
No-Regret Learning with High-Probability in Adversarial Markov Decision Processes

Mahsa Ghasemi^{1,2}

Abolfazl Hashemi²

Haris Vikalo¹

Ufuk Topcu^{2,3}

¹Department of Electrical and Computer Engineering, University of Texas at Austin

²Oden Institute for Computational Engineering and Sciences, University of Texas at Austin

³Department of Aerospace Engineering and Engineering Mechanics, University of Texas at Austin

Abstract

In a variety of problems, a decision-maker is unaware of the loss function associated with a task, yet it has to minimize this unknown loss in order to accomplish the task. Furthermore, the decision-maker’s task may evolve, resulting in a varying loss function. In this setting, we explore sequential decision-making problems modeled by adversarial Markov decision processes, where the loss function may arbitrarily change at every time step. We consider the bandit feedback scenario, where the agent observes only the loss corresponding to its actions. We propose an algorithm, called *online relative-entropy policy search with implicit exploration*, that achieves a sublinear regret not only in expectation but, more importantly, with high probability. In particular, we prove that by employing an *optimistically biased* loss estimator, the proposed algorithm achieves a regret of $\tilde{O}((T|\mathcal{A}||\mathcal{S}|)^{\frac{2}{3}}\sqrt{\tau})$, where $|\mathcal{S}|$ is the number of states, $|\mathcal{A}|$ is the number of actions, τ is the mixing time, and T is the time horizon. To our knowledge, the proposed algorithm is the first scheme that enjoys such high-probability regret bounds for general adversarial Markov decision processes under the presence of bandit feedback.

1 INTRODUCTION

A central notion in the analysis of online and sequential decision-making systems is that of Markov decision processes (MDPs). MDPs enable modeling decision-makers (learners) that need to make a sequence of decisions in the presence of uncertainty in the decision-maker’s environment. In this scenario, a loss (or reward) function captures the task expected from the learner. Therefore, the decision-maker’s goal is to design a learning algorithm that, despite operating

under uncertainty, learns a policy with the lowest cumulative loss (or the highest cumulative reward). In a traditional MDP problem, one assumes that the environment’s dynamics and the losses are stationary (i.e., time-invariant) throughout the time horizon of the interaction between the learner and the environment. However, the stationarity assumption does not hold in scenarios where the agent’s task evolves over time.

The so-called adversarial MDP (A-MDP) [10] is a new paradigm that enables the study of sequential decision-making problems with evolving tasks. In particular, in an A-MDP, the environment’s dynamics remain invariant while the loss function changes arbitrarily over time.¹ The learner aims to follow a policy that minimizes its loss in expectation over the time horizon. A standard metric for evaluating the learner’s performance is regret, i.e., the difference between the learner’s loss and the loss occurred by the best (stationary stochastic) policy in hindsight.

We focus on the setting of A-MDPs with bandit feedback, i.e., after each action, the learner observes only the corresponding loss but not the entire loss function. We consider the general class of uniformly ergodic adversarial MDPs for which the loss may change at every time step and provide the first no-regret algorithm that achieves a high-probability regret bound in this setting.

The state-of-the-art algorithms for learning in A-MDPs with bandit feedback guarantee sublinear regret of $\tilde{O}(\sqrt{T})$ in expectation², where T denotes the time horizon. However, it remains a challenge to establish algorithms that attain sublinear regrets with high probability. Since in many practical settings, e.g., robotics and recommender systems, a learner may operate only once in an environment, a high-probability regret bound is more desirable than a bound in expectation. Nevertheless, high-probability guarantees are considered to be significantly more difficult to obtain than expected guarantees in online tasks with bandit feedback [18].

¹We assume a setting where the adversarial changes in the loss are *oblivious* to the past actions taken by the learner.

²The notation \tilde{O} hides $\log(\cdot)$ factors.

Table 1: Overview of theoretical regret guarantees of learning algorithms for uniformly ergodic A-MDPs. Expected stands for in expectation while High Probability stands for with high probability. $\beta \in (0, 1]$ in [20, 21] is a lower bound on all stationary distributions. The result in [21] only holds for $T = \tilde{\Omega}(|\mathcal{A}|\tau\beta^{-3})$.

Reference	Related Setting	Regret Bound	Regret Type
[10]	full feedback	$\mathcal{O}(\tau^2 \sqrt{T \log \mathcal{A} })$	Expected
[20]	bandit feedback	$\mathcal{O}(T^{2/3} \tau \sqrt{\log(T) \mathcal{A} \log(\mathcal{A})/\beta})$	Expected
[21]	bandit feedback	$\mathcal{O}(\sqrt{\tau^3 T \log(T) \mathcal{A} \log(\mathcal{A})/\beta})$	Expected, Only for large T
[6]	full feedback	$\mathcal{O}(\sqrt{\tau (\log \mathcal{A} \mathcal{S}) T \log T})$	Expected
This work	bandit feedback	$\mathcal{O}((T \mathcal{A} \mathcal{S})^{2/3} \sqrt{\tau \log(\mathcal{S} \mathcal{A}) \log T \log 1/\delta})$	Expected, High Probability

Contribution. We propose a learning algorithm for A-MDPs with bandit feedback that achieves a sublinear regret guarantee with high probability. We consider the linear programming formulation of MDPs where the decision variables are the occupancy measure of state-action pairs. This adaptation enables us to employ online mirror descent as a building block of our proposed scheme to learn low-regret occupancy measures. Furthermore, inspired by the idea of implicit exploration [14, 18] for adversarial bandits, we design a new *optimistically biased* estimator of the loss function for A-MDPs with bandit feedback to establish regret guarantees that hold with arbitrarily high probability. Our specific contributions are as follows:

- We propose a new online learning algorithm, called *online relative-entropy policy search with implicit exploration*, for A-MDPs under bandit feedback that employs a novel optimistically biased loss estimator.
- We design a novel optimistically biased loss estimator which implicitly promotes the learner to explore the action space and learn a sequence of randomized policies by relying on a variant of online mirror descent.
- We prove for uniformly ergodic MDPs that the proposed algorithm achieves a regret bound of $\tilde{\mathcal{O}}((T|\mathcal{A}||\mathcal{S}|)^{2/3} \sqrt{\tau})$, both in expectation and with high probability. Here, $|\mathcal{S}|$ is the number of states, $|\mathcal{A}|$ is the number of actions, and τ is the mixing time of the MDP. To our knowledge, the proposed scheme is the first to achieve the above high probability regret bound for uniformly ergodic MDPs with bandit feedback.

