
Online Episodic Convex Reinforcement Learning

Anonymous Authors¹

Abstract

We study online learning in episodic finite-horizon Markov decision processes (MDPs) with convex objective functions, known as the concave utility reinforcement learning (CURL) problem. This setting generalizes RL from linear to convex losses on the state-action distribution induced by the agent’s policy. The non-linearity of CURL invalidates classical Bellman equations and requires new algorithmic approaches. We introduce the first algorithm achieving near-optimal regret bounds for online CURL without any prior knowledge on the transition function. To achieve this, we use an online mirror descent algorithm with varying constraint sets and a carefully designed exploration bonus. We then address for the first time a bandit version of CURL, where the only feedback is the value of the objective function on the state-action distribution induced by the agent’s policy. We achieve a sub-linear regret bound for this more challenging problem by adapting techniques from bandit convex optimization to the MDP setting.

1. Introduction

Reinforcement learning (RL) studies the problem where an agent interacts with an environment over time, adhering to a probabilistic policy that maps states to actions and aiming to minimize the cumulative expected losses. The environment’s dynamics are represented by a Markov decision process (MDP), assumed here to be episodic, with episodes of length N , a finite state space \mathcal{X} , a finite action space \mathcal{A} , and a sequence of probability transition kernels $p := (p_n)_{n \in [N]}$, such that for each $(x, a) \in \mathcal{X} \times \mathcal{A}$, $p_n(\cdot | x, a) \in \Delta_{\mathcal{X}}$, the simplex over the state space. Formally, the RL problem involves finding a policy π that, under a transition kernel p , induces a state-action distribution sequence $\mu^{\pi, p} \in (\Delta_{\mathcal{X} \times \mathcal{A}})^N$ mini-

mizing the inner product with a loss vector $\ell := (\ell_n)_{n \in [N]}$, with $\ell_n \in \mathbb{R}^{\mathcal{X} \times \mathcal{A}}$, i.e.: $\min_{\pi \in (\Delta_{\mathcal{A}})^{\mathcal{X} \times N}} \langle \ell, \mu^{\pi, p} \rangle$. A large body of literature is devoted to solving the RL problem efficiently and with theoretical guarantees in many challenging environments (Bertsekas, 2019; Sutton & Barto, 2018).

However, numerous practical problems entail more intricate objectives, such as those encountered within the Concave Utility Reinforcement Learning (CURL) framework (Hazan et al., 2019; Zahavy et al., 2021) (also known as convex RL). The CURL problem consists in minimizing a convex function (or maximizing a concave function) on the state-action distributions induced by an agent’s policy:

$$\min_{\pi \in (\Delta_{\mathcal{A}})^{\mathcal{X} \times N}} F(\mu^{\pi, p}). \quad (1)$$

In addition to RL, other examples of machine learning problems that can be written as CURL are pure exploration (Hazan et al., 2019; Mutti et al., 2021; 2022), where $F(\mu^{\pi, p}) = \langle \mu^{\pi, p}, \log(\mu^{\pi, p}) \rangle$; imitation learning (Ghasemipour et al., 2020; Lavington et al., 2022) and apprenticeship learning (Zahavy et al., 2019; Abbeel & Ng, 2004), where $F(\mu^{\pi, p}) = D_g(\mu^{\pi, p}, \mu^*)$, with D_g representing a Bregman divergence induced by a function g and μ^* being a behavior to be imitated; certain instances of mean-field control (Bensoussan et al., 2013), where $F(\mu^{\pi, p}) = \langle \ell(\mu^{\pi, p}), \mu^{\pi, p} \rangle$; mean-field games with potential rewards (Lavigne & Pfeiffer, 2023); risk-averse RL (García et al., 2015; Pan et al., 2019; Greenberg et al., 2022), among others. The non-linearity of CURL alters the additive structure inherent in standard RL, invalidating the classical Bellman equations. Consequently, dynamic programming approaches become infeasible, necessitating the development of novel methodologies.

A natural extension of CURL is the online scenario, wherein a sequence of policies $(\pi^t)_{t \in [T]}$ is computed over T episodes, aimed at minimizing a cumulative loss $L_T := \sum_{t=1}^T F^t(\mu^{\pi^t, p})$, where the objective F^t can change arbitrarily (known as the adversarial scenario (Even-Dar et al., 2009)), and the MDP probability kernel p is unknown. Most existing approaches to CURL fail to address the challenges of the online setting (adversarial losses and unknown dynamics). The few methods that attempt to tackle this problem rely on strong assumptions about the probability transition kernel (Moreno et al., 2024), which can be overly restrictive in real-world scenarios. To overcome this, we need

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

an approach capable of optimizing the objective function while simultaneously learning the environment, effectively balancing the exploration-exploitation dilemma.

Contribution 1. In the full-information feedback setting, where the objective function F^t is fully revealed to the learner at the end of episode t , we propose the first method achieving sub-linear regret for online CURL with adversarial losses and unknown transition kernels, without relying on additional model assumptions. Our algorithm uses an Online Mirror Descent (OMD) variant incorporating well-designed exploration bonuses into the sub-gradient of the objective function to handle the exploration-exploitation trade-off. It achieves a regret of $\tilde{O}(\sqrt{T})$, matching the state-of-the-art (SoTA) in more restricted settings (Moreno et al., 2024), while obtaining a closed-form solution.

Contribution 2. We extend our approach to incorporate bandit feedback on the objective function. We first consider the RL case where $F^t(\mu) := \langle \ell^t, \mu \rangle$. Bandit feedback in this setting means that the agent only observes the loss function in the state-action pairs they visit during each episode, i.e. $(\ell_n^t(x_n^t, a_n^t))_{n \in [N]}$ where $(x_n^t, a_n^t)_{n \in [N]}$ is the agent’s trajectory. We obtain the optimal regret of $\tilde{O}(\sqrt{T})$ in this setting. We then address for the first time the general CURL problem under more strict bandit feedback. In this setting, the learner only has access to the value of the objective function evaluated on the state-action distribution sequence induced by the agent’s policy, i.e., $F^t(\mu^{\pi^t, p})$. We propose two algorithms for this setting and show that they achieve sub-linear regret. One algorithm requires that the MDP is known, while the other, under an additional assumption on the structure of the MDP, operates in the setting where the MDP is estimated progressively from observed trajectories. We rely on gradient estimation techniques from the bandit convex optimization literature, even as the peculiar structure of our constraint set and uncertainty regarding the true transition kernel present some unique challenges.

1.1. Related Work

Offline CURL. An extensive line of work focus on the offline version of CURL (Problem (1)), where the objective function is known and fixed. The methodologies proposed by (Zhang et al., 2020; 2021; Barakat et al., 2023) rely on policy gradient techniques, requiring the estimation of F ’s gradient concerning the policy π , a task often complex. Taking a different approach, Zahavy et al. (2021) cast the CURL problem as a min-max game using Fenchel duality, demonstrating that conventional RL algorithms can be tailored to fit the CURL framework. Recently, Geist et al. (2022) established that CURL is a specific instance of mean-field games. Moreover, Moreno et al. (2024) undertake a convexification of Problem (1) and propose a mirror descent algorithm with a non-standard Bregman divergence. Mutti et al. (2023a;b)

study the gap between evaluating agent performance over infinite realizations versus finite trials and question the classic CURL formulation in Eq. (1). To align with prior work, we adopt the classic CURL formulation.

Online CURL. To the best of our knowledge, Greedy MD-CURL from (Moreno et al., 2024) is the only regret minimization algorithm designed for online CURL. However, it only achieves sublinear regret when the system dynamics follow the form $x_{n+1} = g_n(x_n, a_n, \varepsilon_n)$, where g_n is a known deterministic function, and ε_n is an external noise with an unknown distribution independent of (x_n, a_n) , which significantly limits its applicability, as we empirically show in Sec. 5. This assumption simplifies the problem, as the algorithm only needs to learn the noise distribution, which can be done independently of the policy, eliminating the need for exploration. In contrast, our approach does not assume any specific form for the dynamics, which introduces the challenge of developing a policy that minimizes total loss while simultaneously enabling sufficient exploration to improve estimates of the transition kernels. The technical novelty we introduce to overcome this challenge are well-designed exploration bonuses detailed in Sec. 3.

RL approaches. *Model-optimistic methods* construct a set of plausible MDPs by forming confidence bounds around the empirical transition kernels, then select the policy that maximizes the expected reward in the best feasible MDP. A key example of this approach is UCRL (Upper Confidence RL) methods (Jaksch et al., 2008; Zimin & Neu, 2013; Rosenberg & Mansour, 2019b; Jin et al., 2020). While these methods offer strong theoretical guarantees, they are often difficult to implement due to the complexity of optimizing over all plausible MDPs. While we believe these approaches could be generalized to CURL, their computational complexity has led us to propose an alternative method. *Value-optimistic methods* are value-based approaches that compute optimistic value functions, rather than optimistic models, using dynamic programming. An example is UCB-VI (Azar et al., 2017). However, these methods are limited to stochastic losses. *Policy-optimization (PO) methods* directly optimize the policy and are widely used in RL due to their faster performance and closed-form solutions. Recently, Luo et al. (2021) achieved SoTA regret for PO methods with adversarial losses and bandit feedback by introducing *dilated bonuses*, which satisfy a *dilated* Bellman equation and are added to the Q -function. However, their approach cannot be applied here due to CURL’s non-linearity (the expectation of the trajectory appears inside the objective function) which invalidates the Bellman’s equations.

We achieve our results by computing local bonuses and adding them to the (sub-)gradient of the objective function in each OMD instance as exploration bonuses. This is more computationally efficient than model-optimistic ap-

Table 1. Comparisons of SoTA finite-horizon tabular MDPs methods. MD stands for Mirror Descent, KL for Kullback-Leibler divergence and Γ is defined in Eq. (4). MD + (\cdot) indicates the regularization added to the MD iteration. MD on π indicates a policy optimization approach in which MD iterations are performed on policies instead of state-action distributions (occupancy-measures).

	Algorithm	Optimal regret in T	CURL	Closed-form	Exploration	No model assumption	Adversarial Losses	Bandit feedback
(Jin et al., 2020)	MD + KL	✓	✗	✗	UCRL	✓	✓	✓
(Moreno et al., 2024)	MD + Γ	✓	✓	✓	None	✗	✓	✗
(ours)	MD + Γ	✓	✓	✓	Bonus	✓	✓	✓
(Luo et al., 2021)	MD on π	✓	✗	✓	Bonus	✓	✓	✓

proaches and addresses the exploration issues in previous online CURL methods. We believe our analysis is of independent interest, as it also offers a new way to study RL approaches over occupancy measures, while providing closed-form solutions. See Table 1 for comparisons.

