
Online Learning in MDPs with Linear Function Approximation and Bandit Feedback

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We consider the problem of online learning in an episodic Markov decision process,
2 where the reward function is allowed to change between episodes in an adversarial
3 manner and the learner only observes the rewards associated with its actions.
4 We assume that rewards and the transition function can be represented as linear
5 functions in terms of a known low-dimensional feature map, which allows us to
6 consider the setting where the state space is arbitrarily large. We also assume that
7 the learner has a perfect knowledge of the MDP dynamics. Our main contribution
8 is developing an algorithm whose expected regret after T episodes is bounded
9 by $\tilde{\mathcal{O}}(\sqrt{dHT})$, where H is the number of steps in each episode and d is the
10 dimensionality of the feature map.

11 1 Introduction

12 We study the problem of online learning in episodic Markov Decision Processes (MDP), modeling
13 a sequential decision making problem where the interaction between a learner and its environment
14 is divided into T episodes of fixed length H . At each time step of the episode, the learner observes
15 the current state of the environment, chooses one of the available actions, and earns a reward.
16 Consequently, the state of the environment changes according to the transition function of the
17 underlying MDP, as a function of the previous state and the action taken by the learner. A key
18 distinguishing feature of our setting is that we assume that the reward function can change arbitrarily
19 between episodes, and the learner only has access to bandit feedback: instead of being able to observe
20 the reward function at the end of the episode, the learner only gets to observe the rewards that it
21 actually received. As traditional in this line of work, we aim to design algorithms for the learner with
22 theoretical guarantees on her regret, which is the difference between the total reward accumulated by
23 the learner and the total reward of the best stationary policy fixed in hindsight.

24 Unlike most previous work on this problem, we allow the state space to be very large and aim to
25 prove performance guarantees that do not depend on the size of the state space, bringing theory one
26 step closer to practical scenarios where assuming finite state spaces is unrealistic. To address the
27 challenge of learning in large state spaces, we adopt the classic RL technique of using *linear function*
28 *approximation* and suppose that we have access to a relatively low-dimensional feature map that can
29 be used to represent policies and value functions. We will assume that the feature map is expressive
30 enough so that all action-value functions can be expressed as linear functions of the features, and that
31 the learner has full knowledge of the transition function of the MDP.

32 Our main contribution is designing a computationally efficient algorithm called ONLINE Q-REPS,
33 and prove that in the setting described above, its regret is at most $\mathcal{O}(\sqrt{dHTD}(\mu^*||\mu_0))$, where d is
34 the dimensionality of the feature map and $D(\mu^*||\mu_0)$ is the relative entropy between the state-action
35 distribution μ^* induced by the optimal policy and an initial distribution μ_0 given as input to the

36 algorithm. Notably, our results do not require the likelihood ratio between these distributions to be
 37 uniformly bounded, and the bound shows no dependence on the eigenvalues of the feature covariance
 38 matrices. Our algorithm itself requires solving a d^2 -dimensional convex optimization problem at the
 39 beginning of each episode, which can be solved to arbitrary precision ε in time polynomial in d and
 40 $1/\varepsilon$, independently of the size of the state-action space.

41 Our work fits into a long line of research considering online learning in Markov decision processes.
 42 The problem of regret minimization in stationary MDPs with a *fixed* reward function has been studied
 43 extensively since the work of Burnetas and Katehakis [6], Auer and Ortner [2], Tewari and Bartlett
 44 [31], Jaksch et al. [14], with several important advances made in the past decade [9, 10, 4, 13, 15].
 45 While most of these works considered small finite state spaces, the same techniques have been very
 46 recently extended to accommodate infinite state spaces under the assumption of realizable function
 47 approximation by Jin et al. [17] and Yang and Wang [33]. In particular, the notion of *linear MDPs*
 48 introduced by Jin et al. [17] has become a standard model for linear function approximation and has
 49 been used in several recent works (e.g., 22, 32, 1).

50 Even more relevant is the line of work considering adversarial rewards, initiated by Even-Dar et al.
 51 [12], who consider online learning in continuing MDPs with full feedback about the rewards. They
 52 proposed a MDP-E algorithm, that achieves $\tilde{O}(\tau^2 \sqrt{T} \log K)$ regret, where τ is an upper bound
 53 on the mixing time of the MDP. Later, Neu et al. [25] proposed an algorithm which guarantees
 54 $\tilde{O}(\sqrt{\tau^3 KT/\alpha})$ regret with bandit feedback, essentially assuming that all states are reachable with
 55 probability $\alpha > 0$ under all policies. In our work, we focus on episodic MDPs with a fixed
 56 episode length H . The setting was first considered in the bandit setting by Neu et al. [23], who
 57 proposed an algorithm with a regret bound of $\mathcal{O}(H^2 \sqrt{TK}/\alpha)$. Although the number of states
 58 does not appear explicitly in the bound, the regret scales at least linearly with the size of the state
 59 space \mathcal{X} , since $|\mathcal{X}| \leq H/\alpha$. Later work by Zimin and Neu [35], Dick et al. [11] eliminated the
 60 dependence on α and proposed an algorithm achieving $\tilde{O}(\sqrt{TH|\mathcal{X}|K})$ regret. Regret bounds
 61 for the full-information case without prior knowledge of the MDP were achieved by Neu et al.
 62 [24] and Rosenberg and Mansour [30], of order $\tilde{O}(H|\mathcal{X}|K\sqrt{T})$ and $\tilde{O}(H|\mathcal{X}|\sqrt{KT})$, respectively.
 63 These results were recently extended to handle bandit feedback about the rewards by Jin et al. [16],
 64 ultimately resulting in a regret bound of $\tilde{O}(H|\mathcal{X}|\sqrt{KT})$.

65 As apparent from the above discussion, all work on online learning in MDPs with adversarial rewards
 66 considers finite state spaces. The only exception we are aware of is the recent work of Cai et al.
 67 [7], whose algorithm OPPO is guaranteed to achieve $\tilde{O}(\sqrt{d^3 H^3 T})$, assuming that the learner has
 68 access to d -dimensional features that can perfectly represent all action-value functions. While Cai,
 69 Yang, Jin, and Wang [7] remarkably assumed no prior knowledge of the MDP parameters, their
 70 guarantees are only achieved in the full-information case. This is to be contrasted with our results
 71 that are achieved for the much more restrictive bandit setting, albeit with the stronger assumption of
 72 having full knowledge of the underlying MDP, as required by virtually all prior work in the bandit
 73 setting, with the exception of Jin et al. [16].

74 Our results are made possible by a careful combination of recently proposed techniques for contextual
 75 bandit problems and optimal control in Markov decision processes. In particular, a core component
 76 of our algorithm is a regularized linear programming formulation of optimal control in MDPs due
 77 to Bas-Serrano et al. [5], which allows us to reduce the task of computing near-optimal policies
 78 in linear MDPs to a low-dimensional convex optimization problem. A similar algorithm design
 79 has been previously used for tabular MDPs by Zimin and Neu [35], Dick et al. [11], with the
 80 purpose of removing factors of $1/\alpha$ from the previous state-of-the-art bounds of Neu et al. [23].
 81 Analogously to this improvement, our methodology enables us to make strong assumptions on
 82 problem-dependent constants like likelihood ratios between μ^* and μ_0 or eigenvalues of the feature
 83 covariance matrices. Another important building block of our method is a version of the recently
 84 proposed Matrix Geometric Resampling procedure of Neu and Olkhovskaya [21] that enables us to
 85 efficiently estimate the reward functions. Incorporating these estimators in the algorithmic template
 86 of Bas-Serrano et al. [5] is far from straightforward and requires several subtle adjustments.

87 **Notation.** We use $\langle \cdot, \cdot \rangle$ to denote inner products in Euclidean space and by $\|\cdot\|$ we denote the
 88 Euclidean norm for vectors and the operator norm for matrices. For a symmetric positive definite
 89 matrix A , we use $\lambda_{\min}(A)$ to denote its smallest eigenvalue. We write $\text{tr}(A)$ for the trace of a matrix
 90 A and use $A \succcurlyeq 0$ to denote that an operator A is positive semi-definite, and we use $A \succcurlyeq B$ to denote

91 $A - B \succcurlyeq 0$. For a d -dimensional vector v , we denote the corresponding $d \times d$ diagonal matrix by
 92 $\text{diag}(v)$. For a positive integer N , we use $[N]$ to denote the set of positive integers $\{1, 2, \dots, N\}$.
 93 Finally, we will denote the set of all probability distributions over any set \mathcal{X} by $\Delta_{\mathcal{X}}$.

94 2 Preliminaries

95 An episodic Markovian Decision Process (MDP), denoted by $M = (\mathcal{X}, \mathcal{A}, H, P, r)$ is defined by a
 96 state space \mathcal{X} , action space \mathcal{A} , episode length $H \in \mathbb{Z}_+$, transition function $P : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_{\mathcal{X}}$ and a
 97 reward function $r : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$. For convenience, we will assume that both \mathcal{X} and \mathcal{A} are finite
 98 sets, although we allow the state space \mathcal{X} to be arbitrarily large. Without significant loss of generality,
 99 we will assume that the set of available actions is the same \mathcal{A} in each state, with cardinality $|\mathcal{A}| = K$.
 100 Furthermore, without any loss of generality, we will assume that the MDP has a layered structure,
 101 satisfying the following conditions:

- 102 • The state set \mathcal{X} can be decomposed into H disjoint sets: $\mathcal{X} = \cup_{h=1}^H \mathcal{X}_h$,
- 103 • $\mathcal{X}_1 = \{x_1\}$ and $\mathcal{X}_H = \{x_H\}$ are singletons,
- 104 • transitions are only possible between consecutive layers, that is, for any $x_h \in \mathcal{X}_h$, the
 105 distribution $P(\cdot|x, a)$ is supported on \mathcal{X}_{h+1} for all a and $h \in [H - 1]$.

106 These assumptions are common in the related literature (e.g., 23, 35, 30) and are not essential to our
 107 analysis; their primary role is simplifying our notation.

108 In the present paper, we consider an *online learning* problem where the learner interacts with its envi-
 109 ronment in a sequence of episodes $t = 1, 2, \dots, T$, facing a different *reward functions* $r_{t,1}, \dots, r_{t,H+1}$
 110 selected by a (possibly adaptive) adversary at the beginning of each episode t . Oblivious to the reward
 111 function chosen by the adversary, the learner starts interacting with the MDP in each episode from the
 112 initial state $X_{t,1} = x_1$. At each consecutive step $h \in [H - 1]$ within the episode, the learner observes
 113 the state $X_{t,h}$, picks an action $A_{t,h}$ and observes the reward $r_{t,h}(X_{t,h}, A_{t,h})$. Then, unless $h = H$,
 114 the learner moves to the next state $X_{t,h+1}$, which is generated from the distribution $P(\cdot|X_{t,h}, A_{t,h})$.
 115 At the end of step H , the episode terminates and a new one begins. The aim of the learner is to select
 116 its actions so that the cumulative sum of rewards is as large as possible.