2 RELATED WORK

A number of exact and approximate solutions to the problem of learning optimal policies in MDPs have been proposed in the literature. These include value iteration, policy iteration, and policy gradient techniques (see, e.g. [4, 5, 23] for a detailed discussion). Diverging from these methods, a linear programming (LP) approach has recently gained attention for MDPs with stationary loss functions [1, 25] as well as A-MDPs [28, 24, 6]. Our proposed algorithm relies on the LP formulation for computing an occupancy measure

corresponding to the optimal policy. However, we propose a new loss estimator to deal with the bandit feedback setting.

Learning with A-MDPs can be categorized into two cases: episodic MDPs and uniformly ergodic MDPs where the latter is considered more general and challenging [21].

2.1 EPISODIC A-MDPS

In this setting, the loss may change from episode to episode. Jaksch et al. [11] introduced the UCRL-2 algorithm for MDPs with stochastic rewards under the setting of unknown transition functions and full information feedback. UCRL-2 keeps track of confidence sets that, with high probability, contain the true transition function and shrink over time. For episodes of length L , they showed a regret of $\tilde{\mathcal{O}}(L|\mathcal{S}|\sqrt{|\mathcal{A}|T})$ compared to the optimal policy and provided a min-max lower bound, which can be achieved for a sufficiently large T (see the recent work of Azar et al. [3]). Non-adversarial episodic MDPs are studied in [7, 27] and recently [9] establish the state-of-the-art lowerbound for this setting. For the adversarial episodic MDP model, Neu et al. [19] proposed the *follow-the-perturbed-optimistic-policy* algorithm which relies on the follow-the-perturbed-leader method [13] and provides a regret of $\tilde{\mathcal{O}}(L|\mathcal{S}||\mathcal{A}|\sqrt{T})$. Recently, Rosenberg and Mansour [24] extended the results of Jaksch et al. [11] to MDPs with convex loss functions by employing online convex optimization and provided an algorithm achieving a regret of $\tilde{\mathcal{O}}(L|\mathcal{S}|\sqrt{|\mathcal{A}|T})$. However, these results still rely on the availability of full information feedback.

Under bandit feedback, Zimin and Neu [28] introduced the O-REPS algorithm, which employs the online mirror descent algorithm to achieve an expected regret of $\tilde{\mathcal{O}}(\sqrt{LT|\mathcal{A}||\mathcal{S}|})$. Similarly, Dick et al. [8] provided a regret of $\tilde{\mathcal{O}}(\sqrt{LT|\mathcal{S}||\mathcal{A}|})$ but for a computationally improved algorithm. Jin et al. [12] explored the use of an implicit exploration in the loss estimation for the case of *episodic A-MDPs* with unknown transition functions and obtained a high-probability regret of $\tilde{\mathcal{O}}(L|\mathcal{S}|\sqrt{T|\mathcal{A}|})$. In this work, we employ a different loss estimator with implicit exploration property for the general case of *ergodic A-MDPs* where the

transition functions are known. We also analyze the regret in terms of both the expectation and high-probability bounds.

2.2 UNIFORMLY ERGODIC A-MDPS

Learning with uniformly ergodic MDPs in the adversarial setting is considerably more complicated than the episodic case. In the former, which is also the focus of this paper, the loss may change in every round, as opposed to the episodic setting where the loss in each episode is fixed. This difficulty renders the task of deriving learning algorithms with sublinear regrets more challenging.

Early work by Even-Dar et al. [10] on A-MDPS is on uniformly ergodic MDPs with known transition function and full information feedback. Considering each action to be an expert, the MDP-E algorithm in [10] employs the weighted majority method [17] in each state and achieves a regret of $\mathcal{O}(\tau^2 \sqrt{T \log |\mathcal{A}|})$. Yu et al. [26] improved the computational efficiency utilizing the so-called *follow-the-perturbed-leader* method [13] with a regret of $\mathcal{O}(|\mathcal{A}|^2 |\mathcal{S}| \tau T^{3/4})$. More recently, Cardoso et al. [6] provided a regret of $\mathcal{O}(\sqrt{\tau} (\log |\mathcal{A}| |\mathcal{S}|) T \log T)$ for the same problem of uniformly ergodic MDPs with known transition function and full information feedback. Compared to this line of work, we consider the bandit setting and propose a new algorithm achieving high probability regret bounds.

For uniformly ergodic MDPs and under bandit feedback, Neu et al. [20] developed an algorithm that obtains $\tilde{\mathcal{O}}(\tau T^{2/3} |\mathcal{A}|^{1/3} \beta^{-1/3})$ regret in expectation; here, β is a lower bound on all stationary distributions, typically satisfying $\beta^{-1} = \mathcal{O}(|\mathcal{S}|)$. The expected regret of this algorithm is further improved with respect to T to $\tilde{\mathcal{O}}(\sqrt{\tau^3 T} |\mathcal{A}| \beta^{-1})$ [21]; however, this result is semi-asymptotic, i.e., it holds only for very long time horizons satisfying $T = \tilde{\Omega}(|\mathcal{A}| \tau \beta^{-3})$. We note that all of the aforementioned episodic schemes, as well as the uniformly ergodic results in [10, 20, 21], are guaranteed to achieve certain *expected regrets*. In contrast, in this paper, we propose a new scheme that achieves regret bound of $\tilde{\mathcal{O}}((T |\mathcal{A}| |\mathcal{S}|)^{\frac{2}{3}} \sqrt{\tau})$ not only in expectation but also with high probability.

Table 1 summarizes the differences between our approach with the most relevant existing methods.

3 BACKGROUND AND PRELIMINARY

We briefly overview the definitions of Markov decision processes, uniformly ergodicity assumption, random and expected regret, and the occupancy measures.

3.1 MARKOV DECISION PROCESS

Definition 1. A Markov decision process (MDP) is a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \ell)$, where \mathcal{S} is a finite discrete state space, \mathcal{A} is a finite discrete action space, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is a probabilistic transition function, and $\ell : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is a loss function.

