

# Policy Optimization via Adv2: Adversarial Learning on Advantage Functions

Anonymous authors

Paper under double-blind review

## Abstract

We revisit the reduction of learning in adversarial Markov decision processes [MDPs] to adversarial learning based on  $Q$ -values; this reduction has been considered in a number of recent articles as one building block to perform policy optimization. Namely, we first restate this reduction in a greater generality: it may involve any adversarial learning strategy (not just exponential weights) and it may be based indifferently on  $Q$ -values or on advantage functions. Both twists may be leveraged to obtain improved practical performance. We then present two extensions: convergence of the last iterate for a vast class of adversarial learning strategies (again, not just exponential weights); stronger regret criteria for learning in MDPs, inherited from the stronger regret criteria of adversarial learning called strongly adaptive regret and tracking regret. Finally, we demonstrate how adversarial learning, also referred to as aggregation of experts, relates to orchestration of expert policies (also known as imitation learning): we obtain stronger forms of performance guarantees in this setting than existing ones, via yet another, simple reduction.

A word intended to reviewers: we included three brief formal proofs in the main body, as well as many other detailed arguments. This is why the appendix is short (less than 6 pages, possibly being even shortened to 4 pages if the proof of the performance difference lemma is omitted) but the main body features 17 pages.

## 1 Introduction

In this article, we revisit a specific approach in policy optimization for adversarial Markov decision processes [MDPs] in the episodic setting, namely, the closed-form design of the policies output over time (which change in an incremental way) based on estimated value functions. In virtually all previous work, these policies are computed thanks to the same adversarial-learning strategy, referred to under possibly different names: exponential weights, weighted majority, Boltzmann reweighting, or online mirror descent, to name a few. It turns out that different adversarial-learning strategies may be used, which may have important consequences in practice: this choice can, for instance, influence the computational efficiency or the robustness to estimation errors.

### 1.1 Brief literature review

Before reviewing in detail our contributions, we first provide a concise overview of the related literature and justify some claims contained in the previous paragraph.

**Adversarial MDPs / Reduction to adversarial learning.** The setting of adversarial MDPs was introduced by Even-Dar et al. (2009) and Yu et al. (2009). As in the standard episodic setup, the transition kernels dictating the evolution of the states are unknown and constant over time. However, the reward functions vary over time and may be chosen by some adversary; they are possibly revealed at the end of an episode. Both references were also the first ones to introduce a reduction of the control of adversarial MDPs to standard adversarial learning (a setting also called expert prediction; see Cesa-Bianchi & Lugosi, 2006 for an overview thereof). In this article, we will be interested in closed-form policy optimization, and not in approaches relying on so-called occupancy measures (introduced by Zimin & Neu, 2013), which solve

a complex convex optimization problem at each episode and do not result in closed-form expressions for the policies output (see, e.g., Rosenberg & Mansour, 2019).

**Policy optimization.** Policy optimization refers to designing policies to be used at each episode, often obtained by sequential incremental updates, and may be opposed to value-based learning in MDPs, which focuses on estimating and improving value functions rather than directly constructing policies. Several approaches were considered in policy optimization, for instance, (natural) policy gradient (Sutton et al., 2000, Kakade & Langford, 2002), and variants like Trust Region Optimization or Proximal Policy Optimization (TRPO and PPO, respectively; see Schulman et al., 2015, Schulman et al., 2017). We will rather be interested in the closed-form policy design relying on estimates of  $Q$ -value functions. This vein of research includes the works by Shani et al. (2020), Cai et al. (2020), He et al. (2022), Zhao et al. (2023), Tiapkin et al. (2024) (see also Abbasi-Yadkori et al., 2019) to name a few contributions illustrating well the angle used. The settings differ in these articles depending, among others, on the feedback on the reward functions (full monitoring or bandit feedback) and on the structural assumptions, or lack thereof, on the transition kernels.

However, all cited references have one thing in common: they rely on the same adversarial-learning strategy (except Tiapkin et al., 2024, which points to the present article).

**A single adversarial-learning strategy, based on exponential weights.** This same adversarial-learning strategy is known under different names and relies on exponential weights; namely, Agarwal et al. (2021, Section 5.3) refers to it as multiplicative weights updates, Abbasi-Yadkori et al. (2019), as the Boltzmann policy<sup>1</sup>, Shani et al. (2020) and Zhao et al. (2023), as online mirror descent (with a Kullback-Leibler regularization), while Cai et al. (2020) and He et al. (2022) do not write any explicit name but obtain its expression by some follow-the-regularized-leader approach with a Kullback-Leibler regularization (referring to the same closed-form update obtained by earlier references).

Interestingly, this strategy based on exponential weights aligns with the concept of natural policy gradient for non-adversarial MDPs when the policy parametrization is softmax: both approaches involve the same update rule on the weights (this explicit update rule was, for instance derived, in Agarwal et al., 2021, Section 5.3, see also Kakade, 2001). This specific case, as the intersection of two optimisation paradigms, leads to remarkable theoretical guarantees in non-adversarial MDPs; see, in particular, the recent work by Müller & Montúfar (2024) and references therein.

Two exceptions to the use of the exponential-weight strategy are provided by Even-Dar et al. (2009) and Yu et al. (2009), which resort to a strategy called follow-the-perturbed-leader (Kalai & Vempala, 2005); but their setting and objectives are somewhat different to the ones considered in this article and in the references of the previous paragraph.

**Previous reductions of learning in MDPs to adversarial learning.** We provide a specific analysis of the strategy based on exponential weights in Section 6, obtaining improved regret bounds compared to the analyses provided in the mentioned references. These analyses range from a few-line-long proof by direct reduction to adversarial learning in Shani et al. (2020), that we copy in Section 3.2 (but that can be improved in the specific case of exponential weights), to longer proofs (possibly several pages, see, e.g., Zhao et al., 2023, Appendix A.1). The typical proofs are one-page-long, do not clearly identify a reduction, and consist of ad hoc adaptations of the proof for exponential weights based on telescoping Kullback-Leibler terms<sup>2</sup> à la Freund & Schapire (1999), as in Agarwal et al. (2021, Section 5.3) or Cai et al. (2020). We note that the cited references actually run the exponential-weight strategy on estimated  $Q$ -values: more details are provided in Section 7.

In a nutshell, among all cited references, Shani et al. (2020) already clearly identified how to reduce learning MDPs to adversarial learning, but only leveraged this fact for one specific adversarial learning strategy (and provided suboptimal bounds for this strategy).

<sup>1</sup>Abbasi-Yadkori et al. (2019) even states that “the choice of the Boltzmann policy is not arbitrary”, but one of the point of the present article is to actually show that it is, as many other choices of adversarial-learning strategies are suitable.

<sup>2</sup>Simpler proofs of performance exist for exponential weights in the adversarial setting (based on Hoeffding’s lemma, see Cesa-Bianchi & Lugosi, 2006, Section 2.2).

## 1.2 Contributions and outline of this article

In Section 2, we define formally the setting of episodic adversarial Markov decision processes [MDPs] and state our objective: the minimization of a cumulative regret defined as the sum of the differences between the value functions of the best stationary policy and of the policies output.

Section 3 recalls the reduction of learning in MDPs to adversarial learning as clearly stated by Shani et al. (2020). We essentially replicate their proof, based on the performance difference lemma, up to a single, straightforward extension: the consideration of a vast family of possible adversarial learning strategies, not just exponential weights with a constant learning rate. For instance, the ML-Prod and ML-Poly strategies (Gaillard et al., 2014, Gaillard et al., 2021) are suitable adversarial-learning strategies that exhibit in general much better empirical performance than exponential weights. Another, immediate, remark is that the theoretical guarantees hold when adversarial-learning strategies are fed with advantage functions instead of  $Q$ -functions, which is a second source of improved empirical performance.

We then present three extensions and discuss two twists: a special-case analysis for exponential weights with improved bounds, and how to use the general theory developed in real-case scenarios where advantage functions need to be estimated (in particular due to the transition kernels being unknown).

**Extension 1: convergence of the last iterate.** Section 4 focuses on a simple regret instead of a cumulative regret, in the case reward functions are constant over time: the difference between the last policy output and the best policy. Agarwal et al. (2021, Section 5.3) controlled this quantity for exponential weights with a constant learning rate (in the discounted setting). We show how to extend their argument to a large class of adversarial strategies satisfying a natural property that we call “monotonicity of weights”.

**Extension 2: Stronger forms of regret.** Section 5 shows that the general reduction studied also works for a stronger notion of regret called strongly adaptive regret and consisting of studying the sums of differences in value functions over subintervals of time. As a consequence, the so-called tracking regret may also be controlled: therein, the comparison is made not to the best stationary policy but to the best sequence of policies with few shifts. To the best of our knowledge, the control of such improved forms of regret for MDPs is an original contribution.

**The special case of exponential weights.** Section 6 leverages elements from Extensions 1 and 2 to show that when the adversarial learning strategy consists of exponential weights with a constant learning rate, the (cumulative) regret may be bounded by the number of shifts in the reward sequence. This provides yet another generalization of the results of Agarwal et al. (2021, Section 5.3). In addition, the proof technique of Agarwal et al. (2021, Section 5.3) seemed highly specific to the discounted setting: we provide instead a treatment for the episodic setting.

**Practical versions with estimated advantage functions.** Section 7 puts in perspective the design of policies studied in this article: in practice, advantage functions are unknown but may be estimated, so that the strategies studied earlier in this article should be run on these estimates. We review the literature to explain how the actual regret relates to the regret in terms of estimated value functions. The black-box reduction recalled and emphasized in Section 3 should be helpful in the core of new proofs, to make the latter more modular. (Compare, for instance, the modular approach of Tiapkin et al., 2024 based on the present article to more typical proofs re-deriving regret guarantees in terms of value functions based on mimicking proofs of adversarial regret bounds for exponential weights.)

**Extension 3: Aggregation (orchestration) of expert policies, a.k.a. imitation learning.** Adversarial learning is sometimes called prediction with experts (see Cesa-Bianchi & Lugosi, 2006). Section 8 considers the case where policies output over time are not learned anymore in a direct tabular setting, but are obtained by (state-by-state and stage-by-stage) convex combinations of some expert policies. The aim is to mimic the performance of the overall best such convex combination. This methodology relates to orchestration of expert policies, also called imitation learning (see, for instance, Cheng et al., 2020 and Liu et al., 2023). We show that to deal with this problem, it suffices to consider expert policies as actions in a lifted

MDP and apply all results described earlier in this article. We obtain stronger performance guarantees than in the cited references.

## 2 Setting and aims

**Notation.** We denote by  $\mathcal{P}(\mathcal{X})$  the set of probability distributions over some set  $\mathcal{X}$ , either finite or given by an interval of  $\mathbb{R}$  in the sequel. For an integer  $n \geq 1$ , let  $[n] = \{1, \dots, n\}$  denote the set of the first  $n$  integers.

**Setting.** We consider an  $H$ -episodic and (obviously) adversarial Markov decision process [MDP] with finite state and action spaces  $\mathcal{S}$  and  $\mathcal{A}$ , of respective cardinalities  $S$  and  $A$ : each episode  $t \geq 1$  is of length  $H \geq 1$  and is governed by transition kernels  $\mathcal{T} = (\mathcal{T}_h)_{h \in [H-1]}$ , where  $\mathcal{T}_h : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$ , and by reward functions  $\mathcal{R}_t = (\mathcal{R}_{t,h})_{h \in [H]}$ , where  $\mathcal{R}_{t,h} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}([0, 1])$ . The transition kernels are constant over episodes, while the reward functions  $\mathcal{R}_t$  vary between episodes; they may actually be picked by an adversary in some oblivious way, i.e., the entire sequence  $(\mathcal{R}_t)_{t \geq 1}$  is determined by the adversary before the first episode takes place.

We denote by  $r_{t,h} : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  the mean-payoff function associated with some  $\mathcal{R}_{t,h}$ , i.e.,  $r_{t,h}(s, a)$  is the expectation of the distribution  $\mathcal{R}_{t,h}(s, a)$ , for each  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ .

A (stationary, or one-shot) policy  $\pi = (\pi_h)_{h \in [H]}$  is a sequence of mappings  $\pi_h : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ ; we denote by  $\pi_h(\cdot | s)$  the probability distribution over actions that it uses in stage  $h$  and state  $s$ . The learner should determine a policy  $\pi_t$  at the beginning of each episode  $t \geq 1$ , based on the information gained at rounds  $\tau \leq t - 1$ ; that information consists at least of the states observed and actions played therein, as well as the rewards obtained. In some scenarios, additional observations may be performed, which we will explicitly detail; for instance, in Section 3.2, the learning system may observe, among others, the mean-payoff functions  $\mathbf{r}_\tau = (r_{\tau,h})_{h \in [H]}$  at the end of episode  $\tau$ .