117 Our algorithm and analysis will make use of the concept of (stationary stochastic) *policies* $\pi : \mathcal{X} \rightarrow$
 118 $\Delta_{\mathcal{A}}$. A policy π prescribes a behaviour rule to the learner by assigning probability $\pi(a|x)$ to taking
 119 action a at state x . Let $\tau^\pi = ((X_1, A_1), (X_2, A_2), \dots, (X_H, A_H))$ be a trajectory generated by
 120 following the policy π through the MDP. Then, for any $x_h \in \mathcal{X}_h, a_h \in \mathcal{A}$ we define the occupancy
 121 measure $\mu_h^\pi(x, a) = \mathbb{P}_\pi[(x, a) \in \tau^\pi]$. We will refer to the collection of these distributions across all
 122 layers h as the *occupancy measure* induced by π and denote it as $\mu^\pi = (\mu_1^\pi, \mu_2^\pi, \dots, \mu_H^\pi)$. We will
 123 denote the set of all valid occupancy measures by \mathcal{U} and note that this is a convex set, such that for
 124 every element $\mu \in \mathcal{U}$ the following set of linear constraints is satisfied:

$$\sum_{a \in \mathcal{A}} \mu_{h+1}(x, a) = \sum_{x', a' \in \mathcal{X}_h \times \mathcal{A}} P(x|x', a') \mu_h(x', a'), \quad \forall x \in \mathcal{X}_{h+1}, h \in [H - 1], \quad (1)$$

125 as well as $\sum_a \mu_1(x_1, a) = 1$. From every valid occupancy measure μ , a stationary stochastic
 126 policy $\pi = \pi_1, \dots, \pi_{H-1}$ can be derived as $\pi_{\mu,h}(a|x) = \mu_h(x, a) / \sum_{a'} \mu_h(x, a')$. For each h ,
 127 introducing the linear operators E and P through their action on a set state-action distribution u_h as
 128 $(E^\top u_h)(x) = \sum_{a \in \mathcal{A}} u_h(x, a)$ and $(P_h^\top u_h)(x) = \sum_{x', a' \in \mathcal{X}_h, \mathcal{A}} P(x|x', a') u_h(x', a')$, the constraints
 129 can be simply written as $E^\top \mu_{h+1} = P_h^\top \mu_h$ for each h . We will use the inner product notation for
 130 the sum over the set of states and actions: $\langle \mu_h, r_h \rangle = \sum_{(x,a) \in (\mathcal{X}_h \times \mathcal{A})} \mu_h(x, a) r_{t,h}(x, a)$. Using this
 131 notation, we formulate our objective as selecting a sequence of policies π_t for each episode t in a
 132 way that it minimizes the *total expected regret* defined as

$$\mathfrak{R}_T = \sup_{\pi^*} \sum_{t=1}^T \sum_{h=1}^H (\mathbb{E}_{\pi^*} [r_{t,h}(X_h^*, A_h^*)] - \mathbb{E}_{\pi_t} [r_t(X_{t,h}, A_{t,h})]) = \sup_{\mu^* \in \mathcal{U}} \sum_{t=1}^T \sum_{h=1}^H \langle \mu_h^* - \mu_h^{\pi_t}, r_{t,h} \rangle,$$

133 where the notations $\mathbb{E}_{\pi^*} [\cdot]$ and $\mathbb{E}_{\pi_t} [\cdot]$ emphasize that the state-action trajectories are generated by
 134 following policies π^* and π_t , respectively. As the above expression suggests, we can reformulate
 135 our online learning problem as an instance of online linear optimization where in each episode t , the

136 learner selects an occupancy measure $\mu_t \in \mathcal{U}$ (with $\mu_t = \mu^{\pi_t}$) and gains reward $\sum_{h=1}^H \langle \mu_{t,h}, r_{t,h} \rangle$.
 137 Intuitively, the regret measures the gap between the total reward gained by the learner and that of the
 138 best stationary policy fixed in hindsight, with full knowledge of the sequence of rewards chosen by
 139 the adversary. This performance measure is standard in the related literature on online learning in
 140 MDPs, see, for example Neu et al. [23], Zimin and Neu [35], Neu et al. [24], Rosenberg and Mansour
 141 [30], Cai et al. [7].

142 In this paper, we focus on MDPs with potentially enormous state spaces, which makes it difficult
 143 to design computationally tractable algorithms with nontrivial guarantees, unless we make some as-
 144 sumptions. We particularly focus on the classic technique of relying on *linear function approximation*
 145 and assuming that the reward functions occurring during the learning process can be written as a
 146 linear function of a low-dimensional feature map. We specify the form of function approximation
 147 and the conditions our analysis requires as follows:

148 **Assumption 1** (Linear MDP with adversarial rewards). *There exists a feature map $\varphi : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}^d$*
 149 *and a collection of d signed measures $m = (m_1, \dots, m_d)$ on \mathcal{X} , such that for any $(x, a) \in \mathcal{X} \times \mathcal{A}$*
 150 *the transition function can be written as*

$$P(\cdot|x, a) = \langle m(\cdot), \varphi(x, a) \rangle.$$

151 *Furthermore, the reward function chosen by the adversary in each episode t can be written as*

$$r_{t,h}(x, a) = \langle \theta_{t,h}, \varphi(x, a) \rangle$$

152 *for some $\theta_{t,h} \in \mathbb{R}^d$. We assume that the features and the parameter vectors satisfy $\|\varphi(x, a)\| \leq \sigma$*
 153 *and that the first coordinate $\varphi_1(x, a) = 1$ for all $(x, a) \in \mathcal{X} \times \mathcal{A}$. Also we assume that $\|\theta_{t,h}\| \leq R$.*

154 Online learning under this assumption, but with a fixed reward function, has received substantial
 155 attention in the recent literature, particularly since the work of Jin et al. [17] who popularized the
 156 term ‘‘Linear MDP’’ to refer to this class of MDPs. This has quickly become a common assumption
 157 for studying reinforcement learning algorithms (Cai et al. [7], Jin et al. [17], Neu and Pike-Burke
 158 [22], Agarwal et al. [1]). This is also a special case of *factored linear models* (Yao et al. [34], Pires
 159 and Szepesvari [29]).

160 Linear MDPs come with several attractive properties that allow efficient optimization and learning.
 161 In this work, we will exploit the useful property shown by Neu and Pike-Burke [22] and Bas-Serrano
 162 et al. [5] that all occupancy measures in a linear MDP can be seen to satisfy a relaxed version of the
 163 constraints in Equation (1). Specifically, for all h , defining the feature matrix $\Phi_h \in \mathbb{R}^{(\mathcal{X}_h \times \mathcal{A}) \times d}$ with
 164 its action on the distribution u as $\Phi_h^\top u = \sum_{x,a \in \mathcal{X}_h \times \mathcal{A}} u_h(x, a) \varphi(x, a)$, we define \mathcal{U}_Φ as the set of
 165 state-action distributions $(\mu, u) = ((\mu_1, \dots, \mu_H), (u_1, \dots, u_H))$ satisfying the following constraints:
 166

$$E^\top u_{h+1} = P_h^\top \mu_h \quad (\forall h), \quad \Phi_h^\top u_h = \Phi_h^\top \mu_h \quad (\forall h), \quad E^\top u_1 = 1. \quad (2)$$

167 It is easy to see that for all feasible (μ, u) pairs, u satisfies the original constraints (1) if the MDP
 168 satisfies Assumption 1: since the transition operator can be written as $P_h = \Phi_h M_h$ for some matrix
 169 M_h . In this case, we clearly have

$$E^\top u_{h+1} = P_h^\top \mu_h = M_h^\top \Phi_h^\top \mu_h = M_h^\top \Phi_h^\top u_h = P_h^\top u_h, \quad (3)$$

170 showing that any feasible u is indeed a valid occupancy measure. Furthermore, due to linearity
 171 of the rewards in Φ , we also have $\langle u_h, r_{t,h} \rangle = \langle \mu_h, r_{t,h} \rangle$ for all feasible $(\mu, u) \in \mathcal{U}_\Phi$. While the
 172 number of variables and constraints in Equation (2) is still very large, it has been recently shown that
 173 approximate linear optimization over this set can be performed tractably [22, 5]. Our own algorithm
 174 design described in the next section will heavily build on these recent results.

175 3 Algorithm and main results

176 This section presents our main contributions: a new efficient algorithm for the setting described above,
 177 along with its performance guarantees. Our algorithm design is based on a reduction to online linear
 178 optimization, exploiting the structural results established in the previous section. In particular, we will
 179 heavily rely on the algorithmic ideas established by Bas-Serrano et al. [5], who proposed an efficient
 180 reduction of approximate linear optimization over the high-dimensional set \mathcal{U}_Φ to a low-dimensional
 181 convex optimization problem. Another key component of our algorithm is an efficient estimator of
 182 the reward vectors $\theta_{t,h}$ based on the work of Neu and Olkhovskaya [21]. For reasons that we will
 183 clarify in Section 4, accommodating these reward estimators into the framework of Bas-Serrano et al.
 184 [5] is not straightforward and necessitates some subtle changes.

185 **3.1 The policy update rule**

186 Our algorithm is an instantiation of the well-known ‘‘Follow the Regularized Leader’’ (FTRL) template
 187 commonly used in the design of modern online learning methods (see, e.g., 26). We will make the
 188 following design choices:

- 189 • The decision variables will be the vector $(\mu, u) \in \mathbb{R}^{2(\mathcal{X} \times \mathcal{A})}$, with the feasible set \mathcal{U}_Φ^2 defined
 190 through the constraints

$$E^\top u_h = P_h^\top \mu_h \quad (\forall h), \quad \Phi_h^\top \text{diag}(u_h) \Phi_h = \Phi_h^\top \text{diag}(\mu_h) \Phi_h \quad (\forall h). \quad (4)$$

191 These latter constraints ensure that the feature covariance matrices under u and μ will be
 192 identical, which is necessary for technical reasons that will be clarified in Section 4. Notice
 193 that, due to our assumption that $\varphi_1(x, a) = 1$, we have $\mathcal{U}_\Phi^2 \subseteq \mathcal{U}_\Phi$, so all feasible u ’s continue
 194 to be feasible for the original constraints (1).

- 195 • The regularization function will be chosen as $\frac{1}{\eta} D(\mu \| \mu_0) + \frac{1}{\alpha} D_C(u \| \mu_0)$ for some positive
 196 regularization parameters η and α , where μ_0 is the occupancy measure induced by the
 197 uniform π_0 with $\pi_0(a|x) = \frac{1}{K}$ for all x, a , and D and D_C are the marginal and conditional
 198 relative entropy functions respectively defined as $D(\mu \| \mu_0) = \sum_{h=1}^H D(\mu_h \| \mu_{0,h})$ and
 199 $D_C(\mu \| \mu_0) = \sum_{h=1}^H D_C(\mu_h \| \mu_{0,h})$ with

$$D(\mu_h \| \mu_{0,h}) = \sum_{(x,a) \in (\mathcal{X}_h \times \mathcal{A})} \mu_h(x, a) \log \frac{\mu_h(x, a)}{\mu_{0,h}(x, a)}, \quad \text{and}$$

$$D_C(\mu_h \| \mu_{0,h}) = \sum_{(x,a) \in (\mathcal{X}_h \times \mathcal{A})} \mu_h(x, a) \log \frac{\pi_{\mu,h}(a|x)}{\pi_{0,h}(a|x)}.$$

200 With these choices, the updates of our algorithm in each episode will be given by

$$(\mu_t, u_t) = \arg \max_{(\mu, u) \in \mathcal{U}_\Phi^2} \left\{ \sum_{s=1}^{t-1} \sum_{h=1}^{H-1} \langle \mu_h, \widehat{r}_{s,h} \rangle - \frac{1}{\eta} D(\mu \| \mu_0) - \frac{1}{\alpha} D_C(u \| \mu_0) \right\} \quad (5)$$

201 where $\widehat{r}_{t,h} \in \mathbb{R}^{\mathcal{X} \times \mathcal{A}}$ is an estimator of the reward function $r_{t,h}$ that will be defined shortly.

