



**HAL**  
open science

# Online Markov Decision Processes Under Bandit Feedback

Gergely Neu, András György, Csaba Szepesvári, András Antos

► **To cite this version:**

Gergely Neu, András György, Csaba Szepesvári, András Antos. Online Markov Decision Processes Under Bandit Feedback. IEEE Transactions on Automatic Control, 2014, 59, pp.676 - 691. 10.1109/TAC.2013.2292137 . hal-01079422

**HAL Id: hal-01079422**

**<https://hal.science/hal-01079422v1>**

Submitted on 1 Nov 2014

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

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

# Online Markov Decision Processes under Bandit Feedback

Gergely Neu, András György, Csaba Szepesvári, András Antos

**Abstract**—We consider online learning in finite stochastic Markovian environments where in each time step a new reward function is chosen by an oblivious adversary. The goal of the learning agent is to compete with the best stationary policy in hindsight in terms of the total reward received. Specifically, in each time step the agent observes the current state and the reward associated with the last transition, however, the agent does not observe the rewards associated with other state-action pairs. The agent is assumed to know the transition probabilities. The state of the art result for this setting is an algorithm with an expected regret of  $O(T^{2/3} \ln T)$ . In this paper, assuming that stationary policies mix uniformly fast, we show that after  $T$  time steps, the expected regret of this algorithm (more precisely, a slightly modified version thereof) is  $O(T^{1/2} \ln T)$ , giving the first rigorously proven, essentially tight regret bound for the problem.

## I. INTRODUCTION

In this paper we consider online learning in finite stochastic Markovian environments where in each time step a new reward function may be chosen by an oblivious adversary. The interaction between the learner and the environment is shown in Figure 1. The environment is split into two parts: One part has a controlled Markovian dynamics, while another one has an unrestricted, uncontrolled (autonomous) dynamics. In each discrete time step  $t$ , the learning agent receives the state of the Markovian environment ( $\mathbf{x}_t \in \mathcal{X}$ ) and some information ( $y_{t-1} \in \mathcal{Y}$ ) about the previous state of the autonomous dynamics. The learner then makes a decision about the next action ( $\mathbf{a}_t \in \mathcal{A}$ ), which is sent to the environment. In response, the environment makes a transition: the next state  $\mathbf{x}_{t+1}$  of the Markovian part is drawn from a transition probability kernel  $P(\cdot|\mathbf{x}_t, \mathbf{a}_t)$ , while the other part makes a transition in an autonomous fashion. In the meanwhile, the agent incurs a reward  $\mathbf{r}_t = r(\mathbf{x}_t, \mathbf{a}_t, y_t) \in [0, 1]$  that depends on the *complete* state of the environment and the chosen action; then the process continues with the next step. The goal of the learner is to collect as much reward as possible. The agent knows the transition probability kernel  $P$  and the reward function  $r$ , however, he does not know the sequence  $y_t$  in advance. We call this problem *online learning in Markov Decision Processes (MDPs)*.

We take the viewpoint that the uncontrolled dynamics might be very complex and thus modeling it based on the available limited in-

This research was supported in part by the National Development Agency of Hungary from the Research and Technological Innovation Fund (KTIA-OTKA CNK 77782), the Alberta Innovates Technology Futures, and the Natural Sciences and Engineering Research Council (NSERC) of Canada. Parts of this work have been published at the Twenty-Fourth Annual Conference on Neural Information Processing Systems (NIPS 2010) [16].

G. Neu is with the Department of Computer Science and Information Theory, Budapest University of Technology and Economics, Budapest, Hungary, and with the Computer and Automation Research Institute of the Hungarian Academy of Sciences, Budapest, Hungary (email: neu.gergely@gmail.com).

A. György is with the Department of Computing Science, University of Alberta, Edmonton, Canada (email: gya@cs.bme.hu). During parts of this work he was with the Machine Learning Research Group of the Computer and Automation Research Institute of the Hungarian Academy of Sciences.

Cs. Szepesvári is with the Department of Computing Science, University of Alberta, Edmonton, Canada (email: szepesva@ualberta.ca).

A. Antos is with the Budapest University of Technology and Economics, Budapest, Hungary (email: antos@cs.bme.hu). During parts of this work he was with the Machine Learning Research Group of the Computer and Automation Research Institute of the Hungarian Academy of Sciences, Budapest, Hungary.

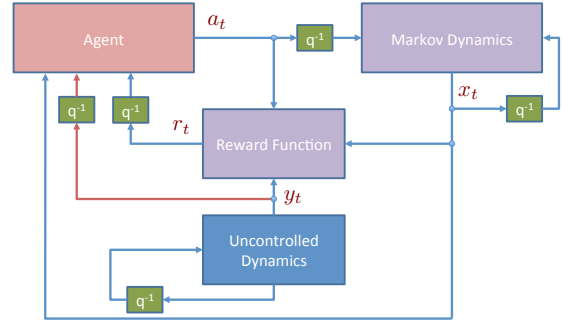


Fig. 1. The interaction between the learning agent and the environment. Here  $q^{-1}$  denotes a unit delay, that is, any information sent through such a box is received at the beginning of the next time step.

formation might be hopeless. Equivalently, we assume that whatever can be modeled about the environment is modeled in the Markovian, controlled part. As a result, when evaluating the performance of the learner, the total reward of the learner will be compared to that of the *best stochastic stationary policy in hindsight* that assigns actions to the states of the Markovian part in a random manner. This stationary policy is thus selected as the policy that maximizes the total reward given the sequence of reward functions  $r_t(\cdot, \cdot) \equiv r(\cdot, \cdot, y_t)$ ,  $t = 1, 2, \dots$ .<sup>1</sup> Given a horizon  $T > 0$ , any policy  $\pi$  and initial distribution uniquely determines a distribution over the sequence space  $(\mathcal{X} \times \mathcal{A})^T$ . Noting that the expected total reward of  $\pi$  is then a linear function of the distribution of  $\pi$  and that the space of distributions is a convex polytope with vertices corresponding to distributions of deterministic policies, we see that there will always exist a deterministic policy that maximizes the total expected reward in  $T$  time steps. Hence, it is enough to consider deterministic policies only as a reference. To make the objective more precise, for a given stationary deterministic policy  $\pi : \mathcal{X} \rightarrow \mathcal{A}$  let  $(\mathbf{x}_t^\pi, \mathbf{a}_t^\pi)$  denote the state-action pair that would have been visited in time step  $t$  had one used policy  $\pi$  from the beginning of time (the initial state being fixed). Then, the goal can be expressed as keeping the (*expected*) *regret*,

$$\hat{L}_T = \max_{\pi} \mathbb{E} \left[ \sum_{t=1}^T r_t(\mathbf{x}_t^\pi, \mathbf{a}_t^\pi) \right] - \mathbb{E} \left[ \sum_{t=1}^T \mathbf{r}_t \right]$$

small, *regardless of the sequence of reward functions*  $\{r_t\}_{t=1}^T$ . In particular, a sublinear regret-growth,  $\hat{L}_T = o(T)$  ( $T \rightarrow \infty$ ) means that the *average* reward collected by the learning agent approaches that of the best policy in hindsight. Naturally, a smaller growth-rate is more desirable.<sup>2</sup>

The motivation to study this problem is manifold. One viewpoint

<sup>1</sup>It is worth noting that the problem can be defined without referring to the uncontrolled, unmodelled dynamics by starting with an arbitrary sequence of reward functions  $\{r_t\}$ . That the two problems are equivalent follows because there is no restriction on the range of  $\{y_t\}$  or its dynamics.

<sup>2</sup>Following previous works in the area, in this paper we only consider regret relative to a fixed stationary policy. However, as usual in online learning, our results and algorithms can also be extended to less restricted sets of reference policies, such as the class of sequences of stationary policies with a restricted number of switches. We discuss such extensions in Section IV-D.

is that a learning agent achieving sublinear regret growth shows *robustness* in the face of arbitrarily assigned rewards, thus, the model provides a useful generalization of learning and acting in Markov Decision Processes. Some examples where the need for such robustness arises naturally are discussed below. Another viewpoint is that this problem is a useful generalization of online learning problems studied in the machine learning literature (e.g., [5]). In particular, in this literature, the problems studied are so-called prediction problems that involve an (oblivious) environment that chooses a sequence of loss functions. The learner’s predictions are elements in the common domain of these loss functions and the goal is to keep the regret small as compared with the best fixed prediction in hindsight. Identifying losses with negative rewards we may notice that this problem coincides exactly with our model with  $|\mathcal{X}|=1$ , that is, our problem is indeed a generalization of this problem where the reward functions have memory represented by multiple states subject to the Markovian control.

Let us now consider some examples that fit the above model. Generally, since our approach assumes that the hard-to-model, uncontrolled part influences the rewards only, the examples concern cases where the reward is difficult to model. This is the case, for example, in various production- and resource-allocation problems, where the major source of difficulty is to model the prices that influence the rewards. Indeed, the prices in these problems tend to depend on external, generally unobserved factors and thus dynamics of the prices might be hard to model. Other examples include problems coming from computer science, such as the  $k$ -server problem, paging problems, or web-optimization (e.g., ad-allocation problems with delayed information) [see, e.g., 7, 22].

Previous results that concern online learning in MDPs (with known transition probability kernels) are summarized in Table I. In

paper	algorithm	feedback	loops	regret bound
Even-Dar et al. [6, 7]	MDP-E	full information	yes	$\tilde{O}(T^{1/2})$
Yu et al. [22]	LAZY-FPL <sup>1</sup>	full information	yes	$\tilde{O}(T^{3/4+\epsilon})$ , $\epsilon > 0$
Yu et al. [22]	Q-FPL <sup>2</sup>	bandit	yes	$o(T)$
Neu et al. [13]		bandit	no	$O(T^{1/2})$
Neu et al. [16]	MDP-EXP3	bandit	yes	$\tilde{O}(T^{2/3})$
this paper	MDP-EXP3	bandit	yes	$\tilde{O}(T^{1/2})$

TABLE I

SUMMARY OF PREVIOUS RESULTS. PREVIOUS WORKS CONCERNED PROBLEMS WITH EITHER FULL-INFORMATION OR BANDIT FEEDBACK, PROBLEMS WHEN THE MDP DYNAMICS MAY OR MAY NOT HAVE LOOPS (TO BE MORE PRECISE, IN NEU ET AL. [13] WE CONSIDERED EPISODIC MDPs WITH RESTARTS). FOR EACH PAPER, THE ORDER OF THE OBTAINED REGRET BOUND IN TERMS OF THE TIME HORIZON  $T$  IS GIVEN.

<sup>1</sup>The Lazy-FPL algorithm has smaller computational complexity than MDP-E.

<sup>2</sup>The stochastic regret of Q-FPL was shown to be sublinear almost surely (not only in expectation).

the current paper we study the problem with recurrent Markovian dynamics while assuming that the *only information received about the uncontrolled part is in the form of the actual reward*  $\mathbf{r}_t$ . In particular, in our model the agent does not receive  $y_t$ , while in most previous works it was assumed that  $y_t$  is observed [6, 7, 22]. Following the terminology used in the online learning literature [2], when  $y_t$  is available (equivalently, the agent receives the reward function  $r_t : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$  in every time step), we say that learning happens under *full information*, while in our case we say that learning happens under *bandit feedback* (note that Even-Dar et al. [7] suggested as an open problem to address the bandit situation studied here). In an

earlier version of this paper [16], we provided an algorithm, MDP-EXP3, for learning in MDPs with recurrent dynamics under bandit feedback, and showed that it achieves a regret of order  $\tilde{O}(T^{2/3})$ .<sup>3</sup> In this paper we improve upon the analysis of [16] and prove an  $\tilde{O}(T^{1/2})$ -regret bound for the same algorithm. As it follows from a lower bound proven by Auer et al. [2] for bandit problems, apart from logarithmic and constant terms the rate obtained is unimprovable. The improvement compared to [16] is achieved by a more elaborate proof technique that builds on a (perhaps) novel observation that the so-called exponential weights technique (that our algorithm builds upon) changes its weights “slowly”. As in previous works where “loopy” Markovian dynamics were considered, our main assumptions on the MDP transition probability kernel will be that stationary policies mix uniformly fast. In addition, we shall assume that the stationary distributions of these policies are bounded away from zero. These assumptions will be discussed later.

We also mention here that Yu and Mannor [20, 21] considered the related problem of online learning in MDPs where the transition probabilities may also change arbitrarily after each transition. This problem is significantly more difficult than the case where only the reward function is allowed to change. Accordingly, the algorithms proposed in these papers do not achieve sublinear regret. Unfortunately, these papers have also gaps in the proofs, as discussed in detail in [13].

Finally, we note in passing that the contextual bandit problem considered by Lazaric and Munos [12] can also be regarded as a simplified version of our online learning problem where the states are generated in an i.i.d. fashion (though we do not consider the problem of competing with the best policy in a restricted subset of stationary policies). For regret bounds concerning learning in purely stochastic *unknown* MDPs, see the work of Jaksch et al. [10] and the references therein. Learning in adversarial MDPs without loops was also considered by György et al. [8] for deterministic transitions under bandit feedback, and under full information but with *unknown* transition probability kernels in our recent paper [14].

The rest of the paper is organized as follows: The problem is laid out in Section II, which is followed by a section that makes our assumptions precise (Section III). The algorithm and the main result are given and discussed in Section IV, with the proofs presented in Section V.

## II. NOTATION AND PROBLEM DEFINITION

The purpose of this section is to provide the formal definition of our problem and to set the goals. We start with some preliminaries, in particular by reviewing the language we use in connection to Markov Decision Processes (MDPs). This will be followed by the definition of the online learning problem. We assume that the reader is familiar with the concepts necessary to study MDPs, our purpose here is to introduce the notation only. For more background about MDPs, consult Puterman [17].

We define a finite Markov Decision Process (MDP)  $M$  by a finite state space  $\mathcal{X}$ , a finite action set  $\mathcal{A}$ , a transition probability kernel  $P : \mathcal{X} \times \mathcal{A} \times \mathcal{X} \rightarrow [0, 1]$ , and a reward function  $r : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$ . At time  $t \in \{1, 2, \dots\}$ , based on the sequence of past states, observed rewards, and actions,  $(\mathbf{x}_1, \mathbf{a}_1, r(\mathbf{x}_1, \mathbf{a}_1), \mathbf{x}_2, \dots, \mathbf{x}_{t-1}, \mathbf{a}_{t-1}, r(\mathbf{x}_{t-1}, \mathbf{a}_{t-1}), \mathbf{x}_t) \in (\mathcal{X} \times \mathcal{A} \times \mathbb{R})^{t-1} \times \mathcal{X}$ , an agent acting in the MDP  $M$  chooses an action  $\mathbf{a}_t \in \mathcal{A}$  to be executed.<sup>4</sup> As a result, the process moves to state  $\mathbf{x}_{t+1} \in \mathcal{X}$  with probability

<sup>3</sup>Here,  $\tilde{O}(g(s))$  denotes the class of functions  $f : \mathbb{N} \rightarrow \mathbb{R}^+$  satisfying  $\sup_{s \in \mathbb{N}} \frac{f(s)}{g(s) \ln^\alpha(g(s))} < \infty$  for some  $\alpha \geq 0$ .

<sup>4</sup>Throughout the paper we will use boldface letters to denote random variables.

$P(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{a}_t)$  and the agent incurs the reward  $r(\mathbf{x}_t, \mathbf{a}_t)$ . We note in passing that at the price of increased notational load, but with essentially no change to the contents, we could consider the case where the set of actions available at time step  $t$  is restricted to a non-empty subset  $\mathcal{A}(\mathbf{x}_t)$  of all actions, where the set-system,  $(\mathcal{A}(x))_{x \in \mathcal{X}}$ , is known to the agent. However, for simplicity, in the rest of the paper we stick to the case  $\mathcal{A}(x) = \mathcal{A}$ . In an MDP the goal of the agent is to maximize the long-term reward. In particular, in the so-called *average-reward problem*, the goal of the agent is to maximize the long-run average reward. In what follows, the symbols  $x, x', \dots$  will be reserved to denote a state in  $\mathcal{X}$ , while  $a, a', b$  will be reserved to denote an action in  $\mathcal{A}$ . In expressions involving sums over  $\mathcal{X}$ , the domain of  $x, x', \dots$  will be suppressed to avoid clutter. The same holds for sums involving actions.

Before defining the learning problem, let us introduce some more notation. We use  $\|v\|_p$  to denote the  $L_p$ -norm of a function or a vector. In particular, for  $p = \infty$  the supremum norm of a function  $v : S \rightarrow \mathbb{R}$  is defined as  $\|v\|_\infty = \sup_{s \in S} |v(s)|$ , and for  $1 \leq p < \infty$  and for any vector  $u = (u_1, \dots, u_d) \in \mathbb{R}^d$ ,  $\|u\|_p = \left(\sum_{i=1}^d |u_i|^p\right)^{1/p}$ . We use  $e_1, \dots, e_d$  to denote the *row* vectors of the canonical basis of the Euclidean space  $\mathbb{R}^d$ . Since we will identify  $\mathcal{X}$  with the integers  $\{1, \dots, |\mathcal{X}|\}$ , we will also use the notation  $e_x$  for  $x \in \mathcal{X}$ . We will use  $\ln$  to denote the natural logarithm function.

#### A. Online learning in MDPs

In this paper we consider a so-called *online learning problem* when the reward function is allowed to change arbitrarily in every time step. That is, instead of a single reward function  $r$ , a sequence of reward functions  $\{r_t\}$  is given. This sequence is assumed to be fixed ahead of time, and, for simplicity, we assume that  $r_t(x, a) \in [0, 1]$  for all  $(x, a) \in \mathcal{X} \times \mathcal{A}$  and  $t \in \{1, 2, \dots\}$ . No other assumptions are made about this sequence.

The learning agent is assumed to know the transition probabilities  $P$ , but is not given the sequence  $\{r_t\}$ . The protocol of interaction with the environment is unchanged: At time step  $t$  the agent selects an action  $\mathbf{a}_t$  based on the information available to it, which is sent to the environment. In response, the reward  $r_t(\mathbf{x}_t, \mathbf{a}_t)$  and the next state  $\mathbf{x}_{t+1}$  are communicated to the agent. The initial state  $\mathbf{x}_1$  is generated from a fixed distribution  $P_1$ , which may or may not be known.

Let the *expected total reward* collected by the agent up to time  $T$  be denoted by

$$\widehat{R}_T = \mathbb{E} \left[ \sum_{t=1}^T r_t(\mathbf{x}_t, \mathbf{a}_t) \right].$$

As before, the goal of the agent is to make this sum as large as possible. In classical approaches to learning one would assume some kind of regularity of  $r_t$  and then derive bounds on how much reward the learning agent loses as compared to the agent that knew about the regularity of the rewards and who acted optimally from the beginning of time. The loss or *regret*, measured in terms of the difference of total expected rewards of the two agents, quantifies the learner's efficiency. In this paper, following the recent trend in the machine learning literature [5], while keeping the regret criterion, we will avoid making any assumption on how the reward sequence is generated, and take a worst-case viewpoint. The potential benefit is that the results will be more generally applicable and the algorithms will enjoy added robustness, while, generalizing from results available for supervised learning [4, 11, 18], the algorithms can also be shown to avoid being too pessimistic.