2. Problem Formulation

2.1. Setting

For a finite set \mathcal{S} , $|\mathcal{S}|$ represents its cardinality, while $\Delta_{\mathcal{S}}$ denotes the $|\mathcal{S}|$ -dimensional simplex. For all $d \in \mathbb{N}$ we denote $[d] := \{1, \dots, d\}$. We let $\|\cdot\|_1$ be the L_1 norm, and for all $v := (v_n)_{n \in [N]}$, such that $v_n \in \mathbb{R}^{\mathcal{X} \times \mathcal{A}}$ we define $\|v\|_{\infty, 1} := \sup_{1 \leq n \leq N} \|v_n\|_1$. We denote by $\|\cdot\|_{1, \infty}$ its dual. Let $\Pi := (\Delta_{\mathcal{A}})^{\mathcal{X} \times N}$ denote the set of policies. We consider an episodic MDP as introduced in Sec. 1. We assume that the initial state-action pair of an agent is sampled from a fixed distribution $\mu_0 \in \Delta_{\mathcal{X} \times \mathcal{A}}$ at the beginning of each episode. At time step $n \in [N]$, the agent moves to a state $x_n \sim p_n(\cdot | x_{n-1}, a_{n-1})$, and chooses an action $a_n \sim \pi_n(\cdot | x_n)$ by means of a policy $\pi_n : \mathcal{X} \rightarrow \Delta_{\mathcal{A}}$. When the agent follows a policy $\pi := (\pi_n)_{n \in [N]}$ for an episode in an environment described by the MDP with a transition kernel p , this induces a state-action distribution, which we denote by $\mu^{\pi, p} := (\mu_n^{\pi, p})_{n \in [N]}$, that can be calculated recursively for all $(n, x, a) \in [N] \times \mathcal{X} \times \mathcal{A}$, by

$$\begin{aligned} \mu_0^{\pi, p}(x, a) &= \mu_0(x, a) \\ \mu_n^{\pi, p}(x, a) &= \sum_{(x', a')} \mu_{n-1}^{\pi, p}(x', a') p_n(x | x', a') \pi_n(a | x). \end{aligned} \quad (2)$$

We define the set of all state-action distribution sequences satisfying the dynamics of the MDP as

$$\begin{aligned} \mathcal{M}_{\mu_0}^p := \left\{ \mu \in (\Delta_{\mathcal{X} \times \mathcal{A}})^N \mid \sum_{a' \in \mathcal{A}} \mu_n(x', a') = \right. \\ \left. \sum_{x \in \mathcal{X}, a \in \mathcal{A}} p_n(x' | x, a) \mu_{n-1}(x, a), \forall x' \in \mathcal{X}, \forall n \in [N] \right\}. \end{aligned} \quad (3)$$

For any $\mu \in \mathcal{M}_{\mu_0}^p$, there is a strategy π such that $\mu^{\pi, p} = \mu$. It suffices to take $\pi_n(a | x) \propto \mu_n(x, a)$ when the normal-

ization factor is non-zero, and arbitrarily defined otherwise. Let $\mathcal{M}_{\mu_0}^{p, *}$ be the subset of $\mathcal{M}_{\mu_0}^p$ where the corresponding policies π satisfy $\pi_n(a | x) \neq 0$ for all (x, a) . For any two probability transition kernels p, q , we define $\Gamma : \mathcal{M}_{\mu_0}^p \times \mathcal{M}_{\mu_0}^{q, *} \rightarrow \mathbb{R}$ such that, for all $\mu, \mu' \in \mathcal{M}_{\mu_0}^p \times \mathcal{M}_{\mu_0}^{q, *}$ with policies π, π' ,

$$\Gamma(\mu, \mu') := \sum_{n=1}^N \mathbb{E}_{(x, a) \sim \mu_n(\cdot)} \left[\log \left(\frac{\pi_n(a | x)}{\pi'_n(a | x)} \right) \right]. \quad (4)$$

In the online extension of CURL, the objective function for episode t is denoted as $F^t := \sum_{n=1}^N f_n^t$, where $f_n^t : \Delta_{\mathcal{X} \times \mathcal{A}} \rightarrow \mathbb{R}$ is convex and L -Lipschitz with respect to the $\|\cdot\|_1$ norm (hence F^t is L_F -Lipschitz with respect to the norm $\|\cdot\|_{\infty, 1}$ with $L_F := LN$). The objective function F^t is unknown to the learner in the start of episode t . In this paper, we examine three types of objective function feedback: *Full-information*: In this case, F^t is fully disclosed to the learner at the end of episode t , and is treated in Sec. 3.2. *Bandit in RL*: Here, $F^t(\mu) := \langle \ell^t, \mu \rangle$, and the learner observes the loss function only for the state-action pairs visited, i.e., $(\ell_n^t(x_n^t, a_n^t))_{n \in [N]}$, which is covered in Sec. 4.1. *Bandit in CURL*: In this scenario, the learner only has access to the objective function evaluated on the state-action distribution sequence induced by the agent's policy, i.e., $F^t(\mu^{\pi^t, p})$, and is treated in Sec. 4.2.

The learner's goal is to compute a sequence of strategies $(\pi^t)_{t \in [T]}$, where T represents the total number of episodes, that minimizes their total loss $L_T := \sum_{t=1}^T F^t(\mu^{\pi^t, p})$. The learner's performance is evaluated by comparing it to any policy $\pi \in (\Delta_{\mathcal{A}})^{\mathcal{X} \times N}$ using the static regret:

$$R_T(\pi) := \sum_{t=1}^T F^t(\mu^{\pi^t, p}) - F^t(\mu^{\pi, p}). \quad (5)$$

We assume the probability transition kernel p is unknown to the learner. Hence, to minimize its total loss, the learner must optimize the objective function while simultaneously learn the environment dynamics, facing an exploration-exploitation dilemma. The interaction between the learner and the environment proceeds in episodes. At each episode t , the learner selects a policy π^t , sends it to the agent, and observes its trajectory $o^t := (x_0^t, a_0^t, \dots, x_N^t, a_N^t)$. The learner uses this observation to compute an estimation of

the probability transition kernel \hat{p}^{t+1} . At the end of episode t , the learner receives one of the three feedbacks described above for the objective function F^t , and then calculates the policy for the next episode, π^{t+1} , based on π^t, \hat{p}^{t+1} , and the feedback on F^t .

2.2. Preliminary Results

The results in this section are either known or extensions of existing results needed for the analysis.

Since the probability transition kernel is unknown, we propose an online mirror descent (OMD) instance that optimizes over the state-action distributions induced by the estimated MDP as if it was the true model. This approach differs from the model-optimistic methods for RL discussed in Sec. 1.1 where each iteration is performed over the union of all state-action distribution sets induced by MDPs within a confidence set around the estimated model, which results in a computationally expensive optimization problem per iteration. Lemma 2.1 presents an auxiliary result concerning the quality of the state-action distribution sequence $(\mu^t)_{t \in [T]}$ when μ^t is the solution of Eq. (6), an OMD instance on the set of state-action distributions induced by a transition kernel q^t . It extends the upper bound result from (Moreno et al., 2024) for OMD with smoothly varying constraint sets to any sequence of bounded vectors $(z^t)_{t \in [T]}$ and any sequence of smoothly varying transitions $(q^t)_{t \in [T]}$.

Lemma 2.1. *Let $(q^t)_{t \in [T]}$ be a sequence of probability transition kernels and $(z^t)_{t \in [T]}$ a sequence of vectors in $\mathbb{R}^{N \times |\mathcal{X}| \times |\mathcal{A}|}$, such that $\max_{t \in [T]} \|z^t\|_{1, \infty} \leq \zeta$. Initialize $\pi_n^1(a|x) := 1/|\mathcal{A}|$. For $t \in [T]$, let $\tilde{\pi}^t := \frac{t}{t+1}\pi^t + \frac{1}{t+1}|\mathcal{A}|^{-1}$ be a smoothed version of the policy and compute iteratively*

$$\mu^{t+1} \in \arg \min_{\mu \in \mathcal{M}_{\mu_0}^{q^{t+1}}} \tau \langle z^t, \mu \rangle + \Gamma(\mu, \mu^{\tilde{\pi}^t, q^t}). \quad (6)$$

Then, there is a $\tau > 0$ such that, for any sequence $(\nu^t)_{t \in [T]}$, with $\nu^t := \nu^{\pi, q^t}$ for a common policy π ,

$$\sum_{t=1}^T \langle z^t, \mu^t - \nu^t \rangle \leq O(\zeta N \sqrt{V_T |\mathcal{X}| \log(|\mathcal{A}|) T \log(T)}),$$

$$\text{where } V_T \geq 1 + \max_{(n, x, a)} \sum_{t=1}^{T-1} \|q_n^t(\cdot|x, a) - q_n^{t+1}(\cdot|x, a)\|_1.$$

This lemma is proved in App. C. It is known (Moreno et al., 2024) that for the divergence Γ defined in Eq. (4), Eq. (6) has a closed-form solution for the policy (see App. A.2).

Learning the model. Since the learner does not know the probability transition kernel, it must estimate p from the agents' trajectories. Below we present the empirical way for estimating the transition and a well-known result (Lem. 2.2) on its quality using Hoeffding's inequality. Let $N_n^t(x, a) = \sum_{s=1}^{t-1} \mathbb{1}_{\{x_n^s = x, a_n^s = a\}}$, $M_n^t(x'|x, a) =$

$\sum_{s=1}^{t-1} \mathbb{1}_{\{x_{n+1}^s = x', x_n^s = x, a_n^s = a\}}$. The learner's estimate for the transition kernel at the end of episode $t-1$, to be used in episode t , is as follows

$$\hat{p}_{n+1}^t(x'|x, a) := \frac{M_n^t(x'|x, a)}{\max\{1, N_n^t(x, a)\}}. \quad (7)$$

Lemma 2.2 (Lem. 17 of Jaksch et al., 2008). *For any $0 < \delta < 1$, with a probability of at least $1 - \delta$,*

$$\|p_n(\cdot|x, a) - \hat{p}_n^t(\cdot|x, a)\|_1 \leq \sqrt{\frac{2|\mathcal{X}| \log\left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta}\right)}{\max\{1, N_{n-1}^t(x, a)\}}}$$

holds simultaneously for all $(t, n, x, a) \in [T] \times [N] \times \mathcal{X} \times \mathcal{A}$.