202 As written above, it is far from obvious if these updates can be calculated efficiently. The following
 203 result shows that, despite the apparent intractability of the maximization problem, it is possible to
 204 reduce the above problem into a d^2 -dimensional unconstrained convex optimization problem:

205 **Proposition 1.** *Define for each $h \in [H-1]$, a matrix $Z_h \in \mathbb{R}^{d \times d}$ and let matrix $Z \in \mathbb{R}^{d \times d(H-1)}$
 206 be defined as $Z = (Z_1, \dots, Z_{H-1})$. We will write $h(x) = h$, if $x \in \mathcal{X}_h$. Define the Q -function taking
 207 values $Q_Z(x, a) = \varphi(x, a)^\top Z_{h(x)} \varphi(x, a)$ and define the value function*

$$V_Z(x) = \frac{1}{\alpha} \log \left(\sum_{a \in A(x)} \pi_0(a|x) e^{\alpha Q_Z(x, a)} \right)$$

208 *For any $h \in [H-1]$ and for any $x \in \mathcal{X}_h$, $a \in A(x)$, denote $P_{x,a} V_Z = \sum_{x' \in \mathcal{X}_{h(x)+1}} P(x'|x, a) V_Z(x')$
 209 and $\Delta_{t,Z}(x, a) = \sum_{s=1}^{t-1} \widehat{r}_{s,h(x)}(x, a) + P_{x,a} V_Z - Q_Z(x, a)$. Then, the optimal solution of the
 210 optimization problem (5) is given as*

$$\widehat{\pi}_{t,h}(a|x) = \pi_0(a|x) e^{\alpha(Q_{Z_t^*}(x, a) - V_{Z_t^*}(x))},$$

$$\widehat{\mu}_{t,h}(x, a) \propto \mu_0(x, a) e^{\eta \Delta_{t,Z_t^*}(x, a)},$$

211 *where $Z_t^* = (Z_{t,1}^*, \dots, Z_{t,H-1}^*)$ is the minimizer of the convex function*

$$\mathcal{G}_t(Z) = \frac{1}{\eta} \sum_{h=1}^{H-1} \log \left(\sum_{x \in \mathcal{X}_h, a \in A(x)} \mu_0(x, a) e^{\eta \Delta_{t,Z}(x, a)} \right) + V_Z(x_1). \quad (6)$$

212 A particular merit of this result is that it gives an explicit formula for the policy π_t that induces the
 213 optimal occupancy measure u_t , and that $\pi_t(a|x)$ can be evaluated straightforwardly as a function of
 214 the features $\varphi(x, a)$ and the parameters Z_t^* . The proof of the result is based on Lagrangian duality,
 215 and mainly follows the proof of Proposition 1 in Bas-Serrano et al. [5], with some subtle differences
 216 due to the episodic setting we consider and the appearance of the constraints $\Phi_h^\top \text{diag}(u_h) \Phi_h =$
 217 $\Phi_h^\top \text{diag}(\mu_h) \Phi_h$. The proof is presented in Appendix A.1.

218 The proposition above inspires a very straightforward implementation that is presented as Algorithm 1.
 219 Due to the direct relation with the algorithm of Bas-Serrano et al. [5], we refer to this method as
 220 ONLINE Q-REPS, where Q-REPS stands for ‘‘Relative Entropy Policy Search with Q-functions’’.
 221 ONLINE Q-REPS adapts the general idea of Q-REPS to the online setting in a similar way as the
 222 O-REPS algorithm of Zimin and Neu [35] adapted the Relative Entropy Policy Search method of
 223 Peters et al. [28] to regret minimization in tabular MDPs with adversarial rewards. While O-REPS
 224 would in principle be still applicable to the large-scale setting we study in this paper and would
 225 plausibly achieve similar regret guarantees, its implementation would be nearly impossible due to the
 226 lack of the structural properties enjoyed by ONLINE Q-REPS, as established in Proposition 1.

Algorithm 1 ONLINE Q-REPS

Parameters: $\eta, \alpha > 0$, exploration parameter $\gamma \in (0, 1)$,

Initialization: Set $\hat{\theta}_{1,h} = 0$ for all h , compute Z_1 .

For $t = 1, \dots, T$, **repeat:**

- Draw $Y_t \sim \text{Ber}(\gamma)$,
- **For** $h = 1, \dots, H$, **do:**
 - Observe $X_{t,h}$ and, for all $a \in \mathcal{A}(X_{t,h})$, set

$$\pi_{t,h}(a|X_{t,h}) = \pi_{0,h}(a|X_{t,h}) e^{\alpha(Q_{Z_t}(X_{t,h}, a) - V_{Z_t}(X_{t,h}))},$$

- if $Y = 0$, draw $A_{t,h} \sim \pi_{t,h}(\cdot|X_{t,h})$, otherwise draw $A_{t,h} \sim \pi_{0,h}(\cdot|X_{t,h})$,
 - observe the reward $r_{t,h}(X_{t,h}, A_{t,h})$.
 - Compute $\hat{\theta}_{t,1}, \dots, \hat{\theta}_{t,H-1}, Z_{t+1}$.
-

227 **3.2 The reward estimator**

228 We now turn to describing the reward estimators $\hat{r}_{t,h}$, which will require several further definitions.
 229 Specifically, a concept of key importance will be the following *feature covariance matrix*:

$$\Sigma_{t,h} = \mathbb{E}_{\pi_t} [\varphi(X_{t,h}, A_{t,h}) \varphi(X_{t,h}, A_{t,h})^\top].$$

230 Making sure that $\Sigma_{t,h}$ is invertible, we can define the estimator

$$\tilde{\theta}_{t,h} = \Sigma_{t,h}^{-1} \varphi(X_{t,h}, A_{t,h}) r_{t,h}(X_{t,h}, A_{t,h}). \quad (7)$$

231 This estimate shares many similarities with the estimates that are broadly used in the literature on
 232 adversarial linear bandits [18, 3, 8]. It is easy to see that $\tilde{\theta}_{t,h}$ is an unbiased estimate of $\theta_{t,h}$:

$$\mathbb{E}_t [\tilde{\theta}_{t,h}] = \mathbb{E}_t \left[\Sigma_{t,h}^{-1} \varphi(X_{t,h}, A_{t,h}) \varphi(X_{t,h}, A_{t,h})^\top \theta_{t,h} \right] = \Sigma_{t,h}^{-1} \Sigma_{t,h} \theta_{t,h} = \theta_{t,h}.$$

233 Unfortunately, exact computation of $\Sigma_{t,h}$ is intractable. To address this issue, we propose a method
 234 to directly estimate the inverse of the covariance matrix $\Sigma_{t,h}$ by adapting the Matrix Geometric
 235 Resampling method of Neu and Olkhovskaya [21] (which itself is originally inspired by the Geometric
 236 Resampling method of [19, 20]). Our adaptation has two parameters $\beta > 0$ and $M \in \mathbb{Z}_+$, and generates
 237 an estimate of the inverse covariance matrix through the following procedure¹:

¹The version we present here is a naïve implementation, optimized for readability. We present a more practical variant in Appendix B

Matrix Geometric Resampling

Input: simulator of P , policy $\tilde{\pi}_t = (\tilde{\pi}_{t,1}, \dots, \tilde{\pi}_{t,H-1})$.

For $i = 1, \dots, M$, **repeat:**

1. Simulate a trajectory

$\tau(i) = \{(X_1(i), A_1(i)), \dots, (X_{H-1}(i), A_{H-1}(i))\}$, following the policy $\tilde{\pi}_t$ in P ,

2. **For** $h = 1, \dots, H - 1$, **repeat:**

Compute

(a) $B_{i,h} = \varphi(X_h(i), A_h(i))\varphi(X_h(i), A_h(i))^\top$,

(b) $C_{i,h} = \prod_{j=1}^i (I - \beta B_{j,h})$.

Return $\hat{\Sigma}_{t,h}^+ = \beta I + \beta \sum_{i=1}^M C_{i,h}$ for all $h \in [H - 1]$.

238 Based on the above procedure, we finally define our estimator as

$$\hat{\theta}_{t,h} = \hat{\Sigma}_{t,h}^+ \varphi(X_{t,h}, A_{t,h}) r_{t,h}(X_{t,h}, A_{t,h}).$$

239 The idea of the estimate is based on the truncation of the Neumann-series expansion of the matrix
 240 $\Sigma_{t,h}^{-1}$ at the M th order term. Then, for large enough M , the matrix $\hat{\Sigma}_{t,h}^+$ is a good estimator of the
 241 inverse covariance matrix, which will be quantified formally in the analysis. For more intuition on
 242 the estimate, see section 3.2. in Neu and Olkhovskaya [21]. With a careful implementation explained
 243 in Appendix B, $\hat{\theta}_{t,h}$ can be computed in $O(MHKd)$ time, using M calls to the simulator.
 244

245 3.3 The regret bound

246 We are now ready to state our main result: a bound on the expected regret of ONLINE Q-REPS.
 247 During the analysis, we will suppose that all the optimization problems solved by the algorithm are
 248 solved up to an additive error of $\varepsilon \geq 0$. Furthermore, we will denote the covariance matrix generated
 249 by the uniform policy at layer h as $\Sigma_{0,h}$, and make the following assumption:

250 **Assumption 2.** The eigenvalues of $\Sigma_{0,h}$ for all h are lower bounded by $\lambda_{\min} > 0$.

251 Our main result is the following guarantee regarding the performance of ONLINE Q-REPS:

252 **Theorem 1.** Suppose that the MDP satisfies Assumptions 1 and 2 and $\lambda_{\min} > 0$. Furthermore,
 253 suppose that, for all t , Z_t satisfies $\mathcal{G}_t(Z_t) \leq \min_Z \mathcal{G}_t(Z) + \varepsilon$ for some $\varepsilon \geq 0$. Then, for $\gamma \in (0, 1)$,
 254 $M \geq 0$, any positive $\eta \leq \frac{2}{(M+2)H}$ and any positive $\beta \leq \frac{1}{2\sigma^2}$, the expected regret of ONLINE
 255 Q-REPS over T episodes satisfies

$$\begin{aligned} \mathfrak{R}_T \leq & 2T\sigma RH \cdot \exp(-\gamma\beta\lambda_{\min}M) + \gamma HT + \eta HdT \frac{4}{3} + \frac{1}{\eta} D(\mu^* \|\mu_0) + \frac{1}{\alpha} D_C(u^* \|\mu_0) \\ & + \sqrt{\alpha\varepsilon}(M+2)HT. \end{aligned}$$

256 Furthermore, letting $\beta = \frac{1}{2\sigma^2}$, $M = \left\lceil \frac{2\sigma^2 \log(TH\sigma R)}{\gamma\lambda_{\min}} \right\rceil$, $\eta = \frac{1}{\sqrt{TdH}}$, $\alpha = \frac{1}{\sqrt{TdH}}$ and $\gamma = \frac{1}{\sqrt{TH}}$ and
 257 supposing that T is large enough so that the above constraints on M, γ, η and β are satisfied, we
 258 also have

$$\mathfrak{R}_T \leq \sqrt{dHT} (2 + D(\mu^* \|\mu_0) + D_C(u^* \|\mu_0)) + \sqrt{HT} + \sqrt{\varepsilon} T^{5/4} (Hd)^{1/4} + 2.$$

259 Thus, when all optimization problems are solved up to precision $\varepsilon = T^{-3/2}$, the regret of ONLINE
 260 Q-REPS is guaranteed to be of $\mathcal{O}(\sqrt{dHTD}(\mu^* \|\mu_0))$.

261 3.4 Implementation

262 While Proposition 1 establishes the form of the ideal policy updates π_t through the solution of an
 263 unconstrained convex optimization problem, it is not obvious that this optimization problem can be
 264 solved efficiently. Indeed, one immediate challenge in optimizing \mathcal{G}_t is that its gradient takes the form

$$\nabla \mathcal{G}_t(Z) = \sum_{x,a} \tilde{\mu}_Z(x,a) \left(\varphi(x,a)\varphi(x,a)^\top - \sum_{x',a'} P(x'|x,a)\pi_Z(a'|x')\varphi(x',a')\varphi(x',a')^\top \right),$$

265 where $\tilde{\mu}_Z(x, a) = \frac{\mu_0(x, a) \exp(\eta \Delta_Z(x, a))}{\sum_{x', a'} \mu_0(x', a') \exp(\eta \Delta_Z(x', a'))}$. Sampling from this latter distribution (and thus ob-
 266 taining unbiased estimators of $\nabla \mathcal{G}_t(Z)$) is problematic due to the intractable normalization constant.