The concept of regret in our case is defined as follows: We shall consider algorithms which are competitive with stochastic stationary policies. Fix a (stochastic) stationary policy  $\pi : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$  and

let  $\{(\mathbf{x}'_t, \mathbf{a}'_t)\}$  be the trajectory that results from following policy  $\pi$  from  $\mathbf{x}'_1 \sim P_1$  (in particular,  $\mathbf{a}'_t \sim \pi(\cdot|\mathbf{x}'_t) \stackrel{\text{def}}{=} \pi(\mathbf{x}'_t, \cdot)$ ). The *expected total reward* of  $\pi$  over the first  $T$  time steps is defined as

$$R_T^\pi = \mathbb{E} \left[ \sum_{t=1}^T r_t(\mathbf{x}'_t, \mathbf{a}'_t) \right].$$

Now, the (*expected*) *regret* (or *expected relative loss*) of the learning agent relative to the class of stationary policies is defined as

$$\widehat{L}_T = \sup_{\pi} R_T^\pi - \widehat{R}_T,$$

where the supremum is taken over all stochastic stationary policies in  $\mathcal{M}$ . Note that the policy maximizing the total expected reward is chosen in hindsight, that is, based on the knowledge of the reward functions  $r_1, \dots, r_T$ . Thus, the regret measures how well the learning agent is able to generalize from its moment to moment knowledge of the rewards to the sequence  $r_1, \dots, r_T$ . If the regret of an agent grows sublinearly with  $T$  then it can be said to act as well as the best (stochastic stationary) policy in the long run (i.e., the average expected reward of the agent in the limit is equal to that of the best policy). In this paper our main result will show that there exists an algorithm such that if that algorithm is followed by the learning agent, then the learning agent's regret will be bounded by  $C\sqrt{T} \ln T$ , where  $C > 0$  is a constant that depends on the transition probability kernel, but is independent of the sequence of rewards  $\{r_t\}$ .

### III. ASSUMPTIONS ON THE TRANSITION PROBABILITY KERNEL

Before describing our assumptions, a few more definitions are needed: First of all, for brevity, in what follows we will call stochastic stationary policies just policies. Further, without loss of generality, we shall identify the states with the first  $|\mathcal{X}|$  integers and assume that  $\mathcal{X} = \{1, 2, \dots, |\mathcal{X}|\}$ . Now, take a policy  $\pi$  and define the Markov kernel  $P^\pi(x'|x) = \sum_a \pi(a|x)P(x'|x, a)$ . The identification of  $\mathcal{X}$  with the first  $|\mathcal{X}|$  integers makes it possible to view  $P^\pi$  as a matrix:  $(P^\pi)_{x,x'} = P^\pi(x'|x)$ . In what follows, we will also take this view when convenient.

In general, distributions will also be treated as *row* vectors. Hence, for a distribution  $\mu$  over  $\mathcal{X}$ ,  $\mu P^\pi$  is the distribution over  $\mathcal{X}$  that results from using policy  $\pi$  for one step after a state is sampled from  $\mu$  (i.e., the “next-state distribution” under  $\pi$ ). Finally, a stationary distribution of a policy  $\pi$  is a distribution  $\mu_{\text{st}}$  that satisfies  $\mu_{\text{st}} P^\pi = \mu_{\text{st}}$ .

In what follows we assume that every (stochastic stationary) policy  $\pi$  has a well-defined unique stationary distribution  $\mu_{\text{st}}^\pi$ . This ensures that the average reward underlying any stationary policy is a well-defined single real number. It is well-known that in this case the convergence to the stationary distribution is exponentially fast. Following Even-Dar et al. [7], we consider the following stronger, “uniform mixing condition” (which implies the existence of the unique stationary distributions):

**Assumption A1:** There exists a number  $\tau \geq 0$  such that for any policy  $\pi$  and any pair of distributions  $\mu$  and  $\mu'$  over  $\mathcal{X}$ ,

$$\|(\mu - \mu')P^\pi\|_1 \leq e^{-1/\tau} \|\mu - \mu'\|_1. \quad (1)$$

As Even-Dar et al. [7], we call the smallest  $\tau$  satisfying this assumption the *mixing time* of the transition probability kernel  $P$ . Together with the existence and uniqueness of the stationary policy, the next assumption ensures that every state is visited eventually no matter what policy is chosen:

**Assumption A2:** The stationary distributions are uniformly bounded away from zero:

$$\inf_{\pi, x} \mu_{\text{st}}^\pi(x) \geq \beta > 0$$

for some  $\beta \in \mathbb{R}$ .

Note that  $e^{-1/\tau}$  is the supremum over all policy  $\pi$  of the Markov-Dobrushin coefficient of ergodicity, defined as  $m_{P^\pi} = \sup_{\mu \neq \mu'} \frac{\|(\mu - \mu')P^\pi\|_1}{\|\mu - \mu'\|_1}$  for the transition probability kernel  $P^\pi$ , see, e.g., [9]. It is also known that  $m_{P^\pi} = 1 - \min_{x, x' \in \mathcal{X}} \sum_{y \in \mathcal{X}} \min\{P^\pi(y|x), P^\pi(y|x')\}$  [9]. Since  $m_{P^\pi}$  is a continuous function of  $\pi$  and the set of policies is compact, there is a policy  $\pi'$  with  $m_{P^{\pi'}} = \sup_{\pi} m_{P^\pi}$ . These facts imply that Assumption A1 is satisfied, that is,  $\sup_{\pi} m_{P^\pi} < 1$ , if and only if for every  $\pi$ ,  $m_{P^\pi} < 1$ , that is,  $P^\pi$  is a *scrambling* matrix ( $P^\pi$  is a scrambling matrix if any two rows of  $P^\pi$  share some column in which they both have a positive element). Furthermore, if  $P^\pi$  is a scrambling matrix for any deterministic policy  $\pi$  then it is also a scrambling matrix for any stochastic policy. Thus, to guarantee Assumption A1 it is enough to verify mixing for deterministic policies only. The assumptions will be further discussed in Section IV-D.

#### IV. LEARNING IN ONLINE MDPs UNDER BANDIT FEEDBACK

In this section we shall first introduce some additional, standard MDP concepts that we will need. That these concepts are well-defined follows from our assumptions on  $P$  and from standard results to be found, for example, in the book by Puterman [17]. After the definitions, we specify our algorithm. The section is finished by the statement of our main result concerning the performance of the proposed algorithm.

##### A. Preliminaries

Fix an arbitrary policy  $\pi$  and  $t \geq 1$ . Let  $\{(\mathbf{x}'_s, \mathbf{a}'_s)\}$  be a random trajectory generated by  $\pi$  and the transition probability kernel  $P$  and an arbitrary everywhere positive initial distribution over the states. We will use  $q_t^\pi$  to denote the *action-value function* underlying  $\pi$  and the immediate reward  $r_t$ , while we will use  $v_t^\pi$  to denote the corresponding (*state*) *value function*.<sup>5</sup> That is, for  $(x, a) \in \mathcal{X} \times \mathcal{A}$ ,

$$q_t^\pi(x, a) = \mathbb{E} \left[ \sum_{s=1}^{\infty} (r_t(\mathbf{x}'_s, \mathbf{a}'_s) - \rho_t^\pi) \middle| \mathbf{x}'_1 = x, \mathbf{a}'_1 = a \right],$$

$$v_t^\pi(x) = \mathbb{E} \left[ \sum_{s=1}^{\infty} (r_t(\mathbf{x}'_s, \mathbf{a}'_s) - \rho_t^\pi) \middle| \mathbf{x}'_1 = x \right],$$

where  $\rho_t^\pi$  is the *average reward per stage* corresponding to  $\pi$ :

$$\rho_t^\pi = \lim_{S \rightarrow \infty} \frac{1}{S} \sum_{s=1}^S \mathbb{E}[r_t(\mathbf{x}'_s, \mathbf{a}'_s)].$$

The average reward per stage can be expressed as

$$\rho_t^\pi = \sum_x \mu_{\text{st}}^\pi(x) \sum_a \pi(a|x) r_t(x, a),$$

where  $\mu_{\text{st}}^\pi$  is the stationary distribution underlying policy  $\pi$ . Under our assumptions stated in the previous section, up to a shift by a constant function, the value functions  $q_t^\pi, v_t^\pi$  are the unique solutions to the *Bellman equations*

$$q_t^\pi(x, a) = r_t(x, a) - \rho_t^\pi + \sum_{x'} P(x'|x, a) v_t^\pi(x'),$$

$$v_t^\pi(x) = \sum_a \pi(a|x) q_t^\pi(x, a),$$
(2)

which hold simultaneously for all  $(x, a) \in \mathcal{X} \times \mathcal{A}$  (Corollary 8.2.7 of [17]). We will use  $q_t^*$  to denote the *optimal action-value function*,

<sup>5</sup>Most sources would call these functions *differential* action- and state-value functions. We omit this adjective for brevity.

that is, the action-value function underlying a policy that maximizes the average-reward in the MDP specified by  $(P, r_t)$ . We will also need these concepts for an arbitrary reward function  $r: \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ . In such a case, we will use  $v^\pi, q^\pi$ , and  $\rho^\pi$  to denote the respective value function, action-value function, and average reward of a policy  $\pi$ .

Now, consider the trajectory  $\{(\mathbf{x}_t, \mathbf{a}_t)\}$  followed by a learning agent with  $\mathbf{x}_1 \sim P_1$ . For any  $t \geq 1$ , define

$$\mathbf{u}_t = (\mathbf{x}_1, \mathbf{a}_1, r_1(\mathbf{x}_1, \mathbf{a}_1), \dots, \mathbf{x}_t, \mathbf{a}_t, r_t(\mathbf{x}_t, \mathbf{a}_t)) \quad (3)$$

and introduce the policy followed in time step  $t$ ,  $\pi_t(a|x) = \mathbb{P}[\mathbf{a}_t = a | \mathbf{u}_{t-1}, \mathbf{x}_t = x]$ , where  $\mathbf{u}_0$  and, more generally  $\mathbf{u}_s$  for all  $s \leq 0$  is defined to be the empty sequence. Note that  $\pi_t$  is computed based on past information and is therefore random. We introduce the following notation:

$$\mathbf{q}_t = q_t^{\pi_t}, \quad \mathbf{v}_t = v_t^{\pi_t}, \quad \boldsymbol{\rho}_t = \rho_t^{\pi_t}.$$

With this, we see that the following equations hold simultaneously for all  $(x, a) \in \mathcal{X} \times \mathcal{A}$ :

$$\mathbf{q}_t(x, a) = r_t(x, a) - \boldsymbol{\rho}_t + \sum_{x'} P(x'|x, a) \mathbf{v}_t(x'),$$

$$\mathbf{v}_t(x) = \sum_a \pi_t(a|x) \mathbf{q}_t(x, a).$$
(4)

##### B. The algorithm

Our algorithm, MDP-EXP3, shown as Algorithm 1, is inspired by that of Even-Dar et al. [7], while also borrowing ideas from the EXP3 algorithm (exponential weights algorithm for exploration and exploitation) of Auer et al. [2]. The main idea of the algorithm is to

---

#### Algorithm 1 MDP-EXP3: an algorithm for online learning in MDPs

---

Set  $N \geq 1$ ,  $\mathbf{w}_1(x, a) = \mathbf{w}_2(x, a) = \dots = \mathbf{w}_{2N-1}(x, a) = 1$ ,  
 $\gamma \in (0, 1)$ ,  $\eta \in (0, \gamma]$ .  
 For  $t = 1, 2, \dots$  repeat:

1) Set

$$\pi_t(a|x) = (1 - \gamma) \frac{\mathbf{w}_t(x, a)}{\sum_b \mathbf{w}_t(x, b)} + \frac{\gamma}{|\mathcal{A}|}$$

for all  $(x, a) \in \mathcal{X} \times \mathcal{A}$ .

2) Draw an action  $\mathbf{a}_t \sim \pi_t(\cdot | \mathbf{x}_t)$ .

3) Receive reward  $r_t(\mathbf{x}_t, \mathbf{a}_t)$  and observe  $\mathbf{x}_{t+1}$ .

4) If  $t \geq N$

a) Compute  $\boldsymbol{\mu}_t^N$  for all  $x \in \mathcal{X}$  using (8).

b) Construct estimates  $\hat{\mathbf{r}}_t$  using (6) and compute  $\hat{\mathbf{q}}_t$  using (5).

c) Set  $\mathbf{w}_{t+N}(x, a) = \mathbf{w}_{t+N-1}(x, a) e^{\eta \hat{\mathbf{q}}_t(x, a)}$  for all  $(x, a) \in \mathcal{X} \times \mathcal{A}$ .

---

construct estimates  $\{\hat{\mathbf{q}}_t\}$  of the action-value functions  $\{\mathbf{q}_t\}$ , which are then used to determine the action-selection probabilities  $\pi_t(\cdot | x)$  in each state  $x$  in each time step  $t$ . In particular, the probability of selecting action  $a$  in state  $x$  at time step  $t$  is computed as the mixture of the uniform distribution (which encourages exploring actions irrespective of what the algorithm has learned about the action-values) and a Gibbs distribution, the mixture parameter being  $\gamma > 0$ . Given a state  $x$ , the Gibbs distribution defines the probability of choosing action  $a$  at time step  $t$  to be proportional to  $\exp(\eta \sum_{s=N}^{t-N} \hat{\mathbf{q}}_s(x, a))$ .<sup>6</sup>

<sup>6</sup>In the algorithm the Gibbs action-selection probabilities are computed in an incremental fashion with the help of the “weights”  $\mathbf{w}_t(x, a)$ . Note that a numerically stable implementation would calculate the action-selection probabilities based on the relative value differences,  $\sum_{s=N}^{t-N} \hat{\mathbf{q}}_s(x, \cdot) - \max_{a \in \mathcal{A}} \sum_{s=N}^{t-N} \hat{\mathbf{q}}_s(x, a)$ . These relative value differences can also be updated incrementally. The form shown in Algorithm 1 is preferred for mathematical clarity.

Here,  $\eta > 0$ ,  $N > 0$  are further parameters of the algorithm. Note that for the single-state setting with  $N = 1$ , MDP-EXP3 is equivalent to the EXP3 algorithm of Auer et al. [2].

It is interesting to discuss how the Gibbs policy (i.e.,  $\frac{\mathbf{w}_t(x, \cdot)}{\sum_b \mathbf{w}_t(x, b)}$ ) is related to what is known as the Boltzmann-exploration policy in the reinforcement learning literature [e.g., 19]. Remember that given a state  $x$ , the Boltzmann-exploration policy would select action  $a$  at time step  $t$  with probability proportional to  $\exp(\eta \hat{\mathbf{q}}_{t-1}^*(x, a))$  for some estimate  $\hat{\mathbf{q}}_{t-1}^*$  of the optimal action-value function in the MDP  $(P, \hat{\mathbf{r}}_{t-1})$ , where  $\{\hat{\mathbf{r}}_t\}$  is the sequence of estimated reward functions. Thus, we can see a couple of differences between the Boltzmann exploration and our Gibbs policy. The first difference is that the Gibbs policy in our algorithm uses the cumulated sum of the estimates of action-values, while the Boltzmann policy uses only the last estimate. By depending on the sum, the Gibbs policy will rely less on the last estimate. This reduces how fast the policies can change, making the learning “smoother”. Another difference is that in our Gibbs policy the sum of previous action-values runs only up to step  $t-N$  instead of using the sum that runs up to the last step  $t-1$ . The reasons for doing this will be explained below. Finally, the Gibbs policy uses the action-value function estimates (in the MDPs  $\{(P, \hat{\mathbf{r}}_s)\}$ ) of the policies  $\{\pi_s\}$  selected by the algorithm, as opposed to using an estimate of the optimal action-value function. This makes our algorithm closer in spirit to (modified) policy iteration than to value iteration and is again expected to reduce the variance of the learning process.

The reason the Gibbs policy does not use the last  $N$  estimates is to allow the construction of a reasonable estimate  $\hat{\mathbf{q}}_t$  of the action-value function  $\mathbf{q}_t$ . If  $r_t$  was available, one could compute  $\mathbf{q}_t$  based on  $r_t$  (cf. (4)) and the sum could then run up to  $t-1$ , resulting in the algorithm of Even-Dar et al. [7]. Since in our problem  $r_t$  is not available, we estimate it using an importance sampling estimator  $\hat{r}_t$  below (from now on,  $t \geq N$ ). Given this  $\hat{r}_t$ , the estimate  $\hat{\mathbf{q}}_t$  of the action-value function  $\mathbf{q}_t$  is defined as the action-value function underlying policy  $\pi_t$  in the average-reward MDP given by the transition probability kernel  $P$  and reward function  $\hat{\mathbf{r}}_t$ . Thus,  $\hat{\mathbf{q}}_t$ , up to a shift by a constant function, can be computed as the solution to the Bellman equations corresponding to  $(P, \hat{\mathbf{r}}_t)$  (cf. (4)):

$$\begin{aligned} \hat{\mathbf{q}}_t(x, a) &= \hat{r}_t(x, a) - \hat{\rho}_t + \sum_{x'} P(x'|x, a) \hat{v}_t(x'), \\ \hat{v}_t(x) &= \sum_{a'} \pi_t(a'|x) \hat{\mathbf{q}}_t(x, a'), \\ \hat{\rho}_t &= \sum_{x', a'} \mu_{\text{st}}^{\pi_t}(x') \pi_t(a'|x') \hat{r}_t(x', a'), \end{aligned} \quad (5)$$

which hold simultaneously for all  $(x, a) \in \mathcal{X} \times \mathcal{A}$ . Since  $\pi_t$  is invariant to constant shifts of  $\hat{\mathbf{q}}_t$ , any of the solutions of these equations leads to the same sequence of policies. Hence, in what follows, without loss of generality we assume that the algorithm uses  $\hat{\mathbf{q}}_t$ , i.e., the value function of  $\pi_t$  in the average-reward MDP defined by  $(P, \hat{\mathbf{r}}_t)$ .

To define the estimator  $\hat{r}_t$  define  $\mu_t^N(x)$  as the probability of visiting state  $x$  at time step  $t$ , conditioned on the history  $\mathbf{u}_{t-N}$  up to time step  $t-N$ , including  $\mathbf{x}_{t-N}$  and  $\mathbf{a}_{t-N}$  (cf. (3) for the definition of  $\{\mathbf{u}_t\}$ ):

$$\mu_t^N(x) \stackrel{\text{def}}{=} \mathbb{P}[\mathbf{x}_t = x | \mathbf{u}_{t-N}], \quad x \in \mathcal{X}.$$

Then, the estimate of  $r_t$  is constructed using

$$\hat{r}_t(x, a) = \begin{cases} \frac{r_t(x, a)}{\pi_t(a|x) \mu_t^N(x)}, & \text{if } (x, a) = (\mathbf{x}_t, \mathbf{a}_t); \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

The importance sampling estimator (6) is well-defined only if for

$x = \mathbf{x}_t$ ,

$$\mu_t^N(x) > 0 \quad (7)$$

holds almost surely (by construction  $\pi_t(\cdot | \mathbf{x}_t) > \gamma / |\mathcal{A}| > 0$ ). To see the intuitive reason of why (7) holds, it is instructive to look into how the distribution  $\mu_t^N$  can be computed.