A sequence of actions by the agent generates a state trajectory over an MDP in the following manner. The agent starts in an initial state $s_{init} \in \mathcal{S}$. At time t , the agent is in state s_t . Upon taking an action $\mathbf{a}_t \in \mathcal{A}$, the environment stochastically selects a next state s_{t+1} according to $\mathcal{P}(\cdot | s_t, \mathbf{a}_t)$ and the agent receives a loss $\ell(s_t, \mathbf{a}_t)$.

A policy π is a mapping from the history of states and actions, i.e., $\mathbf{h}_t = (s_1, \mathbf{a}_1, s_2, \mathbf{a}_2, \dots, s_t, \mathbf{a}_t)$, to the action space. The goal is to find a policy that minimizes the expected cumulative loss

$$\mathbb{E} \left[\sum_{t=1}^T \ell(s_t, \mathbf{a}_t) \right]$$

over a horizon of length T . Since MDPs admit optimal stochastic stationary policies [23], it suffices to search for an optimal policy in the family of stochastic stationary policies, i.e., $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$. Throughout this paper, we use $\pi(a|s)$ to denote the probability of selecting action a in state s and $\mathcal{P}^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ to denote the transition probabilities between the states under policy π . Let $\nu_t \in \Delta_{\mathcal{S}}$ denote the state distribution at time t , where $\Delta_{\mathcal{S}}$ is a simplex in $\mathbb{R}^{|\mathcal{S}|}$. Further, let ν_{t+1}^π denote the state distribution at time t under policy π conforming to

$$\nu_{t+1}^\pi = \nu_t \mathcal{P}^\pi.$$

We use $\bar{\nu}^\pi$ to indicate the stationary distribution of policy π satisfying

$$\bar{\nu}^\pi = \bar{\nu}^\pi \mathcal{P}^\pi.$$

We assume that the MDP satisfies the so-called uniformly ergodic property, stated formally next.

Definition 2. Let $\nu_1, \nu_2 \in \Delta_{\mathcal{S}}$ denote a pair of state distributions. A uniformly ergodic MDP is an MDP for which there exists $\tau \geq 1$ such that, for any policy π it holds that

$$\|\nu_1 \mathcal{P}^\pi - \nu_2 \mathcal{P}^\pi\|_1 \leq e^{-\frac{1}{\tau}} \|\nu_1 - \nu_2\|_1.$$

Intuitively, for every policy over a uniformly ergodic MDP, the convergence rate of state distributions to a unique stationary distribution is exponentially fast. Similar to [21], we assume that the dynamics of the AMDP, i.e., the probabilistic transition function \mathcal{P} is known. The important and more practical setting of dealing with unknown dynamic is left to future work.

In an A-MDP, the loss function, denoted by ℓ_t , varies over time. We assume that $\ell_t \in [0, 1]$ to remove the dependence of the analysis on the magnitude of the loss. We indicate the long-time average loss of a fixed policy π with respect to a fixed loss function ℓ by

$$\xi_\ell^\pi = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\mathbf{s}_t^\pi, \mathbf{a}_t^\pi).$$

One can show that for a fixed loss function ℓ , an optimal policy minimizing ξ_ℓ^π can be computed by solving a linear program (see (2) in [6]). We use ξ_t^π and $\xi_t^{\pi_t}$ to respectively denote the long-time average loss of policy π and policy π_t , with respect to the loss function ℓ_t .

We use the formulation of the regret minimization where the regret is defined with respect to the best policy in hindsight. Define

$$\mathcal{L}_T = \mathbb{E} \left[\sum_{t=1}^T \ell_t(\mathbf{s}_t, \mathbf{a}_t) \right]$$

as the expected cumulative loss of the learner, where the expectation is with respect to the randomness of the trajectories over the MDP. Also, define

$$\mathcal{L}_T(\pi) = \mathbb{E} \left[\sum_{t=1}^T \ell_t(\mathbf{s}_t, \mathbf{a}_t) | \pi \right]$$

as the expected cumulative loss under a fixed policy π . Then, the main goal of this paper is to achieve a low random regret

$$\mathcal{R}_T := \max_{\pi} \mathcal{L}_T - \mathcal{L}_T(\pi), \quad (1)$$

on which we seek a high-probability bound. Note that the randomness of \mathcal{R}_T is injected by the learner and is different from the randomness in the objective function \mathcal{L}_T . The expected regret is defined as

$$\bar{\mathcal{R}}_T = \mathbb{E}[\mathcal{R}_T] = \mathbb{E} \left[\max_{\pi} \mathcal{L}_T - \mathcal{L}_T(\pi) \right].$$

3.2 OCCUPANCY MEASURE

A (stochastic stationary) policy can equivalently be represented using occupancy measures. The occupancy measure of a policy is defined as the distribution induced by the execution of that policy over the state-action pairs, asymptotically, i.e.,

$$\rho^\pi(s, a) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \Pr(\mathbf{s}_t = s, \mathbf{a}_t = a | \pi).$$

A key property of occupancy measures is that the inflow into a state should be balanced by the outflow from that state. Formally, for every state $s \in \mathcal{S}$,

$$\sum_{a \in \mathcal{A}} \rho^\pi(s, a) = \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} \mathcal{P}(s | s', a') \rho^\pi(s', a').$$

Additionally, the occupancy measures are normalized over the entire state-action space, i.e.,

$$\sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \rho(s, a) = 1.$$

There is a one-to-one mapping between stochastic stationary policies and occupancy measures. The policy π^ρ corresponding to an occupancy measure ρ can be computed according to

$$\pi^\rho(a | s) = \frac{\rho(s, a)}{\sum_{a' \in \mathcal{A}} \rho(s, a')}, \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}. \quad (2)$$

The mapping between policies and occupancy measures allows one to reformulate a search over the policy space as a search over the occupancy measure space.

4 ONLINE RELATIVE-ENTROPY POLICY SEARCH WITH IMPLICIT EXPLORATION

Our proposed solution to the regret minimization problem is outlined in Algorithm 1. The algorithm builds upon the relative-entropy policy search of Peters et al. [22] and its online variant O-REPS by [28] while employing a novel loss estimator. In particular, Algorithm 1, which we refer to as *online relative-entropy policy search with implicit exploration (O-REPS-IX)*, employs an online mirror descent (OMD) optimization approach in conjunction to an optimistically biased loss estimator. We study these components of O-REPS-IX next.