267 This challenge can be addressed in a variety of ways. First, one can estimate the gradients via weighted
 268 importance sampling from the distribution $\tilde{\mu}_Z$ and using these in a stochastic optimization procedure.
 269 This approach has been recently proposed and analyzed for an approximate implementation of REPS
 270 by Pacchiano et al. [27], who showed that it results in ε -optimal policy updates given polynomially
 271 many samples in $1/\varepsilon$. Alternatively, one can consider an empirical counterpart of the loss function
 272 replacing the expectation with respect to μ_0 with an empirical average over a number of i.i.d. samples
 273 drawn from the same distribution. The resulting loss function can then be optimized via standard
 274 stochastic optimization methods. This approach has been proposed and analyzed by Bas-Serrano
 275 et al. [5]. We describe the specifics of this latter approach in Appendix C.

276 4 Analysis

277 This section gives the proof of Theorem 1 by stating the main technical results as lemmas and putting
 278 them together to obtain the final bound. In the first part of the proof, we show the upper bound on the
 279 auxiliary regret minimization game with general reward inputs and ideal updates. Then, we relate
 280 this quantity to the true expected regret by taking into account the properties of our reward estimates
 281 and the optimization errors incurred when calculating the updates. The proofs of all the lemmas are
 282 deferred to Appendix A.

283 We start by defining the idealized updates $(\hat{\mu}_t, \hat{u}_t)$ obtained by solving the update steps in Equation (5)
 284 exactly, and we let u_t be the occupancy measure induced by policy π_t that is based on the near-optimal
 285 parameters Z_t satisfying $\mathcal{G}_t(Z_t) \leq \min_Z \mathcal{G}_t(Z) + \varepsilon$. We will also let μ_t be the occupancy measure
 286 resulting from mixing u_t with the exploratory distribution μ_0 and note that $\mu_{t,h} = (1 - \gamma)u_{t,h} + \gamma\mu_{t,h}$.
 287 Using this notation, we will consider an auxiliary online learning problem with the sequence of
 288 reward functions given as $\hat{r}_{t,h}(x, a) = \langle \varphi(x, a), \hat{\theta}_{t,h} \rangle$, and study the performance of the idealized
 289 sequence $(\hat{\mu}_t, \hat{u}_t)$ therein:

$$\hat{\mathfrak{R}}_T = \sum_{t=1}^T \sum_{h=1}^{H-1} \langle \mu_h^* - \hat{u}_{t,h}, \hat{r}_{t,h} \rangle.$$

290 Our first lemma bounds the above quantity:

291 **Lemma 1.** *Suppose that $\hat{\theta}_{t,h}$ is such that $|\eta \cdot \langle \varphi(x, a), \hat{\theta}_{t,h} \rangle| < 1$ holds for all x, a . Then, the
 292 auxiliary regret satisfies*

$$\hat{\mathfrak{R}}_T \leq \eta \sum_{t=1}^T \sum_{h=1}^{H-1} \langle \hat{\mu}_{t,h}, \hat{r}_{t,h}^2 \rangle + \frac{1}{\eta} D(\mu^* \parallel \mu_0) + \frac{1}{\alpha} D_C(u^* \parallel \mu_0).$$

293 While the proof makes use of a general potential-based argument commonly used for analyzing FTRL-
 294 style algorithms, it involves several nontrivial elements exploiting the structural results concerning
 295 ONLINE Q-REPS proved in Proposition 1. In particular, these properties enable us to upper bound
 296 the potential differences in a particularly simple way. The main term on contributing to the regret $\hat{\mathfrak{R}}_T$
 297 can be bounded as follows:

298 **Lemma 2.** *Suppose that $\varphi(X_{t,h}, a)$ is satisfying $\|\varphi(X_{t,h}, a)\|_2 \leq \sigma$ for any a , $0 < \beta \leq \frac{1}{2\sigma^2}$ and
 299 $M > 0$. Then for each t and h ,*

$$\mathbb{E}_t [\langle \hat{\mu}_{t,h}, \hat{r}_{t,h}^2 \rangle] \leq \frac{4d}{3(1 - \gamma)} + (M + 1)^2 \|\hat{u}_{t,h} - u_{t,h}\|_1.$$

300 The proof of this claim makes heavy use of the fact that $\langle \hat{\mu}_{t,h}, \hat{r}_{t,h}^2 \rangle = \langle \hat{u}_{t,h}, \hat{r}_{t,h}^2 \rangle$, which is ensured
 301 by the construction of the reward estimator $\hat{r}_{t,h}$ and the constraints on the feature covariance matrices
 302 in Equation (4). This property is not guaranteed to hold under the first-order constraints (2) used in
 303 the previous works of Neu and Pike-Burke [22] and Bas-Serrano et al. [5], which eventually justifies
 304 the higher complexity of our algorithm.

305 It remains to relate the auxiliary regret to the actual regret. The main challenge is accounting for the
 306 mismatch between μ_t and u_t , and the bias of \hat{r}_t , denoted as $b_{t,h}(x, a) = \mathbb{E}_t [\hat{r}_{t,h}(x, a)] - r_{t,h}(x, a)$.

307 To address these issues, we observe that for any t, h , we have

$$\begin{aligned} \langle \mu_{t,h}, r_{t,h} \rangle &= \langle (1-\gamma)u_{t,h} + \gamma\mu_{0,h}, r_{t,h} \rangle = \langle (1-\gamma)\widehat{u}_{t,h} + \gamma\mu_{0,h}, r_{t,h} \rangle + (1-\gamma)\langle u_{t,h} - \widehat{u}_{t,h}, r_{t,h} \rangle \\ &\geq \mathbb{E}_t[\langle (1-\gamma)\widehat{u}_{t,h} + \gamma\mu_{0,h}, \widehat{r}_{t,h} \rangle] + \|b_{t,h}\|_\infty + (1-\gamma)\|u_{t,h} - \widehat{u}_{t,h}\|_1, \end{aligned}$$

308 where in the last step we used the fact that $\|r_{t,h}\|_\infty \leq 1$. After straightforward algebraic manipulations, this implies that the regret can be bounded as

$$\mathfrak{R}_T \leq (1-\gamma)\mathbb{E}[\widehat{\mathfrak{R}}_T] + \sum_{t=1}^T \sum_{h=1}^H \mathbb{E}[\gamma\langle \mu_{0,h} - \mu_h^*, r_{t,h} \rangle + \|\widehat{u}_{t,h} - u_{t,h}\|_1 + \|b_{t,h}\|_\infty]. \quad (8)$$

310 In order to proceed, we need to verify the condition $|\eta \cdot \langle \varphi(x, a), \widehat{\theta}_{t,h} \rangle| < 1$ so that we can apply
311 Lemma 1 to bound $\widehat{\mathfrak{R}}_T$. This is done in the following lemma:

312 **Lemma 3.** *Suppose that $\eta \leq \frac{2}{(M+2)}$. Then, for all, t, h , the reward estimates satisfy $\eta\|\widehat{r}_{t,h}\|_\infty < 1$.*

313 Proceeding under the condition $\eta(M+1)$, we can apply Lemma 1 to bound the first term on the
314 right-hand side of Equation (8), giving

$$\mathfrak{R}_T \leq \frac{D(\mu^* \|\mu_0)}{\eta} + \frac{D_C(u^* \|\mu_0)}{\alpha} + \frac{4\eta dHT}{3} + \gamma HT + \sum_{t,h} \mathbb{E}[(M+2)\|\widehat{u}_{t,h} - u_{t,h}\|_1 + \|b_{t,h}\|_\infty].$$

315 It remains to bound the bias of the reward estimators and the effect of the optimization errors that
316 result in the mismatch between u_t and \widehat{u}_t . The following lemma shows that this mismatch can be
317 directly controlled as a function of the optimization error:

318 **Lemma 4.** *The following bound is satisfied for all t and h : $\|\widehat{u}_{t,h} - u_{t,h}\|_1 \leq \sqrt{2\alpha\varepsilon}$.*

319 The final element in the proof is the following lemma that bounds the bias of the estimator:

320 **Lemma 5.** *For $M \geq 0$, $\beta = \frac{1}{2\sigma^2}$, we have $\|b_{t,h}\|_\infty \leq \sigma R \exp(-\gamma\beta\lambda_{\min}M)$.*

321 Putting these bounds together with the above derivations concludes the proof of Theorem 1.

322 5 Discussion

323 This paper studies the problem of online learning in MDPs, merging two important lines of work
324 on this problem concerned with linear function approximation [17, 7] and bandit feedback with
325 adversarial rewards [23, 25, 35]. Our results are the first in this setting and not directly comparable
326 with any previous work, although some favorable comparisons can be made with previous results in
327 related settings. In the tabular setting where $d = |\mathcal{X}||\mathcal{A}|$, our bounds exactly recover the minimax
328 optimal guarantees first achieved by the O-REPS algorithm of Zimin and Neu [35]. For realizable
329 linear function approximation, the work closest to ours is that of Cai et al. [7], who prove bounds of
330 order $\sqrt{d^2 H^3 T}$, which is worse by a factor of $\sqrt{d}H$ than our result. Their setting, however, is not
331 exactly comparable to ours due to the different assumptions about the feedback about the rewards
332 and the knowledge of the transition function.

333 One particular strength of our work is providing a complete analysis of the propagation of optimization
334 errors incurred while performing the updates. This is indeed a unique contribution in the related
335 literature, where the effect of such errors typically go unaddressed. Specifically, the algorithms of
336 Zimin and Neu [35], Rosenberg and Mansour [30], and Jin et al. [16] are all based on solving convex
337 optimization problems similar to ours, the effect of optimization errors or potential methods for
338 solving the optimization problems are not discussed at all. That said, we believe that the methods
339 for calculating the updates discussed in Section 3.4 are far from perfect, and more research will be
340 necessary to find truly practical optimization methods to solve this problem.

341 The most important open question we leave behind concerns the requirement to have full prior
342 knowledge of P . In the tabular case, this challenge has been successfully addressed in the adversarial
343 MDP problem recently by Jin et al. [16], whose technique is based on adjusting the constraints (1)
344 with a confidence set over the transition functions, to account for the uncertainty about the dynamics.
345 We find it plausible that a similar extension of ONLINE Q-REPS is possible by incorporating a
346 confidence set for linear MDPs, as has been done in the case of i.i.d. rewards by Neu and Pike-Burke
347 [22]. Nevertheless, the details of such an extension remain highly non-trivial, and we leave the
348 challenge of working them out open for future work.