When  $t = N$ , it should be clear from the definition of  $\mu_t^N$  that, viewing  $\mu_t^N$  as a row vector,  $\mu_N^N = P_1(P^{\pi_1})^{N-1}$ . Now let  $t > N$ . Denote by  $P^a$  the transition probability matrix of the policy that selects action  $a$  in every state and recall that  $e_x$  denotes the  $x^{\text{th}}$  unit row vector of the canonical basis of the  $|\mathcal{X}|$ -dimensional Euclidean space. We may write

$$\mu_t^N = e_{\mathbf{x}_{t-N}} P^{\mathbf{a}_{t-N}} P^{\pi_{t-N+1}} \dots P^{\pi_{t-1}}, \quad t > N. \quad (8)$$

This holds because for any  $t \geq N$ ,  $\pi_t$  is entirely determined by the history  $\mathbf{u}_{t-N}$ , while for  $t > N$  the history  $\mathbf{u}_{t-N}$  also includes (and thus determines)  $\mathbf{x}_{t-N}$ ,  $\mathbf{a}_{t-N}$ . Using the notation  $\mathbf{z} \in \sigma(\mathbf{u}_{t-N})$  to denote that the random variable  $\mathbf{z}$  is measurable with respect to the sigma-algebra generated by the history  $\mathbf{u}_{t-N}$ , the above fact can be stated as

$$\begin{aligned} \mathbf{x}_{t-N}, \mathbf{a}_{t-N} &\in \sigma(\mathbf{u}_{t-N}) \text{ for } t > N, \\ \pi_t &\in \sigma(\mathbf{u}_{t-N}) \text{ for } t \geq N. \end{aligned} \quad (9)$$

Consequently, we also have that  $\pi_{t-1}, \dots, \pi_{t-N+1} \in \sigma(\mathbf{u}_{t-N})$  and therefore (8) follows from the law of total probability. Note also that

$$\mathbb{P}[\mathbf{a}_t = a | \mathbf{x}_t = x, \mathbf{u}_{t-N}] = \mathbb{P}[\mathbf{a}_t = a | \mathbf{x}_t = x, \mathbf{u}_{t-1}] = \pi_t(a|x), \quad (10)$$

where the last equality follows from the definition of  $\pi_t$  and  $\mathbf{a}_t$ .

The algorithm as presented needs to know  $P_1$  to compute  $\mu_t^N$  at step  $t = N$ . When  $P_1$  is unknown, instead of starting the computation of the weights at time step  $t = N$ , we can start the computation at time step  $t = N + 1$  (i.e., change  $t \geq N$  of step 4 to  $t \geq N + 1$ ). Clearly, in the worst-case, the regret can only increase by a constant amount (the magnitude of the largest reward) as a result of this change.

An essential step of the proof of our main result is to show that inequality (7) indeed holds, that is,  $\mu_t^N(x)$  is bounded away from zero. In fact, we will show that this inequality holds almost surely<sup>7</sup> for all  $x \in \mathcal{X}$  provided that  $N$  is large enough, which explains why the sum in the definition of the Gibbs policy runs from time  $N$ . This will be done by first showing that the policies  $\pi_t$  (especially, during the last  $N - 1$  steps) change “sufficiently slowly” (this is where it becomes useful that the Gibbs policy is defined using a sum of previous action values). Consequently,  $\pi_{t-N+1}, \dots, \pi_{t-1}$  will all be “quite close” to the policy of the last time step. Then, the expression on the right-hand side of (8) can be seen to be close to the  $N - 1$ -step state distribution of  $\pi_t$  when starting from  $(\mathbf{x}_{t-N}, \mathbf{a}_{t-N})$ , which, if  $N$  is large enough, will be shown to be close to the stationary distribution of  $\pi_t$  thanks to Assumption A1. Since by Assumption A2,  $\min_{x \in \mathcal{X}} \mu_{\text{st}}^{\pi_t}(x) \geq \beta > 0$  then, by choosing the algorithm’s parameters appropriately, we can show that  $\mu_t^N(x) \geq \beta/2 > 0$  holds for all  $x \in \mathcal{X}$ , that is, inequality (7) follows. This is shown in Lemma 13.

It remains to be seen that the estimate  $\hat{r}_t$  is meaningful. In this regard, we claim that

$$\mathbb{E}[\hat{r}_t(x, a) | \mathbf{u}_{t-N}] = r_t(x, a) \quad (11)$$

<sup>7</sup>In what follows, for the sake of brevity, unless otherwise stated, we will omit the modifier “almost surely” from probabilistic statements. It is worth to mention that the finiteness of  $\mathcal{X}$  and  $\mathcal{A}$  allows several statements concerning conditional expectations to hold always, instead of almost surely.

holds for all  $(x, a) \in \mathcal{X} \times \mathcal{A}$ . First note that

$$\mathbb{E}[\hat{\mathbf{r}}_t(x, a) | \mathbf{u}_{t-N}] = \frac{r_t(x, a)}{\pi_t(a|x)\boldsymbol{\mu}_t^N(x)} \mathbb{E}[\mathbb{I}_{\{(x,a)=(\mathbf{x}_t, \mathbf{a}_t)\}} | \mathbf{u}_{t-N}],$$

where we have exploited that  $\pi_t, \boldsymbol{\mu}_t^N \in \sigma(\mathbf{u}_{t-N})$ . Now,

$$\begin{aligned} \mathbb{E}[\mathbb{I}_{\{(x,a)=(\mathbf{x}_t, \mathbf{a}_t)\}} | \mathbf{u}_{t-N}] \\ = \mathbb{P}[\mathbf{a}_t = a | \mathbf{x}_t = x, \mathbf{u}_{t-N}] \mathbb{P}[\mathbf{x}_t = x | \mathbf{u}_{t-N}]. \end{aligned}$$

By definition,  $\mathbb{P}[\mathbf{x}_t = x | \mathbf{u}_{t-N}] = \boldsymbol{\mu}_t^N(x)$  and by (10),  $\mathbb{P}[\mathbf{a}_t = a | \mathbf{x}_t = x, \mathbf{u}_{t-N}] = \pi_t(a|x)$ . Putting together the equalities obtained, we get (11).

By linearity of expectation and since  $\pi_t, \boldsymbol{\mu}_{\text{st}}^{\pi_t} \in \sigma(\mathbf{u}_{t-N})$ , it then follows from (5) and (11) that  $\mathbb{E}[\hat{\boldsymbol{\rho}}_t | \mathbf{u}_{t-N}] = \boldsymbol{\rho}_t$ , and, hence, by the linearity of the Bellman equations and by our assumption that  $\hat{\mathbf{q}}_t$  is the value function underlying the MDP  $(P, \hat{\mathbf{r}}_t)$  and policy  $\pi_t$ , we have, for all  $(x, a) \in \mathcal{X} \times \mathcal{A}$ ,

$$\begin{aligned} \mathbb{E}[\hat{\mathbf{q}}_t(x, a) | \mathbf{u}_{t-N}] &= \mathbf{q}_t(x, a), \\ \mathbb{E}[\hat{\mathbf{v}}_t(x) | \mathbf{u}_{t-N}] &= \mathbf{v}_t(x). \end{aligned} \quad (12)$$

As a consequence, we also have, for all  $(x, a) \in \mathcal{X} \times \mathcal{A}, t \geq N$ ,

$$\begin{aligned} \mathbb{E}[\hat{\boldsymbol{\rho}}_t] &= \mathbb{E}[\boldsymbol{\rho}_t], \\ \mathbb{E}[\hat{\mathbf{q}}_t(x, a)] &= \mathbb{E}[\mathbf{q}_t(x, a)], \\ \mathbb{E}[\hat{\mathbf{v}}_t(x)] &= \mathbb{E}[\mathbf{v}_t(x)]. \end{aligned} \quad (13)$$

Let us finally comment on the computational complexity of our algorithm. Due to the delay in updating the policies based on the weights, the algorithm needs to store  $N$  policies (or weights, leading to the policies). Thus, the memory requirement of MDP-EXP3 scales with  $N|\mathcal{A}||\mathcal{X}|$  (in the real-number model). The computational complexity of the algorithm is dominated by the cost of computing  $\hat{\mathbf{r}}_t$  and, in particular, by the cost of computing  $\boldsymbol{\mu}_t^N$ , plus the cost of solving the Bellman equations (5). The cost of this is  $O(|\mathcal{X}|^2(N + |\mathcal{X}| + |\mathcal{A}|))$  in the worst case, for each time step, however, it can be much smaller for specific practical cases such as when the number of possible next-states is limited.

### C. Main result

Our main result is the following bound concerning the performance of MDP-EXP3.

**Theorem 1 (Regret under bandit feedback):** Let the transition probability kernel  $P$  satisfy Assumptions A1 and A2. Let  $T > 0$  and let  $N = 1 + \lceil \tau \ln T \rceil$ , and  $h(y) = 2y \ln y$  for  $y > 0$ . Then for an appropriate choice of the parameters  $\eta$  and  $\gamma$  (which depend on  $|\mathcal{A}|, T, \beta, \tau$ ), for any sequence of reward functions  $\{r_t\}$  taking values in  $[0, 1]$ , for

$$T > \max \left\{ c_1 \frac{(|\mathcal{A}|\tau + \frac{\tau^3}{|\mathcal{A}|}) \ln |\mathcal{A}|}{\beta^3}, h \left( c_2 \frac{(\frac{|\mathcal{A}|}{\tau} + \frac{\tau}{|\mathcal{A}|}) \ln |\mathcal{A}|}{\beta} \right) \right\}$$

and  $\tau \geq 1^8$  the regret of the algorithm MDP-EXP3 can be bounded as

$$\hat{L}_T \leq C \sqrt{\frac{\tau^3 T |\mathcal{A}| \ln(|\mathcal{A}|) \ln(T)}{\beta}} + C' \tau^2 \ln T$$

for some universal constants  $c_1, c_2, C, C' > 0$ .

Note that with the specific choice of parameters the total cost of the algorithm for a time horizon of  $T$  is  $O(T|\mathcal{X}|^2(\tau \ln(T) + |\mathcal{X}| + |\mathcal{A}|))$ .

<sup>8</sup>The choice of the lower bound on  $\tau$  is arbitrary, but the constants in the theorem depend on it. Furthermore, with some extra work, our proof also gives rise to a bound for the case when  $\tau \rightarrow 0$ , but for simplicity we decided to leave out this analysis.

The proof is presented in the next section. For comparison, we give now the analogue result for the algorithm of Even-Dar et al. [7] that was developed for the full-information case when the algorithm is given  $r_t$  in each time step. As hinted on before, our algorithm reduces to this algorithm if we set  $N = 1$ ,  $\hat{\mathbf{r}}_t = r_t$  and  $\gamma = 0$ . We call this algorithm MDP-E after Even-Dar et al. [7]. The following regret bound holds for this algorithm:

**Theorem 2 (Regret under full-information feedback):** Fix  $T > 0$ . Let the transition probability kernel  $P$  satisfy Assumption A1. Then, for an appropriate choice of the parameter  $\eta$  (which depends on  $|\mathcal{A}|, T, \tau$ ), for any sequence of reward functions  $\{r_t\}$  taking values in  $[0, 1]$ , the regret of the algorithm MDP-E can be bounded as

$$\hat{L}_T \leq 4(\tau + 1) + \sqrt{2T(2\tau + 3)(2\tau^2 + 6\tau + 5) \ln |\mathcal{A}|}. \quad (14)$$

For pedagogical reasons, we shall present the proof in the next section, too. Note that the constants in this bound are different from those presented in Theorem 5.1 of Even-Dar et al. [7]. In particular, the leading term here is  $2\tau^{3/2} \sqrt{2T \ln |\mathcal{A}|}$ , while their leading term is  $4\tau^2 \sqrt{T \ln |\mathcal{A}|}$ . The above bound both corrects some small mistakes in their calculations and improves the result at the same time.<sup>9</sup>

As Even-Dar et al. [7] note, the regret bound (14) does not depend directly on the number of states,  $|\mathcal{X}|$ , but the dependence appears implicitly through  $\tau$  only. Even-Dar et al. [7] also note that a tighter bound, where only the mixing times of the actual policies chosen appear, can be derived. However, it is unclear whether in the worst-case this could be used to improve the bound. Similarly to (14), our bound depends on  $|\mathcal{X}|$  through other constants. In the bandit case, these are  $\beta$  and  $\tau$ . Comparing the theorems it seems that the main price of not seeing the rewards is the appearance of  $|\mathcal{A}|$  instead of  $\ln |\mathcal{A}|$  (a typical difference between the bandit and full observation cases) and the appearance of a  $\sqrt{1/\beta}$  term in the bound.

### D. Discussion and future work

In this paper, we have presented an online learning algorithm, MDP-EXP3 for adversarial MDPs, that is, finite stochastic Markovian decision environments where the reward function may change after each transition. This is the first algorithm for this setting that has a rigorously proved  $O(\sqrt{T \ln T})$  bound on its regret. We discuss the features of the algorithm, along with future research directions below.

*a) Extensions:* We considered the expected regret relative to the best fixed policy selected in hindsight. A typical extension is to prove a high probability bound on the regret, which we think can be done in a standard way using concentration inequalities. Note, however, that the extension is more complicated than for the bandit problems because the mixing property has to be used together with the martingale reasoning. Another potential extension is to compete with larger policy classes, such as with sequences of policies with a bounded number of policy-switches. Similarly to Neu et al. [13, 15], the MDP-EXP3 algorithm should then be modified by replacing EXP3 with the EXP3.S algorithm of Auer et al. [2], specifically designed to compete with switching experts in place of EXP3. Note that, again, the analysis will be more complicated than in the bandit case, and requires to bound the maximum regret of EXP3.S relative to any fixed policy over any time window. When compared to a policy with  $C$  switches, the resulting regret bound is expected to be  $C$  times

<sup>9</sup>One of the mistakes is in the proof of Theorem 4.1 of Even-Dar et al. [7] where they failed to notice that  $q_t^{\pi_t}$  can take on negative values. Thus, their Assumption 3.1 is not met by  $\{q_t^{\pi_t}\}$  (one needs to extend the upper bound given in their Lemma 2.2 with a lower bound and change Assumption 3.1). As a result, Assumption 3.1 cannot be used to show that the inequality in the proof of Theorem 4.1 holds. This mistake, as well as the others, can easily be corrected, as we show it here.

larger than that of Theorem 1, while the algorithm would not need to know the number of switches  $C$ .

*b) Tuning and complexity:* Setting up and running the algorithm MDP-EXP3 may actually be computationally demanding. Setting the parameters of the algorithm ( $\eta$  and  $\gamma$ ) requires a known lower bound  $\beta^*$  on the visitation probabilities such that  $\beta = \inf_{\pi, x} \mu_{\text{st}}^\pi(x) > \beta^* > 0$  and also the knowledge of an upper bound  $\tau^*$  on the mixing time  $\tau$ . While these quantities can be determined in principle from the transition probability kernel  $P$ , it is not clear how to compute efficiently the minimum over all policies. Computational issues also arise during running the algorithm: as it is discussed in Section IV-B, each step of the MDP-EXP3 algorithm requires  $O(|\mathcal{X}|^2(\tau \ln T + |\mathcal{X}| + |\mathcal{A}|))$  computations, which may be too demanding if, e.g., the size of the state space is large. It is an interesting problem to design a more efficient method that achieves similar performance guarantees.

*c) Assumptions on the Markovian dynamics:* We believe that it should be possible to extend our main result beyond Assumption A1, requiring only the existence of a unique stationary distribution for any policy  $\pi$  (we will refer to this latter assumption as the unichain assumption). Using that the distribution of any unichain Markov chain converges exponentially fast to its stationary distribution, and that it is enough to verify Assumption A1 for deterministic policies only, one can easily show that if  $P$  satisfies the unichain assumption, then there exists an integer  $K > 0$  such that  $(P^\pi)^K$  is a scrambling matrix for any policy  $\pi$ . Then, we conjecture that the MDP-EXP3 algorithm will work as it is, except that the regret will be increased. The key to prove this result is to generalize Lemmas 4 and 5 to this case.

Finally, one may also consider the case when the Markov chains corresponding to  $P^\pi$  are periodic. We speculate that this may be dealt with using *occupancy probabilities* and Cesaro-averages instead of the stationary and state distributions, respectively.

## V. PROOFS

In this section we present the proofs of Theorem 1 and Theorem 2. We start with the proof of Theorem 2 as this is a simpler result. The proof of this result is presented partly for the sake of completeness and partly so that we can be more specific about the corrections required to fix the main result (Theorem 5.2) of Even-Dar et al. [7]. Further, the proof will also serve as a starting point for the proof of our main result, Theorem 1. Nevertheless, the impatient reader may skip this next section and jump immediately to the proof of Theorem 1, which apart from referring to some general lemmas developed in the next subsection, is entirely self-contained.

### A. Proof of Theorem 2

Throughout this section we consider the MDP-E algorithm (given by Algorithm 1 with  $N = 1$ ,  $\hat{\mathbf{r}}_t = \mathbf{r}_t$  and  $\gamma = 0$ ), and we suppose that  $P$  satisfies Assumption A1. Let  $\pi_t$  denote the policy used in step  $t$  of the algorithm. Note that  $\pi_t$  is not random since by assumption the reward function is available at all states (not just the visited ones). Hence, the sequence of policies chosen does not depend on the states visited by the algorithm but is deterministic. Remember that  $\rho_t = \rho_t^{\pi_t}$  denotes the average reward of policy  $\pi_t$  measured with respect to the reward function  $r_t$ . Following Even-Dar et al. [7], fix some policy  $\pi$  and consider the decomposition of the regret relative to  $\pi$ :

$$R_T^\pi - \hat{R}_T = \left( R_T^\pi - \sum_{t=1}^T \rho_t^\pi \right) + \left( \sum_{t=1}^T \rho_t^\pi - \sum_{t=1}^T \rho_t \right) + \left( \sum_{t=1}^T \rho_t - \hat{R}_T \right). \quad (15)$$

The first and the last terms measure the difference between the sum of (asymptotic) average rewards and the actual expected reward. The mixing assumption (Assumption A1) ensures that these differences

are not large. In particular, in the case of a fixed policy, this difference is bounded by a constant of order  $\tau$ :

**Lemma 1:** For any  $T \geq 1$  and any policy  $\pi$ , it holds that

$$R_T^\pi - \sum_{t=1}^T \rho_t^\pi \leq 2\tau + 2. \quad (16)$$

This lemma is also stated in [7]. We give the proof for completeness also to correct slight inaccuracies of the proof given in [7].