At the beginning of each episode  $t \geq 1$ , the same initial state  $s_{t,1} = s_1$  is set. Then, at each stage  $h \in [H-1]$ , the learning system draws an action  $a_{t,h} \sim \pi_{t,h}(\cdot | s_{t,h})$ , after which it obtains and observes a stochastic reward  $r_{t,h} \sim \mathcal{R}_t(s_{t,h}, a_{t,h})$ , while the environment moves to a new state drawn as  $s_{t,h+1} \sim \mathcal{T}_h(\cdot | s_{t,h}, a_{t,h})$ . In the final stage, only an action  $a_{t,H} \sim \pi_{t,H}(\cdot | s_{t,H})$  is drawn, and a reward  $r_{t,H} \sim \mathcal{R}_t(s_{t,H}, a_{t,H})$  is obtained and observed.

By the tower rule, the value function of a given stationary policy  $\pi = (\pi_j)_{j \in [H]}$  at episode  $t \geq 1$  and started at stage  $h \in [H]$  equals, for all  $s \in \mathcal{S}$ ,

$$V_h^{\pi, \mathcal{R}_t}(s) = \mathbb{E}^{\pi, \mathcal{T}} \left[ \sum_{j=h}^H r_{t,j}(s_j, a_j) \mid s_h = s \right], \quad (1)$$

where the piece of notation  $\mathbb{E}^{\pi, \mathcal{T}}$  indicates that actions  $a_h$  and states  $s_h$  in the expectation are governed by the policy  $\pi$  and the transition kernels  $\mathcal{T}$ , as described above.

### 2.1 First aim: direct tabular learning

We evaluate the policies  $\pi_t$  picked over time in terms of their value functions and are interested in mimicking the performance of the best stationary policy in hindsight. More precisely, the learning system aims to control

$$\forall T \geq 1, \quad R_T = \max_{\pi} \sum_{t=1}^T \left( V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) \right), \quad (2)$$

where the maximum is over all stationary policies  $\pi$ . We write “ $\forall T \geq 1$ ” to indicate that either the time horizon  $T$  is unknown or the regret should be controlled for all time horizons. The regret  $R_T$  involves a sum essentially because the reward functions  $\mathcal{R}_t$  evolve over time in a possibly adversarial way; when they are constant over time, then convergence of the last iterate (of the  $T$ -th term in the sum above) may be achieved, see Section 4.

The aim described above is called direct tabular learning as policies  $\pi_t$  are picked by determining, for each stage  $h$  and state  $s$ , the entire probability distribution  $\pi_{t,h}(\cdot|s)$ . The terminology is borrowed from Agarwal et al., 2021, Section 3.

The setting above is summarized in the left part of Box A in Section 8.

The aim described above may be difficult to complete when the number  $A$  of actions is large. In addition, the learning system may sometimes have some prior information given by a finite set of expert policies among which some policies could perform well (the subsets of these good-performing policies could possibly depend on the state). We therefore introduce an alternative aim in Section 8 called aggregation (or orchestration) of expert policies, but actually show that resolving this objective is equivalent to the first aim described above.

## 2.2 Additional notation

For later use, we define  $Q$ -values and advantage functions, and use the same notation as in (1) to that end. For any pair of stationary policy  $\pi$  and reward functions  $\mathcal{R}$ , we define its  $Q$ -value function at episode  $t \in [T]$ , and started from stage  $h \in [H]$ , as

$$Q_h^{\pi, \mathcal{R}_t} : (s, a) \in \mathcal{S} \times \mathcal{A} \mapsto \mathbb{E}^{\pi, \mathcal{T}} \left[ \sum_{j=h}^H r_{t,j}(s_j, a_j) \mid s_h = s, a_h = a \right],$$

and its advantage function as

$$A_h^{\pi, \mathcal{R}_t} : (s, a) \in \mathcal{S} \times \mathcal{A} \mapsto Q_h^{\pi, \mathcal{R}_t}(s, a) - V_h^{\pi, \mathcal{R}_t}(s). \quad (3)$$

We only keep in the notation  $V_h$ ,  $Q_h$ , and  $A_h$  the parameters  $\pi$  and  $\mathcal{R}_t$  that vary, and omit the transition kernels  $\mathcal{T}$ . We use the short-hand notation

$$A_h^{\pi_t, \mathcal{R}_t}(s, \cdot) = (A_h^{\pi_t, \mathcal{R}_t}(s, a))_{a \in \mathcal{A}} \quad (4)$$

to denote the vector of advantages for a given episode  $t$  and a given stage  $h$ .

## 3 Methodology and core result: adversarial learning on advantage functions

**Contributions of this section.** We recall how strategies designed to control the regret in the so-called adversarial setting, i.e., satisfying guarantees as described in Definition 1 below, may be used to construct policies so as to control the regret in terms of value functions. This observation was essentially already made in the literature, at least for exponential weights; see, for instance, how Shani et al. (2020, Section 6) handles its Term (ii).

Before formally stating our main result, we briefly recall what the adversarial setting consists in; see the monograph by Cesa-Bianchi & Lugosi (2006) for a more detailed exposition.

### 3.1 Reminder on adversarial learning

At each round  $t \geq 1$ , based on the information collected during past rounds, a learning strategy picks a convex combination  $w_t = (w_{t,1}, \dots, w_{t,K}) \in \mathcal{P}([K])$  while an opponent player simultaneously picks, possibly at random, a vector  $g_t = (g_{t,1}, \dots, g_{t,K})$  of signed rewards. Both  $w_t$  and  $g_t$  are revealed at the end of the round. More formally, we mean that a learning strategy is a sequence  $\varphi = (\varphi_t)_{t \geq 1}$  of functions  $\varphi_t : \mathbb{R}^{K(t-1)}$  and that  $w_t = \varphi_t((g_\tau)_{\tau \leq t-1})$  for  $t \geq 1$ . This formula means in particular that the initial vector  $w_1 = \varphi_1(\emptyset)$  is constant.

**Definition 1** (adversarial-learning regret bound). *A sequential strategy controls the regret in the adversarial setting with rewards bounded by  $M > 0$  if there exists a sequence  $(B_{T,K})_{T \geq 1}$  of positive numbers with  $B_{T,K}/T \rightarrow 0$  and such that, against all opponent players sequentially picking reward vectors in  $[-M, M]^K$ ,*

$$\forall T \geq 1, \quad \max_{k \in [K]} \sum_{t=1}^T g_{t,k} - \sum_{t=1}^T \sum_{j \in [K]} w_{t,j} g_{t,j} \leq 2M B_{T,K}.$$

The optimal orders of magnitude of  $B_{T,K}$  are  $\sqrt{T \ln K}$  (see Cesa-Bianchi & Lugosi, 2006). In Definition 1, the strategy may know  $M$  and rely on its value. On the contrary, the number  $T$  of rounds is unknown and actually, for the sake of exposition, Definition 1 requires a control of the adversarial regret for all  $T \geq 1$ , which is a mild restriction.

Our main examples of strategies abiding by the constraints of Definition 1 are the potential-based strategies by Cesa-Bianchi & Lugosi (2003). They are defined based on a sequence of non-decreasing functions  $\Phi_t : \mathbb{R} \rightarrow [0, +\infty)$ ; they resort to  $w_{1,k} = 1/K$  and

$$\forall t \geq 2, \quad w_{t,k} = \frac{v_{t,k}}{\sum_{j \in [K]} v_{t,j}}, \quad \text{where} \quad v_{t,k} = \Phi_t \left( \sum_{\tau=1}^{t-1} g_{\tau,k} - \sum_{\tau=1}^{t-1} \sum_{j \in [K]} w_{\tau,j} g_{\tau,j} \right). \quad (5)$$

**Example 1.** *Cesa-Bianchi & Lugosi (2003, Section 2) show that the strategy based on the constant polynomial potentials  $\Phi_t \equiv \Phi : x \mapsto (\max\{x, 0\})^{2 \ln K}$  provides the control  $B_{T,K} = \sqrt{6T \ln K}$  for the regret in the adversarial setting.*

**Example 2.** *Auer et al. (2002) studied exponential potentials  $\Phi_t(x) = \exp(\eta_t x)$  with time-varying learning rates  $\eta_t = (1/M)\sqrt{(\ln K)/t}$ . This sequential strategy controls the regret with  $B_{T,K} = \sqrt{T \ln K}$  in the adversarial setting.*

A final example is of a different, not potential-based, nature.

**Example 3.** *The greedy projection algorithm of Zinkevich (2003) relies on a sequence  $(\eta_t)_{t \geq 1}$  of positive step sizes and sets  $w_{t+1} = \text{proj}(w_t + \eta_t g_t)$  for  $t \geq 1$ , where  $w_1 = (1/K, \dots, 1/K)$  and where  $\text{proj}$  is the convex projection onto  $\mathcal{P}([K])$  in Euclidean norm. For the choices  $\eta_t = (1/M)\sqrt{1/(2Kt)}$ , this strategy controls the regret with  $B_{T,K} = \sqrt{2KT}$  in the adversarial setting.*

### 3.2 Policy optimization via adversarial learning on advantage functions

This section presents rather standard material and must be read accordingly. Indeed, what follows is a reduction that was essentially known, though only applied with exponential weights and on  $Q$ -values rather than on advantage functions. The proof follows the one by Shani et al. (2020, Section 6) (see also Agarwal et al., 2021, proof of Theorem 16), i.e., is based on the performance difference lemma.

For each stage  $h \in [H]$ , we fix a sequential strategy  $\varphi_h = (\varphi_{t,h})_{t \geq 1}$  in the adversarial setting, relying on reward vectors bounded by  $M_h = H - h + 1$  and of dimension  $K = A$ , i.e., indexed by  $\mathcal{A}$ . We run these strategies on the advantage functions, in a stage-by-stage and state-by-state manner, as follows: for all  $t \geq 1$ ,

$$\forall h \in [H], \quad \forall s \in \mathcal{S}, \quad \pi_{t,h}(\cdot | s) = \varphi_{t,h} \left( (A_h^{\pi_{\tau}, \mathcal{R}_{\tau}}(s, \cdot))_{\tau \leq t-1} \right), \quad (6)$$

where we used the notation defined in (4). We refer to this strategy as  $(\varphi_h)_{h \in [H]} \text{-Adv2}$ , for  $(\varphi_h)_{h \in [H]} \text{-adv}$  adversarial learning on advantage functions.

It constitutes a “theoretical” strategy, as it relies on the oracle knowledge of the advantage functions—an issue that we discuss and mitigate later in Section 7. The strategy could be run instead on  $Q$ -values, see Remark 1 below.

**Theorem 1.** *In the setting of Section 2 where rewards lie in  $[0, 1]$ , if, for all  $h \in [H]$ , the sequential strategies  $\varphi_h$  control the regret in the adversarial setting (Definition 1) by  $B_{T,A}$  for  $A$ -dimensional reward vectors bounded by  $H - h + 1$ , then the  $(\varphi_h)_{h \in [H]} \text{-Adv2}$  strategy defined in (6) controls the regret as:*

$$\forall T \geq 1, \quad \max_{\pi} \sum_{t=1}^T \left( V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) \right) \leq H(H+1) B_{T,A}.$$

As indicated above, following Shani et al. (2020, Section 6), the (short) proof of Theorem 1 relies on the so-called performance difference lemma, which we recall next. For the sake of completeness, references for this lemma and a proof thereof are provided in Appendix C.

**Lemma 1** (Performance difference lemma). *Let  $\mu_h^{s_1, \pi, \mathcal{T}}$  be the distribution of the state  $s_{h'}$  of the  $h'$ -th stage, starting from the state  $s_1$  in the first stage, following the stationary policy  $\pi$  and the transition kernels  $\mathcal{T}$ . In a MDP with transition kernels  $\mathcal{T}$ , for all pairs  $\pi, \pi'$  of stationary policies, for all reward functions  $\mathcal{R}$ , and for all stages  $h \in [H]$ ,*

$$\sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \left( V_h^{\pi, \mathcal{R}}(s) - V_h^{\pi', \mathcal{R}}(s) \right) = \sum_{h'=h}^H \sum_{s \in \mathcal{S}} \mu_{h'}^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_{h'}(a|s) A_{h'}^{\pi', \mathcal{R}}(s, a).$$

In particular, for  $h = 1$ ,

$$V_1^{\pi, \mathcal{R}}(s_1) - V_1^{\pi', \mathcal{R}}(s_1) = \sum_{h'=1}^H \sum_{s \in \mathcal{S}} \mu_{h'}^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_{h'}(a|s) A_{h'}^{\pi', \mathcal{R}}(s, a).$$

*Proof of Theorem 1.* We fix a stationary policy  $\pi$  throughout the proof and control the regret with respect to this  $\pi$ .