4.1 ESTIMATING THE LOSS

In the bandit feedback setting, the agent observes only the loss corresponding to its current state and action and consequently has to construct an estimate of the overall loss function. In the episodic setting, the O-REPS algorithm [28] uses the following unbiased estimator to estimate the unobserved part of the loss in each episode and achieves the optimal expected regret:

$$\hat{\ell}_t^{\text{O-REPS}}(x, a) = \frac{\ell_t(s, a) \mathbb{I}\{(s, a) \in \mathbf{e}_t\}}{\rho_t(s, a)}, \quad (3)$$

where \mathbf{e}_t is the t^{th} episode. A similar unbiased loss estimator is further adopted in the uniformly ergodic setting [20, 21]. However, different from O-REPS [28] there is no notion of episode in the setting of uniformly ergodic A-MDPs that we consider here. One implication of this non-episodic aspect is that we can no longer limit the position of the agent to a specific layer at a time step. Furthermore, the agent does not restart from the initial state after every episode. Another implication of having no episode is that the agent updates

Algorithm 1 Online Relative-Entropy Policy Search with Implicit Exploration (O-REPS-IX)

- 1: **Input:** An MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P})$, time horizon T , estimation window N , exploration parameter γ , learning rate η
- 2: **Output:** Occupancy measures $\rho_1, \rho_2, \dots, \rho_T$ at each time step
- 3: Initialize the occupancy measures for the first $2N - 1$ time steps

$$\rho_t(s, a) = \frac{1}{|\mathcal{S}||\mathcal{A}|}, \quad \forall t \in [2N - 1], \forall (s, a) \in \mathcal{S} \times \mathcal{A}$$

- 4: Set the initial state $s_1 = s_{init}$ and the initial history $h_1 = (s_1)$
- 5: **for** $t = 1, \dots, T$ **do**
- 6: Compute the current policy at the current state

$$\pi_t(a|s_t) = \frac{\rho_t(s_t, a)}{\sum_{a' \in \mathcal{A}} \rho_t(s_t, a')}$$

- 7: Draw an action a_t randomly from the distribution $\pi_t(a|s_t)$
- 8: Observe the loss value $\ell_t(s_t, a_t)$ and the next state s_{t+1}
- 9: Update the history,

$$\mathbf{h}_t \leftarrow \mathbf{h}_{t-1} + (a_t, \ell_t(s_t, a_t), s_{t+1})$$

- 10: **if** $t \geq N$ **then**
- 11: Compute $\nu_{t|t-N}$
- 12: Construct the loss estimator $\hat{\ell}_t$ from the current history h_t

$$\hat{\ell}_t(s, a) = \frac{\ell_t(s, a)}{\nu_{t|t-N}(s)\pi_t(a|s) + \gamma} \mathbb{I}\{s_t = s, a_t = a\}$$

- 13: Compute the optimal value function

$$\hat{v}_t = \arg \min_v \ln Z_t(v)$$

- 14: Compute the solution ρ_{t+N} to the projection step (11)

$$\rho_{t+N}(s, a) = \frac{\rho_{t+N-1}(s, a)e^{\delta(s, a|\hat{v}_t, \hat{\ell}_t)}}{Z_t(\hat{v}_t)}$$

- 15: **end if**
 - 16: **end for**
-

the policy every time step given that the loss function may change in every time step. Furthermore, a major limitation of

the estimators in [28, 20, 21] is that they suffer from a high degree of variance. Consequently, although these schemes achieve the optimal expected regret, it can be shown, using the arguments presented in Remark 1 in Section 11.5 and Exercise 11.5 in [16] and Section 3 in [2], that the random regret will be linear with a nonzero probability due to the high variance of the loss estimator.

Given above challenges, designing a loss estimator for the uniformly ergodic setting requires further considerations.

Let

$$\nu_{t|t-N}(s) = \Pr(\mathbf{s}_t = s | \mathbf{h}_{t-N})$$

denote the probability of being in state s at time t given the history at time $t - N$, $t \geq N + 1$. Also, let $\vec{\nu}_{t|t-N}$ represent a vector of dimension $|\mathcal{S}|$, concatenating $\nu_{t|t-N}(s)$ for all $s \in \mathcal{S}$, and $e_{\mathbf{h}_{t-N}}$ represent a unit vector in $\mathbb{R}^{|\mathcal{S}|}$ such that

$$e_{\mathbf{h}_{t-N}}(s) = \begin{cases} 1 & \text{if } \mathbf{s}_{t-N} = s \\ 0 & \text{otherwise.} \end{cases}$$

Then, one can obtain $\vec{\nu}_{t|t-N}$ according to:

$$\vec{\nu}_{t|t-N} = e_{\mathbf{h}_{t-N}} \mathcal{P}^{\mathbf{a}_{t-N}} \mathcal{P}^{\pi_{t-N+1}} \mathcal{P}^{\pi_{t-N+2}} \dots \mathcal{P}^{\pi_{t-1}},$$

where $\mathcal{P}^{\mathbf{a}_{t-N}}$ denotes the transition probabilities between the states upon taking action \mathbf{a}_{t-N} .

We propose the following loss estimator:

$$\hat{\ell}_t(s, a) := \frac{\ell_t(s, a)}{\nu_{t|t-N}(s)\pi_t(a|s) + \gamma} \mathbb{I}\{s_t = s, a_t = a\}, \quad (4)$$

which exploits the bandit observation of the loss in the current time step; here, $\gamma > 0$ is an exploration parameter that induces exploration whose value will be determined in Theorem 1. Intuitively, $\nu_{t|t-N}(s)\pi_t(a|s)$ can be thought of as some form of occupancy measure. Looking at (3), one might be tempted to set $N = 1$ in (4), justified by the fact that in the episodic setting

$$\mathbb{E}[\mathbb{I}\{(s, a) \in \mathbf{e}_t | \mathbf{e}_1, \dots, \mathbf{e}_{t-1}\}] = \rho_t(s, a).$$

However, this does not hold in the uniformly ergodic settings as

$$\mathbb{E}[\mathbb{I}\{s_t = s, a_t = a\} | t - 1] = \mathcal{P}(s, a | s_{t-1}, a_{t-1}) \neq \nu_t(s)\pi_t(a|s). \quad (5)$$

Due to this discrepancy, which is discussed originally by Neu et al. [20], we will consider a sufficiently large N . Larger values of N , which we henceforth refer to as the estimation window, results in a better estimate of the loss function. Intuitively, delaying the policy update leads to lower variance of the random regret, enabling a high-probability analysis since the estimation window N helps to robustify the estimator against the learner's randomness.