These results suffice for analyzing CURL with full-information feedback (Sec. 3). For bandit feedback, more refined tools are needed. In bandit RL, we need Bernstein's inequality to bound the L_1 distance (Lem. D.2). In bandit CURL, we also need a bound on the Kullback-Leibler (KL) divergence (Lem. E.3), which requires the Laplace (add-one) estimator (Eq. (51)), as the KL of the empirical one can be unbounded.

3. Exploration Bonus in CURL

We now present our novel approach for online CURL with adversarial losses and unknown dynamics.

3.1. Limitations of previous approaches

The performance measure of a learner playing a sequence of strategies $(\pi^t)_{t \in [T]}$ is given by the static regret defined in Eq. (5). Using the estimate of the probability transition kernel \hat{p}^t computed by the learner, the static regret can be further decomposed as follows

$$\begin{aligned} R_T(\pi) &= \sum_{t=1}^T F^t(\mu^{\pi^t, p}) - F^t(\mu^{\pi^t, \hat{p}^t}) \\ &\quad + \sum_{t=1}^T F^t(\mu^{\pi^t, \hat{p}^t}) - F^t(\mu^{\pi, p}) \\ &\leq \underbrace{\sum_{t=1}^T \langle \nabla F^t(\mu^{\pi^t, p}), \mu^{\pi^t, p} - \mu^{\pi^t, \hat{p}^t} \rangle}_{R_T^{\text{MDP}}} \\ &\quad + \underbrace{\sum_{t=1}^T \langle \nabla F^t(\mu^{\pi^t, \hat{p}^t}), \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, p} \rangle}_{R_T^{\text{policy}}}, \end{aligned} \quad (8)$$

where the inequality comes from the convexity of F^t . Let $\xi_n^t(x, a) := \|p_n(\cdot|x, a) - \hat{p}_n^t(\cdot|x, a)\|_1$. The term R_T^{MDP} , accounts for the error in estimating the MDP, and satisfies

$R_T^{\text{MDP}} = \tilde{O}(\sqrt{T})$ with high probability. This is a classic result (see Neu et al., 2012). We first show that

$$R_T^{\text{MDP}} \leq L \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} \mu_i^{\pi^t, p}(x, a) \xi_{i+1}^t(x, a). \quad (9)$$

Then, using Lem. 2.2 and that $N_n^t(x, a)$ increases with the empirical version of the state-action distribution $\mu_n^{\pi^t, p}(x, a)$ we achieve the final bound (see App. B.2). The second term, R_T^{policy} , depends on the algorithm used to derive the policies. As mentioned in Sec. 1.1, model-optimistic approaches could be adapted to CURL, but they are computationally expensive. To achieve low complexity, we explore potential problems that might arise from the absence of explicit exploration. We decompose this regret term as follows:

$$\begin{aligned} R_T^{\text{policy}} &= \underbrace{\sum_{t=1}^T \langle \nabla F^t(\mu^{\pi^t, \hat{p}^t}), \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle}_{R_T^{\text{policy/MD}}} \\ &+ \underbrace{\sum_{t=1}^T \langle \nabla F^t(\mu^{\pi^t, \hat{p}^t}), \mu^{\pi, \hat{p}^t} - \mu^{\pi, p} \rangle}_{R_T^{\text{policy/MDP}}}. \end{aligned}$$

Assume the learner computes its policy sequence $(\pi^t)_{t \in [T]}$ by solving Eq. (6) with $q^{t+1} := \hat{p}^{t+1}$ and $z^t := \nabla F^t(\mu^{\pi^t, \hat{p}^t})$. Hence, from Lem. 2.1, $R_T^{\text{policy/MD}} = \tilde{O}(\sqrt{T})$ (Lemmas A.3 and A.4 in the Appendix demonstrate that $\sum_{t=1}^T \|\hat{p}^{t+1}(\cdot|x, a) - \hat{p}^t(\cdot|x, a)\|_1 \leq e \log(T)$. By hypothesis, $\|\nabla F^t(\mu^{\pi^t, \hat{p}^t})\|_{1, \infty} \leq L_F$. Hence, we meet all the assumptions from Lem. 2.1). But the term $R_T^{\text{policy/MDP}}$ poses a challenge. It can be decomposed as R_T^{MDP} in Eq. (9). However, the state-action distribution multiplying $\xi_{i+1}^t(x, a)$ would either be $\mu_i^{\pi, p}(x, a)$ or $\mu_i^{\pi, \hat{p}^{t+1}}(x, a)$, and neither is related to $N_i^t(x, a)$. Consequently, we do not have the same convergence effect as R_T^{MDP} . In fact, this term can become prohibitively large. Without exploration, previous work using similar analysis (Moreno et al., 2024) only achieved optimal regret under strong model assumptions, limiting its applicability in realistic scenarios.

3.2. CURL with full-information feedback

We outline our idea to overcome previous limitations presented in Subsec. 3.1. Let $b^t := (b_n^t)_{n \in [N]}$ be a sequence of vectors, to be properly defined later, such that $b_n^t \in \mathbb{R}^{\mathcal{X} \times \mathcal{A}}$. We assume that π^t is the policy inducing μ^t computed as in Eq. (6) with $q^t := \hat{p}^t$, but instead of considering $z^t = \nabla F^t(\mu^{\pi^t, \hat{p}^t})$ as the (sub-)gradient of MD to be used in episode $t + 1$, we let $z^t := \nabla F^t(\mu^{\pi^t, \hat{p}^t}) - b^t$, i.e.,

$$\mu^{t+1} := \arg \min_{\mu \in \mathcal{M}_{\mu_0}^{\hat{p}^{t+1}}} \left\{ \tau \langle \nabla F^t(\mu^{\pi^t, \hat{p}^t}) - b^t, \mu \rangle + \Gamma(\mu, \tilde{\mu}^t) \right\}.$$

If we assume that b^t is such that, for all $t \in [T]$ and for some $\zeta > 0$, $\|\nabla F^t(\mu^{\pi^t, \hat{p}^t}) - b^t\|_{1, \infty} \leq \zeta$, then by Lem. 2.1 and by adding and subtracting the bonus vector, we would have that R_T^{policy} is bounded by

$$\tilde{O}(\sqrt{T}) + \sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle + R_T^{\text{policy/MDP}}. \quad (10)$$

Let $C_\delta := \sqrt{2|\mathcal{X}| \log(|\mathcal{X}||\mathcal{A}|NT/\delta)}$, and for all $n \in \{0, [N]\}$, $(x, a) \in \mathcal{X} \times \mathcal{A}$, let

$$b_n^t(x, a) := L(N - n) \frac{C_\delta}{\sqrt{\max\{1, N_n^t(x, a)\}}}. \quad (11)$$

Note that $\|b_n^t\|_\infty \leq LNC_\delta$, ensuring that the hypothesis of Lem. 2.1 remains valid for this sequence. Decomposing $R_T^{\text{policy/MDP}}$ as we do for R_T^{MDP} in Eq. (9), and then applying Lem. 2.2, we get that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, $R_T^{\text{policy/MDP}}$ is bounded by

$$\begin{aligned} LC_\delta \sum_{t=1}^T \sum_{n=0}^{N-1} (N - n) \sum_{x,a} \frac{\mu_n^{\pi, \hat{p}^t}(x, a)}{\sqrt{\max\{1, N_n^t(x, a)\}}} \\ = \sum_{t=1}^T \langle \mu^{\pi, \hat{p}^t}, b^t \rangle. \end{aligned} \quad (12)$$

By replacing Eq. (12) in Eq. (10), the additive property in the decomposition allows us to cancel out the problematic regret term $R_T^{\text{policy/MDP}}$. As a result, we obtain that $R_T^{\text{policy}} \leq \tilde{O}(\sqrt{T}) + \sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} \rangle$. All that remains is to analyze the new term due to the added bonus, $\sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} \rangle$, which we do in Prop. 3.1.

Proposition 3.1. *Let $(b^t)_{t \in [T]}$ be the bonus vector in Eq. (11). For any $\delta' \in (0, 1)$, with probability $1 - 3\delta'$,*

$$\sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} \rangle = \tilde{O}(LN^3 |\mathcal{X}|^{3/2} \sqrt{|\mathcal{A}|T}).$$

With all the ingredients in place, we introduce our new method, *Bonus O-MD-CURL*, in Alg. 1. The main result is in Thm. 3.2 and its proof is in App. B.2. In terms of T and $|\mathcal{A}|$, our result matches the optimal one in RL from (Jin et al., 2020), but we have additional factors of N and $\sqrt{|\mathcal{X}|}$ that are due to using bonuses and dealing with convex RL.

Theorem 3.2. *Running Alg. 1 for online CURL with unknown transition kernel, full-information feedback, where $F^t := \sum_{n=1}^N f_n^t$ is convex and each f_n^t is L -Lipschitz under $\|\cdot\|_1$, ensures that, with probability at least $1 - 6\delta$ for any $\delta \in (0, 1)$, the optimal choice of τ achieves, for any $\pi \in \Pi$,*

$$R_T(\pi) = \tilde{O}(LN^3 |\mathcal{X}|^{3/2} \sqrt{|\mathcal{A}|T}).$$

4. Bandit Feedback

4.1. Bandit feedback with bonus in RL

We generalize Alg. 1 to handle the RL case with bandit feedback. Our aim is not to improve the existing algorithms

Algorithm 1 Bonus O-MD-CURL (Full-information)

- 1: **Input:** number of episodes T , initial policy $\pi^1 \in \Pi$, initial state-action distribution μ_0 and state-action distribution sequence $\mu^1 = \tilde{\mu}^1 = \mu^{\pi^1, \hat{p}^1}$ with $\hat{p}_n^1(\cdot|x, a) = 1/|\mathcal{X}|$, learning rate $\tau > 0$.
- 2: **Init.:** $\forall(n, x, a, x'), N_n^1(x, a) = M_n^1(x'|x, a) = 0$
- 3: **for** $t = 1, \dots, T$ **do**
- 4: agent starts at $(x_0^t, a_0^t) \sim \mu_0(\cdot)$
- 5: **for** $n = 1, \dots, N$ **do**
- 6: Env. draws new state $x_n^t \sim p_n(\cdot|x_{n-1}^t, a_{n-1}^t)$
- 7: Update counts

$$N_{n-1}^{t+1}(x_{n-1}^t, a_{n-1}^t) = N_{n-1}^t(x_{n-1}^t, a_{n-1}^t) + 1$$

$$M_{n-1}^{t+1}(x_n^t|x_{n-1}^t, a_{n-1}^t) = M_{n-1}^t(x_n^t|x_{n-1}^t, a_{n-1}^t) + 1$$
- 8: Agent chooses an action $a_n^t \sim \pi_n^t(\cdot|x_n^t)$
- 9: **end for**
- 10: Compute bonus sequence as in Eq. (11)
- 11: Observe objective function F^t
- 12: Compute μ^{π^t, \hat{p}^t} as in Eq. (2)
- 13: Update transition estimate as in Eq. (7)
- 14: Compute the π^{t+1} associated to the solution of Eq. 6 with $z^t := -\nabla F^t(\mu^{\pi^t, \hat{p}^t}) + b^t$ and $q^{t+1} = \hat{p}^{t+1}$
- 15: Compute $\tilde{\pi}^{t+1}$ (Lem. 2.1), and $\mu^{t+1} := \mu^{\tilde{\pi}^{t+1}, \hat{p}^{t+1}}$
- 16: **end for**

for bandit RL; rather, we show that our new methodology and analysis for CURL achieves comparable results to the SoTA in bandit RL. In this case, an adversary selects a sequence of loss functions $(\ell^t)_{t \in [T]}$, with $\ell^t := (\ell_n^t)_{n \in [N]}$, where $\ell_n^t : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$, and the objective function is given by $F^t(\mu) := \langle \ell^t, \mu \rangle = \sum_{n=1}^N \langle \ell_n^t, \mu_n \rangle$. Note that now the gradient of F^t with respect to μ is always equal to ℓ^t due to the linearity of the objective function. Bandit feedback in this setting implies that the learner observes the loss function only for the state-action pairs visited by the agent during each episode, i.e., $(\ell_n^t(x_n^t, a_n^t))_{n \in [N]}$ where $(x_n^t, a_n^t)_{n \in [N]}$ is the agent's trajectory.