The first part consists of applying the adversarial-learning regret upper bound for each  $h \in [H]$ . As the rewards take values in  $[0, 1]$ , we have that  $|A_h^{\pi_\tau, \mathcal{R}_\tau}(s, a)| \leq H - h + 1$  for all  $\tau, s, a$ . By the definition of advantage functions (for the equality to 0) and by Definition 1 and the design of the  $(\varphi_h)_{h \in [H]}$ -Adv2 strategy (for the upper bound), we have, for all  $s \in \mathcal{S}$ ,

$$\max_{a \in \mathcal{A}} \sum_{t=1}^T A_h^{\pi_t, \mathcal{R}_t}(s, a) - \sum_{t=1}^T \sum_{a \in \mathcal{A}} \overbrace{\pi_{t,h}(a|s) A_h^{\pi_t, \mathcal{R}_t}(s, a)}^{=0} \leq 2(H - h + 1) B_{T,A}. \quad (7)$$

The second part consists of applying the performance difference lemma, i.e., Lemma 1 above with  $h = 1$ , which guarantees that

$$V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) = \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_h(a|s) A_h^{\pi_t, \mathcal{R}_t}(s, a).$$

Summing this equality over  $t$  and rearranging, we get

$$\begin{aligned} \sum_{t=1}^T \left( V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) \right) &= \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_h(a|s) \sum_{t=1}^T A_h^{\pi_t, \mathcal{R}_t}(s, a) \\ &\leq \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \underbrace{\max_{a \in \mathcal{A}} \sum_{t=1}^T A_h^{\pi_t, \mathcal{R}_t}(s, a)}_{\leq 2(H-h+1) B_{T,A}} \underbrace{\sum_{h=1}^H \pi_h(a|s)}_{=H(H+1)} \leq 2 \sum_{h=1}^H (H - h + 1) B_{T,A}, \end{aligned} \quad (8)$$

where we substituted (7). Here, we crucially used that the weights  $\mu_h^{s_1, \pi, \mathcal{T}}(s)$  are independent of  $t$  as they only depend on the fixed benchmark policy  $\pi$ , on the common transition kernels  $\mathcal{T}$ , and on the initial state  $s_1$  (identical for all  $t$ ).  $\square$

### 3.3 Comments

In this section, we comment and discuss the Adv2 strategy (6) and its bound.

We first note that the regret bound of Theorem 1 is independent of the size  $S$  of the state space; it only depends on the size  $A$  of the action space, on the number  $T$  of episodes, and on the length  $H$  of the episodes. Given that adversarial-learning strategies have a per-round computational complexity typically proportional to  $K$  (with the notation of Section 3.1), the per-round computational complexity of the Adv2 strategies (6) are typically proportional to  $SAH$  as far as the weight updates are concerned. The main computational issue relies in computing (or estimating, see Section 7) the advantage functions  $A_h^{\pi_\tau, \mathcal{R}_\tau}$ .

Second, for potential-based strategies (5), we note that the original definition (6) of Adv2 and the alternative definition based on  $Q$ -values,

$$\pi_{t,h}(\cdot | s) = \varphi_{t,h} \left( (Q_h^{\pi_\tau, \mathcal{R}_\tau}(s, \cdot))_{\tau \leq t-1} \right), \quad (9)$$

lead to the exact same strategies. This may be shown by induction, based on the fact that for all  $h \in [H]$  and  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , due to (7) for the first equality and due to the definitions of value functions for the second equality,

$$\begin{aligned} & \sum_{\tau=1}^{t-1} A_h^{\pi_\tau, \mathcal{R}_\tau}(s, a) - \sum_{\tau=1}^{t-1} \overbrace{\sum_{a \in \mathcal{A}} \pi_{\tau,h}(a|s) A_h^{\pi_\tau, \mathcal{R}_\tau}(s, a)}^{=0} = \sum_{\tau=1}^{t-1} A_h^{\pi_\tau, \mathcal{R}_\tau}(s, a) \\ \text{and} \quad & \sum_{\tau=1}^{t-1} Q_h^{\pi_\tau, \mathcal{R}_\tau}(s, a) - \sum_{\tau=1}^{t-1} \underbrace{\sum_{a \in \mathcal{A}} \pi_{\tau,h}(a|s) Q_h^{\pi_\tau, \mathcal{R}_\tau}(s, a)}_{=V_h^{\pi_\tau, \mathcal{R}_\tau}(s)} = \sum_{\tau=1}^{t-1} A_h^{\pi_\tau, \mathcal{R}_\tau}(s, a). \end{aligned}$$

For general adversarial-learning strategies, the induced strategies (6) and (9) may differ, though they achieve the same regret guarantees, as detailed by the following remark.

**Remark 1.** *An inspection of the proof of Theorem 1 shows that it would also work for the strategies of the form (9). Indeed, the inequality (7) therein would be replaced equivalently by*

$$2(H - h + 1) B_{T,A} \geq \max_{a \in \mathcal{A}} \sum_{t=1}^T Q_h^{\pi_t, \mathcal{R}_t}(s, a) - \sum_{t=1}^T \overbrace{\sum_{a \in \mathcal{A}} \pi_{t,h}(a|s) Q_h^{\pi_t, \mathcal{R}_t}(s, a)}^{=V_h^{\pi_t, \mathcal{R}_t}(s)} = \max_{a \in \mathcal{A}} \sum_{t=1}^T A_h^{\pi_t, \mathcal{R}_t}(s, a),$$

while the rest of the proof would be unaffected. However, using the advantage functions is preferred in practice, as it provides a greater numerical stability, as well as a possibly a lower variance when the value function are estimated (see Section 7).

## 4 Extension 1:

### Convergence of the last iterate for some adversarial learning strategies

**Contributions of this section.** We generalize an argument of Agarwal et al. (2021, Section 5.3), which was provided for exponential weights only (in the discounted setting): the aim is to control the convergence of the last iterate, i.e., to upper bound  $\max_{\pi} V_1^{\pi, \mathcal{R}}(s_1) - V_1^{\pi^T, \mathcal{R}}(s_1)$ , when (mean) rewards functions are constant over time.

More precisely, for some adversarial-learning strategies  $\varphi$ , satisfying some property which we call monotonicity of weights, and in case reward functions do not vary over time (or even just mean reward functions, see Remark 2) the result of Theorem 1 may be strengthened into a convergence result of the last iterate, at a rate faster by a factor of  $1/T$  compared to the convergence of the cumulative regret (2).

**Definition 2** (monotonicity of weights). *A sequential strategy  $\varphi = (\varphi_t)_{t \geq 1}$  in the adversarial setting satisfies monotonicity of weights if against all opponent players sequentially picking  $K$ -dimensional reward vectors  $g_\tau = (g_{\tau,k})_{k \in [K]}$ , the convex weights output by  $\varphi$  are such that*

$$\forall t \geq 1, \quad \sum_{k \in [K]} w_{t+1,k} \left( g_{t,k} - \sum_{j \in [K]} w_{t,j} g_{t,j} \right) \geq 0,$$

where we recall the notation  $(w_{t,k})_{k \in [K]} = \varphi_t(g_1, \dots, g_{t-1})$  and  $(w_{t+1,k})_{k \in [K]} = \varphi_t(g_1, \dots, g_{t-1}, g_t)$ .

The reason behind the terminology of monotonicity of weights, as well as the proof of the lemma below, may be found in Appendix A.



**Lemma 2.** *The potential-based strategies (5) of Cesa-Bianchi & Lugosi (2003) with constant, non-decreasing potential functions  $\Phi_t \equiv \Phi$  (like in Example 1) and the greedy projection algorithm (Example 3) of Zinkevich (2003) satisfy monotonicity of weights.*

**Theorem 2.** *Assume reward functions do not vary over time and are all equal to some  $\mathcal{R}$ . If, for all  $h \in [H]$ , the sequential strategies  $\varphi_h$  satisfy monotonicity of weights (Definition 2) and control the regret in the adversarial setting (Definition 1) by  $B_{T,A}$  for  $A$ -dimensional reward vectors bounded by  $H - h + 1$ , then the last iterate of the  $(\varphi_h)_{h \in [H]}$ -Adv2 strategy defined in (6) satisfies*

$$\forall T \geq 1, \quad \max_{\pi} V_1^{\pi, \mathcal{R}}(s_1) - V_1^{\pi_T, \mathcal{R}}(s_1) \leq \frac{H(H+1) B_{T,A}}{T}.$$

The bound by Agarwal et al. (2021, Section 5.3), where the exponential weights with a constant learning rate are considered, corresponds to this theorem but is stated separately in Corollary 2, for reasons that will be made clear in Section 6. As the proof of Theorem 2 is concise, we provide it in the main body of this article.

*Proof.* Given the definition (6), the monotonicity of weights (Definition 2), and the definition of advantage functions, we have that, for all  $t \geq 1$ , for all  $h \in [H]$ , and  $s \in \mathcal{S}$ ,

$$\sum_{a \in \mathcal{A}} \pi_{t+1,h}(a|s) A_h^{\pi_t, \mathcal{R}}(s, a) \geq \sum_{a \in \mathcal{A}} \pi_{t,h}(a|s) A_h^{\pi_t, \mathcal{R}}(s, a) = 0.$$

Therefore, the performance difference lemma, i.e., Lemma 1 above with  $h = 1$ , shows that

$$V_1^{\pi_{t+1}, \mathcal{R}}(s_1) - V_1^{\pi_t, \mathcal{R}}(s_1) = \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi_{t+1}, \mathcal{T}}(s) \underbrace{\sum_{a \in \mathcal{A}} \pi_{t+1,h}(a|s) A_h^{\pi_t, \mathcal{R}}(s, a)}_{\geq 0} \geq 0.$$

(This is the part of the proof where crucially use that reward functions do not vary over time.) Thus,

$$\max_{\pi} V_1^{\pi, \mathcal{R}}(s_1) - V_1^{\pi_T, \mathcal{R}}(s_1) \leq \max_{\pi} V_1^{\pi, \mathcal{R}}(s_1) - \frac{1}{T} \sum_{t=1}^T V_1^{\pi_t, \mathcal{R}}(s_1) \leq \frac{H(H+1) B_{T,A}}{T},$$

where we applied Theorem 1 for the final bound.  $\square$

**Remark 2.** *An inspection of the proof above shows that what matters actually is only that mean reward functions  $\mathbf{r}_t = (r_{t,h})_{h \in [H]}$  be constant over time. Indeed, the value and advantage functions only depend on the  $\mathcal{R}_t$  through the  $\mathbf{r}_t$ ; this fact is also illustrated in the proof of the performance difference lemma which only requires identical mean reward functions, not the identity of reward functions.*

## 5 Extension 2: Stronger forms of regret

**Contributions of this section.** *We push the logic of the reduction of the control of MDPs to adversarial learning, and leverage stronger forms of regret in adversarial learning. This section thus presents new regret criteria for learning MDPs.*

Definition 1 considers the simplest definition of adversarial regret. However, several stronger notions of regrets were proposed by the literature. The proof of Theorem 1 shows that the vanilla notion of adversarial regret of Definition 1 regret may be transferred into the vanilla regret (2) in terms of value functions. Actually, this proof may be mimicked to transfer stronger notions of adversarial regret. We illustrate this possibility with two notions of adversarial regrets that replace the comparison to a single global policy by local comparisons (strongly adaptive regret) or by global comparisons to sequences of policies (tracking regret).

## 5.1 Strongly adaptive regret and tracking regret in adversarial learning

We use again the notation for adversarial learning introduced at the beginning of Section 2.1. The first extended notion of regret, called strongly adaptive regret, measures performance simultaneously over each given sub-interval of time with respect to the best component over that sub-interval. It was introduced by Daniely et al. (2015), based on the concept of adaptive regret from Hazan & Seshadhri (2009), itself based on the work by Littlestone & Warmuth (1994).

**Definition 3** (strongly adaptive regret in adversarial learning). *A sequential strategy controls the strongly adaptive regret in the adversarial setting with rewards bounded by  $M > 0$  if there exist positive numbers  $B_{T,K,\tau}$ , where  $T \geq 1$  and  $\tau \in [T]$ , such that, against all opponent players sequentially picking reward vectors in  $[-M, M]^K$ ,*

$$\forall T \geq 1, \quad \forall \tau \in [T], \quad \max_{t_0 \in [T-\tau+1]} \left\{ \max_{k \in [K]} \sum_{t=t_0}^{t_0+\tau-1} g_{t,k} - \sum_{t=t_0}^{t_0+\tau-1} \sum_{j \in [K]} w_{t,j} g_{t,j} \right\} \leq 2M B_{T,K,\tau},$$

$$\text{and} \quad \sup_{\tau \in [T]} \frac{B_{T,K,\tau}}{T} \rightarrow 0 \quad \text{as } T \rightarrow \infty.$$

It follows from Daniely et al. (2015, Theorem 1) that the strongly adaptive regret can be controlled with bounds  $B_{T,K,\tau}$  of order  $\sqrt{\tau}$  up to logarithmic factors.