*Proof:* Let  $\{(\mathbf{x}_t, \mathbf{a}_t)\}$  be the trajectory when  $\pi$  is followed. Note that the difference between  $R_T^\pi$  and  $\sum_{t=1}^T \rho_t^\pi$  is caused by the difference between the initial distribution of  $\mathbf{x}_1$  and the stationary distribution of  $\pi$ . To quantify the difference, write

$$R_T^\pi - \sum_{t=1}^T \rho_t^\pi = \sum_{t=1}^T \sum_x (\nu_t^\pi(x) - \mu_{\text{st}}^\pi(x)) \sum_a \pi(a|x) r_t(x, a),$$

where  $\nu_t^\pi(x) = \mathbb{P}[\mathbf{x}_t = x]$  is the state distribution at time step  $t$ . Viewing  $\nu_t^\pi$  as a row vector, we have  $\nu_t^\pi = \nu_{t-1}^\pi P^\pi$ . Consider the  $t^{\text{th}}$  term of the above difference. Then, using  $r_t(x, a) \in [0, 1]$  and Assumption A1 we get<sup>10</sup>

$$\begin{aligned} & \sum_x (\nu_t^\pi(x) - \mu_{\text{st}}^\pi(x)) \sum_a \pi(a|x) r_t(x, a) \\ & \leq \|\nu_t^\pi - \mu_{\text{st}}^\pi\|_1 = \|\nu_{t-1}^\pi P^\pi - \mu_{\text{st}}^\pi P^\pi\|_1 \\ & \leq e^{-1/\tau} \|\nu_{t-1}^\pi - \mu_{\text{st}}^\pi\|_1 \leq \dots \leq e^{-(t-1)/\tau} \|\nu_1^\pi - \mu_{\text{st}}^\pi\|_1 \\ & \leq 2e^{-(t-1)/\tau}. \end{aligned}$$

This, together with the elementary inequality  $\sum_{t=1}^T e^{-(t-1)/\tau} \leq 1 + \int_0^\infty e^{-t/\tau} dt = 1 + \tau$  gives the desired bound. ■

Consider now the second term of (15) and in particular its  $t^{\text{th}}$  term  $\rho_t^\pi - \rho_t = \rho_t^\pi - \rho_t^{\hat{\pi}_t}$ . This term is the difference of the average reward obtained by  $\pi$  and  $\hat{\pi}_t$ . The following lemma shows that this difference can be rewritten in terms of the state-wise action-disadvantages underlying  $\hat{\pi}_t$ :

**Lemma 2 (Performance difference lemma):** Consider an MDP specified by the transition probability kernel  $P$  and reward function  $r$ . Let  $\pi, \hat{\pi}$  be two (stochastic stationary) policies in the MDP. Assume that  $\mu_{\text{st}}^\pi, \rho^{\hat{\pi}}$  and  $q^{\hat{\pi}}$  are well-defined.<sup>11</sup> Then,

$$\rho^\pi - \rho^{\hat{\pi}} = \sum_{x,a} \mu_{\text{st}}^\pi(x) \pi(a|x) \left[ q^{\hat{\pi}}(x, a) - v^{\hat{\pi}}(x) \right].$$

This lemma appeared as Lemma 4.1 in [7], but similar statements have been known for a while. For example, the book of Cao [3] also puts performance difference statements in the center of the theory of MDPs. For the sake of completeness, we include the easy proof. Note that the statement of the lemma continues to hold even when  $q^{\hat{\pi}}$  and  $v^{\hat{\pi}}$  are shifted by the same constant function.

*Proof:* We have

$$\begin{aligned} & \sum_{x,a} \mu_{\text{st}}^\pi(x) \pi(a|x) q^{\hat{\pi}}(x, a) \\ & = \sum_{x,a} \mu_{\text{st}}^\pi(x) \pi(a|x) \left[ r(x, a) - \rho^{\hat{\pi}} + \sum_{x'} P(x'|x, a) v^{\hat{\pi}}(x') \right] \\ & = \rho^\pi - \rho^{\hat{\pi}} + \sum_x \mu_{\text{st}}^\pi(x) v^{\hat{\pi}}(x), \end{aligned}$$

where the second equality holds since  $\sum_{x,a} \mu_{\text{st}}^\pi(x) \pi(a|x) P(x'|x, a) = \mu_{\text{st}}^\pi(x')$ . Reordering the terms gives the desired result. ■

<sup>10</sup>Even-Dar et al. [7] mistakenly uses  $\|\nu_t^\pi - \mu_{\text{st}}^\pi\|_1 \leq e^{-t/\tau} \|\nu_1^\pi - \mu_{\text{st}}^\pi\|_1$  in their paper ( $t = 1$  immediately shows that this can be false). See, e.g., the proofs of their Lemmas 2.2 and 5.2.

<sup>11</sup>This lemma does not need Assumption A1 and in fact the assumptions we make could be further relaxed with a slight change to the claim.



Because of this lemma,  $\rho_t^\pi - \rho_t = \sum_{x,a} \mu_{st}^\pi(x) \pi(a|x) (q_t^{\pi_t}(x,a) - v_t^{\pi_t}(x))$ . Thus, by flipping the sum that runs over time with the one that runs over the state-action pairs, we get:  $\sum_{t=1}^T \rho_t^\pi - \sum_{t=1}^T \rho_t = \sum_{x,a} \mu_{st}^\pi(x) \pi(a|x) \sum_{t=1}^T (q_t^{\pi_t}(x,a) - v_t^{\pi_t}(x))$ . Thus, it suffices to bound, for a fixed state-action pair  $(x, a)$ , the sum

$$\sum_{t=1}^T (q_t^{\pi_t}(x, a) - v_t^{\pi_t}(x)) = \sum_{t=1}^T \left( q_t^{\pi_t}(x, a) - \sum_{a'} \pi_t(a'|x) q_t^{\pi_t}(x, a') \right).$$

By construction,  $\pi_t(a|x) \propto \exp(\eta \sum_{s=1}^{t-1} q_s^{\pi_s}(x, a))$  (recall that  $\gamma = 0$  in this version of the algorithm), which means that the sum is the regret of the so-called exponential weights algorithm (EWA) against action  $a$  when the algorithm is used on the sequence  $\{q_t^{\pi_t}(x, \cdot)\}$ . Assume for a moment that  $K > 0$  is such that  $\|q_t^{\pi_t}\|_\infty \leq K$  holds for  $1 \leq t \leq T$ . Then, since  $q_t^{\pi_t}$  takes its values from an interval of length  $2K$ , Theorem 2.2 in [5] implies that the regret of EWA can be bounded by

$$\frac{\ln|\mathcal{A}|}{\eta} + \frac{K^2 \eta T}{2}. \quad (17)$$

Notice that  $\{q_t^{\pi_t}\}$  is a sequence that is sequentially generated from  $\{r_t\}$ . It is Lemma 4.1 of [5] that shows that the bound of Theorem 2.2 of [5] continues to hold for such sequentially generated functions. Putting the inequalities together, we obtain

$$\sum_{t=1}^T \rho_t^\pi - \sum_{t=1}^T \rho_t \leq \frac{\ln|\mathcal{A}|}{\eta} + \frac{K^2 \eta T}{2}. \quad (18)$$

According to the next lemma an appropriate value for  $K$  is  $2\tau + 3$ . The lemma is stated in a greater generality than what is needed here because the more general form will be used later.

**Lemma 3:** Pick any policy  $\pi$  in an MDP  $(P, r)$ . Assume that the mixing time of  $\pi$  is  $\tau$  in the sense of (1). If  $|\sum_a \pi(a|x) r(x, a)| \leq R \leq \|r\|_\infty$  holds for any  $x \in \mathcal{X}$ , then  $|v^\pi(x)| \leq 2R(\tau + 1)$  holds for all  $x \in \mathcal{X}$ . Furthermore, for any  $(x, a) \in \mathcal{X} \times \mathcal{A}$ ,  $|q^\pi(x, a)| \leq R(2\tau + 3) + |r(x, a)|$  and, if, in addition,  $r(x, a) \geq 0$  for any  $(x, a) \in \mathcal{X} \times \mathcal{A}$ , then  $|q^\pi(x, a)| \leq (2\tau + 3) \|r\|_\infty$ .

*Proof:* As it is well known and is easy to see from the definitions, the (differential) value of policy  $\pi$  at state  $x$  can be written as

$$v^\pi(x) = \sum_{s=1}^{\infty} \sum_{x'} (\nu_{s,x}^\pi(x') - \mu_{st}^\pi(x')) \sum_a \pi(a|x') r(x', a),$$

where  $\nu_{s,x}^\pi = e_x(P^\pi)^{s-1}$  is the state distribution when following  $\pi$  for  $s-1$  steps starting from state  $x$ . The triangle inequality and then the bound on  $\sum_a \pi(a|x') r(x', a)$  gives

$$|v^\pi(x)| \leq R \sum_{s=1}^{\infty} \sum_{x'} |\nu_{s,x}^\pi(x') - \mu_{st}^\pi(x')| \leq 2R(\tau + 1),$$

where in the second inequality we used  $\|\nu_{s,x}^\pi - \mu_{st}^\pi\|_1 \leq 2e^{-(s-1)/\tau}$  and that  $\sum_{s=1}^{\infty} e^{-(s-1)/\tau} \leq \tau + 1$  (cf. the proof of Lemma 1). This proves the first inequality. The inequalities on  $|q^\pi(x, a)|$  follow from the first part and the Bellman equation:

$$\begin{aligned} |q^\pi(x, a)| &\leq |r(x, a)| + |\rho^\pi| + \sum_{x'} P(x'|x, a) |v^\pi(x')| \\ &\leq R(2\tau + 3) + |r(x, a)|, \\ |q^\pi(x, a)| &\leq |r(x, a) - \rho^\pi| + \sum_{x'} P(x'|x, a) |v^\pi(x')| \\ &\leq (2\tau + 3) \|r\|_\infty. \end{aligned}$$

Here, in the first inequality we used that  $|\rho^\pi| \leq \sum_x \mu_{st}^\pi(x) |\sum_a \pi(a|x) r(x, a)| \leq R$ , while the second inequality holds since  $|r(x, a) - \rho^\pi|, R \in [0, \|r\|_\infty]$ . ■

Let us now consider the third term of (15),  $\sum_{t=1}^T \rho_t - \widehat{R}_T$ . The  $t^{\text{th}}$  term of this difference is the difference between the average reward of  $\pi_t$  and the expected reward obtained in step  $\pi_t$ . If  $\nu_t(x)$  is the distribution of states in time step  $t$ ,  $\sum_{t=1}^T \rho_t - \widehat{R}_T = \sum_{t=1}^T \sum_x (\mu_{st}^\pi(x) - \nu_t(x)) \sum_a \pi(a|x) r_t(x, a)$ . Thus,

$$\sum_{t=1}^T \rho_t - \widehat{R}_T \leq \sum_{t=1}^T \|\mu_{st}^\pi - \nu_t\|_1 \quad (19)$$

and so remains to bound the  $\ell^1$  distances between the distributions  $\mu_{st}^\pi$  and  $\nu_t$ . For this, we will use two general lemmas that will again come useful later. For  $f : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ , introduce the mixed norm  $\|f\|_{1,\infty} = \max_x \sum_a |f(a|x)|$ , where  $f(a|x)$  is identified with  $f(x, a)$ . Clearly,  $\|\nu P^\pi - \nu P^{\hat{\pi}}\|_1 \leq \|\pi - \hat{\pi}\|_{1,\infty}$  holds for any two policies  $\pi, \hat{\pi}$  and any distribution  $\nu$  (cf. Lemma 5.1 in [7]). The first lemma shows that the map  $\pi \mapsto \mu_{st}^\pi$  as a map from the space of stationary policies equipped with the mixed norm  $\|\cdot\|_{1,\infty}$  to the space of distributions equipped with the  $\ell^1$ -norm is  $(\tau+1)$ -Lipschitz:

**Lemma 4:** Let  $P$  be a transition probability kernel over  $\mathcal{X} \times \mathcal{A}$  such that the mixing time of  $P$  is  $\tau < \infty$ . For any two policies,  $\pi, \hat{\pi}$ , it holds that

$$\|\mu_{st}^\pi - \mu_{st}^{\hat{\pi}}\|_1 \leq (\tau + 1) \|\pi - \hat{\pi}\|_{1,\infty}.$$

*Proof:* The statement follows from solving

$$\begin{aligned} \|\mu_{st}^\pi - \mu_{st}^{\hat{\pi}}\|_1 &\leq \left\| \mu_{st}^\pi P^\pi - \mu_{st}^{\hat{\pi}} P^\pi \right\|_1 + \left\| \mu_{st}^{\hat{\pi}} P^\pi - \mu_{st}^{\hat{\pi}} P^{\hat{\pi}} \right\|_1 \\ &\leq e^{-1/\tau} \|\mu_{st}^\pi - \mu_{st}^{\hat{\pi}}\|_1 + \|\pi - \hat{\pi}\|_{1,\infty} \end{aligned}$$

for  $\|\mu_{st}^\pi - \mu_{st}^{\hat{\pi}}\|_1$  and using

$$1/(1 - e^{-1/\tau}) \leq \tau + 1. \quad (20)$$

The next lemma allows us to compare an  $n$ -step distribution under a policy sequence with the stationary distribution of the sequence's last policy:

**Lemma 5:** Let  $P$  be a transition probability kernel over  $\mathcal{X} \times \mathcal{A}$  such that the mixing time of  $P$  is  $\tau < \infty$ . Take any probability distribution  $\nu_1$  over  $\mathcal{X}$ , integer  $n \geq 1$  and policies  $\pi_1, \dots, \pi_n$ . Consider the distribution  $\nu_n = \nu_1 P^{\pi_1} \dots P^{\pi_{n-1}}$ . Then, it holds that

$$\|\nu_n - \mu_{st}^{\pi_n}\|_1 \leq 2e^{-(n-1)/\tau} + (\tau + 1)^2 \max_{1 \leq t \leq n} \|\pi_t - \pi_{t-1}\|_{1,\infty},$$

where, for convenience, we have introduced  $\pi_0 = \pi_1$ .

*Proof:* If  $n = 1$  the result is obtained from  $\|\nu_1 - \mu_{st}^{\pi_1}\|_1 \leq 2$ . Thus, in what follows we assume  $n \geq 2$ . Let  $c = \max_{1 \leq t \leq n} \|\pi_t - \pi_{t-1}\|_{1,\infty}$ . By the triangle inequality,

$$\begin{aligned} \|\nu_n - \mu_{st}^{\pi_n}\|_1 &\leq \|\nu_n - \mu_{st}^{\pi_{n-1}}\|_1 + \|\mu_{st}^{\pi_{n-1}} - \mu_{st}^{\pi_n}\|_1 \\ &\leq e^{-1/\tau} \|\nu_{n-1} - \mu_{st}^{\pi_{n-1}}\|_1 + (\tau + 1)c, \end{aligned}$$

where we used that by the previous lemma  $\|\mu_{st}^{\pi_{n-1}} - \mu_{st}^{\pi_n}\|_1 \leq (\tau + 1) \|\pi_{n-1} - \pi_n\|_{1,\infty} \leq (\tau + 1)c$ . Continuing recursively, we get

$$\begin{aligned} \|\nu_n - \mu_{st}^{\pi_n}\|_1 &\leq e^{-1/\tau} \left( e^{-1/\tau} \|\nu_{n-2} - \mu_{st}^{\pi_{n-2}}\|_1 + (\tau + 1)c \right) + (\tau + 1)c \\ &\vdots \\ &\leq e^{-\frac{n-1}{\tau}} \|\nu_1 - \mu_{st}^{\pi_1}\|_1 + (\tau + 1)c \left( 1 + e^{-\frac{1}{\tau}} + \dots + e^{-\frac{n-2}{\tau}} \right) \\ &\leq 2e^{-(n-1)/\tau} + (\tau + 1)^2 c, \end{aligned}$$

where we bounded the geometrical series by  $1/(1 - e^{-1/\tau})$  and used (20). ■

Applying this lemma to  $\|\nu_t - \mu_{\text{st}}^{\pi_t}\|_1$  we get

$$\|\nu_t - \mu_{\text{st}}^{\pi_t}\|_1 \leq 2e^{-(t-1)/\tau} + (\tau + 1)^2 K',$$

where  $K'$  is a bound on  $\max_{2 \leq t \leq n} \|\pi_t - \pi_{t-1}\|_{1,\infty}$ .<sup>12</sup> Therefore, by (19), we have

$$\sum_{t=1}^T \rho_t - \widehat{R}_T \leq 2 \sum_{t=1}^T e^{-t/\tau} + (\tau + 1)^2 K' T \leq 2\tau + (\tau + 1)^2 K' T.$$

Thus, it remains to find an appropriate value for  $K'$ . It is a well known property of EWA that  $\|\pi_t(\cdot|x) - \pi_{t-1}(\cdot|x)\|_1 \leq \eta \|q_{t-1}^{\pi_{t-1}}(x, \cdot)\|_\infty$ . Indeed, applying Pinsker's inequality and Hoeffding's lemma (see Section A.2 and Lemma A.6 in Cesa-Bianchi and Lugosi 5), we get for any  $x \in \mathcal{X}$

$$\begin{aligned} \|\pi_t(\cdot|x) - \pi_{t-1}(\cdot|x)\|_1 &\leq \sqrt{2D(\pi_{t-1}(\cdot|x) \|\pi_t(\cdot|x))} \\ &= \sqrt{2 \left[ \ln \left( \sum_b \pi_{t-1}(b|x) e^{\eta q_{t-1}^{\pi_{t-1}}(b,x)} \right) - \sum_a \eta \pi_{t-1}(a|x) q_{t-1}^{\pi_{t-1}}(b,x) \right]} \\ &\leq \eta \|q_{t-1}^{\pi_{t-1}}(x, \cdot)\|_\infty \end{aligned}$$

where, for two distributions  $D(v\|v') = \sum_i v_i \ln(v_i/v'_i)$  denotes the Kullback-Leibler divergence of the distributions  $v$  and  $v'$ . Thus,  $\|\pi_t - \pi_{t-1}\|_{1,\infty} \leq \eta \|q_{t-1}^{\pi_{t-1}}\|_\infty$ . Now, by Lemma 3,  $\|q_t^{\pi_t}\|_\infty \leq 2\tau + 3$ , showing that  $K' = \eta(2\tau + 3)$  is suitable. Putting together the inequalities obtained, we get

$$\sum_{t=1}^T \rho_t - \widehat{R}_T \leq 2\tau + (2\tau + 3)(\tau + 1)^2 \eta T.$$

Combining (16), (18) and this last bound, we obtain

$$R_T^\pi - \widehat{R}_T = 4\tau + 2 + \frac{\ln|\mathcal{A}|}{\eta} + \frac{\eta T(2\tau + 3)(2\tau^2 + 6\tau + 5)}{2}.$$

Setting

$$\eta = \sqrt{\frac{2 \ln|\mathcal{A}|}{T(2\tau + 3)(2\tau^2 + 6\tau + 5)}},$$

we get the bound stated in Theorem 2.

### B. Proof of Theorem 1

Throughout this section we consider the MDP-EXP3 algorithm and suppose that both Assumptions A1 and A2 hold for  $P$ . We start from the decomposition (15), which is repeated to emphasize the difference that some of the terms are random now:

$$R_T^\pi - \widehat{R}_T = \left( R_T^\pi - \sum_{t=1}^T \rho_t^\pi \right) + \left( \sum_{t=1}^T \rho_t^\pi - \sum_{t=1}^T \rho_t \right) + \left( \sum_{t=1}^T \rho_t - \widehat{R}_T \right) \quad (21)$$

As before, Lemma 1 shows that the first term is bounded by  $2(\tau + 1)$ . Thus, it remains to bound the expectation of the other two terms. This is done in the following two propositions whose proofs are deferred to the next subsections:

<sup>12</sup>Lemma 5.2 of Even-Dar et al. [7] gives a bound on  $\|\nu_t - \mu_{\text{st}}^{\pi_t}\|_1$  with a slightly different technique. However, there are multiple mistakes in the proof. Once the mistakes are removed, their bounding technique gives the same result as ours. One of the mistakes is that Assumption 3.1 states that  $K' = \sqrt{\ln|\mathcal{A}|/T}$ , whereas since the range of the action-value functions scales with  $\tau$ ,  $K'$  should also scale with  $\tau$ . Unfortunately, in [16] we committed the same mistake, which we correct here. We choose to present an alternate proof, as we find it somewhat cleaner and it also gave us the opportunity to present Lemma 4.