We define Alg. 2 in App. D, a version of *Bonus O-MD-CURL* where for each OMD update we take $z^t := \hat{\ell}^t - b^t$, with $\hat{\ell}^t$ an importance-weighted estimator of ℓ^t defined in Eq. (40) and b^t the bonus vector defined in Eq. (11). Thm. 4.1 states that Alg. 2 achieves the regret bound of $\tilde{O}(\sqrt{T})$ known to be the optimal for RL with bandit feedback (Jin et al., 2020). For the proof and for an overview of approaches for bandit RL see App. D.

Theorem 4.1. *Playing Alg. 2 for RL with adversarial losses $(\ell^t)_{t \in [T]}$, unknown transition kernel, and bandit feedback, obtains with high probability for any policy $\pi \in \Pi$,*

$$R_T(\pi) = \tilde{O}(N^3 |\mathcal{X}|^{3/2} \sqrt{|\mathcal{A}| T} + N^{3/2} |\mathcal{X}|^{5/4} |\mathcal{A}| \sqrt{T}).$$

4.2. CURL with bandit feedback

Returning back to the CURL framework, we now assume that $F^t : \Delta_{\mathcal{X} \times \mathcal{A}} \rightarrow [0, N]$ can be any convex, L -Lipschitz function with respect to $\|\cdot\|_1$. In contrast to Sec. 3, we assume here that after executing a policy π^t we observe $F^t(\mu^{\pi^t, p})$ instead of $\nabla F^t(\mu^{\pi^t, p})$. We will consider both the case when the MDP is known in advance and when it needs (as in previous sections) to be estimated progressively from observed trajectories.

Main challenges. This problem can be broadly categorized as a bandit convex optimization (BCO) problem. This places us in a more challenging domain compared to the bandit feedback setting in the standard RL problem, where the gradient of the loss function is identical for any point in $(\Delta_{\mathcal{X} \times \mathcal{A}})^N$ and is easier to estimate. Moreover, as a BCO problem, the present setting still exhibits distinctive challenges. One being the peculiar nature of our decision set $\mathcal{M}_{\mu_0}^p$ and how it impedes the efficacy of some standard gradient estimation techniques as we explain below. Another issue arises when the MDP is not known as that induces uncertainty over the true set of permissible occupancy measures. This incomplete knowledge of the decision set is atypical in the BCO literature and introduces multiple sources of bias for any adopted method.

4.2.1. ENTROPIC REGULARIZATION METHOD

Our first approach is to extend our MD-based algorithm from Sec. 3, supposing still that the MDP is not known. Since the algorithm required knowledge of the gradient $\nabla F^t(\mu^{\pi^t, p})$, we propose to estimate it by querying F^t at a random perturbation of $\mu^{\pi^t, p}$, a standard approach in the convex bandit literature popularized by Flaxman et al. (2005). This method yields $T^{3/4}$ regret under convex and Lipschitz conditions, and is incapable of doing better (Hu et al., 2016). Although more advanced algorithms and analyses achieve \sqrt{T} regret (Hazan & Li, 2016; Bubeck et al., 2021; Fokkema et al., 2024), they are arguably less practical, more complicated, and have worse dimension dependence. For $d \in \mathbb{Z}_+$, we denote by \mathbb{B}^d and \mathbb{S}^d the unit ball and sphere respectively in \mathbb{R}^d , and by $\mathbb{1}_d \in \mathbb{R}^d$ the vector with all entries equal to one. Let $k : \mathcal{S} \rightarrow \mathbb{R}$ be a convex function, where $\mathcal{S} \subseteq \mathbb{R}^d$ is a convex set satisfying $\mathbb{B}^d \subseteq \mathcal{S}$. Fix some $\delta \in (0, 1)$. The approach of Flaxman et al. (2005) relies on the observation that $\frac{(1-\delta)d}{\delta} \mathbb{E}_{\mathbf{u} \in \mathbb{S}^d} [k((1-\delta)x + \delta \mathbf{u}) \mathbf{u}] \approx \nabla k(x)$. Hence, $\frac{(1-\delta)d}{\delta} k((1-\delta)x + \delta \mathbf{u}) \mathbf{u}$ (for some \mathbf{u} uniformly sampled from \mathbb{S}^d) can be used as a one-point stochastic surrogate for the gradient. Applying this idea to our problem presents several challenges. Mainly, $\mathcal{M}_{\mu_0}^p$ has an empty interior in $\mathbb{R}^{N|\mathcal{X}||\mathcal{A}|}$. This can be addressed, assuming for the moment that the kernel p is known, by defining a bijection $\Lambda_p : (\mathcal{M}_{\mu_0}^p)^- \rightarrow \mathcal{M}_{\mu_0}^p$, where $(\mathcal{M}_{\mu_0}^p)^- \subseteq \mathbb{R}^{N|\mathcal{X}|(|\mathcal{A}|-1)}$ is

a representation of the constraint set in a lower-dimensional space where it is possible for its interior to be non-empty, see App. E.1 for more details. Next, we need to specify a (hyper)sphere that is contained in $(\mathcal{M}_{\mu_0}^p)^-$, which would allow us to use the aforementioned spherical estimation technique while remaining inside the feasible set of occupancy measures. To guarantee the existence of such an object, we rely on the following assumption (discussed further below).

Assumption 4.2. There exists a value $\varepsilon > 0$ such that $p_n(x'|x, a) \geq \varepsilon$ for all $x, x' \in \mathcal{X}^2$, $a \in \mathcal{A}$, and $n \in [N]$.

Under this assumption, we show in App. E.2.1 that for $\kappa := \varepsilon / (|\mathcal{A}| - 1 + \sqrt{|\mathcal{A}| - 1})$, it holds that $\kappa \mathbb{1}_{N|\mathcal{X}|(|\mathcal{A}|-1)} + \kappa \mathbb{B}^{N|\mathcal{X}|(|\mathcal{A}|-1)} \subseteq (\mathcal{M}_{\mu_0}^p)^-$. For any $\mathbf{v} \in \mathbb{B}^{N|\mathcal{X}|(|\mathcal{A}|-1)}$, define $\zeta^{\mathbf{v}, p} := \Lambda_p(\kappa \mathbb{1}_{N|\mathcal{X}|(|\mathcal{A}|-1)} + \kappa \mathbf{v})$. Motivated by the preceding discussion, we use (a simple transformation of)

$$\frac{1-\delta}{\delta\kappa} N|\mathcal{X}|(|\mathcal{A}|-1) F^t((1-\delta)\mu^t + \delta\zeta^{\mathbf{u}^t, p}) \mathbf{u}^t$$

as a surrogate for $\nabla F^t(\mu^t)$, where \mathbf{u}^t is sampled uniformly from $\mathbb{S}^{N|\mathcal{X}|(|\mathcal{A}|-1)}$. What remains is to address the issue that the true kernel p is unknown. Similarly to the full information case, we compute an estimate \hat{p}^t at each round to be used in place of the true kernel, and we employ bonuses to explore. One difference is that we rely on a slightly altered transition kernel estimator (see App. E.2.2) to ensure that \hat{p}^t too satisfies the condition of Asm. 4.2. Another discrepancy to be accounted for in the analysis is that although we compute π^t relying on \hat{p}^t (in particular, π^t is the policy induced by $(1-\delta)\mu^t + \delta\zeta^{\mathbf{u}^t, \hat{p}^t} \in \mathcal{M}_{\mu_0}^{\hat{p}^t}$), we observe $F^t(\mu^{\pi^t, p})$, the evaluation of π^t in the true environment. This induces an extra source of bias in the gradient estimator. We summarize our approach in Alg. 3 in App. E.2.3, and prove the following result in App. E.2.5:

Theorem 4.3. Under Asm. 4.2, Alg. 3 with a suitable tuning of τ , δ , and $(\alpha_t)_{t \in [T]}$ satisfies for any policy $\pi \in \Pi$ that

$$\mathbb{E}[R_T(\pi)] = \tilde{O}(\sqrt{L(L+1)/\varepsilon} |\mathcal{X}|^{5/4} |\mathcal{A}|^{5/4} N^3 T^{3/4}).$$

The main shortcoming of this method is its reliance on the restrictive Asm. 4.2, which also affects the regret guarantee through its dependence on ε . This assumption is not necessary to guarantee that $(\mathcal{M}_{\mu_0}^p)^-$ has a non-empty interior; it suffices instead to assume that every state is reachable at every step, as we do later. Enforcing Asm. 4.2 only serves as a simple way to enable the construction of a sampling sphere with a certain radius. One can construct a different sampling sphere (or ellipsoid) without this assumption; nevertheless, the magnitude of the gradient estimator (which is featured in the current regret bound) would still scale with the reciprocal of the radius of that sphere, the permissible values for which depend on the structure of the MDP and can be arbitrarily small. It seems then that the current approach leads to an inevitable degradation of the bound subject to the structure of the MDP.