A closely related notion is the tracking regret, introduced by Herbster & Warmuth (1998) (see also Cesa-Bianchi & Lugosi, 2006, Chapter 5.2), where the comparison is taken over all time steps but against sequences  $k_{1:T} = (k_1, k_2, \dots, k_T)$  with values in  $[K]$ , with  $C$  shifts (i.e.,  $C$  time steps such that  $k_t \neq k_{t-1}$ ). The tracking regret involves

$$\sum_{t=1}^T g_{t,k_t} - \sum_{t=1}^T \sum_{j \in [K]} w_{t,j} g_{t,j}.$$

There are strong links between strongly adaptive and tracking regret, see Adamskiy et al. (2016). In particular, we explain, in the context of regret with value functions, how strongly adaptive regret with  $B_{T,K,\tau}$  of order  $\sqrt{\tau}$  up to logarithmic factors entails tracking regret of order  $\sqrt{CT}$ ; see Corollary 1.

## 5.2 Transfer to strongly adaptive regret bounds for value functions and policies

Based on Definition 3, we obtain the following regret bound in terms of value functions and policies.

**Theorem 3.** *In the setting of Section 2 where rewards lie in  $[0, 1]$ , if, for all  $h \in [H]$ , the sequential strategies  $\varphi_h$  control the strongly adaptive regret in the adversarial setting (Definition 3) by  $B_{T,A,\tau}$  for  $A$ -dimensional reward vectors bounded by  $H - h + 1$ , then the  $(\varphi_h)_{h \in [H]}$ -Adv2 strategy defined in (6) ensures that*

$$\forall T \geq 1, \quad \forall \tau \in [T], \quad \max_{t_0 \in [T-\tau+1]} \left\{ \max_{\pi} \sum_{t=t_0}^{t_0+\tau-1} \left( V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) \right) \right\} \leq H(H+1) B_{T,A,\tau}.$$

The proof of Theorem 3 is obtained by a direct adaptation of the proof of Theorem 1, which basically consists of considering sums over subintervals only instead of sums over all time periods. Again, since the proof is concise, we provide it here.

*Proof of Theorem 3.* We fix a stationary policy  $\pi$  throughout the proof and control some adaptive regret with respect to this  $\pi$ . By the design (6) of the Adv2 strategy, which operates stage by stage and state by state, we have that for all  $h \in [H]$  and  $s \in \mathcal{S}$ , the following holds, by Definition 3: for all  $T \geq 1$  and  $\tau \in [T]$ ,

$$\max_{t_0 \in [T-\tau+1]} \left\{ \max_{a \in \mathcal{A}} \sum_{t=t_0}^{t_0+\tau-1} A_h^{\pi_t, \mathcal{R}_t}(s, a) - \sum_{t=t_0}^{t_0+\tau-1} \overbrace{\sum_{a \in \mathcal{A}} \pi_{t,h}(a|s) A_h^{\pi_t, \mathcal{R}_t}(s, a)}^{=0} \right\} \leq 2(H-h+1) B_{T,A,\tau}.$$

The same application of the performance difference lemma as in the proof of Theorem 1 entails that for all  $T \geq 1$ ,  $\tau \in [T]$ , and  $t_0 \in [T - \tau + 1]$ ,

$$\begin{aligned} \sum_{t=t_0}^{t_0+\tau-1} \left( V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) \right) &= \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_h(a|s) \sum_{t=t_0}^{t_0+\tau-1} A_h^{\pi_t, \mathcal{R}_t}(s, a) \\ &\leq \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \underbrace{\max_{a \in \mathcal{A}} \sum_{t=t_0}^{t_0+\tau-1} A_h^{\pi_t, \mathcal{R}_t}(s, a)}_{\leq 2(H-h+1) B_{T, A, \tau}} \leq H(H+1) B_{T, A, \tau}. \end{aligned}$$

Here again, we crucially used that the weights  $\mu_h^{s_1, \pi, \mathcal{T}}(s)$  are independent of  $t$ . The claimed bound follows by taking the maximum over  $\pi$  and over  $t_0 \in [T - \tau + 1]$ .  $\square$

### 5.3 Tracking regret bounds for value functions and policies

We now detail a consequence of the bound of Theorem 3 for tracking regret.

We consider sequences  $\pi^{(1:T)} = (\pi^{(1)}, \pi^{(2)}, \dots, \pi^{(T)})$  of stationary policies and define their numbers of shifts  $c(\pi^{(1:T)})$  as follows: the smallest integer  $c'$  such that there exist  $c' - 1$  integers  $\tau_2, \dots, \tau_{c'}$  with values in  $[T]$  such that, denoting  $\tau_1 = 1$  and  $\tau_{c'+1} = T + 1$ ,

$$\forall i \in \{2, \dots, c' + 1\}, \quad \forall t \in \{\tau_{i-1}, \dots, \tau_i - 1\}, \quad \pi^{(t)} = \pi^{(\tau_{i-1})}. \quad (10)$$

The tracking regret against sequences  $\pi^{(1:T)}$  of stationary policies with at most  $C$  shifts is defined as

$$\max_{\substack{\pi^{(1:T)} \text{ such that} \\ c(\pi^{(1:T)}) \leq C}} \sum_{t=1}^T V_1^{\pi^{(t)}, \mathcal{R}_t}(s_1) - \sum_{t=1}^T V_1^{\pi_t, \mathcal{R}_t}(s_1).$$

We fix  $T, C \geq 1$  and a sequence  $\pi^{(1:T)}$  of stationary policies, with  $C$  shifts, occurring at  $\tau_1, \tau_2, \dots, \tau_C$  (where we recall that  $\tau_1 = 1$ ). We introduce  $\tau_{C+1} = T + 1$  and partition time into the  $C$  intervals  $[\tau_i, \tau_{i+1} - 1]$ , for  $i \in [C]$ . The  $C$  values successively taken by the sequence  $\pi^{(1:T)}$  consist of the  $\pi^{(\tau_i)}$ , where  $i \in [C]$ . By applying the bound of Theorem 3 on each of the  $C$  intervals  $[\tau_i, \tau_{i+1} - 1]$ , we obtain the following corollary.

**Corollary 1.** *Under the assumptions of Theorem 3, the  $(\varphi_h)_{h \in [H]}$ -Adv2 strategy defined in (6) also ensures that  $\forall T \geq 1, \forall C \in [T]$ ,*

$$\max_{\substack{\pi^{(1:T)} \text{ such that} \\ c(\pi^{(1:T)}) \leq C}} \sum_{t=1}^T V_1^{\pi^{(t)}, \mathcal{R}_t}(s_1) - \sum_{t=1}^T V_1^{\pi_t, \mathcal{R}_t}(s_1) \leq H(H+1) \max_{\substack{\tau_1, \dots, \tau_C \geq 0: \\ \tau_1 + \tau_2 + \dots + \tau_C = T}} \sum_{i=1}^C B_{T, A, \tau_i}.$$

In particular, if  $B_{T, A, \tau} \leq \ell(T, K) \sqrt{\tau}$ , where  $\ell(T, K)$  is logarithmic in  $T$  and  $K$ , which is a standard bound, then by Jensen's inequality for  $\sqrt{\cdot}$ ,

$$\max_{\substack{\tau_1, \dots, \tau_C \geq 0: \\ \tau_1 + \tau_2 + \dots + \tau_C = T}} \sum_{i=1}^C B_{T, A, \tau_i} \leq \ell(T, K) \max_{\substack{\tau_1, \dots, \tau_C \geq 0: \\ \tau_1 + \tau_2 + \dots + \tau_C = T}} \underbrace{\sum_{i=1}^C \sqrt{\tau_i}}_{\leq \sqrt{C(\tau_1 + \dots + \tau_C)}} = \ell(T, K) \sqrt{CT}.$$

## 6 The special case of exponential weights: improved regret bounds

**Contributions of this section.** *The literature (see Section 1.1) essentially focuses on the adversarial learning strategy given by exponential weights with a constant learning rate  $\eta$ . It turns out that this strategy does not satisfy the requirement of Definition 1 because of a tuning issue: the adversarial regret bound is of the form  $\ln N/\eta + \eta MT/2$  and cannot be simultaneously optimized for all values of  $T$ . The literature*

typically assumes that  $T$  is known and obtains a  $\sqrt{T}$  regret bound for MDPs by taking  $\eta$  of order  $1/\sqrt{T}$ ; see, for instance, among many others, Cai et al. (2020) and Shani et al. (2020). A notable exception, in the discounted setting and for a constant reward function, can be extracted from the proof of Agarwal et al. (2021, Section 5.3)—they handle convergence of the last iterate but their proof technique also applies to cumulative regret. We extend their result to the episodic setting and show that it is not essential that the reward functions be constant over time: we provide an upper bound in terms of the numbers of shifts in the sequence of reward functions.

We study in this section the strategy (6) of Section 3.2 where the adversarial learning strategies are given by the strategy (5) based on a constant exponential potential  $\Phi_t \equiv \Phi : x \mapsto \exp(\eta x)$ . This strategy takes the following simple form: for all  $t \geq 1$ ,

$$\forall h \in [H], \quad \forall s \in \mathcal{S}, \quad \forall a \in \mathcal{A}, \quad \pi_{t,h}(a|s) = \frac{\exp\left(\eta \sum_{\tau=1}^{t-1} A_h^{\pi_\tau, \mathcal{R}_\tau}(s, a)\right)}{\sum_{a' \in \mathcal{A}} \exp\left(\eta \sum_{\tau=1}^{t-1} A_h^{\pi_\tau, \mathcal{R}_\tau}(s, a')\right)}, \quad (11)$$

with the understanding that a sum over no term is null, i.e.,  $\pi_{1,h}(a|s) = 1/A$ .

Agarwal et al. (2021, Section 5.3) showed that the strategy above corresponds to the natural policy gradient strategy based on a softmax parametrization. They proposed a direct analysis (in the discounted setting) with reward functions constant over time. We adapt and extend this analysis to (obviously) adversarial sequences of reward functions. We also claim a more transparent proof scheme, consisting of a suitable adversarial bound (finer than the uniform bounds considered in Definition 1, which in this case would be linear in  $T$ , as recalled in the introduction of this section) applied to policy learning along the lines of the proof of Theorem 1.

Our result is stated in terms of the number  $R$  of regimes shifts in the sequence  $\mathcal{R}_1, \dots, \mathcal{R}_T$  of payoff functions. More formally,  $R$  is the smallest integer such that there exist  $R - 1$  integers  $\tau_2, \dots, \tau_R$  with values in  $[T]$  such that, denoting  $\tau_1 = 1$  and  $\tau_{R+1} = T + 1$ ,

$$\forall k \in \{2, \dots, R + 1\}, \quad \forall t \in \{\tau_{k-1}, \dots, \tau_k - 1\}, \quad \mathcal{R}_t = \mathcal{R}_{\tau_{k-1}}. \quad (12)$$

(The case  $R = 1$  corresponds to a single regime, i.e., the reward functions  $\mathcal{R}_t$  are independent of time.)

The proof of Theorem 4 below may be found in Appendix B. It is more complex than the proof by Agarwal et al. (2021, Section 5.3), which could use a simple argument specific to the discounted setting, with discount factor  $\gamma$ : that distributions over states induced by a starting state  $s_0$ , a policy, and a transition function, put a probability mass at least  $1 - \gamma$  on  $s_0$ , no matter the policy and the transition function. See Remark 7 for more details.

**Theorem 4.** *In the setting of Section 2 where rewards lie in  $[0, 1]$ , the policy learning strategy (11) controls the regret as*

$$\max_{\pi} \sum_{t=1}^T \left( V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) \right) \leq \frac{H \ln A}{\eta} + RH(H + 1),$$

where  $R$  is the number of regimes shifts in the sequence  $\mathcal{R}_1, \dots, \mathcal{R}_T$  of payoff functions.

The bound of Theorem 4 has a smaller order of magnitude than the one of Theorem 1, which is typically of order  $\sqrt{T}$ , as soon as the number of regime shifts satisfies  $R \ll \sqrt{T}$ . (In general, up to  $T - 1$  regime shifts may occur.) In particular, the regret upper bound of Theorem 4 is smaller than a constant when the reward functions do not vary over time.

By Lemma 2 and (the proof of) Theorem 2, we have the following corollary to Theorem 4, in case of a constant sequence of payoff functions. It corresponds to the bound of Agarwal et al. (2021, Section 5.3) with  $H$  playing the role of  $1/(1 - \gamma)$  therein.