Now given a sufficiently large N and an exploration parameter $\gamma > 0$, by taking the expectation of (4) we observe

$$\begin{aligned}\mathbb{E}[\hat{\ell}_t(s, a)|t - N] &= \frac{\ell_t(s, a)\mathbb{E}[\mathbb{I}\{\mathbf{s}_t = s, \mathbf{a}_t = a\}|t - N]}{\boldsymbol{\nu}_{t|t-N}(s)\boldsymbol{\pi}_t(a|s) + \gamma} \\ &= \frac{\ell_t(s, a)\boldsymbol{\nu}_{t|t-N}(s)\boldsymbol{\pi}_t(a|s)}{\boldsymbol{\nu}_{t|t-N}(s)\boldsymbol{\pi}_t(a|s) + \gamma} \leq \ell_t(s, a).\end{aligned}\quad (6)$$

That is, our proposed loss estimator in (4) is *optimistically biased*. This aspect is inline with the optimism principle in online learning [15]. Intuitively, since for a given state and action pair (s, a) the proposed estimator underestimates the true loss, as the agent interacts with the environment the estimated loss of any sub-optimal action will eventually become larger than that of the optimal ones. Furthermore, as we will show in Section 5, the proposed estimator in (4) achieves a variance reducing effect compared to typical unbiased estimators, e.g., (3), thereby enabling a high probability sublinear regret for O-REPS-IX.

4.2 POLICY UPDATE VIA OMD

Given an occupancy measure ρ , we use the unnormalized negative entropy as the potential function of OMD, i.e.,

$$R(\rho) = \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \rho(s, a) \log \rho(s, a) - \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \rho(s, a). \quad (7)$$

Given this potential function, the Bregman divergence $D(\rho||\rho')$ between two occupancy measures ρ and ρ' is the unnormalized Kullback–Leibler divergence:

$$\begin{aligned}D(\rho||\rho') &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \rho(s, a) \log \frac{\rho(s, a)}{\rho'(s, a)} \\ &\quad - \sum_{s \in \mathcal{S}, a \in \mathcal{A}} (\rho(s, a) - \rho'(s, a)).\end{aligned}\quad (8)$$

In t^{th} time step, the agent selects an occupancy measure $\boldsymbol{\rho}_{t+N}$ which minimizes a linear combination of the estimated loss $\hat{\ell}_t$ and the divergence from the previous occupancy measure $\boldsymbol{\rho}_{t+N-1}$. Formally, the agent finds a solution to the constrained optimization problem

$$\boldsymbol{\rho}_{t+N} = \arg \min_{\rho \in \Delta(\mathcal{M})} \left\{ \eta \langle \rho, \hat{\ell}_t \rangle + D(\rho||\boldsymbol{\rho}_{t+N-1}) \right\}, \quad (9)$$

where $\Delta(\mathcal{M})$ denotes the set of all occupancy measures over an MDP \mathcal{M} that satisfy the inflow-outflow balancing and normalization constraints, and $\langle \cdot, \cdot \rangle$ is the inner product in the space of $\mathcal{S} \times \mathcal{A}$. Note that as we discussed $\boldsymbol{\rho}_{t+N}$ is entirely determined by the history \mathbf{h}_t . Hence, in the first $2N - 1$ rounds of learning (see step 3 of Algorithm 1) we initialize the occupancy measure (and consequently the policy $\boldsymbol{\pi}_t$) uniformly, i.e., $\boldsymbol{\rho}_t(s, a) = 1/|\mathcal{S}||\mathcal{A}|$.

Similar to the standard mirror descent techniques, the constrained optimization in (9) can be efficiently solved through a two-step procedure. First, an unconstrained version of the problem is solved, i.e., we find

$$\tilde{\boldsymbol{\rho}}_{t+N} = \arg \min_{\rho} \left\{ \eta \langle \rho, \hat{\ell}_t \rangle + D(\rho||\boldsymbol{\rho}_{t+N-1}) \right\}, \quad (10)$$

which admits a closed form solution

$$\tilde{\boldsymbol{\rho}}_{t+N}(s, a) = \boldsymbol{\rho}_{t+N-1}(s, a) e^{-\eta \hat{\ell}_t(s, a)}.$$

Then, $\tilde{\boldsymbol{\rho}}_{t+N}(s, a)$ is projected to the constraint set $\Delta(\mathcal{M})$, i.e., we find

$$\boldsymbol{\rho}_{t+N} = \arg \min_{\rho \in \Delta(\mathcal{M})} \left\{ D(\rho||\tilde{\boldsymbol{\rho}}_{t+N-1}) \right\}. \quad (11)$$

By enforcing constraints of inflow-outflow balancing and normalization on the occupancy measures, the following constrained optimization yields the solution to the projection step:

$$\begin{aligned}\min_{\rho} \quad & D(\rho||\tilde{\boldsymbol{\rho}}_{t+N-1}) \\ \text{s.t.} \quad & \sum_{a \in \mathcal{A}} \rho(s, a) = \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} \mathcal{P}(s|s', a') \rho(s', a') \quad \forall s \in \mathcal{S}, \\ & \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \rho(s, a) = 1.\end{aligned}\quad (12)$$

We show in Proposition 1 that the above optimization problem using ideas from [28] can equivalently be written as an unconstrained convex optimization problem.

Proposition 1. *Let $v : \mathcal{S} \rightarrow \mathbb{R}$ denote a value function for each state and $\ell : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ denote a loss function. Define*

$$\delta(s, a|v, \ell) = v(s) - \eta \ell(s, a) - \sum_{s' \in \mathcal{S}} v(s') \mathcal{P}(s'|s, a),$$

a function capturing the notion of Bellman error for the value function v . Furthermore, for $t > 1$, define a partition function

$$\mathbf{Z}_t(v) = \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \boldsymbol{\rho}_t(s, a) e^{\delta(s, a|v, \hat{\ell}_t)}.$$

The optimal value function (corresponding to the dual problem) is the solution to an unconstrained optimization problem

$$\hat{v}_t = \arg \min_v \ln \mathbf{Z}_t(v).$$

With these definitions in place, the solution to the projection step (11) is

$$\boldsymbol{\rho}_{t+N}(s, a) = \frac{\boldsymbol{\rho}_{t+N-1}(s, a) e^{\delta(s, a|\hat{v}_t, \hat{\ell}_t)}}{\mathbf{Z}_t(\hat{v}_t)}.$$

Proof. The projection step of online mirror descent has a closed-form solution for episodic A-MDPs as shown in [28]; it is readily derived by differentiating the Lagrangian with respect to $\rho(s, a)$ and setting the gradient to zero. The solution follows by using the second constraint and solving the dual maximization problem.