**Proposition 1:** Let  $L = \frac{2}{\beta}(2\tau + 3)$ ,  $V_{\bar{q}} = \frac{2}{\beta} \left( \frac{|\mathcal{A}|}{\gamma} + 2\tau + 2 \right)$ ,  $U_{\bar{v}} = \frac{4}{\beta}(\tau + 1)$ ,  $U_{\pi_{\bar{q}}} = \frac{4}{\beta}(\tau + 2)$ ,  $U_q = 2\tau + 3$ ,  $U_{\bar{q}} = 2\tau + 4$ ,  $e' = e - 1$ ,  $e'' = e - 2$ ,

$$\begin{aligned} c &= \eta \frac{e(U_{\bar{v}} + L + \gamma V_{\bar{q}})}{1 - \gamma - \eta e N V_{\bar{q}}}, \\ c' &= \eta \frac{e'(L + \gamma V_{\bar{q}}) + (e' U_{\pi_{\bar{q}}} + U_{\bar{v}}) U_{\bar{q}} |\mathcal{A}|}{1 - \gamma - \eta e(N + 1) V_{\bar{q}}}, \end{aligned}$$

and assume that  $\gamma \in (0, 1)$ ,  $c(\tau + 1)^2 < \beta/2$ ,  $N \geq 1 + \left\lceil \tau \ln \left( \frac{4}{\beta - 2c(\tau + 1)^2} \right) \right\rceil$ ,  $0 < \eta < \frac{\beta(1-\gamma)}{2e(N+1)(|\mathcal{A}|/\gamma + 2\tau + 2)}$ . Then, for any policy  $\pi$ , we have

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[\rho_t^\pi - \rho_t] &\leq \frac{\ln|\mathcal{A}|}{\eta} + (N - 1)(U_{\bar{q}} + U_q + 1) \\ &\quad + (T - 2N + 2) \left( c'(N - 1)(1 + \eta e'' V_{\bar{q}}) \right. \\ &\quad \left. + \gamma U_q + \eta e'' |\mathcal{A}| U_{\pi_{\bar{q}}} U_{\bar{q}} \right). \end{aligned}$$

**Proposition 2:** Assume that the conditions of Proposition 1 hold. Then,

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[\rho_t] - \widehat{R}_T &\leq N - 1 + (T - N + 1) c(\tau + 1)^2 \\ &\quad + 2(T - N + 1) e^{-(N-1)/\tau}. \end{aligned} \quad (22)$$

Note that setting  $N \geq 1 + \lceil \tau \ln T \rceil$ , as suggested in Theorem 1, the last term in the right-hand side of (22) becomes  $O(1)$ , while for  $T$  sufficiently large all the conditions of the last two propositions will be satisfied. This leads to the proof of Theorem 1:

*Proof of Theorem 1:* If  $|\mathcal{A}| = 1$  then, due to  $\hat{L}_T = 0$ , the statement is trivial, so we assume  $|\mathcal{A}| \geq 2$  from now on. Define  $\alpha = \frac{\beta}{2} e(U_{\bar{v}} + L + \gamma V_{\bar{q}})$ ,  $\alpha' = \frac{\beta}{2} \{e'(L + \gamma V_{\bar{q}}) + (e' U_{\pi_{\bar{q}}} + U_{\bar{v}}) U_{\bar{q}} |\mathcal{A}|\}$  so that  $c = 2\eta \frac{\alpha}{\beta(1-\gamma) - 2\eta e N V_{\bar{q}}}$  and  $c' = 2\eta \frac{\alpha'}{\beta(1-\gamma) - 2\eta e(N+1) V_{\bar{q}}}$ , where  $V_{\bar{q}} = \frac{\beta}{2} V_{\bar{q}} = |\mathcal{A}|/\gamma + 2\tau + 2$ . In the following we will use the notation  $f \sim g$  for two positive-valued functions  $f, g: D \rightarrow \mathbb{R}^+$  defined on the same domain  $D$  to denote that they are equivalent up to a constant factor, that is,  $\sup_{x \in D} \max\{f(x)/g(x), g(x)/f(x)\} < \infty$ . With this notation, on  $|\mathcal{A}| \geq 2$ ,  $\tau \geq 1$  and as long as  $\gamma \leq 1$ , we have

$$\alpha \sim |\mathcal{A}| + \tau \quad \text{and} \quad \alpha' \sim |\mathcal{A}| \tau^2 \quad (23)$$

independently of the value of  $\beta$  and of the choice of  $\eta, \gamma, N$ . In what follows all the equivalences will be stated for the domain  $|\mathcal{A}| \geq 2, \tau \geq 1$ .

We now show how to choose  $\eta, \gamma$  and  $N$  so as to achieve a small regret bound. In order to do so we will choose these constants so that the conditions of Propositions 1 and 2 are satisfied. For simplicity, we add the constraint  $\gamma \leq 1/2$  that we will also show to hold. Under this additional constraint, the inequality

$$\eta < \frac{\beta(1-\gamma)}{2e(N+1)(|\mathcal{A}|/\gamma + 2\tau + 2)} \quad (24)$$

will be satisfied if we choose  $\gamma = 8e\eta(N+1)(|\mathcal{A}| + \tau + 1)/\beta$ . Indeed, the said inequality holds since it is equivalent to  $D = \beta(1-\gamma) - 2\eta e(N+1)(|\mathcal{A}|/\gamma + 2\tau + 2) > 0$  and

$$\begin{aligned} D &= \beta(1-\gamma) - \frac{2\eta e(N+1)(|\mathcal{A}| + \gamma(2\tau + 2))}{\gamma} \\ &\geq \frac{\beta}{2} - \frac{2\eta e(N+1)(|\mathcal{A}| + \tau + 1)}{\gamma} = \frac{\beta}{4} > 0, \end{aligned}$$

where the first inequality holds because  $\gamma \leq 1/2$  and the second equality holds by the definition of  $\gamma$ . Since  $c \leq 2\eta\alpha/D$  and  $c' =$

$2\eta\alpha'/D$ , this also implies

$$c \leq 8\eta\alpha/\beta \quad \text{and} \quad c' \leq 8\eta\alpha'/\beta. \quad (25)$$

Due to this upper bound on  $c$ ,  $c(\tau+1)^2 < \beta/2$  will be satisfied if

$$\eta < \frac{\beta^2}{16\alpha(\tau+1)^2}. \quad (26)$$

To satisfy  $\gamma \leq 1/2$ , the inequality

$$\eta \leq \frac{\beta}{16e(N+1)(|\mathcal{A}|+\tau+1)} \quad (27)$$

has to be satisfied, too. Before proving (26) and (27), we derive the regret bound they imply.

Taking expectation in (21) and using the bounds of Lemma 1 and Propositions 1 and 2, we get

$$\begin{aligned} \widehat{L}_T &\leq 2\tau + 2 + \frac{\ln|\mathcal{A}|}{\eta} \\ &+ T \left[ c'(N-1)(1 + \eta e''V_{\bar{q}}) + \gamma U_q + \eta e''|\mathcal{A}| U_{\pi_{\bar{q}}} U_{\bar{q}} + c(\tau+1)^2 \right] \\ &+ (N-1)(U_{\bar{q}} + U_q + 2) + 2Te^{-(N-1)/\tau}. \end{aligned}$$

Choosing  $N = 1 + \lceil \tau \ln T \rceil$ , we have  $2Te^{-(N-1)/\tau} \leq 2$ . Furthermore, (24) implies  $\eta e''V_{\bar{q}} \leq \frac{e''}{e(N+1)}$ . This, together with the definition of the different constants above and the bound (25) on  $c$  and  $c'$  gives

$$\begin{aligned} \widehat{L}_T &\leq 2\tau + 2 + \frac{\ln|\mathcal{A}|}{\eta} \\ &+ \frac{\eta T}{\beta} \left[ 8\alpha'(N-1) \left( 1 + \frac{e''}{e(N+1)} \right) \right. \\ &\quad \left. + 8e\beta(N+1)(|\mathcal{A}|+\tau+1)(2\tau+3) \right. \\ &\quad \left. + 8e''|\mathcal{A}|(\tau+2)^2 + 8\alpha(\tau+1)^2 \right] \\ &+ (N-1)(U_{\bar{q}} + U_q + 2) + 2 \\ &\leq \frac{\ln|\mathcal{A}|}{\eta} + \frac{\eta T}{\beta} B + C_1\tau N, \end{aligned}$$

where we introduced  $B$  to denote the expression in the squared brackets and used the fact that  $2\tau+4+(N-1)(U_{\bar{q}}+U_q+2) \leq C_1\tau N$  for some constant  $C_1 > 0$ . Note that (23) implies that

$$B \sim N|\mathcal{A}|\tau^2 + \tau^3 \sim |\mathcal{A}|\tau^3 \ln T \quad (28)$$

since  $N \sim \tau \ln T$  for  $\tau \geq 1, T \geq 2$ . Now, choose  $\eta = \sqrt{\frac{\beta \ln|\mathcal{A}|}{TB}}$ . Then,

$$\begin{aligned} \widehat{L}_T &= 2\sqrt{\frac{TB \ln|\mathcal{A}|}{\beta}} + C_1\tau N \\ &\leq C_2\sqrt{\frac{\tau^3 T |\mathcal{A}| \ln(|\mathcal{A}|) \ln(T)}{\beta}} + C_3\tau^2 \ln T \end{aligned}$$

for some appropriate constants  $C_2, C_3 > 0$ .

It remains to show that for  $T$  large enough, inequalities (26) and (27) will hold, and also the lower bound on  $N$  in the propositions will be satisfied. Instead of (26) we will choose a lower bound on  $T$  to guarantee the stronger condition

$$\eta \leq \frac{\beta^2}{32\alpha(\tau+1)^2}, \quad (29)$$

which, together with (25), also ensures  $c(\tau+1)^2 \leq \beta/4$ . The latter inequality implies that the lower bound on  $N$  in the propositions is satisfied for  $T \geq 8/\beta$ . Using the choice of  $\eta$  and the respective equivalent forms (23) and (28) for  $\alpha$  and  $B$ , one can see that

condition (29) is satisfied if  $T \ln T \geq C_4 \frac{(|\mathcal{A}|\tau+\tau^3/|\mathcal{A}|) \ln|\mathcal{A}|}{\beta^3}$  for some appropriate constant  $C_4 > 0$ . To keep things simple, notice that selecting  $T \geq C_4 \frac{(|\mathcal{A}|\tau+\tau^3/|\mathcal{A}|) \ln|\mathcal{A}|}{\beta^3}$  implies (29), and also  $T \geq 8/\beta$  if  $C_4 \geq 8.13$ . Furthermore, one can similarly show that (27) is satisfied if  $\frac{T}{\ln T} \geq C_5 \frac{(|\mathcal{A}|/\tau+\tau/|\mathcal{A}|) \ln|\mathcal{A}|}{\beta}$  for some appropriate constant  $C_5 > 0$ . By Proposition 3 of Antos et al. [1], for any  $u > 0$ ,  $t/\ln t > u$  if  $t \geq h(u) \stackrel{\text{def}}{=} 2u \ln u$ . Thus, the last condition on  $T$  is satisfied if  $T > h\left(C_5 \frac{(|\mathcal{A}|/\tau+\tau/|\mathcal{A}|) \ln|\mathcal{A}|}{\beta}\right)$ . This finishes the proof of the theorem.  $\blacksquare$

### C. General tools for the proofs of Propositions 1 and 2

Just like in the previous section, throughout this section we suppose that both Assumptions A1 and A2 hold for  $P$  and the rewards are in the  $[0, 1]$  interval. We proceed with a series of lemmas to bound the rate of change of the policies generated by MDP-EXP3.

**Lemma 6:** Let  $1 \leq t \leq T$  and assume that  $\mu_t^N(x) \geq \beta/2$  holds for all states  $x$ . Then, for any  $(x, a) \in \mathcal{X} \times \mathcal{A}$ , we have  $|\hat{v}_t(x)| \leq \frac{4\tau+4}{\beta}$  and  $-\frac{2}{\beta}(2\tau+3) \leq \hat{q}_t(x, a) \leq \frac{2}{\beta} \left( \frac{|\mathcal{A}|}{\gamma} + 2\tau + 2 \right)$ .

*Proof:* Since  $|\sum_a \pi_t(a|x) \hat{r}_t(x, a)| = \pi_t(\mathbf{a}_t|x) \hat{r}_t(x, \mathbf{a}_t) \mathbb{I}_{\{x=\mathbf{x}_t\}} \leq 1/\mu_t^N(x) \leq 2/\beta$  by assumption and  $\hat{v}_t = v^{\pi_t}$ , the first statement of the lemma follows from Lemma 3.

To prove the bounds on  $\hat{q}_t$ , notice that

$$\begin{aligned} 0 &\leq \hat{\rho}_t = \sum_{x,a} \mu_{st}^{\pi_t}(x) \pi_t(a|x) \hat{r}_t(x, a) \\ &= \sum_x \mu_{st}^{\pi_t}(x) \sum_a \pi_t(a|x) \hat{r}_t(x, a) \leq \frac{2}{\beta}. \end{aligned}$$

Applying the above inequalities to the Bellman equations (5), we obtain  $\hat{q}_t(x, a) = \hat{r}_t(x, a) - \hat{\rho}_t + \sum_{x'} P(x'|x, a) \hat{v}_t(x') \geq -\frac{2}{\beta} - \frac{4\tau+4}{\beta} = -\frac{2}{\beta}(2\tau+3)$ . Since  $\pi_t(a|x) \geq \gamma/|\mathcal{A}|$ , the assumption on  $\mu_t^N$  and the definition of  $\hat{r}_t$  imply  $\hat{r}_t(x, a) \leq \frac{2|\mathcal{A}|}{\gamma\beta}$ . Thus, we get the upper bound  $\hat{q}_t(x, a) \leq \frac{2|\mathcal{A}|}{\gamma\beta} + \frac{4\tau+4}{\beta} = \frac{2}{\beta} \left( \frac{|\mathcal{A}|}{\gamma} + 2\tau + 2 \right)$ .  $\blacksquare$

The previous result can be strengthened if one is interested in a bound on  $\mathbb{E}[|\hat{v}_t(x)| | \mathbf{u}_{t-N}]$ :

**Lemma 7:** Let  $1 \leq t \leq T$  and assume that  $\mu_t^N(x) > 0$  holds for all states  $x$ . Then, for any  $x \in \mathcal{X}$ , we have  $\mathbb{E}[|\hat{v}_t(x)| | \mathbf{u}_{t-N}] \leq 2(\tau+1)$ .

*Proof:* Proceeding as in the proof of Lemma 3 and then taking expectations, we get

$$\begin{aligned} \mathbb{E}[|\hat{v}_t(x)| | \mathbf{u}_{t-N}] &\leq \sum_{s=1}^{\infty} \sum_{x'} |\nu_{s,x}^{\pi_t}(x') - \mu_{st}^{\pi_t}(x')| \mathbb{E} \left[ \sum_a \pi_t(a|x') \hat{r}_t(x', a) \middle| \mathbf{u}_{t-N} \right], \end{aligned}$$

where we have exploited that  $\hat{r}_t$  is well-defined by our assumption on  $\mu_t^N$  and it takes only nonnegative values. Now, by (9) and (11),

$$\begin{aligned} \mathbb{E} \left[ \sum_a \pi_t(a|x') \hat{r}_t(x', a) \middle| \mathbf{u}_{t-N} \right] &= \sum_a \pi_t(a|x') \mathbb{E} [\hat{r}_t(x', a) | \mathbf{u}_{t-N}] \\ &= \sum_a \pi_t(a|x') r_t(x', a), \end{aligned}$$

which is bounded between 0 and 1. Hence,  $\mathbb{E}[|\hat{v}_t(x)| | \mathbf{u}_{t-N}] \leq \sum_{s=1}^{\infty} \sum_{x'} |\nu_{s,x}^{\pi_t}(x') - \mu_{st}^{\pi_t}(x')|$ . Finishing as in the proof of Lemma 1 or 3, we get the statement.  $\blacksquare$

Similarly, we will also need a bound on the expected value of  $\mathbb{E}[|\hat{q}_t(x, a)| | \mathbf{u}_{t-N}]$ :

<sup>13</sup>With some extra work one can show that it is sufficient to choose  $T \geq h_1\left(C_4 \frac{\tau^3 |\mathcal{A}| \ln|\mathcal{A}|}{\beta^3}\right)$  with  $C_4 \geq 64$  and  $h_1(y) = y/(\ln y - \ln \ln y)$ .

**Lemma 8:** Let  $1 \leq t \leq T$  and assume that  $\mu_t^N(x) > 0$  holds for all states  $x$ . Then, for any  $(x, a) \in \mathcal{X} \times \mathcal{A}$ , we have that  $\mathbb{E}[|\hat{\mathbf{q}}_t(x, a)| | \mathbf{u}_{t-N}] \leq 2\tau + 4$  and also  $\mathbb{E}[\sum_a \pi_t(a|x) |\hat{\mathbf{q}}_t(x, a)|] \leq 2\tau + 4$ .

*Proof:* By the Bellman equations (5),

$$\begin{aligned} \mathbb{E}[|\hat{\mathbf{q}}_t(x, a)| | \mathbf{u}_{t-N}] &\leq \mathbb{E}[|\hat{\mathbf{r}}_t(x, a)| | \mathbf{u}_{t-N}] + \mathbb{E}[|\hat{\rho}_t| | \mathbf{u}_{t-N}] \\ &\quad + \sum_{x'} P(x'|x, a) \mathbb{E}[|\hat{\mathbf{v}}_t(x')| | \mathbf{u}_{t-N}]. \end{aligned}$$

As before,  $\mathbb{E}[|\hat{\mathbf{r}}_t(x, a)| | \mathbf{u}_{t-N}] \leq 1$ , and also  $\mathbb{E}[|\hat{\rho}_t| | \mathbf{u}_{t-N}] \leq 1$ . Combining these with the result of the previous lemma, we get the first part of the statement. To get the second part note that

$$\begin{aligned} \mathbb{E}\left[\sum_a \pi_t(a|x) |\hat{\mathbf{q}}_t(x, a)|\right] &= \mathbb{E}\left[\mathbb{E}\left[\sum_a \pi_t(a|x) |\hat{\mathbf{q}}_t(x, a)| \mid \mathbf{u}_{t-N}\right]\right] \\ &= \mathbb{E}\left[\sum_a \pi_t(a|x) \mathbb{E}[|\hat{\mathbf{q}}_t(x, a)| | \mathbf{u}_{t-N}]\right] \\ &\leq \mathbb{E}\left[\sum_a \pi_t(a|x) (2\tau + 4)\right] \leq 2\tau + 4. \end{aligned}$$

The quantity  $\pi_t(x, a) |\hat{\mathbf{q}}_t(x, a)|$  also enjoys a bound which is independent of the exploration rate  $\gamma$ :

**Lemma 9:** Let  $1 \leq t \leq T$  and assume that  $\mu_t^N(x) \geq \beta/2$  holds for all states  $x$ . Then, for any  $(x, a) \in \mathcal{X} \times \mathcal{A}$ , it holds that  $\pi_t(x, a) |\hat{\mathbf{q}}_t(x, a)| \leq \frac{4}{\beta}(\tau + 2)$ .

*Proof:* By our assumption on  $\mu_t^N$  and the construction of  $\hat{\mathbf{r}}_t(x, a)$ ,

$$\pi_t(x, a) |\hat{\mathbf{r}}_t(x, a)| \leq \frac{2}{\beta}. \quad (30)$$

Since  $\hat{\mathbf{r}}_t(x, a) = 0$  unless  $a = \mathbf{a}_t$ , Lemma 3 can be applied with  $R = 2/\beta$  to obtain  $|\hat{\mathbf{q}}_t(x, a)| \leq \frac{2}{\beta}(2\tau + 3) + |\hat{\mathbf{r}}_t(x, a)|$ . Multiplying both sides by  $\pi_t(x, a)$  and using (30) again finishes the proof. ■

Now we show that if the policies that we follow up to time step  $t$  change slowly,  $\mu_t^N$  is “close” to  $\mu_{st}^{\pi_t}$ :

**Lemma 10:** Let  $1 \leq N \leq t \leq T$  and  $c > 0$  be such that  $\|\pi_{s+1} - \pi_s\|_{1, \infty} \leq c$  holds for  $1 \leq s \leq t-1$ . Then we have  $\|\mu_t^N - \mu_{st}^{\pi_t}\|_1 \leq c(\tau + 1)^2 + 2e^{-(N-1)/\tau}$ .