**Corollary 2.** *In the setting of Section 2 where rewards lie in  $[0, 1]$ , if the reward functions do not vary over time and are all equal to some  $\mathcal{R}$ , then the last iterate of the policy learning strategy (11) satisfies*

$$\max_{\pi} V_1^{\pi, \mathcal{R}}(s_1) - V_1^{\pi_T, \mathcal{R}}(s_1) \leq \frac{H \ln A}{\eta T} + \frac{H(H+1)}{T}.$$

As in Agarwal et al. (2021, Section 5.3), the bounds obtained in Theorem 4 and Corollary 2 suggest choosing  $\eta$  as large as possible.

## 7 Practical versions with estimated advantage functions

**Contributions of this section.** *We review in greater detail how the literature resorted or should resort to adversarial learning strategies in practice: value functions are typically not observed and must be estimated.*

To implement the strategy (6), advantage functions of the form  $A_h^{\pi_t, \mathcal{R}_t}$  should be computed. The main issue in doing so is that the transition kernels  $\mathcal{T}$  are unknown; that the reward functions  $\mathcal{R}_t$  are fully revealed (full-information feedback) or not (bandit feedback, where only actual rewards are observed) at the end of an episode may be handled (see, among others, Shani et al., 2020). This is why the literature of policy optimization typically replaces the unknown  $A_h^{\pi_t, \mathcal{R}_t}$  by estimates  $\hat{A}_h^t$ , often based on a principle of optimism, and builds the policies of the form

$$\forall h \in [H], \quad \forall s \in \mathcal{S}, \quad \pi_{t,h}(\cdot | s) = \varphi_{t,h} \left( (\hat{A}_h^{\tau}(s, \cdot))_{\tau \leq t-1} \right),$$

where we used the same notation as in (6); see, among others, Abbasi-Yadkori et al. (2019), Shani et al. (2020), Cai et al. (2020), He et al. (2022), Zhao et al. (2023), Tiapkin et al. (2024). In all these references, the estimates of advantage functions satisfy the following facts, which we could interpret as a linear consistency of estimation. Other approaches, based on learning value functions, would not satisfy these properties, especially when  $\hat{V}_h^t(s)$  is defined as  $\max_{a' \in \mathcal{A}} \hat{Q}_h^t(s, a')$ .

**Assumption 1.** *The estimates  $\hat{A}_h^t$  are defined based on estimates  $\hat{Q}_h^t$  of  $Q$ -value functions: for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ ,*

$$\hat{V}_h^t(s) \stackrel{\text{def}}{=} \sum_{a' \in \mathcal{A}} \pi_{t,h}(a' | s) \hat{Q}_h^t(s, a') \quad \text{and} \quad \hat{A}_h^t(s) \stackrel{\text{def}}{=} \hat{Q}_h^t(s, a) - \hat{V}_h^t(s, a).$$

*In addition, the estimates  $\hat{A}_h^t$  are bounded, i.e.,  $|\hat{A}_h^t| \leq M_H$  for some quantity  $M_H$  (typically depending on  $H$ ).*

**Remark 3.** *The same comments as in Remark 1 apply: the policies could be built based on the  $\hat{Q}_h^{\tau}(s, \cdot)$  instead of the  $\hat{A}_h^t(s, \cdot)$ , but the latter are preferred empirically.*

When the sequential strategies  $\varphi_h$  control the regret in the adversarial setting (Definition 1) by  $B_{T,A}$  for  $A$ -dimensional reward vectors, the same argument as in (7), together with Assumption 1 for the equality to 0, shows that the strategy defined above satisfies: for all  $h \in [H]$  and  $s \in \mathcal{S}$ ,

$$\max_{a \in \mathcal{A}} \sum_{t=1}^T \hat{A}_h^t(s, a) - \sum_{t=1}^T \overbrace{\sum_{a \in \mathcal{A}} \pi_{t,h}(a | s) \hat{A}_h^t(s, a)}^{=0} \leq 2M_H B_{T,A},$$

thus, for all  $h \in [H]$  and  $s \in \mathcal{S}$ ,

$$\hat{R}_T \stackrel{\text{def}}{=} \max_{\pi} \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_h(a | s) \sum_{t=1}^T \hat{A}_h^t(s, a) \leq 2M_H B_{T,A}. \quad (13)$$

Specific arguments (see detail below) then relate the quantity above to the target quantity, stated as in (8):

$$R_T = \max_{\pi} \sum_{t=1}^T \left( V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) \right) = \max_{\pi} \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_h(a | s) \sum_{t=1}^T A_h^{\pi_t, \mathcal{R}_t}(s, a). \quad (14)$$

In general,  $\widehat{R}_T$  is not equal to a sum of differences of value functions. An exception is to be found in Tiapkin et al. (2024): they obtain the estimates  $\widehat{Q}_h^t$  as the  $Q$ -value functions corresponding to the policies  $\pi_t$ , to some reward functions  $\mathcal{R}'_t$  (based on the actual reward function  $\mathcal{R}_t$  revealed at the end of the episode plus some bonus function), and to some estimated transition kernels  $\widehat{\mathcal{T}}_t$  (that are constant over epochs).

We now provide two lines of arguments relating (13) and (14).

**Optimism.** The most popular approach is to build optimistic estimates  $\widehat{Q}_h^t$  of the true  $Q$ -value functions  $Q_h^{\pi_t, \mathcal{R}_t}$ , i.e., estimates that upper bound the true values with high probability. These optimistic estimates may or may not rely on structural assumptions (e.g., Cai et al., 2020 and He et al., 2022 assume some linear representation of the transition kernels). The total regret  $R_T$  is then typically decomposed into three terms, and  $\widehat{R}_T$  is one of these three terms: in Shani et al. (2020), term (ii); in Cai et al. (2020), term (i); in He et al. (2022), term  $I_1$ ; in Zhao et al. (2023), the “OMD regret term”; in Tiapkin et al. (2024), term (B). The two other terms are controlled in ways that are specific to each approach and setting, but our interesting observation is the systematic presence of the exact  $\widehat{R}_T$  term. The adversarial learning strategy  $(\varphi_h)_{h \in [H]}$  considered in all these references (but Tiapkin et al., 2024, which refers to the present article) is the exponential potential with a constant learning rate (see Section 6), possibly seen as an instance of online mirror descent.

**Simulator.** An alternative, but less natural, approach is to assume that some simulator is available, as in Agarwal et al. (2021, Section 6). Then, simulations at the end of each episode  $t$  may be performed to estimate the quantities  $Q_h^{\pi_t, \mathcal{R}_t}(s, a)$ , for each pair  $(s, a)$ . Applications of the Hoeffding-Azuma inequality then relate with high probability, for each pair  $(s, a)$ ,

$$\sum_{t=1}^T \widehat{Q}_h^t(s, a) \quad \text{and} \quad \sum_{t=1}^T Q_h^{\pi_t, \mathcal{R}_t}(s, a)$$

or, equivalently, the advantage functions. The regret  $R_T$  is then smaller, with high probability, than the upper bound on  $\widehat{R}_T$  plus some deviation bound of order  $\sqrt{T}$  up to logarithmic terms. We omit the immediate details.

## 8 Extension 3:

### Aggregation (orchestration) of expert policies, a.k.a. imitation learning

**Contributions of this section.** *Adversarial learning is sometimes called prediction with experts (see Cesa-Bianchi & Lugosi, 2006). We again push the logic of the reduction of the control of MDPs to adversarial learning and now rather aggregate expert policies. The aim is to mimic the performance of the overall best convex combination of expert policies (which is in particular better than the performance of the best policy taken in isolation). This setting was termed imitation learning; see, for instance, Cheng et al., 2020 and Liu et al., 2023. We obtain stronger forms of performance guarantees than in the latter references, see Remark 4. We do so via some reduction to the standard tabular case for a lifted MDP.*

We go back to the considerations of Section 2.1 and consider a finite number  $K$  of stationary policies. We denote by  $\Pi = \{\pi_1, \dots, \pi_K\}$  the set these policies and will refer to them as expert policies. We further denote by  $\Pi_h = \{\pi_{1,h}, \dots, \pi_{K,h}\}$  the policies corresponding to a given stage  $h \in [H]$ .

We combine expert policies over time through state-stage-dependent weights  $\mathbf{p}_t = (p_{t,h})_{h \in [H]} \in \mathcal{P}([K])^{[H] \times S}$ , where  $p_{t,h}(\cdot | s) \in \mathcal{P}([K])$  may be interpreted either as a probability distribution over the policies in  $\Pi_h$  or as providing convex weights for the aggregation of the policies in  $\Pi_h$ . More precisely, for each episode  $t \geq 1$ , we denote by  $\mathbf{p}_t \Pi = (p_{t,h} \Pi_h)_{h \in [H]}$  the stationary policy such that, for all stages  $h \in [H]$ ,

$$p_{t,h} \Pi_h : s \in \mathcal{S} \mapsto p_{t,h} \Pi_h(\cdot | s) = \sum_{k \in [K]} p_{t,h}(k | s) \pi_{k,h}(\cdot | s) \in \mathcal{P}(\mathcal{A}). \quad (15)$$

Picking an action  $a'$  according to  $p_{t,h} \Pi_h(\cdot | s)$  amounts to performing a two-stage randomization: first, drawing a policy index  $k' \sim p_{t,h}(\cdot | s)$ , then drawing  $a' \sim \pi_{k',h}(\cdot | s)$ . This remark is important in the cases

where it is difficult or computationally complex to explicitly write the  $\pi_{k,h}(\cdot|s)$  but where it is easy to simulate them.

As indicated above, the set of all possible state-stage-dependent weights  $\mathbf{q}$  corresponds to  $\mathcal{P}([K])^{[H] \times \mathcal{S}}$ . We consider the class  $\mathcal{C}(\Pi)$  of all possible policies defined according to the model above by using the same weights over time:

$$\mathcal{C}(\Pi) = \left\{ \mathbf{q}\Pi, \mathbf{q} \in \mathcal{P}([K])^{[H] \times \mathcal{S}} \right\},$$

and aim to learn a good policy in this class. To do so, the learning strategies pick weights  $\mathbf{p}_t \in \mathcal{P}([K])^{[H] \times \mathcal{S}}$  over time and output  $\boldsymbol{\pi}_t = \mathbf{p}_t \Pi$ . We will minimize the corresponding regret criterion:

$$\forall T \geq 1, \quad R_T^\Pi = \max_{\mathbf{q}} \sum_{t=1}^T \left( V_1^{\mathbf{q}\Pi, \mathbf{R}_t}(s_1) - V_1^{\mathbf{p}_t \Pi, \mathbf{R}_t}(s_1) \right).$$

**Remark 4.** *To the best of our understanding, Cheng et al. (2020) and Liu et al. (2023) consider a more restrictive setting with a constant reward function and in addition target a weaker notion of regret, corresponding to*

$$\max_{k \in [K]} V_1^{\boldsymbol{\delta}_k \Pi, \mathbf{R}}(s_1) - \max_{t \in [T]} V_1^{\mathbf{p}_t \Pi, \mathbf{R}}(s_1),$$

where each  $\boldsymbol{\delta}_k$  is a collection of state-stage-dependent weights that are all given by Dirac masses on expert  $k$ ; i.e.,  $V_1^{\boldsymbol{\delta}_k \Pi, \mathbf{R}} = V_1^{\pi_k, \mathbf{R}}$ .

Actually, the total regret  $R_T$  defined in Section 2.1 may be decomposed as some approximation error, i.e., how good the policies in  $\mathcal{C}(\Pi)$  are in terms of values, plus the regret with respect to  $\mathcal{C}(\Pi)$ :

$$R_T = \max_{\boldsymbol{\pi}} \sum_{t=1}^T \left( V_1^{\boldsymbol{\pi}, \mathbf{R}_t}(s_1) - V_1^{\boldsymbol{\pi}_t, \mathbf{R}_t}(s_1) \right) = \underbrace{\max_{\boldsymbol{\pi}} \sum_{t=1}^T V_1^{\boldsymbol{\pi}, \mathbf{R}_t}(s_1) - \max_{\mathbf{q}} \sum_{t=1}^T V_1^{\mathbf{q}\Pi, \mathbf{R}_t}(s_1)}_{\text{approximation error}} + R_T^\Pi.$$

In this section, we aim to control  $R_T^\Pi$  only and will basically assume that the approximation error is small due to a proper choice of  $\Pi$ . This situation should arise often, as explained by the following remark.