349 **References**

- 350 [1] A. Agarwal, S. Kakade, A. Krishnamurthy, and W. Sun. FLAMBE: Structural complexity and
 351 representation learning of low rank MDPs. In *Advances in Neural Information Processing
 352 Systems (NeurIPS)*, 2020.
- 353 [2] P. Auer and R. Ortner. Logarithmic online regret bounds for undiscounted reinforcement
 354 learning. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Advances in Neural Information
 355 Processing Systems*, volume 19, pages 49–56. MIT Press, 2007.
- 356 [3] B. Awerbuch and R. D. Kleinberg. Adaptive routing with end-to-end feedback: distributed
 357 learning and geometric approaches. In *STOC 2004*, pages 45–53, 2004.
- 358 [4] M. G. Azar, I. Osband, and R. Munos. Minimax regret bounds for reinforcement learning. In
 359 *Proceedings of the 34th International Conference on Machine Learning*, pages 263–272, 2017.
- 360 [5] J. Bas-Serrano, S. Curi, A. Krause, and G. Neu. Logistic Q-learning. In *AI & Statistics*, pages
 361 3610–3618, 2021.
- 362 [6] A. N. Burnetas and M. N. Katehakis. Optimal adaptive policies for Markov Decision Processes.
 363 *Mathematics of Operations Research*, 22(1):222–255, 1997.
- 364 [7] Q. Cai, Z. Yang, C. Jin, and Z. Wang. Provably efficient exploration in policy optimization.
 365 *arXiv e-prints*, art. arXiv:1912.05830, Dec. 2019.
- 366 [8] V. Dani, S. M. Kakade, and T. P. Hayes. The price of bandit information for online optimization.
 367 In *Advances in Neural Information Processing Systems 20*, pages 345–352. 2008.
- 368 [9] C. Dann and E. Brunskill. Sample complexity of episodic fixed-horizon reinforcement learning.
 369 In *Advances in Neural Information Processing Systems*, pages 2818–2826, 2015.
- 370 [10] C. Dann, T. Lattimore, and E. Brunskill. Unifying PAC and regret: Uniform PAC bounds for
 371 episodic reinforcement learning. In *Advances in Neural Information Processing Systems 30*,
 372 pages 5713–5723. 2017.
- 373 [11] T. Dick, A. György, and Cs. Szepesvári. Online learning in Markov decision processes with
 374 changing cost sequences. In *International Conference on Machine Learning*, pages 512–520,
 375 2014.
- 376 [12] E. Even-Dar, S. M. Kakade, and Y. Mansour. Online Markov decision processes. *Math. Oper.
 377 Res.*, 34(3):726–736, 2009.
- 378 [13] R. Fruit, M. Pirodda, A. Lazaric, and R. Ortner. Efficient bias-span-constrained exploration-
 379 exploitation in reinforcement learning. In *International Conference on Machine Learning*, pages
 380 1573–1581, 2018.
- 381 [14] T. Jaksch, R. Ortner, and P. Auer. Near-optimal regret bounds for reinforcement learning.
 382 *Journal of Machine Learning Research*, 99:1563–1600, August 2010. ISSN 1532-4435.
- 383 [15] C. Jin, Z. Allen-Zhu, S. Bubeck, and M. I. Jordan. Is Q-learning provably efficient? In *Advances
 384 in Neural Information Processing Systems*, pages 4863–4873, 2018.
- 385 [16] C. Jin, T. Jin, H. Luo, S. Sra, and T. Yu. Learning adversarial MDPs with bandit feedback and
 386 unknown transition. In *International Conference on Machine Learning*, 2020.
- 387 [17] C. Jin, Z. Yang, Z. Wang, and M. I. Jordan. Provably efficient reinforcement learning with linear
 388 function approximation. In *Proceedings of the 33rd Annual Conference on Learning Theory
 389 (COLT 2020)*, pages 2137–2143, 2020.
- 390 [18] H. B. McMahan and A. Blum. Online geometric optimization in the bandit setting against an
 391 adaptive adversary. In *COLT 2004*, pages 109–123, 2004.
- 392 [19] G. Neu and G. Bartók. An efficient algorithm for learning with semi-bandit feedback. In
 393 *Proceedings of the 24th International Conference on Algorithmic Learning Theory (ALT 2013)*,
 394 pages 234–248, 2013.
- 395 [20] G. Neu and G. Bartók. Importance weighting without importance weights: An efficient
 396 algorithm for combinatorial semi-bandits. *Journal of Machine Learning Research*, 17:1–21,
 397 2016.
- 398 [21] G. Neu and J. Olkhovskaya. Efficient and robust algorithms for adversarial linear contextual
 399 bandits. In *Proceedings of the 33rd Annual Conference on Learning Theory (COLT 2020)*,
 400 pages 3049–3068, 2020.

- 401 [22] G. Neu and C. Pike-Burke. A unifying view of optimism in episodic reinforcement learning. In
402 *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- 403 [23] G. Neu, A. György, and Cs. Szepesvári. The online loop-free stochastic shortest-path problem.
404 In *Proceedings of the 23rd Annual Conference on Learning Theory (COLT 2010)*, pages
405 231–243, 2010.
- 406 [24] G. Neu, A. György, and Cs. Szepesvári. The adversarial stochastic shortest path problem with
407 unknown transition probabilities. In *Proceedings of the Fifteenth International Conference on*
408 *Artificial Intelligence and Statistics*, pages 805–813, 2012.
- 409 [25] G. Neu, A. György, Cs. Szepesvári, and A. Antos. Online Markov decision processes under
410 bandit feedback. volume 59, pages 1804–1812, 01 2013. doi: 10.1109/TAC.2013.2292137.
- 411 [26] F. Orabona. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*, 2019.
- 412 [27] A. Pacchiano, J. Lee, P. Bartlett, and O. Nachum. Near optimal policy optimization via REPS.
413 *arXiv preprint arXiv:2103.09756*, 2021.
- 414 [28] J. Peters, K. Mülling, and Y. Altun. Relative entropy policy search. In *AAAI 2010*, pages
415 1607–1612, 2010. ISBN 978-1-57735-463-5.
- 416 [29] B. Á. Pires and Cs. Szepesvári. Policy error bounds for model-based reinforcement learning
417 with factored linear models. In *Conference on Learning Theory*, pages 121–151, 2016.
- 418 [30] A. Rosenberg and Y. Mansour. Online convex optimization in adversarial Markov decision
419 processes. In *Proceedings of the 36th International Conference on Machine Learning*, pages
420 5478–5486, 2019.
- 421 [31] A. Tewari and P. L. Bartlett. Optimistic linear programming gives logarithmic regret for
422 irreducible MDPs. In *Advances in Neural Information Processing Systems 20*, pages 1505–
423 1512.
- 424 [32] C.-Y. Wei, M. Jafarnia Jahromi, H. Luo, and R. Jain. Learning infinite-horizon average-reward
425 MDPs with linear function approximation. In *Proceedings of The 24th International Conference*
426 *on Artificial Intelligence and Statistics*, pages 3007–3015, 2021.
- 427 [33] L. F. Yang and M. Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and
428 regret bound. In *Proceedings of the 36th International Conference on Machine Learning*, 2020.
- 429 [34] H. Yao, Cs. Szepesvári, B. Pires, and X. Zhang. Pseudo-MDPs and factored linear action
430 models. 10 2014. doi: 10.1109/ADPRL.2014.7010633.
- 431 [35] A. Zimin and G. Neu. Online learning in episodic markovian decision processes by relative
432 entropy policy search. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q.
433 Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 1583–1591.
434 Curran Associates, Inc., 2013.

435 **Checklist**

- 436 1. For all authors...
- 437 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
438 contributions and scope? [Yes]
- 439 (b) Did you describe the limitations of your work? [Yes] **All limitations are discussed in
440 Section 5.**
- 441 (c) Did you discuss any potential negative societal impacts of your work? [N/A] **Our
442 work is theoretical and does not raise any specific concerns about negative soci-
443 etal impact.**
- 444 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
445 them? [Yes]
- 446 2. If you are including theoretical results...
- 447 (a) Did you state the full set of assumptions of all theoretical results? [Yes] **The main set
448 of assumptions can be found in Section 2 and an additional technical assumption
449 is stated in Assumption 2 in Section 3.3.**
- 450 (b) Did you include complete proofs of all theoretical results? [Yes] **Proofs of all stated
451 results can be found in the appendix.**
- 452 3. If you ran experiments...
- 453 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
454 mental results (either in the supplemental material or as a URL)? [N/A]
- 455 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
456 were chosen)? [N/A]
- 457 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
458 ments multiple times)? [N/A]
- 459 (d) Did you include the total amount of compute and the type of resources used (e.g., type
460 of GPUs, internal cluster, or cloud provider)? [N/A]
- 461 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 462 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 463 (b) Did you mention the license of the assets? [N/A]
- 464 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 465
- 466 (d) Did you discuss whether and how consent was obtained from people whose data you're
467 using/curating? [N/A]
- 468 (e) Did you discuss whether the data you are using/curating contains personally identifiable
469 information or offensive content? [N/A]
- 470 5. If you used crowdsourcing or conducted research with human subjects...
- 471 (a) Did you include the full text of instructions given to participants and screenshots, if
472 applicable? [N/A]
- 473 (b) Did you describe any potential participant risks, with links to Institutional Review
474 Board (IRB) approvals, if applicable? [N/A]
- 475 (c) Did you include the estimated hourly wage paid to participants and the total amount
476 spent on participant compensation? [N/A]

477 **A Omitted proofs**

478 **A.1 The proof of Proposition 1**

479 The proof is based on Lagrangian duality: for each $h \in [H - 1]$, we introduce a set of multipliers
 480 $V_h \in \mathbb{R}^{|\mathcal{X}_h|}$ and $Z_h \in \mathbb{R}^{d \times d}$ corresponding to the two sets of constraints connecting $\mu_{t,h}$ and $u_{t,h}$,
 481 and $\rho_{t,h}$ for the normalization constraint of $\mu_{t,h}$. Then, we can write the Lagrangian of the constrained
 482 optimization problem as

$$\begin{aligned} \mathcal{L}(\mu, u; V, Z, \rho) &= \sum_{h=1}^{H-1} \sum_{s=1}^{t-1} \langle \mu_h, \hat{r}_{s,h} \rangle + \langle Z_h, \Phi_h^\top (\text{diag}(u_h) - \text{diag}(\mu_h)) \Phi_h \rangle \\ &\quad + \sum_{h=1}^{H-1} \left(\rho_h (1 - \langle \mu_h, \mathbf{1} \rangle) - \frac{1}{\eta} D(\mu_h \| \mu_{0,h}) - \frac{1}{\alpha} D_C(u_h \| \mu_{0,h}) \right) \\ &\quad + V_1(x_1) (1 - E^\top u_1) + \sum_{h=1}^{H-1} \langle V_{h+1}, P^\top \mu_h - E^\top u_{h+1} \rangle. \end{aligned}$$

For any $h \in [H - 1]$, for any $x \in \mathcal{X}_h, a \in A(x)$, denote $Q_Z(x, a) = \varphi(x, a)^\top Z_h(x) \varphi(x, a)$,
 $P_{x,a} V_{h+1} = \sum_{x' \in \mathcal{X}_{h+1}} P(x'|x, a) V_{h+1}(x')$ and $\Delta_{t,z}(x, a) = \sum_{s=1}^{t-1} \hat{r}_{s,h(x)}(x, a) + P_{x,a} V_{h(x)+1} -$
 $Q_Z(x, a)$. The above Lagrangian is strictly concave, so the maximum of $\mathcal{L}(\mu, d; V, Z, \rho)$ can be
 found by setting the derivatives with respect to its parameters to zero. This gives the following
 expressions for the choices of π and μ :

$$\begin{aligned} \pi_{t,h}^*(a|x) &= \pi_{0,h}(a|x) e^{\alpha(Q_Z(x,a) - V_h(x))}, \\ \mu_{t,h}^*(x, a) &= \mu_0(x, a) e^{\eta(\Delta_{t,z}(x,a) - \rho_{t,h})}, \end{aligned}$$

From the constraint $\sum_{x \in \mathcal{X}_h, a \in A(x)} \mu_{t,h}^*(x, a) = 1$ for all h , we get that

$$\rho_{t,h}^* = \frac{1}{\eta} \log \left(\sum_{x \in \mathcal{X}_h, a \in A(x)} \mu_0(x, a) e^{\eta \Delta_{t,z}(x,a)} \right)$$

483 and from the constraint $\sum_a \pi_t^*(a|x) = 1$, we get

$$V_h^*(x) = \frac{1}{\alpha} \log \left(\sum_a \pi_0(a|x) e^{\alpha Q_Z(x,a)} \right).$$