The projection step of online mirror descent for uniformly ergodic MDPs in (12) is different from the one for episodic MDPs only in constraints on the occupancy measures. That is, for an episodic MDP, occupancy measures are normalized across each layer while for a uniformly ergodic MDP the normalization is over the entire state space. Therefore, with comparable arguments to those stated in [28], we can derive the presented closed form solution. ■

5 REGRET ANALYSIS

We now present the theoretical analysis of O-REPS-IX for A-MDPs satisfying the uniform ergodicity assumption. The detailed proofs are deferred to the Appendix.

5.1 BOUND ON THE RANDOM REGRET

We start by stating our main theoretical result, establishing a high-probability bound on the random regret \mathcal{R}_T .

Theorem 1. *Let*

$$\begin{aligned}\eta &= (T|\mathcal{S}||\mathcal{A}|)^{-2/3} \sqrt{\log(|\mathcal{S}||\mathcal{A}|)}, \\ \gamma &= (T|\mathcal{S}||\mathcal{A}|)^{-1/3} \sqrt{\tau \log T \log \frac{1}{\delta}}, \\ N &= 1 + \lceil \tau \log T \rceil.\end{aligned}\quad (13)$$

Then, for any $\delta \in (0, 1)$, with probability at least $1 - 4\delta$, it holds that the random regret of Algorithm 1 satisfies

$$\begin{aligned}\mathcal{R}_T &\leq C (T|\mathcal{A}||\mathcal{S}|)^{\frac{2}{3}} \sqrt{\tau \log(|\mathcal{S}||\mathcal{A}|) \log T \log \frac{1}{\delta}} \\ &\quad + C' \tau \log T,\end{aligned}$$

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

Remark 1. *Theorem 1 establishes that Algorithm 1 achieves $\tilde{\mathcal{O}}(T^{\frac{2}{3}})$ regret bound with high probability. The pioneering algorithm in [21] achieves an optimal regret of $\tilde{\mathcal{O}}(\sqrt{T})$ only in expectation for sufficiently large T , and under an extra assumption compared to our result (see Assumption A2 there) that bounds all stationary distributions from zero. Algorithm 1 enjoys a high-probability regret bound which is a much stronger type of guarantee. Additionally, the proposed algorithm has a better dependence on τ compared to the algorithm in [21] (that is, $\mathcal{O}(\sqrt{\tau})$ vs. $\mathcal{O}(\tau\sqrt{\tau})$). In our current analysis we employ a relatively large γ to ensure that the proposed estimator is uniformly bounded with probability one. However, this restriction results in a sub-optimal*

regret bound $\tilde{\mathcal{O}}(T^{\frac{2}{3}})$. Hence, it remains an open problem to see whether the high probability regret of Algorithm 1 can be improved to $\tilde{\mathcal{O}}(\sqrt{T})$, i.e., the optimal regret (with respect to T).

Proof of Theorem 1. Recall step 3 of O-REPS-IX (see Algorithm 1). Given that by assumption $\ell_t(s, a) \leq 1$, the first $2N - 1$ terms in the random regret \mathcal{R}_T can be bounded by $2N$. Hence, we study the regret of O-REPS-IX starting $t = 2N$. To obtain the regret bound, we first decompose the random regret:

$$\begin{aligned}\mathcal{R}_T &= \mathcal{O}(N) + \underbrace{\max_{\rho \in \Delta(\mathcal{M})} \sum_{t=2N}^T \xi_t^\pi - \mathbb{E} \left[\sum_{t=2N}^T \ell_t(s_t^\pi, a_t^\pi) \right]}_{\text{I}} \\ &\quad + \underbrace{\sum_{t=2N}^T \xi_t^{\pi_t} - \sum_{t=2N}^T \xi_t^\pi}_{\text{II}} \\ &\quad + \underbrace{\mathbb{E} \left[\sum_{t=2N}^T \ell_t(s_t^{\pi_t}, a_t^{\pi_t}) \right] - \sum_{t=2N}^T \xi_t^{\pi_t}}_{\text{III}}.\end{aligned}$$

Next, we decompose the second term according to

$$\begin{aligned}\sum_{t=2N}^T \xi_t^{\pi_t} - \sum_{t=2N}^T \xi_t^\pi &= \sum_{t=2N}^T \langle \rho_t - \rho, \ell_t \rangle \\ &= \underbrace{\sum_{t=2N}^T \langle \rho_t, \ell_t - \hat{\ell}_t \rangle}_{\text{II-I}} + \underbrace{\sum_{t=2N}^T \langle \rho, \hat{\ell}_t - \ell_t \rangle}_{\text{II-II}} \\ &\quad + \underbrace{\sum_{t=2N}^T \langle \rho_t - \rho, \hat{\ell}_t \rangle}_{\text{II-III}},\end{aligned}$$

by recalling definition of ξ^π and ξ^{π_t} :

$$\begin{aligned}\xi^\pi &= \lim_{T' \rightarrow \infty} \sum_{t=2N}^{T'} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr(s_{t'}^\pi = s, a_{t'}^\pi = a) \ell_t(s, a) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \rho(s, a) \ell_t(s, a),\end{aligned}\quad (14)$$

and

$$\begin{aligned}\xi^{\pi_t} &= \lim_{T' \rightarrow \infty} \sum_{t=2N}^{T'} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr(s_{t'}^{\pi_t} = s, a_{t'}^{\pi_t} = a) \ell_t(s, a) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \rho_t(s, a) \ell_t(s, a).\end{aligned}\quad (15)$$

Now, we bound each term individually. For a fixed policy π over a finite time horizon T , term I measures the difference between the expected reward starting from the initial

distribution ν_1 and the expected reward starting from the stationary distribution $\bar{\nu}^\pi$. Note that this term is deterministic and we can bound this difference by a factor of τ due to the uniform ergodicity assumption which ensures fast mixing time (Lemma 1 and Appendix B).