**Remark 5.** *Denote by  $\boldsymbol{\pi}^*$  a stationary policy achieving the maximum in the definition of  $R_T$ . Given that expert policies are combined through state-stage-dependent weights, the approximation error defined above is null as soon as*

$$\forall h \in [H], \forall s \in \mathcal{S}, \quad \exists q_h(\cdot|s) \in \mathcal{P}([K]) \quad \text{s.t.} \quad \pi_h^*(\cdot|s) = \sum_{k \in [K]} q_h(k|s) \pi_{k,h}(\cdot|s).$$

*In particular, it suffices that there exists  $k_{h,s} \in [K]$  such that  $\pi_h^*(\cdot|s) = \pi_{k_{h,s},h}(\cdot|s)$ . Put differently, it suffices that at each stage  $h \in [H]$  and for each state  $s \in \mathcal{S}$ , one of the expert policies (but not necessarily always the same) coincides with an optimal policy. This observation motivates the use of expert policies in the cases where finitely many easy-to-identify distributions are candidates to be optimal distributions for each given stage-state pair  $(h, s)$ .*

**Summary.** We provide in Box A a summary of the settings and aims considered, here in Section 8 and earlier in Section 2.1.

## 8.1 Equivalence between direct tabular learning and aggregation of expert policies

We now explain why any learning scheme minimizing the standard regret  $R_T$  induces a learning scheme minimizing the regret  $R_T^\Pi$  with respect to a finite set  $\Pi$  of expert policies, and vice versa. In a nutshell, the equivalence stems from considering the indexes  $k \in [K]$  of expert policies as meta-actions, i.e., actions in a sequences of lifted MDPs.

As a consequence, for the sake of clarity and completeness, we re-state the counterpart of our main result, Theorem 1, in the setting of policy orchestration: see Section 8.2.

## BOX A: POLICY OPTIMIZATION, POSSIBLY BASED ON EXPERT POLICIES

Direct tabular learning (Section 2.1)

Aggregation of expert policies (Section 8)

**MDP parameters:** state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , initial state  $s_1 \in \mathcal{S}$ ,  
transition kernels  $\mathcal{T}$

(No additional parameters)

Set  $\Pi$  of  $K$  expert policiesThe environment picks a sequence  $(\mathcal{R}_t)_{t \geq 1}$  of reward functions**For episodes**  $t = 1, 2, \dots$ :1. The initial state is set to  $s_{t,1} = s_1$ 2. **For stages**  $h = 1, \dots, H$ :(a) The learner picks a policy  $\pi_{t,h} : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ (b) and draws an action  $a_{t,h} \sim \pi_{t,h}(\cdot | s_{t,h})$ (a) The learner picks weights  $p_{t,h} \in \mathcal{P}([K])^{\mathcal{S}}$ ,(b) draws  $k_{t,h} \sim p_{t,h}(\cdot | s_{t,h})$ , the index of the expert policy,(c) and draws an action  $a_{t,h} \sim \pi_{k_{t,h},h}(\cdot | s_{t,h})$  according to expert policy  $k_{t,h}$ 4. The learner receives and observes a reward  $r_{t,h} \sim \mathcal{R}_{t,h}(s_{t,h}, a_{t,h})$ ,  
with conditional expectation  $r_{t,h}(s_{t,h}, a_{t,h})$ 5. If  $h \leq H - 1$ , the next state  $s_{t,h+1} \sim \mathcal{T}_h(\cdot | s_{t,h}, a_{t,h})$  is drawn**Goal:** Minimize the regret

$$R_T = \max_{\pi} \sum_{t=1}^T \left( V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) \right)$$

$$R_T^{\Pi} = \max_{\mathbf{q}} \sum_{t=1}^T \left( V_1^{\mathbf{q}^{\Pi}, \mathcal{R}_t}(s_1) - V_1^{\mathbf{p}_t^{\Pi}, \mathcal{R}_t}(s_1) \right)$$

**Direct tabular learning as aggregation of expert policies.** We set  $K = A$  and take as expert policies the Dirac masses on the arms; more precisely, for each  $a \in \mathcal{A}$ , and for all  $h \in [H]$  and  $s \in \mathcal{S}$ , we set  $\pi_{a,h}(\cdot | s) = \delta_a$ , the Dirac mass at  $a$ . This defines the expert policy  $\Delta_a$ . We consider

$$\Delta = \{\Delta_a : a \in \mathcal{A}\} \quad \text{and} \quad \mathcal{C}(\Delta) = \{\mathbf{p}\Delta, \mathbf{p} \in \mathcal{P}(\mathcal{A})^{[H] \times \mathcal{S}}\};$$

$\mathcal{C}(\Delta)$  is the set of all stationary policies, stated in their direct tabular form.

**From direct tabular learning to aggregation of expert policies.** Conversely, we note that aggregation of expert policies in  $\Pi$  amounts to performing direct tabular learning in the following sequence of (lifted) MDPs: the action space is  $\bar{\mathcal{A}} = [K]$ , the state space is  $\bar{\mathcal{S}} = \mathcal{S}$ , the transition kernels  $\bar{\mathcal{T}}$  and the reward functions  $\bar{\mathcal{R}}_t$  are defined, for all  $t \geq 1$  and  $h \in [H]$ , by

$$\begin{aligned} \bar{\mathcal{T}}_h : (s, k) \in \mathcal{S} \times [K] &\mapsto \sum_{a \in \mathcal{A}} \pi_{k,h}(a | s) \mathcal{T}_h(\cdot | s, a) \\ \text{and} \quad \bar{\mathcal{R}}_{t,h} : (s, k) \in \mathcal{S} \times [K] &\mapsto \sum_{a \in \mathcal{A}} \pi_{k,h}(a | s) \mathcal{R}_{t,h}(s, a). \end{aligned}$$



Direct tabular learning on the sequence of lifted MDPs defined above provides policies  $\bar{\pi}_t$  which correspond to the convex weights  $\mathbf{p}_t$  discussed above: for all  $t \geq 1$ ,  $h \in [H]$ , and  $s \in \mathcal{S}$ , we use  $p_{t,h}(\cdot|s) = \bar{\pi}_{t,h}(\cdot|s)$  to aggregated expert policies in the original MDP. Denoting by  $\bar{R}_T$  the regret suffered with direct tabular learning in the lifted MDP, we have:  $R_T^\Pi = \bar{R}_T$ .

**Remark 6.** In the final part of the proof of Theorem 1, we critically used that the transition kernels  $\mathcal{T}$  do not depend on time. The expression above for  $\bar{\mathcal{T}}$  is indeed independent on time, which would not be the case if the expert policies were evolving over time. This explains why we restricted our attention to constant expert policies.

## 8.2 Adversarial learning on advantage functions for aggregation of expert policies

The counterpart for imitation learning of the strategy defined in Section 3.2 is defined as follows, given the equivalence stated above.

For each stage  $h \in [H]$ , we fix a sequential strategy  $\varphi_h = (\varphi_{t,h})_{t \geq 1}$  in the adversarial setting, relying on reward vectors bounded by  $M_h = H - h + 1$  and of dimension  $K$ .

We run these strategies on the advantage functions of the lifted MDPs described above: for all  $t \geq 1$ ,  $h \in [H]$ , and  $s \in \mathcal{S}$ ,

$$\bar{A}_h^{\mathbf{p}_t, \bar{\mathcal{R}}_t}(s, \cdot) = \left( \bar{A}_h^{\mathbf{p}_t, \bar{\mathcal{R}}_t}(s, k) \right)_{k \in [K]}, \quad \text{where} \quad \bar{A}_h^{\mathbf{p}_t, \bar{\mathcal{R}}_t}(s, k) = \sum_{a \in \mathcal{A}} \pi_{k,h}(a|s) A_h^{\mathbf{p}_t, \bar{\mathcal{R}}_t}(s, a). \quad (16)$$

More precisely, we run the strategies  $(\varphi_h)_{h \in [H]}$  in the following stage-by-stage and state-by-state manner: for all  $t \geq 1$ ,

$$p_{t,h}(\cdot|s) = \varphi_{t,h} \left( \left( \bar{A}_h^{\mathbf{p}_\tau, \bar{\mathcal{R}}_\tau}(s, \cdot) \right)_{\tau \leq t-1} \right). \quad (17)$$

We refer to this strategy as  $(\varphi_h)_{h \in [H]}$ -**Adv2-Aggr**, for  $(\varphi_h)_{h \in [H]}$ -adversarial learning on advantage functions for aggregation of expert policies.

Theorem 1 immediately entails the following performance guarantee.

**Corollary 3.** In the setting of Section 8 where rewards lie in  $[0, 1]$ , if, for all  $h \in [H]$ , the sequential strategies  $\varphi_h$  control the regret in the adversarial setting (Definition 1) by  $B_{T,K}$  for  $K$ -dimensional reward vectors bounded by  $H - h + 1$ , then the  $(\varphi_h)_{h \in [H]}$ -**Adv2-Aggr** strategy defined in (17) over the set  $\Pi$  of  $K$  expert policies controls the regret with respect to  $\mathcal{C}(\Pi)$  as:

$$\forall T \geq 1, \quad R_T^\Pi = \max_q \sum_{t=1}^T \left( V_1^{q^\Pi, \bar{\mathcal{R}}_t}(s_1) - V_1^{\mathbf{p}_t^\Pi, \bar{\mathcal{R}}_t}(s_1) \right) \leq H(H+1) B_{T,K}.$$

## References

- Y. Abbasi-Yadkori, P. Bartlett, K. Bhatia, N. Lazic, C. Szepesvari, and G. Weisz. POLITEX: Regret bounds for policy iteration using expert prediction. In *Proceedings of the Thirty-Sixth International Conference on Machine Learning (ICML'19)*, volume 97 of PMLR, pp. 3692–3702, 2019.
- D. Adamskiy, W.K. Koolen, A. Chernov, and V. Vovk. A closer look at adaptive regret. *Journal of Machine Learning Research*, 17(23):1–21, 2016.
- A. Agarwal, S.M. Kakade, J.D. Lee, and G. Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76, 2021.
- P. Auer, N. Cesa-Bianchi, and C. Gentile. Adaptive and self-confident on-line learning algorithms. *Journal of Computer and System Sciences*, 64(1):48–75, 2002.
- Q. Cai, Z. Yang, C. Jin, and Z. Wang. Provably efficient exploration in policy optimization. In *Proceedings of the Thirty-Seventh International Conference on Machine Learning (ICML'20)*, volume 119 of PMLR, pp. 1283–1294, 2020.

- N. Cesa-Bianchi and G. Lugosi. Potential-based algorithms in on-line prediction and game theory. *Machine Learning*, 51:239–261, 2003.
- N. Cesa-Bianchi and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- C.-A. Cheng, A. Kolobov, and A. Agarwal. Policy improvement via imitation of multiple oracles. In *Advances in Neural Information Processing Systems (Neurips’20)*, volume 33, pp. 5587–5598, 2020.
- A. Daniely, A. Gonen, and S. Shalev-Shwartz. Strongly adaptive online learning. In *Proceedings of the Thirty-Second International Conference on Machine Learning (ICML’15)*, pp. 1405–1411, 2015.
- Steven de Rooij, Tim van Erven, Peter D. Grünwald, and Wouter M. Koolen. Follow the leader if you can, hedge if you must. *Journal of Machine Learning Research*, 15(37):1281–1316, 2014.
- E. Even-Dar, S.M. Kakade, and Y. Mansour. Online Markov decision processes. *Mathematics of Operations Research*, 34(3):726–736, 2009.
- Y. Freund and R.E. Schapire. Adaptive game playing using multiplicative weights. *Games and Economic Behavior*, 29(1):79–103, 1999.
- P. Gaillard, G. Stoltz, and T. van Erven. A second-order bound with excess losses. In *Proceedings of The Twenty-Seventh Conference on Learning Theory (COLT’14)*, volume 35 of *Proceedings of Machine Learning Research*, pp. 176–196, 2014.
- P. Gaillard, Y. Goude, L. Plagne, T. Dubois, and B. Thieurmél. *opera: Online Prediction by Expert Aggregation*, 2021. URL <https://CRAN.R-project.org/package=opera>. R package version 1.2.0.
- E. Hazan and C. Seshadhri. Efficient learning algorithms for changing environments. In *Proceedings of the Twenty-Sixth Annual International Conference on Machine Learning (ICML’09)*, pp. 393–400, 2009.
- J. He, D. Zhou, and Q. Gu. Near-optimal policy optimization algorithms for learning adversarial linear mixture MDPs. In *Proceedings of Twenty-Fifth International Conference on Artificial Intelligence and Statistics (AISTats’22)*, volume 151 of PMLR, pp. 4259–4280, 2022.
- M. Herbster and M.K. Warmuth. Tracking the best expert. *Machine Learning*, 32:151–178, 1998.
- S. Kakade. A natural policy gradient. In *Advances in Neural Information Processing Systems*, volume 14, pp. 1531–1538, 2001.
- S. Kakade and J. Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of the Nineteenth International Conference on Machine Learning (ICML’02)*, pp. 267–274, 2002.
- A. Kalai and S. Vempala. Efficient algorithms for online decision problems. *Journal of Computer and System Sciences*, 71(3):291–307, 2005.
- N. Littlestone and M.K. Warmuth. The weighted majority algorithm. *Information and Computation*, 108: 212–261, 1994.
- X. Liu, T. Yoneda, C. Wang, M. Walter, and Y. Chen. Active policy improvement from multiple black-box oracles. In *Proceedings of the Fourtieth International Conference on Machine Learning (ICML’23)*, volume 202 of PMLR, pp. 22320–22337, 2023.
- J. Müller and G. Montúfar. Geometry and convergence of natural policy gradient methods. *Information Geometry*, 7(1):485–523, 2024.
- A. Rosenberg and Y. Mansour. Online convex optimization in adversarial Markov decision processes. In *Proceedings of the Thirty-Sixth International Conference on Machine Learning (ICML’19)*, volume 97 of PMLR, pp. 5478–5486, 2019.
- J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *Proceedings of the Thirty-First International Conference on Machine Learning (ICML’15)*, pp. 1889–1897, 2015.