4.2.2. SELF-CONCORDANT REGULARIZATION METHOD

Fortunately, we can adopt a more principled approach via the use of self-concordant regularization, which is a common technique in bandit convex (and linear) optimization (see, e.g., Abernethy et al., 2008; Saha & Tewari, 2011; Hazan & Levy, 2014), and has been used for online learning in MDPs in different (linear) settings (Lee et al., 2020; Cohen et al., 2021; Van der Hoeven et al., 2023). We show in App. E.3 that $(\mathcal{M}_{\mu_0}^p)^-$ is a convex polytope specified as the intersection of $N|\mathcal{X}||\mathcal{A}|$ half-spaces. We define $\psi_{\text{lb}}: (\mathcal{M}_{\mu_0}^p)^- \rightarrow \mathbb{R}$ as the standard logarithmic barrier for $(\mathcal{M}_{\mu_0}^p)^-$ (see Nemirovski, 2004, Cor. 3.1.1) which is a ϑ -self-concordant barrier (see Nemirovski, 2004, Def. 3.1.1) for $(\mathcal{M}_{\mu_0}^p)^-$ with $\vartheta = N|\mathcal{X}||\mathcal{A}|$. The second approach we adopt here is to run mirror descent directly on $(\mathcal{M}_{\mu_0}^p)^-$ as the decision set and take ψ_{lb} as the regularizer in place of the entropic regularizer that induces Γ . Let ξ belong to the interior of $(\mathcal{M}_{\mu_0}^p)^-$, which we assume is not empty. Property I in (Nemirovski, 2004, Sec. 2.2) implies that $\xi + (\nabla^2 \psi_{\text{lb}}(\xi))^{-1/2} \mathbb{B}^{NS(A-1)} \subseteq (\mathcal{M}_{\mu_0}^p)^-$. Hence, we can construct an ellipsoid—entirely contained in $(\mathcal{M}_{\mu_0}^p)^-$ —around any point in $\text{int}(\mathcal{M}_{\mu_0}^p)^-$. Let ξ^t be the output of mirror descent at round t and $U_t := (\nabla^2 \psi_{\text{lb}}(\xi_t))^{-1/2}$. We can then use the following as a surrogate for the gradient of $F^t \circ \Lambda_p$ at ξ^t (see also Saha & Tewari, 2011):

$$\frac{(1-\delta)}{\delta} N|\mathcal{X}|(|\mathcal{A}|-1) F^t(\Lambda_p(\xi^t + \delta U_t \mathbf{u}^t)) U_t^{-1} \mathbf{u}^t$$

with \mathbf{u}^t again sampled uniformly from $\mathbb{S}^{N|\mathcal{X}|(|\mathcal{A}|-1)}$. The eigenvalues of U_t correspond to the lengths of the semi-axes of the ellipsoid used at round t , which could be arbitrarily small and lead again to the gradient surrogate having large magnitude. However, thanks to the relationship between ξ^t and U_t , a local norm analysis of mirror descent (see, e.g., Lem. 6.16 in Orabona, 2023) absolves the regret of any dependence on the properties of U_t . Unfortunately, due to technical barriers, this log-barrier-based approach is not readily extendable to the setting where the decision set can change over time (in particular, it is not clear whether an analogue for Lem. 2.1 can hold in this case). Hence, we restrict its application *only to the case when the MDP is known*, see Alg. 4 in App. E.3. We state next a regret bound for this algorithm (proved in App. E.3.2), which requires the following less restrictive assumption in place of Asm. 4.2.

Assumption 4.4. For every state $x \in \mathcal{X}$ and step $n \in [N]$, there exists a policy π such that $\sum_{a \in \mathcal{A}} \mu_n^{\pi, p}(x, a) > 0$.

Note that this can be imposed without loss of generality since the MDP is known; defining $\mathcal{X}_n \subseteq \mathcal{X}$ as the subset of states reachable at step n , one can represent occupancy measures as sequences of distributions in $(\Delta_{\mathcal{X}_n \times \mathcal{A}})_{n \in [N]}$.

Theorem 4.5. Under Asm. 4.4, Alg. 4 with a suitable tuning of τ and δ satisfies for any policy $\pi \in \Pi$ that

$$\mathbb{E}[R_T(\pi)] = \tilde{O}(\sqrt{LN}^{7/4} (|\mathcal{X}||\mathcal{A}|T)^{3/4}).$$

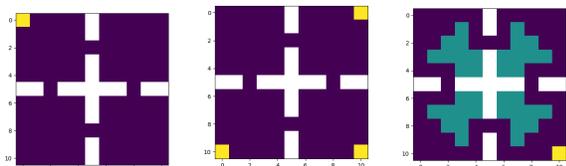


Figure 1. [left] Initial agent distribution; [middle] The three targets from *multi-objectives*; [right] The *constrained MDP* (reward in yellow, constraints in blue).

Though holding only for the known MDP case, this bound maintains the $T^{3/4}$ rate of Thm. 4.3 while eliminating its reliance on Asm. 4.2 and its undesirable dependence on the MDP’s structure. We leave extending this result to unknown MDPs and designing practical approaches enjoying the optimal \sqrt{T} rate for future work.

5. Experiments

We evaluate Bonus O-MD-CURL on the *multi-objective* and *constrained MDP* tasks from (Geist et al., 2022), which use fixed objective functions and fixed probability kernels across time steps. Adversarial and bandit MDPs are harder to implement due to challenges in finding optimal stationary policies, and there is a lack of experimental validation in the literature. We focus on evaluating how well the additive bonus helps the algorithm to learn the environment. We also compare it to Greedy MD-CURL from (Moreno et al., 2024). The state space is an 11×11 four-room grid world, with a single door connecting adjacent rooms. The agent can choose to stay still or move right, left, up, or down, as long as there are no walls blocking the path: $x_{n+1} = x_n + a_n + \varepsilon_n$. The external noise ε_n is a perturbation that can move the agent to a neighboring state with some probability. The initial distribution is a Dirac delta at the upper left corner of the grid, as in Fig. 1 [left]. We take $N = 40$, $\tau = 0.01$, and 5 repetitions per experiment.

Multi-objectives: The goal is to concentrate the distribution on three targets by the final step N , as in Fig. 1 [middle]. The objective function is defined as $f_n(\mu_n^{\pi,p}) := -\sum_{k=1}^3 (1 - \langle \mu_n^{\pi,p}, e^k \rangle)^2$, where $e^k \in \mathbb{R}^{|\mathcal{X}|}$ is a vector with a 1 at the target state and 0 elsewhere. **Constrained MDPs:** The goal is to concentrate the state distribution on the yellow target in Fig. 1 [right] while avoiding the constraint states in blue. The objective function is defined as $f_n(\mu_n^{\pi,p}) := -\langle r, \mu_n^{\pi,p} \rangle + (\langle \mu_n^{\pi,p}, c \rangle)^2$, where $r, c \in \mathbb{R}_+^{|\mathcal{X}| \times |\mathcal{A}|}$. Here, r and c are zero everywhere except at the target and constraint states respectively. For the **Multi-objective** task, Fig. 2 displays the state distribution at the final time step after 50 iterations for Bonus O-MD-CURL [up, left], and Greedy MD-CURL [up, right], and plot the log-loss [down, left] and regret [down, right] after 1000 iterations. We see that Bonus O-MD-CURL reaches

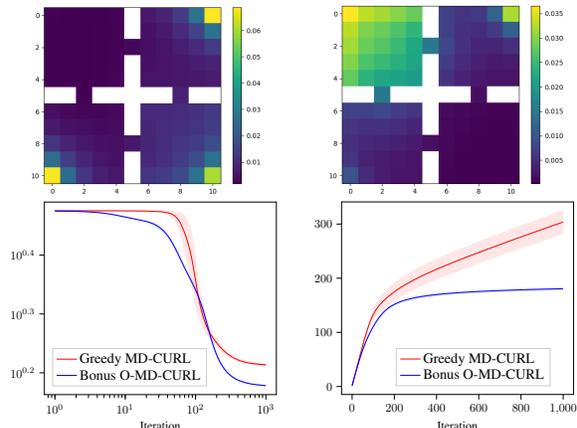


Figure 2. Multi-objective: distribution at $N = 40$ after 50 iters. for Bonus O-MD-CURL [up, left], Greedy MD-CURL [up, right]; log-loss [down, left] and regret [down, right] for 10^3 iters.

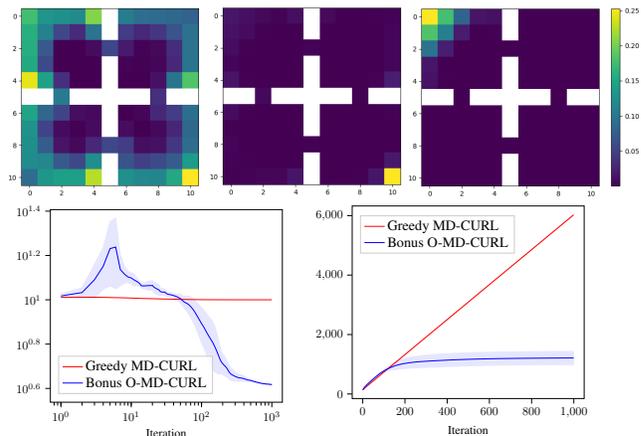


Figure 3. Constrained MDP after 10^3 iters.: sum distributions over all time steps $n \in [40]$ at [up, left]; distribution at the last time step $N = 40$ for Bonus O-MD-CURL [up, center], and Greedy MD-CURL [up, right]; the log-loss [down, left] and regret [down, right].

the targets much faster than Greedy MD-CURL. As for the **Constrained MDP** task, Fig. 3 displays the log-sum of all state distributions for all time steps $n \in [40]$ at iteration 1000 for Bonus O-MD-CURL [up, left]; the state distribution at the last time step $n = 40$ after 1000 iterations for Bonus O-MD-CURL [up, center], and Greedy MD-CURL [up, right]; and the log-loss [down, left] and regret [down, right]. In this case, Greedy MD-CURL fails to reach the target state even after 1000 iterations, while Bonus O-MD-CURL successfully reaches the target state avoiding constrained states to minimize cost thanks to the additive bonuses. These examples empirically demonstrate the value of the additive bonus in tasks requiring exploration.

Impact Statement

This work is of a theoretical nature, we do not foresee any notable societal consequences.