Lemma 1 (Bounding term I [10]). *For any $T \geq 1$ and any policy π , it holds that*

$$\sum_{t=2N}^T \xi_t^\pi - \mathbb{E} \left[\sum_{t=2N}^T \ell_t(s_t^\pi, a_t^\pi) \right] \leq 2(1 + \tau). \quad (16)$$

Term II – which is random – is studied in Lemma 2 (see Appendix C for the unabridged statement). To analyze term II-I, we show in Lemma 9 the fact that the evolving policies from the mirror descent algorithm do not change much between consecutive time steps as long as N satisfies the condition given in Theorem 1. We further need to show that $\hat{\ell}_t$ is a good estimate for ℓ_t , which we prove in two steps. First, we show that ℓ_t is close to $\mathbb{E}[\hat{\ell}_t]$ by a factor directly proportional to γ . Note that with $\gamma = 0$, $\hat{\ell}_t$ becomes an unbiased estimator of ℓ_t . Second, $\hat{\ell}_t$ concentrates around its mean $\mathbb{E}[\hat{\ell}_t]$. We bound term II-II in Lemma 10 by relying on closeness of $\hat{\ell}_t$ to ℓ_t . In particular, the optimistic bias of $\hat{\ell}_t$ allows us to use a concentration result based on the Cramer-Chernoff method (Lemma 4). Analysis of term II-III also has two components, where one of them depends on the regret of the mirror descent algorithm and the other one depends on the fact that the iterates of the OMD do not change too rapidly as long as η and γ satisfy the conditions of Theorem 1.

Lemma 2 (Bounding term II). *With η , γ , and N given in (13), it holds, with probability exceeding $1 - 4\delta$, that*

$$\begin{aligned} & \sum_{t=2N}^T \xi^{\pi_t} - \sum_{t=2N}^T \xi^\pi \\ &= \mathcal{O} \left((T|\mathcal{A}||\mathcal{S}|)^{\frac{2}{3}} \sqrt{\tau \log(|\mathcal{S}||\mathcal{A}|) \log T \log \frac{1}{\delta}} \right). \end{aligned}$$

Term III, similar to term II, is random and captures the difference between the expected reward actually obtained by the agent and the expected reward obtained by the agent had it been in the stationary distribution $\bar{\nu}^{\pi_t}$ at each time step t . By using the fact that the evolving policies from the mirror descent algorithm do not change much between consecutive time steps, in Lemma 3 we establish an upper bound on term III (see Appendix D for the proof).

Lemma 3 (Bounding term III). *For any $\delta \in (0, 1)$, with probability at least $1 - \delta$, it holds that*

$$\begin{aligned} \mathbb{E} \left[\sum_{t=2N}^T \ell_t(s_t, a_t) \right] - \sum_{t=2N}^T \xi_t^{\pi_t} &\leq 2(1 + \tau) \\ &+ 2\eta(1 + \tau) \left(\frac{\log \frac{1}{\delta}}{2\gamma} + T|\mathcal{S}||\mathcal{A}| \right). \end{aligned}$$

Lastly, by adding the bounds from each term and properly selecting the values of γ and η , we obtain the desired result stated in Theorem 1. The details of the regret bound is provided in Appendix E. ■

5.2 BOUND ON THE EXPECTED REGRET

An upper bound on the expected regret of Algorithm 1 can also be obtained by integrating the tail of the high-probability regret bound provided in Theorem 1. The result is formalized in the following theorem.

Theorem 2. *With η , γ , and N given in (13), the expected regret of Algorithm 1 satisfies*

$$\begin{aligned} \bar{\mathcal{R}}_T &\leq C (T|\mathcal{A}||\mathcal{S}|)^{\frac{2}{3}} \sqrt{\tau \log(|\mathcal{S}||\mathcal{A}|) \log T \log \frac{1}{\delta}} \\ &+ C' \tau \log T, \end{aligned} \quad (17)$$

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

Proof of Theorem 2. First, we relate the expected regret to the random regret on which we have derived a high-probability bound. In particular, defining $\mathcal{R}_T^+ := \max\{\mathcal{R}_T, 0\}$, we have

$$\begin{aligned} \bar{\mathcal{R}}_T &\leq \mathbb{E} [\mathcal{R}_T^+] \stackrel{(a)}{=} \int_0^\infty \Pr(\mathcal{R}_T^+ \geq u) du \\ &= \int_0^\infty \Pr(\mathcal{R}_T \geq u) du, \end{aligned}$$

where (a) is due to the fact that for a non-negative random variable X it holds that $\mathbb{E}[X] = \int_0^\infty \Pr(X > x) dx$. We can evaluate the integral using the high-probability bound of Theorem 1 and a change of variables. Assume $\tau > 1$ and let $B = (T|\mathcal{A}||\mathcal{S}|)^{\frac{2}{3}} \sqrt{\tau \log(|\mathcal{S}||\mathcal{A}|) \log T}$. Then, it is apparent that with probability at most 4δ , the following lower bound holds on the random regret:

$$\mathcal{R}_T \geq CB \log \frac{1}{\delta} + C' \tau \log T, \quad (18)$$

for some $C, C' > 0$. Note that the second term is deterministic. Next, let $u = CB \log \left(\frac{1}{\delta}\right)$ and thus $\delta = \exp(-u/CB)$. If $\delta \rightarrow 0^+$, then $u \rightarrow \infty$, while if $\delta \rightarrow \left(\frac{1}{4}\right)^-$, then $u \rightarrow (CB \log 4)^+$. Then,

$$\begin{aligned} \bar{\mathcal{R}}_T &\stackrel{(b)}{\leq} C' \tau \log T + \int_{CB \log 4}^\infty \Pr \left(\mathcal{R}_T \geq CB \log \frac{L}{\delta} \right) du \\ &\stackrel{(c)}{\leq} C' \tau \log T + \int_{CB \log 4}^\infty 4 \exp \left(-\frac{u}{CB} \right) du, \end{aligned}$$

where (b) is due to the nonnegativity of the integrand and (c) corresponds to the simplified high-probability bound in (18). Lastly, a simple integration yields the desired result. ■

6 CONCLUSION AND FUTURE WORK

We considered the general class of uniformly ergodic A-MDPs whose loss functions may change arbitrarily over time. By relying on an optimistically biased loss estimator and online linear optimization techniques, we proposed O-REPS-IX that finds a policy achieving sublinear regret bounds both with high probability and in expectation. In particular, the algorithm achieves the regret of $\tilde{O}((T|\mathcal{A}||\mathcal{S}|)^{\frac{2}{3}}\sqrt{\tau})$ with respect to the best stationary policy in hindsight. The proposed scheme is the first algorithm achieving a high probability sublinear regret bound in the setting of learning with uniformly ergodic A-MDPs and bandit feedback.