- J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms, 2017. Preprint, arXiv:1707.06347.
- L. Shani, Y. Efroni, A. Rosenberg, and S. Mannor. Optimistic policy optimization with bandit feedback. In *Proceedings of the Thirty-Seventh International Conference on Machine Learning (ICML'20)*, volume 119 of PMLR, pp. 8604–8613, 2020.
- R.S. Sutton, D.A. McAllester, S.P. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*, volume 13, pp. 1057–1063, 2000.
- D. Tiapkin, E. Chzhen, and G. Stoltz. Narrowing the gap between adversarial and stochastic MDPs via policy optimization, 2024. Preprint, arXiv:2407.05704.
- J.Y. Yu, S. Mannor, and N. Shimkin. Markov decision processes with arbitrary reward processes. *Mathematics of Operations Research*, 34(3):737–757, 2009.
- C. Zhao, R. Yang, B. Wang, X. Zhang, and S. Li. Learning adversarial low-rank Markov decision processes with unknown transition and full-information feedback. In *Advances in Neural Information Processing Systems (Neurips'23)*, volume 36, pp. 59107–59123, 2023.
- A. Zimin and G. Neu. Online learning in episodic Markovian decision processes by relative entropy policy search. In *Advances in Neural Information Processing Systems*, volume 26, pp. 1583–1591, 2013.
- M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the Twentieth International Conference on Machine Learning (ICML'03)*, pp. 928–936, 2003.

**Appendix.** The appendix provides proofs omitted from the main body of the article.

## A On monotonicity of weights for an adversarial-learning strategy

The proof of Lemma 2 (re-stated below) explains why the property of Definition 2 is termed monotonicity of weights, and why it is a natural property.

**Lemma 2.** *The potential-based strategies (5) of Cesa-Bianchi & Lugosi (2003) with constant, non-decreasing potential functions  $\Phi_t \equiv \Phi$  (like in Example 1) and the greedy projection algorithm (Example 3) of Zinkevich (2003) satisfy monotonicity of weights.*

*Proof.* We start with the potential-based strategies (5), in case of a constant, non-decreasing potential function  $\Phi_t \equiv \Phi$ , and use the notation defined therein. For each  $t \geq 1$ , since  $\Phi$  is non-decreasing, we have, for all  $k \in [K]$ ,

$$v_{t+1,k} \geq v_{t,k} \iff g_{t,k} - \sum_{j \in [K]} w_{t,j} g_{t,j} \geq 0, \quad \text{thus} \quad (v_{t+1,k} - v_{t,k}) \left( g_{t,k} - \sum_{j \in [K]} w_{t,j} g_{t,j} \right) \geq 0$$

in all cases. Therefore,

$$\sum_{k \in [K]} v_{t+1,k} \left( g_{t,k} - \sum_{j \in [K]} w_{t,j} g_{t,j} \right) \geq \sum_{k \in [K]} v_{t,k} \left( g_{t,k} - \sum_{j \in [K]} w_{t,j} g_{t,j} \right) = 0,$$

where the equality to 0 and the final result of Definition 2 are obtained, respectively, by normalizing the  $v_{t+1,k}$  and  $v_{t,k}$  into  $w_{t+1,k}$  and  $w_{t,k}$ .

The calculation above show that the monotonicity of weights is satisfied as soon as weights for components  $k$  associated with a good (respectively, bad) reward  $g_{t,k}$  in the previous round increase (respectively, decrease), where good or bad is determined by the sign of

$$g_{t,k} - \sum_{j \in [K]} w_{t,j} g_{t,j}.$$

This is why we termed this property monotonicity of weights. It looks like a natural property of an adversarial learning strategy.

For the greedy projection algorithm (Example 3) of Zinkevich (2003), we note that by a property of Euclidean projection onto a convex set (here,  $w_{t+1}$  is the projection of  $w_t + \eta_t g_t$  onto the simplex, and  $w_t$  also belongs to the simplex), the following inner product is non-positive:

$$0 \geq \langle w_t - w_{t+1}, (w_t + \eta_t g_t) - w_{t+1} \rangle = \|w_t - w_{t+1}\|^2 + \eta_t \langle w_t - w_{t+1}, g_t \rangle,$$

so that  $\langle w_{t+1} - w_t, g_t \rangle \geq 0$ , which is exactly monotonicity of weights.  $\square$

## B Proof of Theorem 4 (analysis of NPG with softmax parametrization)

For the convenience of the reader, we restate the result to be proved.

**Theorem 4.** *In the setting of Section 2 where rewards lie in  $[0, 1]$ , the policy learning strategy (11) controls the regret as*

$$\max_{\pi} \sum_{t=1}^T \left( V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) \right) \leq \frac{H \ln A}{\eta} + RH(H+1),$$

where  $R$  is the number of regimes shifts in the sequence  $\mathcal{R}_1, \dots, \mathcal{R}_T$  of payoff functions.

As indicated in Section 6, the proof below is based on the analysis of NPG with softmax parametrization proposed by Agarwal et al. (2021, Section 5.3) in the discounted setting with reward functions constant over time. See Remark 7 for an explanation of why the proof in the discounted is significantly simpler than the proof in the episodic setting.

We extend the proof of Agarwal et al. (2021, Section 5.3) to the episodic setting and to (obviously) adversarial sequences of reward functions. We also claim a more transparent proof scheme, consisting of an ad hoc adversarial bound (Lemma 3) which is then applied to policy learning along the lines of the proof of Theorem 1.

The first piece of the proof of Theorem 4 is to replace the uniform regret bounds considered in Definition 1 by some ad hoc, data-based, bound (of the same flavor as the bounds by de Rooij et al., 2014, Section 2 in terms of so-called mixability gaps). Indeed, the uniform regret bound that could be proved (see, e.g., Cesa-Bianchi & Lugosi, 2006, Theorem 2.2) for the adversarial strategy of Lemma 3 is  $B_{T,K} = \ln K / \eta + \eta T / 8$ , which is not sublinear.

**Lemma 3.** *The strategy (5) based on a constant exponential potential  $\Phi_t \equiv \Phi : x \mapsto \exp(\eta x)$ , i.e., picking weights*

$$\forall t \geq 1, \quad w_{t,k} = \frac{v_{t,k}}{\sum_{j \in [K]} v_{t,j}}, \quad \text{where} \quad v_{t,k} = \exp\left(\eta \sum_{\tau=1}^{t-1} g_{\tau,k}\right),$$

with the convention that  $v_{1,k} = 1$  and  $w_{1,k} = 1/K$ , satisfies the following bound: against all opponents sequentially picking reward vectors in  $\mathbb{R}^K$ ,

$$\forall T \geq 1, \quad \max_{k \in [K]} \sum_{t=1}^T g_{t,k} \leq \frac{\ln K}{\eta} + \sum_{t=1}^T \sum_{j \in [K]} w_{t+1,j} g_{t,j}.$$

This lemma is proved at the end of this section and we now apply it to prove Theorem 4.

*Proof of Theorem 4.* We adapt the proof of Theorem 1 by replacing (7) by the ad hoc bound stemming from Lemma 3. Namely,

$$\forall h \in [H], \quad \forall s \in \mathcal{S}, \quad \max_{a \in \mathcal{A}} \sum_{t=1}^T A_h^{\pi_t, \mathcal{R}_t}(s, a) \leq \frac{\ln A}{\eta} + \sum_{t=1}^T \sum_{a \in \mathcal{A}} \pi_{t+1, h}(a|s) A_h^{\pi_t, \mathcal{R}_t}(s, a). \quad (18)$$

We fix a comparator policy  $\pi$ . The combination of the obtained inequality (18) with the application (8) of the performance difference lemma yields

$$\begin{aligned} \sum_{t=1}^T \left( V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) \right) &\leq \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \max_{a \in \mathcal{A}} \sum_{t=1}^T A_h^{\pi_t, \mathcal{R}_t}(s, a) \\ &\leq \frac{H \ln A}{\eta} + \underbrace{\sum_{t=1}^T \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_{t+1, h}(a|s) A_h^{\pi_t, \mathcal{R}_t}(s, a)}_{\text{to be bounded}}. \end{aligned} \quad (19)$$

We fix  $t \in [T]$  and  $h \in [H]$  in what follows. We define a new one-shot policy  $\tilde{\pi}_{t+1}^h = (\tilde{\pi}_{t+1, h'}^h)_{h' \in [H]}$  as follows:

$$\tilde{\pi}_{t+1, h'}^h = \begin{cases} \pi_{h'} & \text{if } h' \leq h-1, \\ \pi_{t+1, h'} & \text{if } h' \geq h. \end{cases}$$

As  $\pi$  and  $\tilde{\pi}_{t+1}^h$  coincide in the first  $h-1$  stages, we have  $\mu_h^{s_1, \pi, \mathcal{T}}(s) = \mu_h^{s_1, \tilde{\pi}_{t+1}^h, \mathcal{T}}(s)$ . In addition, the definition of  $\tilde{\pi}_{t+1}^h$ , the definition of the strategy, Lemma 2, and the definition of advantage functions entail that for all  $s \in \mathcal{S}$  and all  $h' \geq h$ ,

$$\sum_{a \in \mathcal{A}} \tilde{\pi}_{t+1, h'}^h(a|s) A_{h'}^{\pi_t, \mathcal{R}_t}(s, a) = \sum_{a \in \mathcal{A}} \pi_{t+1, h'}(a|s) A_{h'}^{\pi_t, \mathcal{R}_t}(s, a) \geq \sum_{a \in \mathcal{A}} \pi_{t, h'}(a|s) A_{h'}^{\pi_t, \mathcal{R}_t}(s, a) = 0.$$

Therefore, the sum marked as “to be bounded” in (19) can be controlled as

$$\begin{aligned} \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_{t+1, h}(a|s) A_h^{\pi_t, \mathcal{R}_t}(s, a) &= \sum_{s \in \mathcal{S}} \mu_h^{s_1, \tilde{\pi}_{t+1}^h, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \tilde{\pi}_{t+1, h}^h(a|s) A_h^{\pi_t, \mathcal{R}_t}(s, a) \\ &\leq \sum_{h'=h}^H \sum_{s \in \mathcal{S}} \mu_{h'}^{s_1, \tilde{\pi}_{t+1}^h, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \tilde{\pi}_{t+1, h'}^h(a|s) A_{h'}^{\pi_t, \mathcal{R}_t}(s, a) \\ &= \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \left( V_h^{\tilde{\pi}_{t+1}^h, \mathcal{R}_t}(s) - V_h^{\pi_t, \mathcal{R}_t}(s) \right), \end{aligned} \quad (20)$$

where we used the performance difference lemma (Lemma 1) together with  $\mu_h^{s_1, \pi, \mathcal{T}}(s) = \mu_h^{s_1, \tilde{\pi}_{t+1}^h, \mathcal{T}}(s)$  for the final equality.