References

- Abbeel, P. and Ng, A. Y. Apprenticeship learning via inverse reinforcement learning. In *Proceedings of the Twenty-First International Conference on Machine Learning, ICML '04*, pp. 1, New York, NY, USA, 2004. Association for Computing Machinery. ISBN 1581138385. doi: 10.1145/1015330.1015430.
- Abernethy, J., Hazan, E., and Rakhlin, A. Competing in the dark: An efficient algorithm for bandit linear optimization. In *21st Annual Conference on Learning Theory, COLT 2008*, volume 3, pp. 263–274, 2008.
- Azar, M. G., Osband, I., and Munos, R. Minimax regret bounds for reinforcement learning. In Precup, D. and Teh, Y. W. (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 263–272. PMLR, 06–11 Aug 2017.
- Barakat, A., Fatkhullin, I., and He, N. Reinforcement learning with general utilities: simpler variance reduction and large state-action space. In *Proceedings of the 40th International Conference on Machine Learning, ICML'23*. JMLR.org, 2023.
- Beck, A. and Teboulle, M. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Oper. Res. Lett.*, 31(3):167–175, may 2003. ISSN 0167-6377. doi: 10.1016/S0167-6377(02)00231-6.
- Bensoussan, A., Yam, P., and Frehse, J. *Mean Field Games and Mean Field Type Control Theory*. SpringerBriefs in Mathematics. Springer, 2013. ISBN 978-1-4614-8507-0. doi: 10.1007/978-1-4614-8508-7.
- Bertsekas, D. *Reinforcement Learning and Optimal Control*. Athena Scientific optimization and computation series. Athena Scientific, 2019. ISBN 9781886529397.
- Bubeck, S. Convex optimization: Algorithms and complexity. *Found. Trends Mach. Learn.*, 8(3–4):231–357, nov 2015. ISSN 1935-8237. doi: 10.1561/22000000050.
- Bubeck, S., Eldan, R., and Lee, Y. T. Kernel-based methods for bandit convex optimization. *Journal of the ACM (JACM)*, 68(4):1–35, 2021.
- Cammardella, N., Bušić, A., and Meyn, S. P. Kullback–leibler-quadratic optimal control. *SIAM Journal on Control and Optimization*, 61(5):3234–3258, 2023. doi: 10.1137/21M1433654.
- Canonne, C. L., Sun, Z., and Suresh, A. T. Concentration bounds for discrete distribution estimation in kl divergence, 2023.
- Cohen, A., Kaplan, H., Koren, T., and Mansour, Y. Online markov decision processes with aggregate bandit feedback. In *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings of Machine Learning Research*, pp. 1301–1329. PMLR, 2021.
- Cover, T. M. and Thomas, J. A. *Elements of information theory (2. ed.)*. Wiley, 2006. ISBN 978-0-471-24195-9.
- Even-Dar, E., Kakade, S. M., and Mansour, Y. Online markov decision processes. *Mathematics of Operations Research*, 34(3):726–736, 2009. ISSN 0364765X, 15265471.
- Flaxman, A. D., Kalai, A. T., and McMahan, H. B. Online convex optimization in the bandit setting: gradient descent without a gradient. In *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '05*, pp. 385–394. Society for Industrial and Applied Mathematics, 2005.
- Fokkema, H., van der Hoeven, D., Lattimore, T., and Mayo, J. J. Online newton method for bandit convex optimisation, 2024.
- García, J., Fern, and o Fernández. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research*, 16(42):1437–1480, 2015.
- Geist, M., Pérolat, J., Laurière, M., Elie, R., Perrin, S., Bachem, O., Munos, R., and Pietquin, O. Concave utility reinforcement learning: The mean-field game viewpoint. In *International Conference on Autonomous Agents and Multiagent Systems*, pp. 489–497, Richland, SC, 2022. ISBN 9781450392136.
- Ghasemipour, S. K. S., Zemel, R., and Gu, S. A divergence minimization perspective on imitation learning methods. In *Proceedings of the Conference on Robot Learning*, volume 100, pp. 1259–1277, 30 Oct–01 Nov 2020.
- Greenberg, I., Chow, Y., Ghavamzadeh, M., and Mannor, S. Efficient risk-averse reinforcement learning. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 32639–32652. Curran Associates, Inc., 2022.
- Hazan, E. Introduction to Online Convex Optimization, December 2021. arXiv preprint arXiv:1909.05207.
- Hazan, E. and Levy, K. Bandit convex optimization: Towards tight bounds. In *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.

- 495 Hazan, E. and Li, Y. An optimal algorithm for bandit convex
496 optimization. *arXiv preprint arXiv:1603.04350*, 2016.
497
- 498 Hazan, E., Kakade, S., Singh, K., and Van Soest, A. Prov-
499 ably efficient maximum entropy exploration. In *Internat-*
500 *ional Conference on Machine Learning*, volume 97, pp.
501 2681–2691, 09–15 Jun 2019.
- 502
- 503 Hu, X., L.A., P., György, A., and Szepesvari, C. (bandit)
504 convex optimization with biased noisy gradient oracles.
505 In Gretton, A. and Robert, C. C. (eds.), *Proceedings of the*
506 *19th International Conference on Artificial Intelligence*
507 *and Statistics*, volume 51 of *Proceedings of Machine*
508 *Learning Research*, pp. 819–828, Cadiz, Spain, 09–11
509 May 2016. PMLR.
- 510
- 511 Jaksch, T., Ortner, R., and Auer, P. Near-optimal regret
512 bounds for reinforcement learning. *J. Mach. Learn. Res.*,
513 11:1563–1600, 2008.
- 514
- 515 Jin, C., Jin, T., Luo, H., Sra, S., and Yu, T. Learning adver-
516 sarial Markov decision processes with bandit feedback
517 and unknown transition. In III, H. D. and Singh, A. (eds.),
518 *Proceedings of the 37th International Conference on Ma-*
519 *chine Learning*, volume 119 of *Proceedings of Machine*
520 *Learning Research*, pp. 4860–4869. PMLR, 13–18 Jul
521 2020.
- 522
- 523 Jézéquel, R., Ostrovskii, D. M., and Gaillard, P. Efficient
524 and near-optimal online portfolio selection, 2022.
- 525
- 526 Lavigne, P. and Pfeiffer, L. Generalized conditional gradient
527 and learning in potential mean field games, 2023.
- 528
- 529 Lavington, J. W., Vaswani, S., and Schmidt, M. Improved
530 policy optimization for online imitation learning. In Chan-
531 dar, S., Pascanu, R., and Precup, D. (eds.), *Proceedings of*
532 *The 1st Conference on Lifelong Learning Agents*, volume
533 199 of *Proceedings of Machine Learning Research*, pp.
534 1146–1173. PMLR, 22–24 Aug 2022.
- 535
- 536 Lee, C.-W., Luo, H., Wei, C.-Y., and Zhang, M. Bias no
537 more: high-probability data-dependent regret bounds for
538 adversarial bandits and mdps. In *Advances in Neural*
539 *Information Processing Systems*, volume 33, pp. 15522–
540 15533. Curran Associates, Inc., 2020.
- 541
- 542 Luo, H., Wei, C.-Y., and Lee, C.-W. Policy optimization
543 in adversarial mdps: Improved exploration via dilated
544 bonuses. In Ranzato, M., Beygelzimer, A., Dauphin, Y.,
545 Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural*
546 *Information Processing Systems*, volume 34, pp. 22931–
547 22942. Curran Associates, Inc., 2021.
- 548
- 549 Maurer, A. and Pontil, M. Empirical bernstein bounds and
sample variance penalization, 2009.
- Moreno, B. M., Bregere, M., Gaillard, P., and Oudjane,
N. Efficient model-based concave utility reinforcement
learning through greedy mirror descent. In Dasgupta, S.,
Mandt, S., and Li, Y. (eds.), *Proceedings of The 27th*
International Conference on Artificial Intelligence and
Statistics, volume 238 of *Proceedings of Machine Learn-*
ing Research, pp. 2206–2214. PMLR, 02–04 May 2024.
- Mourtada, J. and Gaïffas, S. An improper estimator with
optimal excess risk in misspecified density estimation
and logistic regression. *Journal of Machine Learning*
Research, 23(31):1–49, 2022.
- Mutti, M., Pratisoli, L., and Restelli, M. Task-agnostic
exploration via policy gradient of a non-parametric state
entropy estimate, 2021.
- Mutti, M., Santi, R. D., and Restelli, M. The importance of
non-markovianity in maximum state entropy exploration,
2022.
- Mutti, M., Santi, R. D., Bartolomeis, P. D., and Restelli, M.
Challenging common assumptions in convex reinforce-
ment learning, 2023a.
- Mutti, M., Santi, R. D., Bartolomeis, P. D., and Restelli, M.
Convex reinforcement learning in finite trials. *Journal of*
Machine Learning Research, 24(250):1–42, 2023b.
- Nemirovski, A. Interior point polynomial time methods in
convex programming. *Lecture notes*, 42(16):3215–3224,
2004.
- Neu, G. Explore no more: Improved high-probability regret
bounds for non-stochastic bandits, 2015.
- Neu, G., Gyorgy, A., and Szepesvari, C. The adversarial
stochastic shortest path problem with unknown transition
probabilities. In Lawrence, N. D. and Girolami, M. (eds.),
Proceedings of the Fifteenth International Conference
on Artificial Intelligence and Statistics, volume 22 of
Proceedings of Machine Learning Research, pp. 805–813,
La Palma, Canary Islands, 21–23 Apr 2012. PMLR.
- Orabona, F. A modern introduction to online learning. *arXiv*
preprint, arXiv:1912.13213, 2023.
- Pan, X., Seita, D., Gao, Y., and Canny, J. Risk averse
robust adversarial reinforcement learning. In *2019 Inter-*
national Conference on Robotics and Automation (ICRA),
pp. 8522–8528, 2019. doi: 10.1109/ICRA.2019.8794293.
- Rosenberg, A. and Mansour, Y. Online stochastic shortest
path with bandit feedback and unknown transition func-
tion. In Wallach, H., Larochelle, H., Beygelzimer, A.,
d’Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances*
in Neural Information Processing Systems, volume 32.
Curran Associates, Inc., 2019a.