*Proof:* This follows directly from Lemma 5 since, thanks to the recursive form of  $\mu_t^N$ ,  $\mu_t^N = \mu_1 P^{\pi_t - N + 1} \dots P^{\pi_t - 1}$ , where  $\mu_1 = e_{x_{t-N}} P^{\mathbf{a}_{t-N}}$  for  $t \geq N+1$  and  $\mu_1 = P_1$  if  $t = N$ . ■

In the lemma that follows we compute the rate of change of the policies produced by MDP-EXP3. We will use this lemma for multiple purposes, including showing that for a large enough value of  $N$ ,  $\mu_t^N$  can be uniformly bounded from below by  $\beta/2$ .

To simplify the presentation, we recall some short-hand notation from Proposition 1. In particular, we denote the lower and upper bounds for  $\hat{\mathbf{q}}_t$  by  $-L$  and  $V_{\hat{\mathbf{q}}}$  of Lemma 6, respectively, and the upper bound on  $|\hat{\mathbf{v}}_t|$  from the same lemma by  $U_{\hat{\mathbf{v}}}$ . Thus setting

$$L = \frac{2}{\beta}(2\tau + 3), \quad V_{\hat{\mathbf{q}}} = \frac{2}{\beta} \left( \frac{|\mathcal{A}|}{\gamma} + 2\tau + 2 \right), \quad U_{\hat{\mathbf{v}}} = \frac{4}{\beta}(\tau + 1)$$

we have  $-L \leq \hat{\mathbf{q}}_t(x, a) \leq V_{\hat{\mathbf{q}}}$ , and  $|\hat{\mathbf{v}}_t(x)| \leq U_{\hat{\mathbf{v}}}$  for all state-action pairs  $(x, a)$ .

**Lemma 11:** Assume that  $0 < \eta \leq 1/(V_{\hat{\mathbf{q}}} + L)$ . For  $(x, a) \in \mathcal{X} \times \mathcal{A}$ ,  $1 \leq t \leq T$ ,  $i \in \{0, 1\}$  let

$$\mathbf{d}_{t,i}(x, a) = \hat{\mathbf{q}}_t(x, a) - \sum_{b \in \mathcal{A}} \frac{\pi_{t+N-1+i}(b|x) - \frac{\gamma}{|\mathcal{A}|}}{1-\gamma} \hat{\mathbf{q}}_t(x, b). \quad (31)$$

Then, for all  $N \leq t \leq T$ ,

$$\begin{aligned} |\pi_{t+N-1}(a|x) - \pi_{t+N}(a|x)| \\ \leq \eta \pi_{t+N-1}(a|x) \max\{(e-1)\mathbf{d}_{t,0}(x, a), -\mathbf{d}_{t,1}(x, a)\}. \end{aligned}$$

*Proof:* Fix some state-action pair  $(x, a) \in \mathcal{X} \times \mathcal{A}$  and let  $\mathbf{W}_t(x) = \sum_a \mathbf{w}_t(x, a)$  where  $t = 1, 2, \dots, T$ . Since  $\mathbf{w}_{t+N}$  is computed using the exponential weight update for  $t \geq N$ , we have

$$\begin{aligned} &|\pi_{t+N-1}(a|x) - \pi_{t+N}(a|x)| \\ &= (1-\gamma) \left| \frac{\mathbf{w}_{t+N-1}(x, a)}{\mathbf{W}_{t+N-1}(x)} - \frac{\mathbf{w}_{t+N}(x, a)}{\mathbf{W}_{t+N}(x)} \right| \\ &= (1-\gamma) \frac{\mathbf{w}_{t+N-1}(x, a)}{\mathbf{W}_{t+N-1}(x)} \left| 1 - \frac{\mathbf{w}_{t+N}(x, a)}{\mathbf{w}_{t+N-1}(x, a)} \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} \right| \\ &= (1-\gamma) \frac{\mathbf{w}_{t+N-1}(x, a)}{\mathbf{W}_{t+N-1}(x)} \left| 1 - e^{\eta \hat{\mathbf{q}}_t(x, a)} \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} \right| \\ &\leq \pi_{t+N-1}(a|x) \left| 1 - e^{\eta \hat{\mathbf{q}}_t(x, a)} \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} \right|. \end{aligned} \quad (32)$$

We examine two separate cases depending on the sign of the expression in the absolute value on the right-hand side.

*Case a)*  $1 - e^{\eta \hat{\mathbf{q}}_t(x, a)} \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} \leq 0$ : First notice that the logarithm of the second term is positive by the condition, that is,  $\eta \hat{\mathbf{q}}_t(x, a) + \ln \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} \geq 0$ . Furthermore, it is bounded from above by 1. Indeed, by Jensen’s inequality,

$$\begin{aligned} \ln \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} &= -\ln \sum_b \frac{\mathbf{w}_{t+N-1}(x, b)}{\mathbf{W}_{t+N-1}(x)} e^{\eta \hat{\mathbf{q}}_t(x, b)} \\ &\leq -\eta \sum_b \frac{\mathbf{w}_{t+N-1}(x, b)}{\mathbf{W}_{t+N-1}(x)} \hat{\mathbf{q}}_t(x, b) \end{aligned} \quad (33)$$

and thus

$$\begin{aligned} \eta \hat{\mathbf{q}}_t(x, a) + \ln \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} \\ \leq \eta \hat{\mathbf{q}}_t(x, a) - \eta \sum_b \frac{\mathbf{w}_{t+N-1}(x, b)}{\mathbf{W}_{t+N-1}(x)} \hat{\mathbf{q}}_t(x, b) \\ \leq \eta(V_{\hat{\mathbf{q}}} + L) \leq 1, \end{aligned}$$

where the second inequality holds by our choice of  $V_{\hat{\mathbf{q}}}$  and  $L$ , while the third one holds by our assumption on  $\eta$ . Thus, using  $e^z - 1 \leq (e-1)z$ , which holds for any  $0 \leq z \leq 1$ , we get

$$\begin{aligned} &\left| 1 - e^{\eta \hat{\mathbf{q}}_t(x, a)} \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} \right| \\ &\leq (e-1) \left( \eta \hat{\mathbf{q}}_t(x, a) + \ln \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} \right). \end{aligned}$$

From this inequality, (32) and (33) and using the definition of  $\pi_{t+N-1}$ , we get

$$\begin{aligned} &|\pi_{t+N-1}(a|x) - \pi_{t+N}(a|x)| \\ &\leq \eta(e-1) \pi_{t+N-1}(a|x) \left( \hat{\mathbf{q}}_t(x, a) - \sum_b \frac{\pi_{t+N-1}(b|x) - \frac{\gamma}{|\mathcal{A}|}}{1-\gamma} \hat{\mathbf{q}}_t(x, b) \right). \end{aligned}$$

*Case b)*  $1 - e^{\eta \hat{\mathbf{q}}_t(x, a)} \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} \geq 0$ : Using  $1 - e^z \leq -z$  (which holds for all  $z \in \mathbb{R}$ ), we get

$$\left| 1 - e^{\eta \hat{\mathbf{q}}_t(x, a)} \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} \right| \leq -\eta \hat{\mathbf{q}}_t(x, a) - \ln \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)}.$$

Applying Jensen’s inequality, the second term can be bounded as

$$\begin{aligned} \ln \frac{\mathbf{W}_{t+N-1}(x)}{\mathbf{W}_{t+N}(x)} &= \ln \sum_b \frac{\mathbf{w}_{t+N-1}(x, b)}{\mathbf{W}_{t+N-1}(x)} e^{-\eta \hat{\mathbf{q}}_t(x, b)} \\ &\geq -\eta \sum_b \frac{\pi_{t+N-1}(b|x) - \frac{\gamma}{|\mathcal{A}|}}{1-\gamma} \hat{\mathbf{q}}_t(x, b). \end{aligned}$$

Combining these inequalities with (32), we get

$$\begin{aligned} & |\pi_{t+N-1}(a|x) - \pi_{t+N}(a|x)| \\ & \leq \eta \pi_{t+N-1}(a|x) \left( -\hat{\mathbf{q}}_t(x, a) + \sum_b \frac{\pi_{t+N}(b|x) - \frac{\gamma}{|\mathcal{A}|}}{1-\gamma} \hat{\mathbf{q}}_t(x, b) \right). \end{aligned} \quad \leq \frac{1}{1-\gamma} \left( \pi_t(a|x) \hat{\mathbf{q}}_t(x, a) + \mathbf{F}_t(x, a) V_{\hat{\mathbf{q}}} + \pi_{t+N-1}(a|x) (L + \gamma V_{\hat{\mathbf{q}}}) \right).$$

The two cases together prove the lemma.  $\blacksquare$

By the lemma just proved, the rate of change  $\{\pi_t\}$  is partially governed by  $\{\mathbf{d}_{t,i}\}$ . To further bound the rate of change, we develop lower and upper bounds on  $\mathbf{d}_{t,i}$ . To facilitate this, we rewrite  $\mathbf{d}_{t,i}$  by grouping the terms with identical signs:

$$\begin{aligned} \mathbf{d}_{t,i}(x, a) &= \hat{\mathbf{q}}_t(x, a) + \frac{\gamma}{1-\gamma} \frac{1}{|\mathcal{A}|} \sum_{b \in \mathcal{A}} \hat{\mathbf{q}}_t(x, b) \\ &\quad - \frac{1}{1-\gamma} \sum_{b \in \mathcal{A}} \pi_{t+N-1+i}(b|x) \hat{\mathbf{q}}_t(x, b). \end{aligned} \quad (34)$$

Then, from  $-L \leq \hat{\mathbf{q}}_t(x, a) \leq V_{\hat{\mathbf{q}}}$ , as long as  $0 \leq \gamma < 1$ , we have

$$\begin{aligned} \frac{-L + \sum_b \pi_{t+N-1+i}(b|x) \hat{\mathbf{q}}_t(x, b)}{1-\gamma} &\leq \mathbf{d}_{t,i}(x, a) \\ &\leq \hat{\mathbf{q}}_t(x, a) + \frac{L + \gamma V_{\hat{\mathbf{q}}}}{1-\gamma}. \end{aligned} \quad (35)$$

Notice that since  $V_{\hat{\mathbf{q}}}$  scales with  $1/\gamma$ , we avoided upper bounding  $\hat{\mathbf{q}}_t$  by  $V_{\hat{\mathbf{q}}}$  except when  $\hat{\mathbf{q}}_t$  is multiplied by  $\gamma$ , and so the bounds will not “blow up” as  $\gamma \rightarrow 0$ . In fact, this is one of the main reasons that in this paper we succeed in proving an  $\tilde{O}(T^{1/2})$  regret bound as compared to the  $O(T^{2/3})$  regret bound of [16].

Let us now show that, provided  $\mu_t^N$  is uniformly bounded away from zero, the sequence  $\{\pi_t\}$  changes slowly.

**Lemma 12:** Assume that  $0 \leq \gamma < 1$ ,  $0 < \eta < \min\left(\frac{1}{V_{\hat{\mathbf{q}}+L}, \frac{1-\gamma}{eNV_{\hat{\mathbf{q}}}}}\right) = \frac{\beta(1-\gamma)}{2eN(|\mathcal{A}|/\gamma + 2\tau + 2)}$  and that for all  $N \leq t \leq T$  and states  $x$ ,  $\mu_t^N(x) \geq \beta/2$  holds true. Set  $c = \eta e^{\frac{U_{\hat{\mathbf{v}}} + L + \gamma V_{\hat{\mathbf{q}}}}{1-\gamma - \eta eNV_{\hat{\mathbf{q}}}}}$ . Then, for any  $1 \leq t \leq T$ ,

$$\|\pi_{t+N-1} - \pi_{t+N}\|_{1,\infty} \leq c. \quad (36)$$

*Proof:* We prove the statement by induction on  $t$ . To show the bound for time step  $t$  assume that  $\|\pi_{s+N-1} - \pi_{s+N}\|_{1,\infty} \leq c$  holds for all  $s = 1, 2, \dots, t-1$ . As  $\pi_{s+N-1} = \pi_{s+N}$  for all  $s = 1, 2, \dots, N-1$ , the assumption holds for  $t = 1, \dots, N-1$  and we are left with proving the induction step for  $t \geq N$ .

Fix  $x \in \mathcal{X}$ . For any  $a \in \mathcal{A}$ , by Lemma 11,

$$\begin{aligned} & |\pi_{t+N-1}(a|x) - \pi_{t+N}(a|x)| \\ & \leq \eta \pi_{t+N-1}(a|x) \max\{(e-1)\mathbf{d}_{t,0}(x, a), -\mathbf{d}_{t,1}(x, a)\}. \end{aligned}$$

Our goal is to upper-bound  $\pi_{t+N-1}(a|x)\mathbf{d}_{t,0}(x, a)$  and lower-bound  $\pi_{t+N-1}(a|x)\mathbf{d}_{t,1}(x, a)$ . As before, we make an effort to avoid terms that scale with  $1/\gamma$ , but we allow terms that scale with  $c/\gamma$  as  $c$  will be seen to scale with  $\gamma$  (and  $\eta$ ).

Consider first an upper bound on  $\pi_{t+N-1}(a|x)\mathbf{d}_{t,0}(x, a)$ . From (35), it remains to bound  $\pi_{t+N-1}(a|x)\hat{\mathbf{q}}_t(x, a)$ . By a simple telescoping argument, we bound this by

$$\begin{aligned} & \pi_{t+N-1}(a|x) \hat{\mathbf{q}}_t(x, a) \\ &= \pi_t(a|x) \hat{\mathbf{q}}_t(x, a) + \left[ \sum_{s=t}^{t+N-2} (\pi_{s+1}(a|x) - \pi_s(a|x)) \right] \hat{\mathbf{q}}_t(x, a) \\ &\leq \pi_t(a|x) \hat{\mathbf{q}}_t(x, a) + \mathbf{F}_t(x, a) V_{\hat{\mathbf{q}}}, \end{aligned}$$

where we have introduced

$$\mathbf{F}_t(x, a) = \sum_{s=t}^{t+N-2} |\pi_{s+1}(a|x) - \pi_s(a|x)|.$$

Now, using (34),

$$\pi_{t+N-1}(a|x)\mathbf{d}_{t,0}(x, a) \leq \frac{1}{1-\gamma} \left( \pi_t(a|x) \hat{\mathbf{q}}_t(x, a) + \mathbf{F}_t(x, a) V_{\hat{\mathbf{q}}} + \pi_{t+N-1}(a|x) (L + \gamma V_{\hat{\mathbf{q}}}) \right).$$

Now, let us consider upper bounding  $-\pi_{t+N-1}(a|x)\mathbf{d}_{t,1}(x, a)$ . In this case, we use telescoping for the second term on the left-hand side of (35):

$$\begin{aligned} & \sum_b \pi_{t+N}(b|x) \hat{\mathbf{q}}_t(x, b) \\ &= \sum_b \pi_t(b|x) \hat{\mathbf{q}}_t(x, b) + \sum_{s=t}^{t+N-1} \sum_b (\pi_{s+1}(b|x) - \pi_s(b|x)) \hat{\mathbf{q}}_t(x, b) \\ &\leq \hat{\mathbf{v}}_t(x) + NcV_{\hat{\mathbf{q}}} \leq U_{\hat{\mathbf{v}}} + NcV_{\hat{\mathbf{q}}}, \end{aligned}$$

where we used  $\hat{\mathbf{q}}_t(x, a) \leq V_{\hat{\mathbf{q}}}$ ,  $\hat{\mathbf{v}}_t(x) \leq U_{\hat{\mathbf{v}}}$  and the induction hypothesis. Plugging this into the lower bound in (35), we get

$$-\pi_{t+N-1}(a|x)\mathbf{d}_{t,1}(x, a) \leq \frac{\pi_{t+N-1}(a|x)}{1-\gamma} (L + U_{\hat{\mathbf{v}}} + NcV_{\hat{\mathbf{q}}}).$$

Combining the two cases, we get

$$\begin{aligned} & |\pi_{t+N-1}(a|x) - \pi_{t+N}(a|x)| \\ & \leq \max \left( \frac{\eta(e-1)}{1-\gamma} \left( \pi_t(a|x) \hat{\mathbf{q}}_t(x, a) + \mathbf{F}_t(x, a) V_{\hat{\mathbf{q}}} \right. \right. \\ & \quad \left. \left. + \pi_{t+N-1}(a|x) (L + \gamma V_{\hat{\mathbf{q}}}) \right), \right. \\ & \quad \left. \frac{\eta \pi_{t+N-1}(a|x)}{1-\gamma} (U_{\hat{\mathbf{v}}} + NcV_{\hat{\mathbf{q}}} + L) \right) \\ & \leq \frac{\eta}{1-\gamma} \left( (e-1) (\pi_t(a|x) \hat{\mathbf{q}}_t(x, a) + \mathbf{F}_t(x, a) V_{\hat{\mathbf{q}}}) \right. \\ & \quad \left. + \pi_{t+N-1}(a|x) (U_{\hat{\mathbf{v}}} + NcV_{\hat{\mathbf{q}}}) \right. \\ & \quad \left. + (e-1) \pi_{t+N-1}(a|x) (L + \gamma V_{\hat{\mathbf{q}}}) \right). \end{aligned}$$

Summing these inequalities for all  $a$  and taking the maximum over  $x$  gives

$$\|\pi_{t+N-1} - \pi_{t+N}\|_{1,\infty} \leq \frac{\eta e}{1-\gamma} (U_{\hat{\mathbf{v}}} + NcV_{\hat{\mathbf{q}}} + L + \gamma V_{\hat{\mathbf{q}}}),$$

where we upper bounded  $(e-1)(L + \gamma V_{\hat{\mathbf{q}}})$  by  $e(L + \gamma V_{\hat{\mathbf{q}}})$  and used that the inequality  $\sum_a \mathbf{F}_t(x, a) \leq (N-1)c \leq Nc$  holds by the induction hypothesis. Now, the result follows because, thanks to the definition of  $c$ , the right-hand side equals  $c$  (in fact, this is how the definition of  $c$  is obtained).  $\blacksquare$

**Lemma 13:** Let  $c, \gamma$  be as in Lemma 12. Assume further that  $c(\tau + 1)^2 < \beta/2$ , and let

$$N \geq 1 + \left\lceil \tau \ln \left( \frac{4}{\beta - 2c(\tau + 1)^2} \right) \right\rceil. \quad (37)$$

Then, for all  $N \leq t \leq T$ ,  $x \in \mathcal{X}$ , we have  $\mu_t^N(x) \geq \beta/2$  and  $\|\pi_{t+1} - \pi_t\|_{1,\infty} \leq c$ .

*Proof:* We prove the lemma by induction on  $t$ . The induction hypothesis is that for  $N \leq t \leq T$ ,  $\min_x \mu_s^N(x) \geq \beta/2$  and  $\max_x \sum_a |\pi_{s+1}(a|x) - \pi_s(a|x)| \leq c$  hold for all  $N \leq s \leq t$ .