484 We will further use the notation $V_Z(x) := V_h^*(x)$. Then, by plugging $\pi_{t,h}^*, \mu_{t,h}^*, V_Z(x)$ into the
 485 Lagrangian, we get

$$\mathcal{G}_t(Z) = \mathcal{L}(\mu^*, u^*; V^*, Z, \rho^*) = \frac{1}{\eta} \sum_{h=1}^{H-1} \log \left(\sum_{x \in \mathcal{X}_h, a \in A(x)} \mu_0(x, a) e^{\eta \Delta_{t,z}(x,a)} \right) + V_Z(x_1).$$

486 Then, the solution of the optimization problem can be written as

$$\max_{\mu, u \in U} \min_{V, Z, \rho} \mathcal{L}(\mu, u; V, Z, \rho) = \min_{V, Z, \rho, \mu, u \in U} \max \mathcal{L}(\mu, u; V, Z, \rho) = \min_Z \mathcal{L}(\mu^*, u^*; V^*, Z, \rho^*) = \min_Z \mathcal{G}_t(Z).$$

487 This concludes the proof. ■

488 **A.2 The proof of Lemma 1**

489 The proof is based on a variation of the FTRL analysis that studies the evolution of the potential
 490 function Ψ_t defined for each t as

$$\Psi_t = \max_{(\mu, u) \in \mathcal{U}_\Phi^2} \left\{ \sum_{s=1}^{t-1} \sum_{h=1}^H \langle \mu_h, \hat{r}_{s,h} \rangle - \frac{1}{\eta} D(\mu \| \mu_0) - \frac{1}{\alpha} D_C(u \| u_0) \right\}.$$

491 This definition immediately implies the following bound:

$$\Psi_{T+1} \geq \sum_{s=1}^T \sum_{h=1}^{H-1} \langle \mu_h^*, \widehat{r}_{s,h} \rangle - \frac{1}{\eta} D(\mu^* \|\mu_0) - \frac{1}{\alpha} D_C(u^* \|u_0). \quad (9)$$

492 To proceed, we will heavily exploit the fact that, by Proposition 1, the potential satisfies $\Psi_t = \min_Z \mathcal{G}_t$.
 493 Introducing the notation $Z_t^* = \arg \min_Z \mathcal{G}_t(Z)$, we have

$$\begin{aligned} \Psi_{t+1} - \Psi_t &= \mathcal{G}_{t+1}(Z_{t+1}^*) - \mathcal{G}_t(Z_t^*) \leq \mathcal{G}_{t+1}(Z_t^*) - \mathcal{G}_t(Z_t^*) \\ &= \frac{1}{\eta} \sum_{h=1}^{H-1} \log \frac{\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} \mu_{0,h}(x, a) \exp \left(\eta \left(\sum_{s=1}^t \widehat{r}_{s,h}(x, a) + P_{x,a} V_{Z_t^*} - Q_{Z_t^*}(x, a) \right) \right)}{\sum_{x' \in \mathcal{X}_h, a' \in \mathcal{A}} \mu_{0,h}(x', a') \exp \left(\eta \left(\sum_{s=1}^{t-1} \widehat{r}_{s,h}(x', a') + P_{x',a'} V_{Z_t^*} - Q_{Z_t^*}(x', a') \right) \right)} \\ &= \frac{1}{\eta} \sum_{h=1}^{H-1} \log \left(\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} \mu_{t,h}(x, a) \exp(\eta \widehat{r}_{t,h}(x, a)) \right) \\ &\quad \text{(using the expression of } \mu_{t,h}(x, a) \text{ obtained in Proposition 1)} \\ &\leq \frac{1}{\eta} \sum_{h=1}^{H-1} \log \left(1 + \sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} \mu_{t,h}(x, a) \eta (\widehat{r}_{t,h}(x, a) + \eta \widehat{r}_{t,h}^2(x, a)) \right) \\ &\leq \sum_{h=1}^{H-1} (\langle \mu_{t,h}, \widehat{r}_{t,h} \rangle + \eta \langle \mu_{t,h}, \widehat{r}_{t,h}^2 \rangle), \end{aligned}$$

494 where in the last two lines we have used the inequalities $e^z \leq 1 + z + z^2$, which holds for $z \leq 1$ and
 495 $\log(1 + z) \leq z$, which holds for all $z > -1$, which conditions are verified due to our constraint on η .
 496 Summing up both sides for all t and combining the result with the inequality (9), we obtain

$$\widehat{\mathfrak{R}}_T = \sum_{s=1}^T \sum_{h=1}^{H-1} \langle \mu_h^*, \widehat{r}_{s,h} \rangle - \sum_{t=1}^T \langle \mu_t, \widehat{r}_t \rangle \leq \eta \sum_{t=1}^T \sum_{h=1}^{H-1} \langle \mu_{t,h}, \widehat{r}_{t,h}^2 \rangle + \frac{1}{\eta} D(\mu^* \|\mu_0) + \frac{1}{\alpha} D_C(u^* \|u_0),$$

497 concluding the proof. ■

498 A.3 The proof of Lemma 2

499 We start by using the the definition of $\widehat{\theta}_{t,h}$ to obtain

$$\begin{aligned} \mathbb{E}_t \left[\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} \widehat{\mu}_{t,h}(x, a) \langle \varphi(x, a), \widehat{\theta}_{t,h} \rangle^2 \right] &= \mathbb{E}_t \left[\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} \widehat{\mu}_{t,h}(x, a) \text{tr} \left(\varphi(x, a) \varphi(x, a)^\top \widehat{\theta}_{t,h} \widehat{\theta}_{t,h}^\top \right) \right] \\ &= \mathbb{E}_t \left[\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} \widehat{u}_{t,h}(x, a) \text{tr} \left(\varphi(x, a) \varphi(x, a)^\top \widehat{\theta}_{t,h} \widehat{\theta}_{t,h}^\top \right) \right] \\ &\quad \text{(by the constraint } \Phi_h^\top \text{diag}(\widehat{\mu}_t) \Phi_h = \Phi_h^\top \text{diag}(\widehat{u}_t) \Phi_h) \\ &= \mathbb{E}_t \left[\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} u_{t,h}(x, a) \text{tr} \left(\varphi(x, a) \varphi(x, a)^\top \widehat{\theta}_{t,h} \widehat{\theta}_{t,h}^\top \right) \right] \\ &\quad + \sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} (u_{t,h}(x, a) - \widehat{u}_{t,h}(x, a)) \mathbb{E}_t \left[\langle \varphi(x, a), \widehat{\theta}_{t,h} \rangle^2 \right] \\ &\leq \mathbb{E}_t \left[\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} u_{t,h}(x, a) \text{tr} \left(\varphi(x, a) \varphi(x, a)^\top \widehat{\theta}_{t,h} \widehat{\theta}_{t,h}^\top \right) \right] + \|u_{t,h} - \widehat{u}_{t,h}\|_1 \cdot \|\mathbb{E}_t [\widehat{r}_{t,h}^2]\|_\infty \end{aligned}$$

500 The second term can be bounded straightforwardly by $\|u_{t,h} - \widehat{u}_{t,h}\|_1 (M+1)^2$, using Lemma 3 to
 501 bound $\|\widehat{r}_{t,h}\|_\infty \leq (M+1)$. As for the first term, we have

$$\begin{aligned}
 & (1-\gamma)\mathbb{E}_t \left[\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} u_{t,h}(x,a) \operatorname{tr} \left(\varphi(x,a) \varphi(x,a)^\top \widehat{\theta}_{t,h} \widehat{\theta}_{t,h}^\top \right) \right] \\
 & \leq (1-\gamma)\mathbb{E}_t \left[\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} \operatorname{tr} \left(u_{t,h}(x,a) \varphi(x,a) \varphi(x,a)^\top \widehat{\Sigma}_{t,h}^+ \varphi(X_{t,h}, A_{t,h}) \varphi(X_{t,h}, A_{t,h})^\top \widehat{\Sigma}_{t,h}^+ \right) \right], \\
 & \leq (1-\gamma)\mathbb{E}_t \left[\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} \operatorname{tr} \left(u_{t,h}(x,a) \varphi(x,a) \varphi(x,a)^\top \widehat{\Sigma}_{t,h}^+ \varphi(X_{t,h}, A_{t,h}) \varphi(X_{t,h}, A_{t,h})^\top \widehat{\Sigma}_{t,h}^+ \right) \right] \\
 & \quad + \gamma \mathbb{E}_t \left[\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}} \operatorname{tr} \left(u(x,a) \varphi(x,a) \varphi(x,a)^\top \widehat{\Sigma}_{t,h}^+ \varphi(X_{t,h}, A_{t,h}) \varphi(X_{t,h}, A_{t,h})^\top \widehat{\Sigma}_{t,h}^+ \right) \right] \\
 & = \mathbb{E}_t \left[\operatorname{tr} \left(\Sigma_{t,h} \widehat{\Sigma}_{t,h}^+ \Sigma_{t,h} \widehat{\Sigma}_{t,h}^+ \right) \right],
 \end{aligned}$$

502 where we used $|r_{t,h}(X_{t,h}, A_{t,h})| \leq 1$ in the first inequality. For ease of readability, we will omit the
 503 indices h in the rest of the proof. Using the definition of Σ_t^+ and elementary manipulations, we get

$$\begin{aligned}
 \mathbb{E}_t \left[\operatorname{tr} \left(\Sigma_t \Sigma_t^+ \Sigma_t \Sigma_t^+ \right) \right] &= \beta^2 \cdot \mathbb{E}_t \left[\operatorname{tr} \left(\Sigma^* \left(\sum_{k=0}^M C_k \right) \Sigma_t \left(\sum_{j=0}^M C_j \right) \right) \right] \\
 &= \beta^2 \mathbb{E}_t \left[\sum_{k=0}^M \sum_{j=0}^M \operatorname{tr} \left(\Sigma_t C_k \Sigma_t C_j \right) \right] \\
 &= \beta^2 \mathbb{E}_t \left[\sum_{k=0}^M \operatorname{tr} \left(\Sigma_t C_k \Sigma_t C_k \right) \right] + 2\beta^2 \mathbb{E}_t \left[\sum_{k=0}^M \sum_{j=k+1}^M \operatorname{tr} \left(\Sigma_t C_k \Sigma_t C_j \right) \right].
 \end{aligned}$$

504 Let us first address the first term on the right hand side. To this end, consider any symmetric positive
 505 definite matrix S that commutes with Σ_t and observe that

$$\begin{aligned}
 & \mathbb{E}_t \left[(I - \beta B_k) S (I - \beta B_k) \right] \\
 &= \mathbb{E} \left[(I - \beta \varphi(X(k), A(k)) \varphi(X(k), A(k))^\top) S (I - \beta \varphi(X(k), A(k)) \varphi(X(k), A(k))^\top) \right] \\
 &= S - \beta \mathbb{E} \left[\varphi(X(k), A(k)) \varphi(X(k), A(k))^\top S \right] - \beta \mathbb{E}_t \left[S \varphi(X(k), A(k)) \varphi(X(k), A(k))^\top \right] \\
 & \quad + \beta^2 \mathbb{E}_t \left[\varphi(X(k), A(k)) \varphi(X(k), A(k))^\top S \varphi(X(k), A(k)) \varphi(X(k), A(k))^\top \right] \\
 & \preceq S - 2\beta S \Sigma_t + \beta^2 \sigma^2 S \Sigma_t = S (I - \beta(2 - \beta\sigma^2) \Sigma_t),
 \end{aligned}$$