- 550 Rosenberg, A. and Mansour, Y. Online convex optimization
551 in adversarial Markov decision processes. In Chaudhuri,
552 K. and Salakhutdinov, R. (eds.), *Proceedings of the*
553 *36th International Conference on Machine Learning*, vol-
554 *ume 97 of Proceedings of Machine Learning Research*,
555 pp. 5478–5486. PMLR, 09–15 Jun 2019b.
- 556
557 Saha, A. and Tewari, A. Improved regret guarantees for on-
558 line smooth convex optimization with bandit feedback. In
559 *Proceedings of the Fourteenth International Conference*
560 *on Artificial Intelligence and Statistics*, volume 15 of *Pro-*
561 *ceedings of Machine Learning Research*, pp. 636–642.
562 PMLR, 2011.
- 563
564 Shalev-Shwartz, S. Online learning and online convex op-
565 timization. *Found. Trends Mach. Learn.*, 4(2):107–194,
566 February 2012. ISSN 1935-8237. doi: 10.1561/
567 2200000018.
- 568
569 Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An*
570 *Introduction*. A Bradford Book, Cambridge, MA, USA,
571 2018. ISBN 0262039249.
- 572
573 van der Hoeven, D., Zierahn, L., Lancewicki, T., Rosen-
574 berg, A., and Cesa-Bianchi, N. A unified analysis of
575 nonstochastic delayed feedback for combinatorial semi-
576 bandits, linear bandits, and mdps. In *Proceedings of*
577 *Thirty Sixth Conference on Learning Theory*, volume 195
578 *of Proceedings of Machine Learning Research*, pp. 1285–
579 1321. PMLR, 2023.
- 580
581 Wei, C.-Y. and Luo, H. More adaptive algorithms for ad-
582 versarial bandits. In *Conference On Learning Theory*, pp.
583 1263–1291. PMLR, 2018.
- 584
585 Wolsey, L. and Nemhauser, G. *Integer and Combinatorial*
586 *Optimization*. Wiley Series in Discrete Mathematics and
587 Optimization. Wiley, 1999.
- 588
589 Zahavy, T., Cohen, A., Kaplan, H., and Mansour, Y. Appren-
590 ticeship learning via frank-wolfe. In *AAAI Conference on*
591 *Artificial Intelligence*, 2019.
- 592
593 Zahavy, T., O’ Donoghue, B., Desjardins, G., and Singh,
594 S. Reward is enough for convex mdps. In Ranzato, M.,
595 Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan,
596 J. W. (eds.), *Advances in Neural Information Process-*
597 *ing Systems*, volume 34, pp. 25746–25759. Curran Asso-
598 *ciates, Inc.*, 2021.
- 599
600 Zhang, J., Koppel, A., Bedi, A. S., Szepesvari, C., and Wang,
601 M. Variational policy gradient method for reinforcement
602 learning with general utilities. In Larochelle, H., Ranzato,
603 M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances*
604 *in Neural Information Processing Systems*, volume 33,
pp. 4572–4583. Curran Associates, Inc., 2020.
- Zhang, J., Ni, C., Yu, z., Szepesvari, C., and Wang, M.
On the convergence and sample efficiency of variance-
reduced policy gradient method. In Ranzato, M., Beygelz-
imer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.),
Advances in Neural Information Processing Systems, vol-
ume 34, pp. 2228–2240. Curran Associates, Inc., 2021.
- Zimin, A. and Neu, G. Online learning in episodic marko-
vian decision processes by relative entropy policy search.
In *Advances in Neural Information Processing Systems*,
volume 26, pp. 1583–1591, 2013.

A. Auxiliary results

A.1. Auxiliary lemmas

Lemma A.1. For $0 < \delta < 1$,

$$\begin{aligned} & \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} \mu_i^{\pi^t, p}(x, a) \|p_{i+1}(\cdot|x, a) - \hat{p}_{i+1}^t(\cdot|x, a)\|_1 \\ & \leq 3|\mathcal{X}|N^2 \sqrt{2|\mathcal{A}|T \log\left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta}\right)} + 2|\mathcal{X}|N^2 \sqrt{2T \log\left(\frac{N}{\delta}\right)} \end{aligned}$$

with probability at least $1 - 2\delta$.

Proof. Let $\xi_n^t(x, a) := \|p_n(\cdot|x, a) - \hat{p}_n^t(\cdot|x, a)\|_1$. We denote by $o^t := (x_n^t, a_n^t)_{n \in [N]}$ the trajectory of the agent at episode t when playing policy π^t . Let $\hat{\mu}_n^{\pi^t, p}(x, a) := \mathbf{1}_{\{(x_n^t, a_n^t) = (x, a)\}}$ be the empirical state-action distribution computed from the agent's trajectory. We consider the following decomposition:

$$\begin{aligned} \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} \mu_i^{\pi^t, p}(x, a) \xi_{i+1}^t(x, a) &= \underbrace{\sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} \hat{\mu}_i^{\pi^t, p}(x, a) \xi_{i+1}^t(x, a)}_{(1)} \\ &+ \underbrace{\sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} (\mu_n^{\pi^t, p} - \hat{\mu}_i^{\pi^t, p}(x, a)) \xi_{i+1}^t(x, a)}_{(2)}. \end{aligned}$$

Term (1) analysis. We start by analysing the first term. Using Lem. 2.2, we have that for $\delta \in (0, 1)$, with probability $1 - \delta$,

$$(1) = \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} \hat{\mu}_i^{\pi^t, p}(x, a) \xi_{i+1}^t(x, a) \leq \sqrt{2|\mathcal{X}| \log\left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta}\right)} \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} \frac{\mu_i^{\pi^t, p}(x, a)}{\sqrt{\max\{1, N_i^t(x, a)\}}}.$$

Using Lem. 19 from (Jaksch et al., 2008), we have that for all $i \in [N]$ and $(x, a) \in \mathcal{X} \times \mathcal{A}$,

$$\sum_{t=1}^T \frac{\hat{\mu}_i^{\pi^t, p}(x, a)}{\sqrt{\max\{1, N_i^t(x, a)\}}} \leq (\sqrt{2} + 1) \sqrt{N_i^T(x, a)}.$$

Therefore, using Jensen's inequality and that $\sum_{(x,a)} N_i^T(x, a) = T$ for all $i \in [N]$, we have that

$$\begin{aligned} \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} \frac{\hat{\mu}_i^{\pi^t, p}(x, a)}{\sqrt{\max\{1, N_i^t(x, a)\}}} &\leq 3 \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} \sqrt{N_i^T(x, a)} \\ &\leq 3 \sum_{n=1}^N \sum_{i=0}^{n-1} \sqrt{|\mathcal{X}||\mathcal{A}|T} \\ &\leq 3N^2 \sqrt{|\mathcal{X}||\mathcal{A}|T}. \end{aligned} \tag{13}$$

Substituting this inequality into the upper bound for term (1) yields

$$\begin{aligned} (1) &\leq \sqrt{2|\mathcal{X}| \log\left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta}\right)} 3N^2 \sqrt{|\mathcal{X}||\mathcal{A}|T} \\ &= 3|\mathcal{X}|N^2 \sqrt{2|\mathcal{A}|T \log\left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta}\right)}. \end{aligned} \tag{14}$$

Term (2) analysis. We now analyse the second term. Let $\mathcal{F}^t := \sigma(o^1, \dots, o^{t-1})$ be the filtration generated by the trajectories of the agent from the first episode, up to the end of episode $t - 1$. Note that $\xi_{n+1}^t(x, a)$ is \mathcal{F}^t measurable, as it only depends on observations up to episode $t - 1$. Therefore,

$$\mathbb{E}[\xi_{n+1}^t(x, a) \hat{\mu}_n^{\pi^t, p}(x, a) | \mathcal{F}_n^t] = \xi_{n+1}^t(x, a) \mathbb{E}[\hat{\mu}_n^{\pi^t, p}(x, a) | \mathcal{F}_n^t] = \xi_{n+1}^t(x, a) \mu_n^{\pi^t, p}(x, a).$$

For all $n \in [N]$, let $M_n^0 = 0$ and for all $t \in [T]$,

$$M_n^t := \sum_{s=1}^t \sum_{x,a} (\mu_n^{\pi^s, p}(x, a) - \hat{\mu}_n^{\pi^s, p}(x, a)) \xi_{n+1}^s(x, a).$$

From the observation above, $(M_n^t)_{t \in [T]}$ is a martingale sequence with respect to the filtration \mathcal{F}^t . Furthermore, as by definition $|\xi_{n+1}^t(x, a)| \leq 2$,

$$|M_n^t - M_n^{t-1}| \leq \sum_{x \in \mathcal{X}} \left| \sum_{a \in \mathcal{A}} (\mu_n^{\pi^t, p}(x, a) - \hat{\mu}_n^{\pi^t, p}(x, a)) \xi_{n+1}^t(x, a) \right| \leq 2|\mathcal{X}|.$$

Therefore, by Azuma-Hoeffding, we have that for any $\varepsilon > 0$,

$$\mathbb{P}(M_n^T \geq \varepsilon) \leq \exp\left(\frac{-\varepsilon^2}{8|\mathcal{X}|^2 T}\right).$$

Applying the union bound on all $n \in [N]$, we then have that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$M_n^T \leq 2|\mathcal{X}| \sqrt{2T \log\left(\frac{N}{\delta}\right)}$$

holds simultaneously for all $n \in [N]$.

Substituting this inequality into term (2) and summing over $n \in [N]$ and $i \in [n - 1]$, we obtain, with probability at least $1 - \delta$, that

$$(2) = \sum_{n=1}^N \sum_{i=0}^{n-1} M_i^T \leq 2|\mathcal{X}| N^2 \sqrt{2T \log\left(\frac{N}{\delta}\right)}. \quad (15)$$

Final step. Combining the upper bounds for term (1) from Eq. (14) and term (2) from Eq. (15), we obtain, with probability at least $1 - 2\delta$, that

$$\sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} \mu_i^{\pi^t, p}(x, a) \xi_{i+1}^t(x, a) \leq 3|\mathcal{X}| N^2 \sqrt{2|\mathcal{A}| T \log\left(\frac{|\mathcal{X}||\mathcal{A}| N T}{\delta}\right)} + 2|\mathcal{X}| N^2 \sqrt{2T \log\left(\frac{N}{\delta}\right)},$$

concluding the proof. □

Lemma A.2. For any $0 < \delta < 1$,

$$\sum_{t=1}^T \sum_{n=0}^N (N - n) \sum_{x,a} \frac{\mu_n^{\pi^t, p}(x, a)}{\sqrt{\max\{1, N_n^t(x, a)\}}} \leq 3N^2 \sqrt{|\mathcal{X}||\mathcal{A}| T} + |\mathcal{X}| N^2 \sqrt{2T \log\left(\frac{N}{\delta}\right)},$$

holds with probability at least $1 - \delta$.