Let us first show that this hypothesis holds when  $N \leq t \leq 2N-2$ . By the construction of the policies, we have  $\max_x \sum_a |\pi_{t+1}(a|x) - \pi_t(a|x)| = 0 \leq c$  for all  $1 \leq t \leq 2N-2$ . Thus, by Lemma 10, we get that  $\|\mu_t^N - \mu_{st}^{\pi_t}\|_1 \leq c(\tau + 1)^2 + 2e^{-(N-1)/\tau}$  holds for all  $N \leq t \leq 2N-2$ . By our assumption about  $N$ , we have

$$c(\tau + 1)^2 + 2e^{-(N-1)/\tau} \leq \beta/2, \quad (38)$$

thus for any  $N \leq t \leq 2N - 2$ ,

$$\left\| \boldsymbol{\mu}_t^N - \boldsymbol{\mu}_{st}^{\pi_t} \right\|_{\infty} \leq \left\| \boldsymbol{\mu}_t^N - \boldsymbol{\mu}_{st}^{\pi_t} \right\|_1 \leq \beta/2. \quad (39)$$

Since, by assumption,  $\mu_{st}^{\pi_t}(x) \geq \beta$  holds for any stationary policy  $\pi$ , we also have  $\mu_{st}^{\pi_t}(x) \geq \beta$  ( $x \in \mathcal{X}$ ). This, together with (39) gives that  $\boldsymbol{\mu}_t^N(x) \geq \beta/2$  holds for any  $x \in \mathcal{X}$ .

Now, fix a time index  $2N - 1 \leq t \leq T$  and assume that the induction hypothesis holds for time  $t - 1$ . Then, thanks to  $\min_x \mu_{t-N+1}^N(x) \geq \beta/2$ , Lemma 12 implies  $\|\boldsymbol{\pi}_{t+1} - \boldsymbol{\pi}_t\|_{1,\infty} \leq c$ . Now, by Lemma 10, we have  $\|\boldsymbol{\mu}_t^N - \boldsymbol{\mu}_{st}^{\pi_t}\|_1 \leq c(\tau + 1)^2 + 2e^{-(N-1)/\tau}$ . Using the same reasoning as above, we finish the inductive step and thus the proof.  $\blacksquare$

In our final result we study the weighted sums

$$\Delta_{t,s}(x) \stackrel{\text{def}}{=} \sum_{a \in \mathcal{A}} |\boldsymbol{\pi}_{t+1}(a|x) - \boldsymbol{\pi}_t(a|x)| |\hat{\mathbf{q}}_s(x, a)|$$

for  $x \in \mathcal{X}$ ,  $t, s \geq 1$ . To state this result recall the definitions

$$U_{\pi\hat{\mathbf{q}}} = \frac{4}{\beta}(\tau + 2), \quad U_{\hat{\mathbf{q}}} = 2\tau + 4,$$

from Proposition 1. Note that by Lemma 9,  $\boldsymbol{\pi}_t(x, a) |\hat{\mathbf{q}}_t(x, a)| \leq U_{\pi\hat{\mathbf{q}}}$  and by Lemma 8,  $\mathbb{E} [|\hat{\mathbf{q}}_s(x, a)|] \leq U_{\hat{\mathbf{q}}}$ .

**Lemma 14:** Let  $e' = e - 1$ ,  $C_1 = e'(L + \gamma V_{\hat{\mathbf{q}}}) + (e' U_{\pi\hat{\mathbf{q}}} + U_{\hat{\mathbf{v}}}) U_{\hat{\mathbf{q}}} |\mathcal{A}|$ , and  $c' = \frac{\eta C_1}{1 - \gamma - \eta e^{(N+1)} V_{\hat{\mathbf{q}}}}$ . Assume that  $0 < \gamma \leq 1$ ,  $0 \leq \eta < \min\left(\frac{1}{V_{\hat{\mathbf{q}}+L}, \frac{1-\gamma}{e^{(N+1)} V_{\hat{\mathbf{q}}}}}\right) = \frac{\beta(1-\gamma)}{2e^{(N+1)}(|\mathcal{A}|/\gamma + 2\tau + 2)}$  and that for all  $N \leq t \leq T$  and states  $x$ ,  $\boldsymbol{\mu}_t^N(x) \geq \beta/2$  holds true. Then, for any  $1 \leq t \leq T$ ,  $1 \leq s \leq T$  and  $x \in \mathcal{X}$ , it holds that  $\mathbb{E} [\Delta_{t,s}(x)] \leq c'$ . An observation that will be needed later is that the conditions of this lemma on  $\eta, \gamma$  and  $N$  imply those of Lemma 12.

*Proof:* Fix  $x \in \mathcal{X}$ . We will prove the result by induction. Since in the algorithm the weights are kept fixed for  $1 \leq t \leq 2N - 1$ ,  $\Delta_{t,s} = 0$  when  $1 \leq t \leq 2N - 2$ . This establishes the base case of the induction. Thus, fix  $t \geq N$  and assume that  $\mathbb{E} [\Delta_{t',s}(x)] \leq c'$  holds for all pairs  $(t', s)$  such that  $1 \leq t' \leq t + N - 2$ ,  $1 \leq s \leq T$ .

By Lemma 11,

$$\begin{aligned} \Delta_{t+N-1,s}(x) &\leq \eta \sum_a \boldsymbol{\pi}_{t+N-1}(a|x) \max \{e' \mathbf{d}_{t,0}(x, a), -\mathbf{d}_{t,1}(x, a)\} |\hat{\mathbf{q}}_s(x, a)|. \end{aligned} \quad (40)$$

As in the proof of Lemma 12, we upper bound the two terms resulting from the maximum on the right-hand side of the above expression separately. Considering the first of these, using the upper bound from (35), we get

$$\boldsymbol{\pi}_{t+N-1}(a|x) \mathbf{d}_{t,0}(x, a) \leq \boldsymbol{\pi}_{t+N-1}(a|x) \frac{\hat{\mathbf{q}}_t(x, a) + L + \gamma V_{\hat{\mathbf{q}}}}{1 - \gamma}.$$

We use telescoping to bound the first term in the numerator on the right-hand side:

$$\begin{aligned} \boldsymbol{\pi}_{t+N-1}(a|x) \hat{\mathbf{q}}_t(x, a) &\leq \boldsymbol{\pi}_t(a|x) \hat{\mathbf{q}}_t(x, a) \\ &\quad + \sum_{t'=t}^{t+N-2} |\boldsymbol{\pi}_{t'+1}(a|x) - \boldsymbol{\pi}_{t'}(a|x)| |\hat{\mathbf{q}}_t(x, a)|. \end{aligned}$$

Hence,

$$\begin{aligned} \boldsymbol{\pi}_{t+N-1}(a|x) \mathbf{d}_{t,0}(x, a) &\leq \\ &\frac{1}{1 - \gamma} \left( \boldsymbol{\pi}_{t+N-1}(a|x) (L + \gamma V_{\hat{\mathbf{q}}}) + \boldsymbol{\pi}_t(a|x) \hat{\mathbf{q}}_t(x, a) \right. \\ &\quad \left. + \sum_{t'=t}^{t+N-1} |\boldsymbol{\pi}_{t'+1}(a|x) - \boldsymbol{\pi}_{t'}(a|x)| |\hat{\mathbf{q}}_t(x, a)| \right). \end{aligned}$$

Now, considering the second branch of the maximum, using this time the lower bound from (35),

$$\begin{aligned} & - \boldsymbol{\pi}_{t+N-1}(a|x) \mathbf{d}_{t,1}(x, a) \\ & \leq \boldsymbol{\pi}_{t+N-1}(a|x) \frac{L + \sum_b \boldsymbol{\pi}_{t+N}(b|x) \hat{\mathbf{q}}_t(x, b)}{1 - \gamma} \\ & \leq \frac{\boldsymbol{\pi}_{t+N-1}(a|x)}{1 - \gamma} \left( L + \left| \sum_b \boldsymbol{\pi}_t(b|x) \hat{\mathbf{q}}_t(x, b) \right| \right. \\ & \quad \left. + \sum_b \sum_{t'=t}^{t+N-1} |\boldsymbol{\pi}_{t'+1}(b|x) - \boldsymbol{\pi}_{t'}(b|x)| |\hat{\mathbf{q}}_t(x, b)| \right). \end{aligned}$$

Combining these two inequalities, introducing  $\hat{C} = e' \frac{L + \gamma V_{\hat{\mathbf{q}}}}{1 - \gamma}$ , we get

$$\begin{aligned} & \boldsymbol{\pi}_{t+N-1}(a|x) \max \{e' \mathbf{d}_{t,0}(x, a), -\mathbf{d}_{t,1}(x, a)\} \\ & \leq \boldsymbol{\pi}_{t+N-1}(a|x) \hat{C} + \frac{e'}{1 - \gamma} \boldsymbol{\pi}_t(a|x) |\hat{\mathbf{q}}_t(x, a)| \\ & \quad + \frac{e'}{1 - \gamma} \sum_{t'=t}^{t+N-2} |\boldsymbol{\pi}_{t'+1}(a|x) - \boldsymbol{\pi}_{t'}(a|x)| |\hat{\mathbf{q}}_t(x, a)| \\ & \quad + \frac{\boldsymbol{\pi}_{t+N-1}(a|x)}{1 - \gamma} \left[ \left| \sum_b \boldsymbol{\pi}_t(b|x) \hat{\mathbf{q}}_t(x, b) \right| \right. \\ & \quad \left. + \sum_b \sum_{t'=t}^{t+N-1} |\boldsymbol{\pi}_{t'+1}(b|x) - \boldsymbol{\pi}_{t'}(b|x)| |\hat{\mathbf{q}}_t(x, a)| \right]. \end{aligned}$$

Plugging this into (40), we get

$$\begin{aligned} \Delta_{t+N-1,s}(x) &\leq \eta \hat{C} \\ & \quad + \frac{\eta e'}{1 - \gamma} \sum_a \boldsymbol{\pi}_t(a|x) |\hat{\mathbf{q}}_t(x, a)| |\hat{\mathbf{q}}_s(x, a)| \\ & \quad + \frac{\eta e'}{1 - \gamma} \sum_{t'=t}^{t+N-2} \sum_a |\boldsymbol{\pi}_{t'+1}(a|x) - \boldsymbol{\pi}_{t'}(a|x)| |\hat{\mathbf{q}}_t(x, a)| |\hat{\mathbf{q}}_s(x, a)| \\ & \quad + \eta \sum_a |\hat{\mathbf{q}}_s(x, a)| \frac{\boldsymbol{\pi}_{t+N-1}(a|x)}{1 - \gamma} \left[ \left| \sum_b \boldsymbol{\pi}_t(b|x) \hat{\mathbf{q}}_t(x, b) \right| + \sum_{t'=t}^{t+N-1} \Delta_{t',t}(x) \right]. \end{aligned}$$

Using  $\boldsymbol{\pi}_t(a|x) |\hat{\mathbf{q}}_t(x, a)| \leq U_{\pi\hat{\mathbf{q}}}$ ,  $|\hat{\mathbf{v}}_t(x)| \leq U_{\hat{\mathbf{v}}}$  and  $|\hat{\mathbf{q}}_s(x, a)| \leq V_{\hat{\mathbf{q}}}$  (where the last inequality holds thanks to  $\gamma \leq 1$ ), we obtain

$$\begin{aligned} \Delta_{t+N-1,s}(x) &\leq \eta \hat{C} + \frac{\eta e' U_{\pi\hat{\mathbf{q}}}}{1 - \gamma} \sum_a |\hat{\mathbf{q}}_s(x, a)| + \frac{\eta e' V_{\hat{\mathbf{q}}}}{1 - \gamma} \sum_{t'=t}^{t+N-2} \Delta_{t',t}(x) \\ & \quad + \eta \sum_a |\hat{\mathbf{q}}_s(x, a)| \frac{\boldsymbol{\pi}_{t+N-1}(a|x)}{1 - \gamma} \left[ U_{\hat{\mathbf{v}}} + \sum_{t'=t}^{t+N-1} \Delta_{t',t}(x) \right] \\ &\leq \eta \hat{C} + \frac{\eta e' U_{\pi\hat{\mathbf{q}}}}{1 - \gamma} \sum_a |\hat{\mathbf{q}}_s(x, a)| + \frac{\eta e' V_{\hat{\mathbf{q}}}}{1 - \gamma} \sum_{t'=t}^{t+N-2} \Delta_{t',t}(x) \\ & \quad + \frac{\eta U_{\hat{\mathbf{v}}}}{1 - \gamma} \sum_a |\hat{\mathbf{q}}_s(x, a)| + \frac{\eta V_{\hat{\mathbf{q}}}}{1 - \gamma} \sum_{t'=t}^{t+N-1} \Delta_{t',t}(x) \\ & = \eta \hat{C} + \frac{\eta (e' U_{\pi\hat{\mathbf{q}}} + U_{\hat{\mathbf{v}}})}{1 - \gamma} \sum_a |\hat{\mathbf{q}}_s(x, a)| + \frac{\eta e V_{\hat{\mathbf{q}}}}{1 - \gamma} \sum_{t'=t}^{t+N-1} \Delta_{t',t}(x). \end{aligned}$$

Now, take the expectation of both sides and use that  $\mathbb{E} [|\hat{\mathbf{q}}_s(x, a)|] \leq U_{\hat{\mathbf{q}}}$ . Introducing the constant  $C_1 = \hat{C}(1 - \gamma) + (e' U_{\pi\hat{\mathbf{q}}} + U_{\hat{\mathbf{v}}}) U_{\hat{\mathbf{q}}} |\mathcal{A}| = e'(L + \gamma V_{\hat{\mathbf{q}}}) + (e' U_{\pi\hat{\mathbf{q}}} + U_{\hat{\mathbf{v}}}) U_{\hat{\mathbf{q}}} |\mathcal{A}|$ , we get

$$\mathbb{E} [\Delta_{t+N-1,s}(x)] \leq \frac{\eta C_1}{1 - \gamma} + \frac{\eta e V_{\hat{\mathbf{q}}}}{1 - \gamma} \sum_{t'=t}^{t+N-1} \mathbb{E} [\Delta_{t',t}(x)]. \quad (41)$$

Taking  $s = t$  in (41), using  $1 - \frac{\eta e V_{\hat{q}}}{1-\gamma} > 0$  which holds by our assumption on  $\eta$  and  $\gamma$ , reordering gives

$$\begin{aligned} & \mathbb{E} [\Delta_{t+N-1,t}(x)] \\ & \leq \frac{\eta(1-\gamma)}{1-\gamma-\eta e V_{\hat{q}}} \left\{ \frac{C_1}{1-\gamma} + \frac{e V_{\hat{q}}}{1-\gamma} \sum_{t'=t}^{t+N-2} \mathbb{E} [\Delta_{t',t}(x)] \right\} \\ & \leq \frac{\eta}{1-\gamma-\eta e V_{\hat{q}}} \{C_1 + e V_{\hat{q}} (N-1)c'\} \\ & \leq \frac{\eta}{1-\gamma-\eta e V_{\hat{q}}} \{C_1 + e V_{\hat{q}} N c'\} = c', \end{aligned} \quad (42)$$

where the second inequality follows since, by our induction hypothesis,  $\mathbb{E} [\Delta_{t',t}(x)] \leq c'$  holds for any  $t'$  such that  $1 \leq t' \leq t+N-2$ , the third follows by our assumptions on  $N, \eta$  and  $\gamma$ , while the last equality holds by the definition of  $c'$ . This shows that the induction hypothesis holds for the pair  $(t+N-1, t)$ .

Let us now consider the pairs  $(t+N-1, s)$ , where  $s \neq t$ . We start from (41) again. Note that by our induction hypothesis,  $\mathbb{E} [\Delta_{t',t}(x)] \leq c'$  for  $t \leq t' \leq t+N-2$ . Furthermore, by (42), we also have  $\mathbb{E} [\Delta_{t+N-1,t}(x)] \leq c'$ . Hence,

$$\begin{aligned} \mathbb{E} [\Delta_{t+N-1,s}(x)] & \leq \frac{\eta C_1}{1-\gamma} + \frac{\eta e V_{\hat{q}}}{1-\gamma} N c' \\ & \leq \frac{\eta C_1}{1-\gamma-\eta e V_{\hat{q}}} + \frac{\eta e V_{\hat{q}}}{1-\gamma-\eta e V_{\hat{q}}} N c' = c'. \end{aligned}$$

■

#### D. Proof of Proposition 1

For every  $x, a$  define  $\mathbf{Q}_T(x, a) = \sum_{t=N}^T \mathbf{q}_t(x, a)$  and  $\mathbf{V}_T(x) = \sum_{t=N}^T \mathbf{v}_t(x)$ . Lemma 2 shows that in order to prove Proposition 1, it suffices to prove an upper bound on  $\mathbb{E} [\mathbf{Q}_T(x, a) - \mathbf{V}_T(x)]$ .

**Lemma 15:** Let  $c$  be as in Lemma 12 and  $c'$  be as in Lemma 14. Assume that  $\gamma \in (0, 1)$ ,  $c(\tau+1)^2 < \beta/2$ ,  $N \geq 1 + \lceil \tau \ln \left( \frac{4}{\beta-2c(\tau+1)^2} \right) \rceil$ ,  $0 < \eta < \frac{\beta(1-\gamma)}{2e(N+1)(|\mathcal{A}|/\gamma+2\tau+2)}$ , and  $T \geq N$  hold. Then, for all  $(x, a) \in \mathcal{X} \times \mathcal{A}$ ,

$$\begin{aligned} \mathbb{E} [\mathbf{Q}_T(x, b) - \mathbf{V}_T(x)] & \leq \frac{\ln|\mathcal{A}|}{\eta} + (N-1)(U_{\hat{q}} + U_q) \\ & \quad + (T-2N+2) \left( c'(N-1)(1+\eta e'' V_{\hat{q}}) \right. \\ & \quad \left. + \gamma U_q + \eta e'' |\mathcal{A}| U_{\pi_{\hat{q}}} U_{\hat{q}} \right). \end{aligned}$$

*Proof:* Note that if  $T < 2N$ , the conclusion of the lemma trivially holds. Therefore, in what follows we assume that  $T \geq 2N$ . First, note that the conditions of both Lemmas 13 and 9 are satisfied. Thus, the conclusions of Lemma 9 and therefore also those of Lemmas 6–9 hold. In particular, by Lemma 9,  $\pi_t(a|x)|\hat{q}_t(x, a)| \leq U_{\pi_{\hat{q}}} = \frac{4}{\beta}(\tau+2)$  holds for any  $t \geq N+1$ . Now, using Lemma 6, we have  $\hat{q}_t(x, a) \leq V_{\hat{q}} = \frac{2}{\beta}(|\mathcal{A}|/\gamma+2\tau+2)$ , thus by the constraint on  $\eta$ ,  $\eta \hat{q}_t(x, a) \leq 1$ .