As a future research direction, it is important to establish whether the high-probability regret of O-REPS-IX can be improved to $\tilde{O}(\sqrt{T})$, i.e., the optimal regret. Furthermore, we would like to explore the potential of using the proposed algorithm for learning in safety-critical scenarios. In these scenarios, the high-probability guarantees of O-REPS-IX can be employed to provide desirable safety assurances. Finally, it is valuable to extend our results to the class of risk-aware MDPs.

Author Contributions

The first two authors contributed equally.

Acknowledgements

This work was supported partially by AFRL grant FA9550-19-1-0169, DARPA grant D19AP00004, and NSF ECCS grant 1809327.

References

- [1] ABBASI-YADKORI, Y., BARTLETT, P. L., CHEN, X., AND MALEK, A. Large-scale Markov decision problems via the linear programming dual. *arXiv preprint arXiv:1901.01992* (2019).
- [2] ABERNETHY, J., AND RAKHLIN, A. Beating the adaptive bandit with high probability. In *2009 Information Theory and Applications Workshop* (2009), IEEE, pp. 280–289.
- [3] AZAR, M. G., OSBAND, I., AND MUNOS, R. Minimax regret bounds for reinforcement learning. In *Proceedings of International Conference on Machine Learning (ICML)* (2017), vol. 70, JMLR. org, pp. 263–272.
- [4] BERTSEKAS, D. P., BERTSEKAS, D. P., BERTSEKAS, D. P., AND BERTSEKAS, D. P. *Dynamic programming and optimal control*, vol. 1. Athena scientific Belmont, MA, 1995.
- [5] BERTSEKAS, D. P., AND TSITSIKLIS, J. N. *Neurodynamic programming*, vol. 5. Athena Scientific Belmont, MA, 1996.
- [6] CARDOSO, A. R., WANG, H., AND XU, H. Large scale Markov decision processes with changing rewards. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)* (2019), pp. 2337–2347.
- [7] DANN, C., LI, L., WEI, W., AND BRUNSKILL, E. Policy certificates: Towards accountable reinforcement learning. In *International Conference on Machine Learning* (2019), PMLR, pp. 1507–1516.
- [8] DICK, T., GYORGY, A., AND SZEPESVARI, C. Online learning in Markov decision processes with changing cost sequences. In *Proceedings of International Conference on Machine Learning (ICML)* (2014), pp. 512–520.
- [9] DOMINGUES, O. D., MÉNARD, P., KAUFMANN, E., AND VALKO, M. Episodic reinforcement learning in finite mdps: Minimax lower bounds revisited. In *Algorithmic Learning Theory* (2021), PMLR, pp. 578–598.
- [10] EVEN-DAR, E., KAKADE, S. M., AND MANSOUR, Y. Online Markov decision processes. *Mathematics of Operations Research* 34, 3 (2009), 726–736.
- [11] JAKSCH, T., ORTNER, R., AND AUER, P. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research (JMLR)* 11, Apr (2010), 1563–1600.
- [12] JIN, C., JIN, T., LUO, H., SRA, S., AND YU, T. Learning adversarial MDPs with bandit feedback and unknown transition. *arXiv preprint arXiv:1912.01192* (2019).
- [13] KALAI, A., AND VEMPALA, S. Efficient algorithms for online decision problems. In *Learning Theory and Kernel Machines*. Springer, 2003, pp. 26–40.
- [14] KOCÁK, T., NEU, G., VALKO, M., AND MUNOS, R. Efficient learning by implicit exploration in bandit problems with side observations. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)* (2014), pp. 613–621.
- [15] LAI, T. L., AND ROBBINS, H. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics* 6, 1 (1985), 4–22.
- [16] LATTIMORE, T., AND SZEPESVÁRI, C. Bandit algorithms. *preprint* (2018).

- [17] LITTLESTONE, N., AND WARMUTH, M. K. The weighted majority algorithm. *Information and computation* 108, 2 (1994), 212–261.
- [18] NEU, G. Explore no more: Improved high-probability regret bounds for non-stochastic bandits. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)* (2015), pp. 3168–3176.
- [19] NEU, G., GYORGY, A., AND SZEPESVÁRI, C. The adversarial stochastic shortest path problem with unknown transition probabilities. In *Proceedings of Artificial Intelligence and Statistics (AISTATS)* (2012), pp. 805–813.
- [20] NEU, G., GYÖRGY, A., SZEPESVÁRI, C., AND ANTOS, A. Online markov decision processes under bandit feedback. In *Proceedings of the 23rd International Conference on Neural Information Processing Systems-Volume 2* (2010), pp. 1804–1812.
- [21] NEU, G., GYORGY, A., SZEPESVARI, C., AND ANTOS, A. Online Markov decision processes under bandit feedback. *IEEE Transactions on Automatic Control (TAC)* 3, 59 (2014), 676–691.
- [22] PETERS, J., MÜLLING, K., AND ALTÜN, Y. Relative entropy policy search. In *Proceedings of AAAI Conference on Artificial Intelligence* (2010), pp. 1607–1612.
- [23] PUTERMAN, M. L. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 2014.
- [24] ROSENBERG, A., AND MANSOUR, Y. Online convex optimization in adversarialMarkov decision processes. In *Proceedings of International Conference on Machine Learning (ICML)* (2019), pp. 5478–5486.
- [25] WANG, M. Primal-dual π learning: Sample complexity and sublinear run time for ergodic Markov decision problems. *arXiv preprint arXiv:1710.06100* (2017).
- [26] YU, J. Y., MANNOR, S., AND SHIMKIN, N. Markov decision processes with arbitrary reward processes. *Mathematics of Operations Research* 34, 3 (2009), 737–757.
- [27] ZANETTE, A., AND BRUNSKILL, E. Tighter problem-dependent regret bounds in reinforcement learning without domain knowledge using value function bounds. In *International Conference on Machine Learning* (2019), PMLR, pp. 7304–7312.
- [28] ZIMIN, A., AND NEU, G. Online learning in episodic Markovian decision processes by relative entropy policy search. In *Proceedings of Advances in neural information processing systems (NeurIPS)* (2013), pp. 1583–1591.