b}^\ell - y_b^\ell - \eta/N_B)_+^2}.$$

Therefore, our target is achieved: using first that  $(\cdot)_+$  is subadditive, then applying the Cauchy-Schwarz inequality,

$$\begin{aligned} \sum_{\ell=1}^{N_\eta} \frac{N_\tau(\ell)}{\tau} \left( \|\mathbf{y}^\ell - \bar{\mathbf{y}}_\tau\|_1 - \eta \right)_+ &\leq \sum_{\ell=1}^{N_\eta} \frac{N_\tau(\ell)}{\tau} \sum_{b \in \mathcal{B}} \left( |y_b^\ell - \bar{y}_{\tau,b}^\ell| - \frac{\eta}{N_B} \right)_+ \\ &\leq \sqrt{N_\eta N_B} \sqrt{\sum_{\ell=1}^{N_\eta} \left( \frac{N_\tau(\ell)}{\tau} \right)^2 \sum_{b \in \mathcal{B}} \left( |y_b^\ell - \bar{y}_{\tau,b}^\ell| - \frac{\eta}{N_B} \right)_+^2} \\ &\leq 2M \sqrt{N_\eta N_B} \sqrt{\frac{2}{\delta T}}, \end{aligned}$$

where the last inequality holds, by (16), for all  $\tau \geq T$  with probability at least  $1 - \delta$ . Choosing an integer  $T_\delta$  sufficiently large so that

$$2M \sqrt{N_\eta N_B} \sqrt{\frac{2}{\delta T_\delta}} \leq \delta$$

concludes the proof of the property stated in (4).

### Appendix B. Proof of Lemma 6.3.

PROOF. For all  $(i, j) \in \mathcal{I} \times \mathcal{J}$ , the quantity  $H(i, j)$  is a probability distribution over the set of signals  $\mathcal{H}$ ; we denote by  $H_s(i, j)$  the probability mass that it puts on some signal  $s \in \mathcal{H}$ .

Equation (7) indicates that for each pair  $(i, s) \in \mathcal{I} \times \mathcal{H}$ ,

$$\sum_{t=(n-1)L+1}^{nL} \left( \frac{\mathbb{I}_{\{S_t=s\}} \mathbb{I}_{\{I_t=i\}}}{p_{I_t,n}} - H_s(i, J_t) \right)$$

is a sum of  $L$  elements of a martingale difference sequence, with respect to the filtration whose  $t$ -th element is generated by  $\mathbf{p}_n$ , the pairs  $(I_s, S_s)$  for  $s \leq t$ , and  $J_s$  for  $s \leq t + 1$ . The conditional variances of the increments are bounded by

$$\mathbb{E}_t \left[ \left( \frac{\mathbb{I}_{\{S_t=s\}} \mathbb{I}_{\{I_t=i\}}}{p_{I_t,n}} \right)^2 \right] \leq \frac{1}{p_{i,n}^2} \mathbb{E}_t [\mathbb{I}_{\{I_t=i\}}] = \frac{1}{p_{i,n}};$$

since by definition of the strategy,  $\mathbf{p}_n = (1 - \gamma) \mathbf{x}_n + \gamma \mathbf{u}$ , we have that  $p_{i,n} \geq \gamma/N_{\mathcal{I}}$ , which shows that the sum of the conditional variances is bounded by

$$\sum_{t=(n-1)L+1}^{nL} \text{Var}_t \left( \frac{\mathbb{I}_{\{S_t=s\}} \mathbb{I}_{\{I_t=i\}}}{p_{I_t,n}} \right) \leq \frac{LN_{\mathcal{I}}}{\gamma}.$$

The Bernstein-Freedman inequality (see Freedman [11] or Cesa-Bianchi et al. [6], Lemma A.1) therefore indicates that with probability at least  $1 - \delta$ ,

$$\left| \frac{1}{L} \sum_{t=(n-1)L+1}^{nL} \frac{\mathbb{I}_{\{S_t=s\}} \mathbb{I}_{\{I_t=i\}}}{p_{I_t,n}} - \underbrace{\frac{1}{L} \sum_{t=(n-1)L+1}^{nL} H_s(i, J_t)}_{= H_s(i, \hat{\mathbf{q}}_n)} \right| \leq \sqrt{2 \frac{N_{\mathcal{I}}}{\gamma L} \ln \frac{2}{\delta}} + \frac{1}{3} \frac{N_{\mathcal{I}}}{\gamma L} \ln \frac{2}{\delta}.$$

Therefore, by summing the above inequalities over  $i \in \mathcal{I}$  and  $s \in \mathcal{H}$ , we get (after a union bound) that with probability at least  $1 - N_{\mathcal{I}}N_{\mathcal{H}}\delta$ ,

$$\left\| \tilde{\sigma}_n - \tilde{H}(\hat{\mathbf{q}}_n) \right\|_2 \leq \sqrt{N_{\mathcal{I}}N_{\mathcal{H}}} \left( \sqrt{2 \frac{N_{\mathcal{I}}}{\gamma L} \ln \frac{2}{\delta}} + \frac{1}{3} \frac{N_{\mathcal{I}}}{\gamma L} \ln \frac{2}{\delta} \right).$$