506 where we used our assumption that $\|\varphi(X(k), A(k))\| \leq \sigma$, which implies
 507 $\mathbb{E}_t \left[\|\varphi(X(k), A(k))\|_2^2 \varphi(X(k), A(k)) \varphi(X(k), A(k))^\top \right] \preceq \sigma^2 \Sigma_t$. Now, recalling the defini-
 508 tion $C_k = \prod_{j=1}^k (I - \beta B_j)$ and using the above relation repeatedly, we can obtain

$$\begin{aligned}
 \operatorname{tr} \left(\mathbb{E}_t \left[\Sigma_t C_k \Sigma_t C_k \right] \right) &= \operatorname{tr} \left(\mathbb{E}_t \left[\Sigma_t C_{k-1} \mathbb{E}_t \left[(I - \beta B_k) \Sigma_t (I - \beta B_k) \right] C_{k-1} \right] \right) \\
 &\leq \operatorname{tr} \left(\mathbb{E}_t \left[\Sigma_t C_{k-1} \Sigma_t (I - \beta(2 - \beta\sigma^2) \Sigma_t) C_{k-1} \right] \right) \quad (10) \\
 &\leq \dots \leq \operatorname{tr} \left(\Sigma_t \Sigma_t (I - \beta(2 - \beta\sigma^2) \Sigma_t)^k \right).
 \end{aligned}$$

509 Thus, we can see that

$$\begin{aligned}
 \beta^2 \sum_{k=0}^M \operatorname{tr} \left(\mathbb{E}_t \left[\Sigma_t C_k \Sigma_t C_k \right] \right) &= \beta^2 \sum_{k=0}^M \operatorname{tr} \left(\Sigma_t \Sigma_t (I - \beta(2 - \beta\sigma^2) \Sigma_t)^k \right) \\
 &= \frac{\beta^2}{\beta(2 - \beta\sigma^2)} \operatorname{tr} \left(\Sigma_t \Sigma_t \Sigma_t^{-1} (I - (I - \beta(2 - \beta\sigma^2) \Sigma_t)^M) \right) \leq \frac{\beta \operatorname{tr}(\Sigma_t)}{2 - \beta\sigma^2} \leq \frac{2\beta \operatorname{tr}(\Sigma_t)}{3},
 \end{aligned}$$

510 where we used the condition $\beta \leq \frac{1}{2\sigma^2}$ and the fact that $(I - \beta(2 - \beta\sigma^2)\Sigma_t)^M \succcurlyeq 0$ by the same
 511 condition. We can finally observe that our assumption on the contexts implies $\text{tr}(\Sigma_t) \leq \text{tr}(\sigma^2 I) =$
 512 $\sigma^2 d$, so again by our condition on β we have $\beta \text{tr}(\Sigma_t) \leq \frac{d}{2}$, and the first term is bounded by $\frac{d}{3}$.

513 Moving on to the second term, we first note that for any $j > k$, the conditional expectation of B_j given
 514 $B_{\leq k} = (B_1, B_2, \dots, B_k)$ satisfies $\mathbb{E}[C_j | B_{\leq k}] = C_j(I - \beta\Sigma)^{j-k}$ due to conditional independence
 515 of all B_j given B_k , for $i > k$. We make use of this equality by writing

$$\begin{aligned}
 \beta^2 \sum_{k=0}^M \sum_{j=k+1}^M \mathbb{E}[\text{tr}(\Sigma_t C_k \Sigma_t C_j)] &= \beta^2 \sum_{k=0}^M \mathbb{E} \left[\mathbb{E} \left[\sum_{j=k+1}^M \text{tr}(\Sigma_t C_k \Sigma_t C_j) \middle| B_{\leq k} \right] \right] \\
 &= \beta^2 \sum_{k=0}^M \mathbb{E} \left[\mathbb{E} \left[\sum_{j=k+1}^M \text{tr}(\Sigma_t C_k \Sigma_t C_j (I - \beta\Sigma_t)^{j-k}) \middle| B_{\leq k} \right] \right] \\
 &= \beta \sum_{k=0}^M \mathbb{E} \left[\mathbb{E} \left[\text{tr}(\Sigma_t C_k \Sigma_t C_k \Sigma_t^{-1} (I - (I - \beta\Sigma_t)^{M-k})) \middle| B_{\leq k} \right] \right] \\
 &\leq \beta \sum_{k=0}^M \mathbb{E} \left[\mathbb{E} \left[\text{tr}(\Sigma_t C_k \Sigma_t C_k \Sigma_t^{-1}) \middle| B_{\leq k} \right] \right] \\
 &\quad \text{(due to } (I - \beta\Sigma_t)^{M-k} \succcurlyeq 0 \text{)} \\
 &\leq \beta \sum_{k=0}^M \text{tr}(\Sigma_t \Sigma_t (I - \beta(2 - \beta\sigma^2)\Sigma_t)^k \Sigma_t^{-1}) \\
 &\quad \text{(by the same argument as in Equation (10))} \\
 &\leq \frac{1}{(2 - \beta\sigma^2)} \text{tr}(\Sigma_t \Sigma_t \Sigma_t^{-1} (I - (I - \beta(2 - \beta\sigma^2)\Sigma_t)^M \Sigma_t^{-1})) \\
 &\leq \text{tr}(\Sigma_t \Sigma_t^{-1} \Sigma_t \Sigma_t^{-1}) = d.
 \end{aligned}$$

516 The proof of the lemma is finished by putting everything together. ■

517 A.4 The proof of Lemma 5

518 We first observe that the bias of $\hat{\theta}_{t,h}$ can be easily expressed as

$$\begin{aligned}
 \mathbb{E}_t[\hat{\theta}_{t,h}] &= \mathbb{E}_t \left[\widehat{\Sigma}_{t,h}^+ \varphi(X_{t,h}, A_{t,h}) \varphi(X_{t,h}, A_{t,h})^\top \theta_{t,h} \right] \\
 &= \mathbb{E}_t \left[\widehat{\Sigma}_{t,h}^+ \right] \mathbb{E}_t \left[\varphi(X_{t,h}, A_{t,h}) \varphi(X_{t,h}, A_{t,h})^\top \right] \theta_{t,h} \\
 &= \mathbb{E}_t \left[\widehat{\Sigma}_{t,h}^+ \right] \Sigma_{t,h} \theta_{t,h} = \theta_{t,h} - (I - \beta\Sigma_{t,h})^M \theta_{t,h}.
 \end{aligned}$$

519 Thus, the bias is bounded as

$$\left| \mathbb{E}_t \left[\varphi(X_{t,h}, a)^\top (I - \beta\Sigma_{t,h})^M \theta_{t,h} \right] \right| \leq \|\varphi(X_{t,h}, a)\|_2 \cdot \|\theta_{t,h}\|_2 \|(I - \beta\Sigma_{t,h})^M\|_{\text{op}}.$$

520 In order to bound the last factor above, observe that $\Sigma_{t,h} \succcurlyeq \gamma\Sigma_h$ due to the uniform exploration used
 521 in the first layer by MDP-LINEXP3, which implies that

$$\|(I - \beta\Sigma_{t,h})^M\|_{\text{op}} \leq (1 - \gamma\beta\lambda_{\min})^M \leq \exp(-\gamma\beta\lambda_{\min}M),$$

522 where the second inequality uses $1 - z \leq e^{-z}$ that holds for all z . This concludes the proof. ■

523 A.5 The proof of Lemma 4

524 The proof consists of two main components: proving that the conditional relative entropy between u_t
 525 and \hat{u}_t can be bounded in terms of the optimization error ε , and then using this quantity to bound the
 526 total variation distance between these occupancy measures. For ease of readability, we state these
 527 results as separate lemmas.

528 We will first need the following statement:

529 **Lemma 6.** $D_C(\hat{u}_t \| u_t) \leq \alpha \varepsilon$.

530 The proof follows along similar lines as the proof of Lemma 1 in Bas-Serrano et al. [5]. To preserve
531 clarity, we delegate its proof to Appendix A.6 below. The second lemma lemma bounds the relative
532 entropy between two occupancy measures in terms of their *conditional* relative entropies:

533 **Lemma 7.** For any two occupancy measures u and u' and any h , we have

$$D(u_h \| u'_h) \leq \sum_{k=1}^h D_C(u_k \| u'_k).$$

534 *Proof.* The proof follows from exploiting some basic properties of the relative entropy. Specifically,
535 the result follows from the following chain of inequalities:

$$\begin{aligned} D(u_h \| u'_h) &= D(E^\top u_h \| E^\top u'_h) + D_C(u_h \| u'_h) \\ &\quad \text{(by the chain rule of the relative entropy)} \\ &= D(P^\top u_{h-1} \| P^\top u'_{h-1}) + D_C(u_h \| u'_h) \\ &\quad \text{(by the fact that } u \text{ and } u' \text{ are valid occupancy measures)} \\ &\leq D(u_{h-1} \| u'_{h-1}) + D_C(u_h \| u'_h) \\ &\quad \text{(by the data processing inequality)} \\ &\leq \dots \leq \sum_{k=1}^h D_C(u_k \| u'_k), \end{aligned}$$

536 where the last step follows from iterating the same argument for all layers. ■

537 Putting the above two lemmas together and using Pinsker's inequality, we obtain

$$\|\hat{u}_{t,h} - u_{t,h}\|_1 \leq \sqrt{2D(\hat{u}_{t,h} \| u_{t,h})} \leq \sqrt{2 \sum_{k=1}^h D_C(\hat{u}_{t,k} \| u_{t,k})} \leq \sqrt{2D_C(\hat{u}_t \| u_t)} \leq \sqrt{2\alpha\varepsilon},$$

538 concluding the proof of Lemma 4. ■

539 A.6 The proof of Lemma 6

540 For the proof, let us introduce the notation $\tilde{\mu}_{t,h}$ with

$$\tilde{\mu}_{t,h}(x, a) = \frac{\mu_{0,h}(x, a) e^{\eta \Delta_{t,Z_t}(x,a)}}{\sum_{(x',a') \in (\mathcal{X}_h \times \mathcal{A})} \mu_{0,h}(x', a') e^{\eta \Delta_{t,Z_t}(x',a')}}.$$

541 and also $\mathcal{G}_{t,h}(Z) = \frac{1}{\eta} \log \left(\sum_{x \in \mathcal{X}_h, a \in \mathcal{A}(x)} \mu_0(x, a) e^{\eta \Delta_{t,Z}(x,a)} \right)$ and $Z_t^* = \arg \min_Z \mathcal{G}_t(Z)$. Then,
542 observe that

$$\begin{aligned} D(\hat{\mu}_{t,h} \| \tilde{\mu}_{t,h}) &= \sum_{x,a \in \mathcal{X}_h \times \mathcal{A}} \hat{\mu}_{t,h}(x, a) \log \frac{\hat{\mu}_{t,h}(x, a)}{\tilde{\mu}_{t,h}(x, a)} \\ &= \eta \langle \hat{\mu}_{t,h}, \Delta_{t,Z_t^*} - \mathcal{G}_{t,h}(Z_t^*) \mathbf{1} - \Delta_{t,Z_t} + \mathcal{G}_{t,h}(Z_t) \mathbf{1} \rangle \\ &= \eta \langle \hat{\mu}_{t,h}, P_h V_{Z_t^*} - Q_{Z_t^*} - P_h V_{Z_t} + Q_{Z_t} \rangle + \eta (\mathcal{G}_{t,h}(Z_t^*) - \mathcal{G}_{t,h}(Z_t)) \\ &= \eta \sum_{(x,a) \in (\mathcal{X}_h \times \mathcal{A})} \sum_{x' \in \mathcal{X}_{h+1}} \hat{\mu}_{t,h}(x, a) P(x'|x, a) (V_{Z_t^*}(x') - V_{Z_t}(x')) \\ &\quad + \eta (\mathcal{G}_{t,h}(Z_t^*) - \mathcal{G}_{t,h}(Z_t)) + \eta \sum_{(x,a) \in (\mathcal{X}_h \times \mathcal{A})} \hat{\mu}_{t,h}(x, a) \varphi(x, a)^\top (Z_{t,h} - Z_{t,h}^*) \varphi(x, a) \\ &= \eta \sum_{(x',a') \in (\mathcal{X}_{h+1} \times \mathcal{A})} \hat{u}_{t,h+1}(x', a') (V_{Z_t^*}(x') - V_{Z_t}(x')) + \eta (\mathcal{G}_{t,h}(Z_t) - \mathcal{G}_{t,h}(Z_t^*)). \end{aligned}$$

$$+ \eta \sum_{(x,a) \in (\mathcal{X}_h \times \mathcal{A})} \widehat{u}_{t,h}(x,a) \varphi(x,a)^\top (Z_{t,h} - Z_{t,h}^*) \varphi(x,a).$$