We will follow the steps of the proof in Auer et al. [2]. For  $1 \leq t \leq T$ , define  $\mathbf{W}_t(x) = \sum_a \mathbf{w}_t(x, a)$ . Fix a time step  $t$  such that  $2N-1 \leq t \leq T$  and a state-action pair  $(x, a) \in \mathcal{X} \times \mathcal{A}$ . Recalling  $e'' = e - 2$ , we have the following inequalities:

$$\begin{aligned} \frac{\mathbf{W}_{t+1}(x)}{\mathbf{W}_t(x)} & = \sum_a \frac{\mathbf{w}_{t+1}(x, a)}{\mathbf{W}_t(x)} = \sum_a \frac{\mathbf{w}_t(x, a)}{\mathbf{W}_t(x)} e^{\eta \hat{q}_{t-N+1}(x, a)} \\ & = \sum_a \frac{\pi_t(a|x) - \gamma/|\mathcal{A}|}{1-\gamma} e^{\eta \hat{q}_{t-N+1}(x, a)} \\ & \leq \sum_a \frac{\pi_t(a|x) - \gamma/|\mathcal{A}|}{1-\gamma} \left( 1 + \eta \hat{q}_{t-N+1}(x, a) + e'' (\eta \hat{q}_{t-N+1}(x, a))^2 \right) \end{aligned}$$

(as  $\eta \hat{q}_{t-N+1}(x, a) \leq 1$  and for  $x \leq 1$ ,  $e^x \leq 1 + x + e'' x^2$ ) where  $\hat{\mathbf{Q}}_T^N(x, b) = \sum_{t=N}^{T-N+1} \hat{q}_t(x, b)$ .

$$\begin{aligned} & \leq 1 + \frac{\eta}{1-\gamma} \sum_a \pi_t(a|x) \hat{q}_{t-N+1}(x, a) \\ & \quad + \frac{\eta^2 e''}{1-\gamma} \sum_a \pi_t(a|x) (\hat{q}_{t-N+1}(x, a))^2. \end{aligned}$$

Introduce

$$\mathbf{q}_{2,t-N+1}(x) = \sum_a \pi_t(a|x) (\hat{q}_{t-N+1}(x, a))^2.$$

We now show a bound on the expectation of this quantity that will be useful later. For this, write

$$\begin{aligned} \mathbf{q}_{2,t-N+1}(x) & = \sum_a \pi_t(a|x) (\hat{q}_{t-N+1}(x, a))^2 \\ & = \sum_a \pi_{t-N+1}(a|x) (\hat{q}_{t-N+1}(x, a))^2 \\ & \quad + \sum_a (\pi_t(a|x) - \pi_{t-N+1}(a|x)) (\hat{q}_{t-N+1}(x, a))^2 \end{aligned}$$

By Lemma 9, the first term on the right-hand side can be bounded as follows:

$$\sum_a \pi_{t-N+1}(a|x) (\hat{q}_{t-N+1}(x, a))^2 \leq U_{\pi_{\hat{q}}} \sum_a |\hat{q}_{t-N+1}(x, a)|,$$

while, thanks to Lemma 14, the second one is bounded, in expectation, by

$$\begin{aligned} & \mathbb{E} \left[ \sum_a (\pi_t(a|x) - \pi_{t-N+1}(a|x)) (\hat{q}_{t-N+1}(x, a))^2 \right] \\ & \leq \mathbb{E} \left[ V_{\hat{q}} \sum_{t'=t-N+1}^{t-1} \Delta_{t',t-N+1}(x) \right] \leq (N-1) c' V_{\hat{q}}, \end{aligned}$$

where we have used that  $|\hat{q}_t(x, a)| \leq V_{\hat{q}}$  and also that  $\mathbb{E} [\Delta_{t',s}(x)] \leq c'$ . Combining these inequalities, we get that

$$\mathbb{E} [\mathbf{q}_{2,t-N+1}(x)] \leq |\mathcal{A}| U_{\pi_{\hat{q}}} U_{\hat{q}} + (N-1) c' V_{\hat{q}}. \quad (43)$$

Let us now return to developing an upper-bound on  $\frac{\mathbf{W}_{t+1}(x)}{\mathbf{W}_t(x)}$ . Defining  $\hat{\mathbf{v}}_t^N(x) = \sum_a \pi_t(a|x) \hat{q}_{t-N+1}(x, a)$ , we obtain

$$\frac{\mathbf{W}_{t+1}(x)}{\mathbf{W}_t(x)} \leq 1 + \frac{\eta}{1-\gamma} \hat{\mathbf{v}}_t^N(x) + \frac{\eta^2 e''}{1-\gamma} \mathbf{q}_{2,t-N+1}(x).$$

Using  $1 + x \leq e^x$  and then taking logarithms gives

$$\ln \frac{\mathbf{W}_{t+1}(x)}{\mathbf{W}_t(x)} \leq \frac{\eta}{1-\gamma} \hat{\mathbf{v}}_t^N(x) + \frac{\eta^2 e''}{1-\gamma} \mathbf{q}_{2,t-N+1}(x).$$

Summing over  $t = 2N-1, 2N, \dots, T$ , we get

$$\ln \frac{\mathbf{W}_{T+1}(x)}{\mathbf{W}_{2N-1}(x)} \leq \frac{\eta}{1-\gamma} \hat{\mathbf{V}}_T^N(x) + \frac{\eta^2 e''}{(1-\gamma)} \mathbf{Q}_{2,T}^N(x), \quad (44)$$

where  $\hat{\mathbf{V}}_T^N(x) = \sum_{t=2N-1}^T \hat{\mathbf{v}}_t^N(x)$  and  $\mathbf{Q}_{2,T}^N(x) = \sum_{t=N}^{T-N+1} \mathbf{q}_{2,t}(x)$ .

Now, considering a lower bound on the left-hand side, we have for any action  $b$ ,

$$\ln \frac{\mathbf{W}_{T+1}(x)}{\mathbf{W}_{2N-1}(x)} \geq \ln \frac{\mathbf{w}_{T+1}(x, b)}{\mathbf{w}_{2N-1}(x)} = \eta \sum_{t=N}^{T-N+1} \hat{q}_t(x, b) - \ln|\mathcal{A}|,$$

where we used that  $\mathbf{w}_{2N-1}(x, a) = 1$  holds for all  $a \in \mathcal{A}$ . Combining with (44), we get

$$\hat{\mathbf{V}}_T^N(x) \geq (1-\gamma) \hat{\mathbf{Q}}_T^N(x, b) - \frac{\ln|\mathcal{A}|}{\eta} - \eta e'' \mathbf{Q}_{2,T}^N(x), \quad (45)$$

Let us now bound the difference of  $\widehat{\mathbf{V}}_T^N(x)$  and

$$\widehat{\mathbf{V}}_T(x) = \sum_{t=N}^T \widehat{\mathbf{v}}_t(x) = \sum_{t=N}^T \sum_a \boldsymbol{\pi}_t(a|x) \widehat{\mathbf{q}}_t(x, a).$$

Note that

$$\widehat{\mathbf{V}}_T^N(x) = \sum_{t=N}^{T-N+1} \sum_a \boldsymbol{\pi}_{t+N-1}(a|x) \widehat{\mathbf{q}}_t(x, a).$$

Therefore,

$$\begin{aligned} & \widehat{\mathbf{V}}_T^N(x) - \widehat{\mathbf{V}}_T(x) \\ & \leq \sum_{t=N}^{T-N+1} \sum_a |\widehat{\mathbf{q}}_t(x, a)| \left| \boldsymbol{\pi}_{t+N-1}(a|x) - \boldsymbol{\pi}_t(a|x) \right| \\ & \quad + \sum_{t=T-N+2}^T \sum_a \boldsymbol{\pi}_t(a|x) |\widehat{\mathbf{q}}_t(x, a)| \\ & = \sum_{t=N}^{T-N+1} \sum_{t'=t}^{t+N-2} \boldsymbol{\Delta}_{t',t}(x) + \sum_{t=T-N+2}^T \sum_a \boldsymbol{\pi}_t(a|x) |\widehat{\mathbf{q}}_t(x, a)|, \end{aligned}$$

where we used the definition of  $\boldsymbol{\Delta}_{t',t}(x)$ . Taking the expectation of both sides, using Lemmas 14 and 8, we get

$$\begin{aligned} \mathbb{E} \left[ \widehat{\mathbf{V}}_T^N(x) \right] - \mathbb{E} \left[ \widehat{\mathbf{V}}_T(x) \right] & \leq (T - 2N + 2)(N - 1)c' + (N - 1)U_{\bar{q}} \\ & = (N - 1)\{(T - 2N + 2)c' + U_{\bar{q}}\}. \end{aligned}$$

This, together with (45) gives

$$\begin{aligned} & \mathbb{E} \left[ \widehat{\mathbf{V}}_T(x) \right] + (N - 1) \{(T - 2N + 2)c' + U_{\bar{q}}\} \\ & \geq (1 - \gamma) \mathbb{E} \left[ \widehat{\mathbf{Q}}_T^N(x, b) \right] - \frac{\ln|\mathcal{A}|}{\eta} - \eta e'' \mathbb{E} \left[ \mathbf{Q}_{2,T}^N(x) \right]. \end{aligned} \quad (46)$$

By equation (13), we have  $\mathbb{E} \left[ \widehat{\mathbf{V}}_T(x) \right] = \mathbb{E} \left[ \mathbf{V}_T(x) \right]$  and with the definition  $\mathbf{Q}_T^N(x, b) = \sum_{t=N}^{T-N+1} \mathbf{q}_t(x, b)$ , we also have  $\mathbb{E} \left[ \widehat{\mathbf{Q}}_T^N(x, b) \right] = \mathbb{E} \left[ \mathbf{Q}_T^N(x, b) \right]$ . Thus, using (43) again, we get

$$\begin{aligned} & \mathbb{E} \left[ \mathbf{V}_T(x) \right] + (N - 1) \{(T - 2N + 2)c' + U_{\bar{q}}\} \\ & \geq (1 - \gamma) \mathbb{E} \left[ \mathbf{Q}_T^N(x, b) \right] - \frac{\ln|\mathcal{A}|}{\eta} \\ & \quad - \eta e'' (T - 2N + 2) \{ |\mathcal{A}| U_{\pi_{\bar{q}}} U_{\bar{q}} + (N - 1)c' V_{\bar{q}} \}. \end{aligned}$$

By reordering the terms, this becomes

$$\begin{aligned} & \mathbb{E} \left[ \mathbf{Q}_T^N(x, b) - \mathbf{V}_T(x) \right] \\ & \leq \gamma \mathbb{E} \left[ \mathbf{Q}_T^N(x, b) \right] + (N - 1) \{(T - 2N + 2)c' + U_{\bar{q}}\} + \frac{\ln|\mathcal{A}|}{\eta} \\ & \quad + \eta e'' (T - 2N + 2) \{ |\mathcal{A}| U_{\pi_{\bar{q}}} U_{\bar{q}} + (N - 1)c' V_{\bar{q}} \}. \end{aligned} \quad (47)$$

We now lower bound  $\mathbf{Q}_T^N(x, a)$  by  $\mathbf{Q}_T(x, a)$ . Since the rewards  $r_t(x, a)$  are bounded between 0 and 1, by Lemma 3 we have

$$\mathbf{q}_t(x, b) \leq U_q = 2\tau + 3. \quad (48)$$

Therefore,

$$\mathbf{Q}_T(x, a) - \mathbf{Q}_T^N(x, a) = \sum_{t=T-N+2}^T \mathbf{q}_t(x, a) \leq U_q(N - 1), \quad (49)$$

Moreover, (48) also implies that  $\mathbb{E} \left[ \mathbf{Q}_T^N(x, b) \right] \leq U_q(T - 2N + 2)$ . Combining this with (49) and (47), we obtain the desired bound:

$$\begin{aligned} \mathbb{E} \left[ \mathbf{Q}_T(x, b) - \mathbf{V}_T(x) \right] & \leq \frac{\ln|\mathcal{A}|}{\eta} + (N - 1)(U_{\bar{q}} + U_q) \\ & \quad + (T - 2N + 2) \left( c'(N - 1)(1 + \eta e'' V_{\bar{q}}) + \gamma U_q + \eta e'' |\mathcal{A}| U_{\pi_{\bar{q}}} U_{\bar{q}} \right). \end{aligned}$$

The proof of Proposition 1 is now easy: ■

*Proof of Proposition 1:* Under the conditions of the proposition, combining Lemmas 2 and 15, and using that  $0 \leq \rho_t^\pi, \rho_t \leq 1$  yields

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E} [\rho_t^\pi - \rho_t] \leq N - 1 + \sum_{x,a} \mu_{st}^\pi(x) \pi(a|x) \mathbb{E} [\mathbf{Q}_T(x, a) - \mathbf{V}_T(x)] \\ & \leq (N - 1)(U_{\bar{q}} + U_q + 1) + \frac{\ln|\mathcal{A}|}{\eta} \\ & \quad + (T - 2N + 2) \left( c'(N - 1)(1 + \eta e'' V_{\bar{q}}) + \gamma U_q + \eta e'' |\mathcal{A}| U_{\pi_{\bar{q}}} U_{\bar{q}} \right). \end{aligned}$$

proving Proposition 1. ■

### E. Proof of Proposition 2

Let  $t \geq N$ . First, since  $\boldsymbol{\pi}_t$  is  $\sigma(\mathbf{u}_{t-N})$ -measurable,  $\mathbb{E} [\boldsymbol{\rho}_t] = \mathbb{E} \left[ \sum_x \mu_{st}^{\boldsymbol{\pi}_t}(x) \mathbb{E} [r_t(x, \mathbf{a}_t) | \mathbf{u}_{t-N}] \right]$ . We also have

$$\begin{aligned} \mathbb{E} [r_t(\mathbf{x}_t, \mathbf{a}_t)] & = \mathbb{E} \left[ \mathbb{E} [r_t(\mathbf{x}_t, \mathbf{a}_t) | \mathbf{u}_{t-N}] \right] \\ & = \mathbb{E} \left[ \sum_x \mu_t^N(x) \mathbb{E} [r_t(x, \mathbf{a}_t) | \mathbf{u}_{t-N}] \right]. \end{aligned}$$

Hence,

$$\begin{aligned} \mathbb{E} [\boldsymbol{\rho}_t - r_t(\mathbf{x}_t, \mathbf{a}_t)] & = \mathbb{E} \left[ \sum_x (\mu_{st}^{\boldsymbol{\pi}_t}(x) - \mu_t^N(x)) \mathbb{E} [r_t(x, \mathbf{a}_t) | \mathbf{u}_{t-N}] \right] \\ & \leq \mathbb{E} \left[ \sum_x \left| \mu_{st}^{\boldsymbol{\pi}_t}(x) - \mu_t^N(x) \right| \right], \end{aligned}$$

where we have used that  $r_t(x, a) \in [0, 1]$ .

Thanks to Lemma 13, Lemma 10 is applicable. Hence,  $\sum_x \left| \mu_{st}^{\boldsymbol{\pi}_t}(x) - \mu_t^N(x) \right| \leq c(\tau + 1)^2 + 2e^{-(N-1)/\tau}$ , and thus  $\mathbb{E} [\boldsymbol{\rho}_t - r_t(\mathbf{x}_t, \mathbf{a}_t)] \leq c(\tau + 1)^2 + 2e^{-(N-1)/\tau}$ . Summing up these inequalities for  $t = N, \dots, T$ , and using the trivial bound  $\mathbb{E} [\boldsymbol{\rho}_t - r_t(\mathbf{x}_t, \mathbf{a}_t)] \leq 1$  for the first  $N - 1$  terms, we get the desired result. ■

### ACKNOWLEDGMENTS

This work was supported in part by the Hungarian Scientific Research Fund and the Hungarian National Office for Research and Technology (KTIA-OTKA CNK 77782), the PASCAL2 Network of Excellence under EC grant no. 216886, NSERC, AITF, the Alberta Ingenuity Centre for Machine Learning, the DARPA GALE project (HR0011-08-C-0110) and iCore.

### REFERENCES

- [1] Antos, A., Grover, V., and Szepesvári, C. (2010). Active learning in heteroscedastic noise. *Theoretical Computer Science*, 411(29–30):2712–2728.
- [2] Auer, P., Cesa-Bianchi, N., Freund, Y., and Schapire, R. E. (2002). The nonstochastic multiarmed bandit problem. *SIAM J. Comput.*, 32(1):48–77.
- [3] Cao, X.-R. (2007). *Stochastic Learning and Optimization: A Sensitivity-Based Approach*. Springer, New York.
- [4] Cesa-Bianchi, N., Conconi, A., and Gentile, C. (2004). On the generalization ability of on-line learning algorithms. *IEEE Transactions on Information Theory*, 50:2050–2057.
- [5] Cesa-Bianchi, N. and Lugosi, G. (2006). *Prediction, Learning, and Games*. Cambridge University Press, New York, NY, USA.
- [6] Even-Dar, E., Kakade, S. M., and Mansour, Y. (2005). Experts in a Markov decision process. In Saul, L. K., Weiss, Y., and Bottou, L., editors, *Advances in Neural Information Processing Systems 17*, pages 401–408, Cambridge, MA, USA. MIT Press.



- [7] Even-Dar, E., Kakade, S. M., and Mansour, Y. (2009). Online Markov decision processes. *Mathematics of Operations Research*, 34(3):726–736.
- [8] György, A., Linder, T., Lugosi, G., and Ottucsák, Gy. (2007). The on-line shortest path problem under partial monitoring. *Journal of Machine Learning Research*, 8:2369–2403.
- [9] Ipsen, I. C. F. and Selee, T. M. (2011). Ergodicity coefficients defined by vector norms. *SIAM J. Matrix Analysis Applications*, 32(1):153–200.
- [10] Jaksch, T., Ortner, R., and Auer, P. (2010). Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 99:1563–1600.
- [11] Kakade, S., Sridharan, K., and Tewari, A. (2010). On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In Bengio, Y., Schuurmans, D., Lafferty, J., Williams, C., and Culotta, A., editors, *Advances in Neural Information Processing Systems 22*, pages 793–800. Curran Associates. (December 7–10, 2009).
- [12] Lazaric, A. and Munos, R. (2011). Learning with stochastic inputs and adversarial outputs. *Journal of Computer and System Sciences*. To appear.
- [13] Neu, G., György, A., and Szepesvári, Cs. (2010). The online loop-free stochastic shortest-path problem. In Kalai, A. and Mohri, M., editors, *Proceedings of the 23rd Annual Conference on Learning Theory*, pages 231–243.
- [14] Neu, G., György, A., and Szepesvári, Cs. (2012). The adversarial stochastic shortest path problem with unknown transition probabilities. In Lawrence, N. and Girolami, M., editors, *Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics*, volume 22 of *JMLR Workshop and Conference Proceedings*, pages 805–813.
- [15] Neu, G., György, A., and Szepesvári, Cs. (2013). The online loop-free stochastic shortest-path problem. *In preparation*.
- [16] Neu, G., György, A., Szepesvári, Cs., and Antos, A. (2011). Online Markov decision processes under bandit feedback. In Lafferty, J., Williams, C., Shawe-Taylor, J., Zemel, R., and Culotta, A., editors, *Advances in Neural Information Processing Systems 23*, pages 1804–1812.
- [17] Puterman, M. L. (1994). *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley-Interscience.
- [18] Rakhlin, A., Sridharan, K., and Tewari, A. (2012). Online learning: Stochastic and constrained adversaries. In Shawe-Taylor, J., Zemel, R. S., Bartlett, P. L., Pereira, F. C., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems 24*, page 2752, Cambridge, MA, USA. Curran Associates. (December 12–15, 2011).
- [19] Sutton, R. and Barto, A. (1998). *Reinforcement Learning: An Introduction*. MIT Press.
- [20] Yu, J. Y. and Mannor, S. (2009a). Arbitrarily modulated Markov decision processes. In *Joint 48th IEEE Conference on Decision and Control and 28th Chinese Control Conference*, pages 2946–2953. IEEE Press.
- [21] Yu, J. Y. and Mannor, S. (2009b). Online learning in Markov decision processes with arbitrarily changing rewards and transitions. In *GameNets’09: Proceedings of the First ICST International Conference on Game Theory for Networks*, pages 314–322, Piscataway, NJ, USA. IEEE Press.
- [22] Yu, J. Y., Mannor, S., and Shimkin, N. (2009). Markov decision processes with arbitrary reward processes. *Mathematics of Operations Research*, 34(3):737–757.