Finally, since  $\hat{\sigma}_n$  is the projection in the  $\ell^2$ -norm of  $\tilde{\sigma}_n$  onto the convex set  $\mathcal{F}$ , to which  $\tilde{H}(\hat{\mathbf{q}}_n)$  belongs, we have that

$$\left\| \hat{\sigma}_n - \tilde{H}(\hat{\mathbf{q}}_n) \right\|_2 \leq \left\| \tilde{\sigma}_n - \tilde{H}(\hat{\mathbf{q}}_n) \right\|_2,$$

and this concludes the proof.  $\square$

## References

- [1] J. Abernethy, P. L. Bartlett, and E. Hazan. Blackwell approachability and low-regret learning are equivalent. In *Proceedings of the Twenty-Fourth Annual Conference on Learning Theory (COLT'11)*. Omnipress, 2011.
- [2] D. Blackwell. An analog of the minimax theorem for vector payoffs. *Pacific Journal of Mathematics*, 6:1–8, 1956.
- [3] D. Blackwell. Controlled random walks. In *Proceedings of the International Congress of Mathematicians, 1954, Amsterdam, vol. III*, pages 336–338, 1956.
- [4] A. Blum and Y. Mansour. From external to internal regret. *Journal of Machine Learning Research*, 8:1307–1324, 2007.
- [5] N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- [6] N. Cesa-Bianchi, G. Lugosi, and G. Stoltz. Regret minimization under partial monitoring. *Mathematics of Operations Research*, 31:562–580, 2006.
- [7] X. Chen and H. White. Laws of large numbers for Hilbert space-valued mixingales with applications. *Econometric Theory*, 12:284–304, 1996.
- [8] A.P. Dawid. The well-calibrated Bayesian. *Journal of the American Statistical Association*, 77: 605–613, 1982.
- [9] D. Foster and R. Vohra. Asymptotic calibration. *Biometrika*, 85:379–390, 1998.
- [10] D. Foster and R. Vohra. Regret in the on-line decision problem. *Games and Economic Behavior*, 29:7–36, 1999.
- [11] D.A. Freedman. On tail probabilities for martingales. *Annals of Probability*, 3:100–118, 1975.
- [12] J.E. Goodman and J. O'Rourke, editors. *Handbook of Discrete and Computational Geometry*. Discrete Mathematics and its Applications. Chapman & Hall/CRC, Boca Raton, FL, second edition, 2004.

- [13] S. Hart and A. Mas-Colell. A simple adaptive procedure leading to correlated equilibrium. *Econometrica*, 68:1127–1150, 2000.
- [14] S. Hart and A. Mas-Colell. A general class of adaptive strategies. *Journal of Economic Theory*, 98:26–54, 2001.
- [15] E. Lehrer and E. Solan. Learning to play partially-specified equilibrium. Mimeo, 2007.
- [16] G. Lugosi, S. Mannor, and G. Stoltz. Strategies for prediction under imperfect monitoring. *Mathematics of Operations Research*, 33:513–528, 2008. An extended abstract was presented at COLT’07.
- [17] S. Mannor and N. Shimkin. On-line learning with imperfect monitoring. In *Proceedings of the Sixteenth Annual Conference on Learning Theory (COLT’03)*, pages 552–567. Springer, 2003.
- [18] S. Mannor and N. Shimkin. Regret minimization in repeated matrix games with variable stage duration. *Games and Economic Behavior*, 63(1):227–258, 2008.
- [19] S. Mannor and G. Stoltz. A geometric proof of calibration. *Mathematics of Operations Research*, 35:721–727, 2010.
- [20] S. Mannor, J. Tsitsiklis, and J. Y. Yu. Online learning with sample path constraints. *Journal of Machine Learning Research*, 10(Mar):569–590, 2009.
- [21] J.-F. Mertens, S. Sorin, and S. Zamir. Repeated games. Technical Report no. 9420, 9421, 9422, Université de Louvain-la-Neuve, 1994.
- [22] V. Perchet. Calibration and internal no-regret with random signals. In *Proceedings of the Twentieth International Conference on Algorithmic Learning Theory (ALT’09)*, pages 68–82, 2009.
- [23] V. Perchet. Approachability of convex sets in games with partial monitoring. *Journal of Optimization Theory and Applications*, 149:665–677, 2011.
- [24] V. Perchet. Internal regret with partial monitoring calibration-based optimal algorithms. *Journal of Machine Learning Research*, 2011. In press.
- [25] V. Perchet and M. Quincampoix. On an unified framework for approachability in games with or without signals. Mimeo, 2011.
- [26] A. Piccolboni and C. Schindelhauer. Discrete prediction games with arbitrary feedback and loss. In *Proceedings of the Fourteenth Annual Conference on Computational Learning Theory (COLT’01)*, pages 208–223, 2001.
- [27] A. Rakhlin, K. Sridharan, and A. Tewari. Online learning: Beyond regret. In *Proceedings of the Twenty-Fourth Annual Conference on Learning Theory (COLT’11)*. Omnipress, 2011.
- [28] J. Rambau and G. Ziegler. Projections of polytopes and the generalized Baues conjecture. *Discrete and Computational Geometry*, 16:215–237, 1996.
- [29] A. Rustichini. Minimizing regret: The general case. *Games and Economic Behavior*, 29:224–243, 1999.

**Acknowledgments.** Shie Mannor was partially supported by the ISF under contract 890015 and the Google Inter-university center for Electronic Markets and Auctions. Vianney Perchet benefited from the support of the ANR under grant ANR-10-BLAN 0112. Gilles Stoltz acknowledges support from the French National Research Agency (ANR) under grant EXPLO/RA (“Exploration–exploitation for efficient resource allocation”) and by the PASCAL2 Network of Excellence under EC grant no. 506778.

An extended abstract of this paper appeared in the *Proceedings of the 24th Annual Conference on Learning Theory (COLT’11)*, JMLR Workshop and Conference Proceedings, Volume 19, pages 515–536, 2011.