543 Here, the last equality follows from the fact that $(\widehat{\mu}_t, \widehat{u}_t)$ satisfy the constraints of the optimization
 544 problem (5). On the other hand, we have

$$\begin{aligned} D_C(\widehat{u}_{t,h} \| u_{t,h}) &= \sum_{(x,a) \in (\mathcal{X}_h \times \mathcal{A})} \widehat{u}_{t,h}(x,a) \log \frac{\widehat{\pi}_{t,h}(a|x)}{\pi_{t,h}(a|x)} \\ &= \alpha \sum_{(x,a) \in (\mathcal{X}_h \times \mathcal{A})} \widehat{u}_{t,h}(x,a) \sum_{x' \in \mathcal{X}_{h+1}} P(x'|x,a) (V_{Z_t^*}(x') - V_{Z_t}(x')) \\ &\quad + \alpha \sum_{(x,a) \in (\mathcal{X}_h \times \mathcal{A})} \widehat{u}_{t,h}(x,a) \varphi(x,a)^\top (Z_{t,h} - Z_{t,h}^*) \varphi(x,a) \\ &= \alpha \sum_{(x',a') \in (\mathcal{X}_{h+1} \times \mathcal{A})} \widehat{u}_{t,h+1}(x',a') (V_{Z_t^*}(x') - V_{Z_t}(x')) \\ &\quad + \alpha \sum_{(x,a) \in (\mathcal{X}_h \times \mathcal{A})} \widehat{u}_{t,h}(x,a) \varphi(x,a)^\top (Z_{t,h} - Z_{t,h}^*) \varphi(x,a), \end{aligned}$$

545 where the last equality follows from the fact that \widehat{u}_t is a valid occupancy measure, as shown in
 546 Equation (3). Putting the two equalities together, we get

$$\frac{D(\widehat{\mu}_{t,h} \| \mu_{t,h})}{\eta} - \frac{D_C(\widehat{u}_{t,h} \| u_{t,h})}{\alpha} = \mathcal{G}_{t,h}(Z_t) - \mathcal{G}_{t,h}(Z_t^*).$$

547 Then, summing up over all h gives

$$\frac{D(\widehat{\mu}_t \| \mu_t)}{\eta} - \frac{D_C(\widehat{u}_t \| u_t)}{\alpha} = \sum_{h=1}^H (\mathcal{G}_{t,h}(Z_t) - \mathcal{G}_{t,h}(Z_t^*)) = \mathcal{G}_t(Z_t) - \mathcal{G}_t(Z_t^*) \leq \varepsilon.$$

548 Reordering gives the result. ■

549 A.7 The proof of Lemma 3

550 The claim is proven by the following straightforward calculation:

$$\begin{aligned} \eta \cdot |\langle \varphi(X_{t,h}, a), \widehat{\theta}_t \rangle| &= \eta \cdot |\varphi(X_{t,h}, a)^\top \widehat{\Sigma}_{t,h}^+ \varphi(X_{t,h}, a) \langle \varphi(X_{t,h}, a), \theta_t \rangle| \\ &\leq \eta |\varphi(X_{t,h}, a)^\top \widehat{\Sigma}_{t,h}^+ \varphi(X_{t,h}, a)| \leq \eta \sigma^2 \left\| \widehat{\Sigma}_{t,h}^+ \right\|_{\text{op}} \\ &\leq \eta \sigma^2 \beta \left(1 + \sum_{k=1}^M \|C_{k,h}\|_{\text{op}} \right) \leq \eta(M+1)/2, \end{aligned}$$

551 where we used the fact that our choice of β ensures $\|C_{k,h}\|_{\text{op}} = \left\| \prod_{j=0}^k (I - \beta B_{j,h}) \right\|_{\text{op}} \leq 1$. ■

552 B Fast Matrix Geometric Resampling

553 The naïve implementation of the MGR procedure presented in the main text requires $O(MKHd +$
 554 $MHd^2)$ time due to the matrix-matrix multiplications involved. In this section we explain how to
 555 compute $\widehat{\theta}_t$ in $O(MKHd)$ time, exploiting the fact that the matrices $\widehat{\Sigma}_{t,h}$ never actually need to be
 556 computed, since the algorithm only works with products of the form $\widehat{\Sigma}_{t,h} \varphi(X_{t,h}, A_{t,h})$ for vectors
 557 $X_{t,h}$, $h \in [H]$. This motivates the following procedure:

Fast Matrix Geometric Resampling

Input: simulator of transition function P , policy π_t

Initialization: Compute $Y_{0,h} = \varphi(x_h)$ for all $h \in [H]$.

For $k = 1, \dots, M$, **repeat:**

1. Generate a path $U(i) = \{(X_1(i), A_1(i)), \dots, (X_H(i), A_H(i))\}$, following the policy π_t in the simulator of P ,

2. **For** $h = 1, \dots, H$, **repeat:**

- (a) if $A_h(k) = a_h$, set $Y_{k,h} = Y_{k-1,h} - \beta \langle Y_{k-1,h}, \varphi(X_h(k), A_h(k)) \rangle \varphi(X_h(k), A_h(k))$,
- (b) otherwise, set $Y_{k,h} = Y_{k-1,h}$.

Return $q_{t,h} = \beta Y_{0,h} + \beta \sum_{k=1}^M Y_{k,h}$ for all $h \in [H]$.

560 It is easy to see from the above procedure that each iteration k can be computed using $(K +$
 561 $1)Hd$ vector-vector multiplications: sampling each action $A_h(k)$ takes Kd time due to having to
 562 compute the products $\langle \varphi(X_h(k)), \sum_{s=1}^{t-1} \hat{\theta}_{s,a,h} \rangle$ for each action a , and updating $Y_{k,h}$ can be done by
 563 computing the product $\langle Y_{k-1,h}, \varphi(X_h(k)) \rangle$. Overall, this results in a total runtime of order $MKHd$
 564 as promised above.

565 C Implementation by optimizing the empirical loss

566 This section outlines a possible implementation of the policy update steps based on approximate
 567 minimization of an empirical counterpart of the loss function \mathcal{G}_t . To this end, we define

$$\mathcal{G}_{t,h}(Z) = \frac{1}{\eta} \log \left(\sum_{x,a} \mu_0(x,a) e^{\eta \Delta_{Z,t,h}(x,a)} \right)$$

568 and its empirical counterpart that replaces the expectation by an empirical mean over state-action
 569 pairs sampled from μ_0 . Concretely, for all h , we let $(X_h(i), A_h(i))_{i=1}^N$ be N independent samples
 570 from μ_0 that can be obtained by running policy π_0 in the transition model P . Using these samples,
 571 we define

$$\hat{\mathcal{G}}_{t,h}(Z) = \frac{1}{\eta} \log \left(\sum_{n=1}^N e^{\eta \Delta_{Z,t,h}(X_h(i), A_h(i))} \right). \quad (11)$$

572 This objective function has several desirable properties: it is convex in Z , has bounded gradients,
 573 and is $(\alpha + \eta)$ -smooth. Furthermore, its gradients can be evaluated efficiently in $\mathcal{O}(N)$ time, given
 574 that we can efficiently evaluate expectations of the form $\sum_{x'} P(x'|x,a) V(x')$. As a result, it can be
 575 optimized up to arbitrary precision ε in time polynomial in $1/\varepsilon$ and N .

576 The downside of this estimator is that it is potentially biased. Nevertheless, as the following lemma
 577 shows, it is well-concentrated around the true objective function, under some reasonable conditions:

578 **Lemma 8.** Fix Z and suppose that $|\Delta_Z(x,a)| \leq B$ for all x,a . Then, with probability at least $1 - \delta$,
 579 the following holds:

$$\left| \hat{\mathcal{G}}_{t,h}(Z) - \mathcal{G}_{t,h}(Z) \right| \leq 56 \sqrt{\frac{\log(1/\delta)}{N}}.$$

580 This statement is a variant of Theorem 1 from Bas-Serrano et al. [5], with the key difference being
 581 that being able to exactly calculate expectations with respect to $P(\cdot|x,a)$ enables us to prove a tighter
 582 bound.

583 *Proof.* Let us start by defining the shorthand notations $\hat{S}_i = \Delta_{Z,t}(X_h(i), A_h(i))$ and $W =$
 584 $\frac{1}{N} \sum_{i=1}^N e^{\eta S_i}$. Furthermore, we define the function

$$f(s_1, s_2, \dots, s_N) = \frac{1}{N} \sum_{i=1}^N e^{\eta s_i}$$

585 and notice that it satisfies the bounded-differences property

$$f(s_1, s_2, \dots, s_i, \dots, s_N) - f(s_1, s_2, \dots, s'_i, \dots, s_N) = \frac{1}{N} (e^{\eta s_i} - e^{\eta s'_i}) \leq \frac{\eta e^{2\eta B}}{N}.$$

586 Here, the last step follows from Taylor's theorem that implies that there exists a $\chi \in (0, 1)$ such that

$$e^{\eta s'_i} = e^{\eta s_i} + \eta e^{\eta \chi (s'_i - s_i)}$$

587 holds, so that $e^{\eta s'_i} - e^{\eta s_i} = \eta e^{\eta \chi (s'_i - s_i)} \leq \eta e^{2\eta B}$, where we used the assumption that $|s_i - s'_i| \leq 2B$
 588 in the last step. Notice that our assumption $\eta B \leq 1$ further implies that $e^{2\eta B} \leq e^2$. Thus, also
 589 noticing that $W = f(S_1, \dots, S_N)$, we can apply McDiarmid's inequality that to show that the
 590 following holds with probability at least $1 - \delta'$:

$$|W - \mathbb{E}[W]| \leq \eta e^2 \sqrt{\frac{\log(2/\delta')}{N}}. \quad (12)$$

591 Thus, we can write

$$\begin{aligned} \widehat{\mathcal{G}}_{t,h}(\theta) - \mathcal{G}_{t,h}(\theta) &= \frac{1}{\eta} \log(W) - \frac{1}{\eta} \log(\mathbb{E}[W]) = \frac{1}{\eta} \log\left(\frac{W}{\mathbb{E}[W]}\right) \\ &= \frac{1}{\eta} \log\left(1 + \frac{W - \mathbb{E}[W]}{\mathbb{E}[W]}\right) \leq \frac{W - \mathbb{E}[W]}{\eta \mathbb{E}[W]} \leq e^4 \sqrt{\frac{\log(2/\delta')}{N}}, \end{aligned}$$

592 where the last line follows from the inequality $\log(1 + u) \leq u$ that holds for $u > -1$ and our
 593 assumption on η that implies $W \geq e^{-2}$. Similarly, we can show

$$\mathcal{G}_{t,h}(\theta) - \widehat{\mathcal{G}}_{t,h}(\theta) = \frac{1}{\eta} \log\left(1 + \frac{\mathbb{E}[W] - W}{W}\right) \leq \frac{\mathbb{E}[W] - W}{\eta W} \leq e^4 \sqrt{\frac{\log(2/\delta')}{N}}.$$

594 This concludes the proof. ■