As  $\tilde{\pi}_{t+1}^h$  and  $\pi_{t+1}$  coincide in the last  $h$  stages, we have  $V_h^{\tilde{\pi}_{t+1}^h, \mathcal{R}_t}(s) = V_h^{\pi_{t+1}, \mathcal{R}_t}(s)$  for all  $s \in \mathcal{S}$ . This observation, combined with (20), entails

$$\sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_{t+1, h}(a|s) A_h^{\pi_t, \mathcal{R}_t}(s, a) \leq \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \left( V_h^{\pi_{t+1}, \mathcal{R}_t}(s) - V_h^{\pi_t, \mathcal{R}_t}(s) \right),$$

and we thus get, after substitution into (19),

$$\sum_{t=1}^T \left( V_1^{\pi, \mathcal{R}_t}(s_1) - V_1^{\pi_t, \mathcal{R}_t}(s_1) \right) \leq \frac{H \ln A}{\eta} + \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \sum_{t=1}^T \left( V_h^{\pi_{t+1}, \mathcal{R}_t}(s) - V_h^{\pi_t, \mathcal{R}_t}(s) \right). \quad (21)$$

We obtain telescoping sums on regimes of payoffs. More precisely, with the notation (12),

$$\forall k \in \{2, \dots, R+1\}, \quad \sum_{t=\tau_{k-1}}^{\tau_k-1} \left( V_h^{\pi_{t+1}, \mathcal{R}_t}(s) - V_h^{\pi_t, \mathcal{R}_t}(s) \right) = V_h^{\pi_{\tau_k}, \mathcal{R}_{\tau_{k-1}}}(s) - V_h^{\pi_{\tau_{k-1}}, \mathcal{R}_{\tau_{k-1}}}(s) \leq H - h + 1,$$

where the upper bound follows from the boundedness of rewards in  $[0, 1]$ . Together with (21), we finally obtain

$$\sum_{t=1}^T \left( V_1^{\pi_t, \mathcal{R}_t}(s_1) - V_1^{\pi_{t-1}, \mathcal{R}_t}(s_1) \right) \leq \frac{H \ln A}{\eta} + \sum_{h=1}^H \sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \sum_{k=2}^{R+1} (H - h + 1) = \frac{H \ln A}{\eta} + RH(H + 1),$$

which leads to the claimed regret upper bound after taking the maximum over all policies  $\pi$ .  $\square$

**Remark 7.** *The arguments between (19) and (21) may be bypassed in the discounted setting with discount factor  $\gamma$ ; see Agarwal et al. (2021, Section 5.3). Indeed (with obvious notation, for value functions defined in the standard way for discounted rewards, and for a constant reward function), for each  $s \in \mathcal{S}$ ,*

$$\begin{aligned} \max_{a \in \mathcal{A}} \sum_{t=1}^T A^{\pi_t}(s, a) &\leq \frac{\ln A}{\eta} + \sum_{t=1}^T \overbrace{\sum_{a \in \mathcal{A}} \pi_{t+1}(a|s) A^{\pi_t}(s, a)}^{\geq 0} \\ &\leq \frac{\ln A}{\eta} + \sum_{t=1}^T \underbrace{\frac{1}{1-\gamma} \sum_{s' \in \mathcal{S}} \mu^{s, \pi_{t+1}}(s') \sum_{a \in \mathcal{A}} \pi_{t+1}(a|s') A^{\pi_t}(s', a)}_{= V^{\pi_{t+1}}(s) - V^{\pi_t}(s)} = \frac{\ln A}{\eta} + \frac{V^{\pi_{T+1}}(s) - V^{\pi_1}(s)}{1-\gamma} \leq \frac{1}{(1-\gamma)^2}, \end{aligned}$$

where the first inequality is by Lemma 3, where the non-negativity is guaranteed by monotonicity of weights (see Lemma 2), where the second inequality comes from the fact that distributions induced by a starting state  $s$ , a given policy, and a given transition function put a probability mass at least  $1 - \gamma$  on  $s$ , no matter the policy and transition function (this is the property extremely specific to the discounted setting), where the equality to  $V^{\pi_{t+1}}(s) - V^{\pi_t}(s)$  is by the performance difference lemma, and where the final equality is by telescoping. The inequality obtained above is the key; the rest of the proof merely consists of yet another (now standard) application of the performance difference lemma:

$$\sum_{t=1}^T \left( V^{\pi}(s_1) - V^{\pi_t}(s_1) \right) = \sum_{t=1}^T \frac{1}{1-\gamma} \sum_{s \in \mathcal{S}} \mu^{s_1, \pi}(s) \sum_{a \in \mathcal{A}} \pi(a|s) A^{\pi_t}(s, a) \leq \frac{1}{1-\gamma} \sum_{s \in \mathcal{S}} \mu^{s_1, \pi}(s) \underbrace{\max_{a \in \mathcal{A}} \sum_{t=1}^T A^{\pi_t}(s, a)}_{\leq (\ln A)/\eta + 1/(1-\gamma)^2},$$

which is the bound claimed by Agarwal et al. (2021, Section 5.3).

We conclude this section with a proof of Lemma 3.

*Proof of Lemma 3.* First, a bound “à la Pisier” yields that for all sequences of payoffs  $g_{t,j}$ , possibly signed and unbounded:

$$\begin{aligned} \max_{k \in [K]} \sum_{t=1}^T g_{t,k} &= \frac{1}{\eta} \ln \left( \max_{j \in [K]} \exp \left( \eta \sum_{t=1}^T g_{t,j} \right) \right) \\ &\leq \frac{1}{\eta} \ln \left( \sum_{j \in [K]} \exp \left( \eta \sum_{t=1}^T g_{t,j} \right) \right) = \frac{\ln K}{\eta} + \frac{1}{\eta} \sum_{t=1}^T \ln \left( \sum_{j \in [K]} w_{t,j} \exp(\eta g_{t,j}) \right), \end{aligned}$$

where the equality follows by telescoping: indeed, by definition of the weights,

$$\sum_{j \in [K]} \underbrace{\exp \left( \eta \sum_{t=1}^T g_{t,j} \right)}_{= v_{T+1,j}} = K \prod_{t=1}^T \frac{\sum_{j \in [K]} v_{t+1,j}}{\sum_{j \in [K]} v_{t,j}} = K \prod_{t=1}^T \frac{\sum_{j \in [K]} v_{t,j} \exp(\eta g_{t,j})}{\sum_{j \in [K]} v_{t,j}} = K \prod_{t=1}^T w_{t,j} \exp(\eta g_{t,j}).$$

Second, by the application of Jensen's inequality to the convex function  $x \mapsto x \ln x$ ,

$$\left( \sum_{j \in [K]} w_{t,j} \exp(\eta g_{t,j}) \right) \ln \left( \sum_{j \in [K]} w_{t,j} \exp(\eta g_{t,j}) \right) \leq \sum_{j \in [K]} w_{t,j} \exp(\eta g_{t,j}) \ln(\exp(\eta g_{t,j})),$$

that is, after rearranging and given the definition of the weights  $w_{j,t+1}$ ,

$$\ln \left( \sum_{j \in [K]} w_{t,j} \exp(\eta g_{t,j}) \right) \leq \eta \sum_{j \in [K]} w_{t+1,j} g_{t,j}.$$

The claimed bound follows from combining the two inequalities obtained.  $\square$

## C Proof of the performance difference lemma

A word intended to reviewers: we would be happy to drop this section if required; there is a slight but straightforward generalization compared to earlier statements, lying in the fact that sums over  $h' \geq h$  are considered.

One of the first references stating the performance difference lemma (in the discounted setting) is Kakade & Langford (2002). Statements (possibly of generalizations) of this lemma for  $H$ -episodic MDPs are ubiquitous in the literature (see, e.g., Cai et al., 2020, Lemma 3.2 for a simple statement, and Shani et al., 2020, Lemma 1 for an extension to approximated advantage functions). We state yet another, straightforward, generalization, in terms of advantage and value functions starting at a given stage  $h$ ; this generalization is useful in the proof of Theorem 4 in Appendix B.

**Lemma 1** (Performance difference lemma). *Let  $\mu_{h'}^{s_1, \pi, \mathcal{T}}$  be the distribution of the state  $s_{h'}$  of the  $h'$ -th stage, starting from the state  $s_1$  in the first stage, following the stationary policy  $\pi$  and the transition kernels  $\mathcal{T}$ . In a MDP with transition kernels  $\mathcal{T}$ , for all pairs  $\pi, \pi'$  of stationary policies, for all reward functions  $\mathcal{R}$ , and for all stages  $h \in [H]$ ,*

$$\sum_{s \in \mathcal{S}} \mu_h^{s_1, \pi, \mathcal{T}}(s) \left( V_h^{\pi, \mathcal{R}}(s) - V_h^{\pi', \mathcal{R}}(s) \right) = \sum_{h'=h}^H \sum_{s \in \mathcal{S}} \mu_{h'}^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_{h'}(a|s) A_{h'}^{\pi', \mathcal{R}}(s, a).$$

In particular, for  $h = 1$ ,

$$V_1^{\pi, \mathcal{R}}(s_1) - V_1^{\pi', \mathcal{R}}(s_1) = \sum_{h'=1}^H \sum_{s \in \mathcal{S}} \mu_{h'}^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_{h'}(a|s) A_{h'}^{\pi', \mathcal{R}}(s, a).$$

*Proof.* We denote by  $\mathbb{P}^{s_1, \pi, \mathcal{T}}$  the probability distribution underlying the  $H$ -episodic MDP  $(s_1, a_1, \dots, s_H, a_H)$  starting at  $s_1$ , drawing actions according to  $\pi$ , and subject to the transition kernels  $\mathcal{T}$ . In particular, by definition, for any function  $f : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  and all  $h' \in [H]$ ,

$$\sum_{s \in \mathcal{S}} \mu_{h'}^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_{h'}(a|s) f(s, a) = \mathbb{E}^{s_1, \pi, \mathcal{T}}[f(s_{h'}, a_{h'})].$$

Letting successively  $f$  be  $A_{h'}^{\pi', \mathcal{R}}$  for  $h \leq h' \leq H$  and using the definition  $A_{h'}^{\pi', \mathcal{R}} = Q_{h'}^{\pi', \mathcal{R}} - V_{h'}^{\pi', \mathcal{R}}$ ,

$$\sum_{h'=h}^H \sum_{s \in \mathcal{S}} \mu_{h'}^{s_1, \pi, \mathcal{T}}(s) \sum_{a \in \mathcal{A}} \pi_{h'}(a|s) A_{h'}^{\pi', \mathcal{R}}(s, a) = \mathbb{E}^{s_1, \pi, \mathcal{T}} \left[ \sum_{h'=h}^H (Q_{h'}^{\pi', \mathcal{R}}(s_{h'}, a_{h'}) - V_{h'}^{\pi', \mathcal{R}}(s_{h'})) \right]. \quad (22)$$

Now, by definition of the  $Q$ -values, recalling that  $r$  denotes the mean-payoff functions associated with  $\mathcal{R}$ , we have, for  $h' \leq H - 1$ ,

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \quad Q_{h'}^{\pi', \mathcal{R}}(s, a) = r_{h'}(s, a) + \sum_{s' \in \mathcal{S}} \mathcal{T}_{h'}(s' | s, a) V_{h'+1}^{\pi', \mathcal{R}}(s'). \quad (23)$$

By definition of the MDP, for any function  $g : \mathcal{S} \rightarrow \mathbb{R}$ ,

$$\mathbb{E}^{s_1, \pi, \mathcal{T}} \left[ \sum_{s' \in \mathcal{S}} \mathcal{T}_{h'}(s' | s_{h'}, a_{h'}) g(s') \right] = \mathbb{E}^{s_1, \pi, \mathcal{T}} [g(s_{h'+1})] .$$

Thus, letting  $s = s_{h'}$  and  $a = a_{h'}$  in (23) and taking expectations yields

$$\mathbb{E}^{s_1, \pi, \mathcal{T}} [Q_{h'}^{\pi', \mathcal{R}}(s_{h'}, a_{h'})] = \mathbb{E}^{s_1, \pi, \mathcal{T}} [r_{h'}(s_{h'}, a_{h'})] + \mathbb{E}^{s_1, \pi, \mathcal{T}} [V_{h'+1}^{\pi', \mathcal{R}}(s_{h'+1})] .$$

For  $h' = H$ , we have  $Q_H^{\pi', \mathcal{R}}(s, a) = r_H(s, a)$ . As a consequence of the equalities above, a telescoping sum appears in the right-hand side of (22):

$$\begin{aligned} & \mathbb{E}^{s_1, \pi, \mathcal{T}} \left[ \sum_{h'=h}^H (Q_{h'}^{\pi', \mathcal{R}}(s_{h'}, a_{h'}) - V_{h'}^{\pi', \mathcal{R}}(s_{h'})) \right] \\ &= \mathbb{E}^{s_1, \pi, \mathcal{T}} \left[ r_{h'}(s_H, a_H) + \sum_{h'=h}^{H-1} (r_{h'}(s_{h'}, a_{h'}) + V_{h'+1}^{\pi', \mathcal{R}}(s_{h'+1})) - \sum_{h'=h}^H V_{h'}^{\pi', \mathcal{R}}(s_{h'}) \right] \\ &= \mathbb{E}^{s_1, \pi, \mathcal{T}} \left[ \sum_{h'=h}^H r_{h'}(s_{h'}, a_{h'}) \right] - \mathbb{E}^{s_1, \pi, \mathcal{T}} [V_1^{\pi', \mathcal{R}}(s_h)] . \end{aligned}$$

Finally, the tower rule shows that

$$\mathbb{E}^{s_1, \pi, \mathcal{T}} \left[ \sum_{h'=h}^H r_{h'}(s_{h'}, a_{h'}) \right] = \mathbb{E}^{s_1, \pi, \mathcal{T}} [V_1^{\pi, \mathcal{R}}(s_h)] .$$

The proof is concluded by collecting all the bounds. □