Proof. Recall that we denote by $(x_n^t, a_n^t)_{n \in \{0, \dots, [N]\}}$ the trajectory of the agent during episode t , when playing policy π^t , and that we define by $\hat{\mu}^{\pi^t, p}(x, a) := \mathbb{1}_{\{(x_n^t, a_n^t) = (x, a)\}}$ as the empirical state-action distribution computed from the trajectory of

the agent. We consider the following decomposition:

$$\begin{aligned}
 \sum_{t=1}^T \sum_{n=0}^N (N-n) \sum_{x,a} \frac{\mu_n^{\pi^t \cdot p}(x,a)}{\sqrt{\max\{1, N_n^t(x,a)\}}} &= \underbrace{\sum_{t=1}^T \sum_{n=0}^N (N-n) \sum_{x,a} \frac{\hat{\mu}_n^{\pi^t \cdot p}(x,a)}{\sqrt{\max\{1, N_n^t(x,a)\}}}}_{(1)} \\
 &+ \underbrace{\sum_{t=1}^T \sum_{n=0}^N (N-n) \sum_{x,a} \frac{(\mu_n^{\pi^t \cdot p} - \hat{\mu}_n^{\pi^t \cdot p})(x,a)}{\sqrt{\max\{1, N_n^t(x,a)\}}}}_{(2)}
 \end{aligned} \tag{16}$$

Term (1) analysis. Using the same decomposition of term (1) of Lem. A.1 in Eq. (13) we have that

$$\sum_{t=1}^T \sum_{n=0}^N (N-n) \sum_{x,a} \frac{\hat{\mu}_n^{\pi^t \cdot p}(x,a)}{\sqrt{\max\{1, N_n^t(x,a)\}}} \leq 3N^2 \sqrt{|\mathcal{X}| |\mathcal{A}| T}. \tag{17}$$

Term (2) analysis. The analysis of term (2) follows a similar approach to the analysis of term (2) in Lem. A.1, with the key difference being that, instead of carrying the term related to the difference between the true probability transition and the estimated one, we now have the term $1/\sqrt{\max\{1, N_n^t(x,a)\}}$.

Let $\mathcal{F}^t := \sigma(o^1, \dots, o^{t-1})$ be the filtration generated by the trajectories of the agent from the first episode, up to the end of episode $t-1$. Note that $1/\sqrt{\max\{1, N_n^t(x,a)\}}$ is \mathcal{F}^t measurable, as it only depends from observations of time step n up to episode $t-1$. Therefore,

$$\mathbb{E}[1/\sqrt{\max\{1, N_n^t(x,a)\}} \hat{\mu}_n^{\pi^t \cdot p}(x,a) | \mathcal{F}_n^t] = 1/\sqrt{\max\{1, N_n^t(x,a)\}} \mu_n^{\pi^t \cdot p}(x,a).$$

For all $n \in [N]$, let $M_n^0 = 0$ and for all $t \in [T]$,

$$M_n^t := \sum_{s=1}^t (N-n) \sum_{x,a} (\mu_n^{\pi^s \cdot p}(x,a) - \hat{\mu}_n^{\pi^s \cdot p}(x,a)) 1/\sqrt{\max\{1, N_n^s(x,a)\}}.$$

From the observation above, $(M_n^t)_{t \in [T]}$ is a martingale sequence with respect to the filtration \mathcal{F}^t . Furthermore, as by definition $|1/\sqrt{\max\{1, N_n^t(x,a)\}}| \leq 1$,

$$|M_n^t - M_n^{t-1}| \leq (N-n) \sum_{x \in \mathcal{X}} \left| \sum_{a \in \mathcal{A}} (\mu_n^{\pi^t \cdot p}(x,a) - \hat{\mu}_n^{\pi^t \cdot p}(x,a)) \xi_n^t(x,a) \right| \leq (N-n) |\mathcal{X}|.$$

Therefore, by Azuma-Hoeffding, we have that for any $\varepsilon > 0$,

$$\mathbb{P}(M_n^T \geq \varepsilon) \leq \exp\left(\frac{-\varepsilon^2}{2|\mathcal{X}|^2 (N-n)^2 T}\right).$$

Applying the union bound on all $n \in [N]$, we then have that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$M_n^T \leq |\mathcal{X}| N \sqrt{2T \log\left(\frac{N}{\delta}\right)}.$$

Summing over $n \in [N]$, we have that with probability at least $1 - \delta$,

$$(2) = \sum_{n=0}^N M_n^T \leq |\mathcal{X}| N^2 \sqrt{2T \log\left(\frac{N}{\delta}\right)}. \tag{18}$$

Joining terms (1) and (2). To conclude, we replace the final upper bounds of the terms (1) and (2) of Eq. (17) and (18) respectively in the decomposition of Eq. (16), and we obtain that, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\sum_{t=1}^T \sum_{n=0}^N (N-n) \sum_{x,a} \frac{\mu_n^{\pi^t, p}(x, a)}{\sqrt{\max\{1, N_n^t(x, a)\}}} \leq 3N^2 \sqrt{|\mathcal{X}||\mathcal{A}|T} + |\mathcal{X}|N^2 \sqrt{2T \log\left(\frac{N}{\delta}\right)},$$

concluding the proof. \square

Lemma A.3. For all $n \in [N]$, $(x, a, x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}$, and $t \in [T]$, let $\hat{p}_{n+1}^t(x'|x, a)$ be defined as in Eq. (7). Hence,

$$\|\hat{p}_{n+1}^{t+1}(\cdot|x, a) - \hat{p}_{n+1}^t(\cdot|x, a)\|_1 \leq \frac{\mathbb{1}_{\{x_n^t=x, a_n^t=a\}}}{\max\{1, N_n^{t+1}(x, a)\}}.$$

Proof. From the definition of the estimator \hat{p}^t , we have that

$$\hat{p}_{n+1}^{t+1}(x'|x, a) = \frac{1}{\max\{1, N_n^{t+1}(x, a)\}} \left(N_n^t(x, a) \hat{p}_{n+1}^t(x'|x, a) + \mathbb{1}_{\{x_{n+1}^t=x', x_n^t=x, a_n^t=a\}} \right).$$

Therefore,

$$\begin{aligned} |\hat{p}_{n+1}^{t+1}(x'|x, a) - \hat{p}_{n+1}^t(x'|x, a)| &= \frac{1}{\max\{1, N_n^{t+1}(x, a)\}} \left| \mathbb{1}_{\{x_{n+1}^t=x', x_n^t=x, a_n^t=a\}} - \hat{p}_{n+1}^t(x'|x, a) (N_n^{t+1}(x, a) - N_n^t(x, a)) \right| \\ &= \frac{1}{\max\{1, N_n^{t+1}(x, a)\}} \left| \mathbb{1}_{\{x_{n+1}^t=x', x_n^t=x, a_n^t=a\}} - \hat{p}_{n+1}^t(x'|x, a) \mathbb{1}_{\{x_n^t=x, a_n^t=a\}} \right|. \end{aligned}$$

Summing over $x' \in \mathcal{X}$ we then have that

$$\|\hat{p}_{n+1}^{t+1}(\cdot|x, a) - \hat{p}_{n+1}^t(\cdot|x, a)\|_1 \leq \frac{\mathbb{1}_{\{x_n^t=x, a_n^t=a\}}}{\max\{1, N_n^{t+1}(x, a)\}},$$

concluding the proof. \square

Lemma A.4. For $(n, x, a) \in [N] \times \mathcal{X} \times \mathcal{A}$, let $(q^t)_{t \in [T]}$ be a sequence of probability transition kernels with $q^t := (q_n^t)_{n \in [N]}$ such that

$$\|q_n^{t+1}(\cdot|x, a) - q_n^t(\cdot|x, a)\|_1 \leq \frac{c \mathbb{1}_{\{x_{n-1}^t=x, a_{n-1}^t=a\}}}{\max\{1, N_{n-1}^{t+1}(x, a)\}}$$

for some constant $c > 0$. Then,

$$\sum_{t=1}^T \|q_n^{t+1}(\cdot|x, a) - q_n^t(\cdot|x, a)\|_1 \leq ec \log(T).$$

Proof. We have that

$$\sum_{t=1}^T \|q_n^{t+1}(\cdot|x, a) - q_n^t(\cdot|x, a)\|_1 \leq c \sum_{t=1}^T \frac{\mathbb{1}_{\{x_{n-1}^t=x, a_{n-1}^t=a\}}}{\max\{1, N_{n-1}^{t+1}(x, a)\}} = c \sum_{t=1}^{N_{n-1}^{T+1}(x, a)} \frac{1}{t} \leq c \sum_{t=1}^T \frac{1}{t} \leq c \log(eT) \stackrel{T \geq 2}{\leq} ec \log(T).$$

Lemma A.5. Let $(q^t)_{t \in [T]}$ be a sequence of probability transition kernels, i.e., $q^t := (q_n^t)_{n \in [N]}$ such that for any state-action pair (x, a) and any step $n \in [N]$, $\sum_{t=1}^T \|q_n^{t+1}(\cdot|x, a) - q_n^t(\cdot|x, a)\|_1 \leq c \log(T)$ for some constant $c > 0$. Then, for any sequence of policies $(\pi^t)_{t \in [T]}$,

$$\sum_{t=1}^T \|\mu^{\pi^t, q^{t+1}} - \mu^{\pi^t, q^t}\|_{\infty, 1} \leq c |\mathcal{X}| |\mathcal{A}| N \log(T).$$

While for a fixed policy π ,

$$\sum_{t=1}^T \|\mu^{\pi, q^{t+1}} - \mu^{\pi, q^t}\|_{\infty, 1} \leq c |\mathcal{X}| N \log(T).$$

825 *Proof.* Using Lem. B.1 we obtain that

$$\begin{aligned}
 826 \quad & \sum_{t=1}^T \|\mu^{\pi^t, q^{t+1}} - \mu^{\pi^t, q^t}\|_{\infty, 1} \leq \sum_{t=1}^T \sup_{n \in [N]} \sum_{i=0}^{n-1} \sum_{x, a} \mu_i^{\pi^t, q^t}(x, a) \|q_{i+1}^{t+1}(\cdot|x, a) - q_{i+1}^t(\cdot|x, a)\|_1 \\
 827 \quad & \\
 828 \quad & \\
 829 \quad & \\
 830 \quad & = \sum_{t=1}^T \sum_{n=0}^{N-1} \sum_{x, a} \mu_n^{\pi^t, q^t}(x, a) \|q_{n+1}^{t+1}(\cdot|x, a) - q_{n+1}^t(\cdot|x, a)\|_1 \\
 831 \quad & \\
 832 \quad & \\
 833 \quad & \leq \sum_{n=0}^{N-1} \sum_{x, a} \sum_{t=1}^T \|q_{n+1}^{t+1}(\cdot|x, a) - q_{n+1}^t(\cdot|x, a)\|_1 \leq c|\mathcal{X}||\mathcal{A}|N \log(T). \\
 834 \quad & \\
 835 \quad &
 \end{aligned}$$

836 While for a fixed policy π ,

$$\begin{aligned}
 837 \quad & \\
 838 \quad & \sum_{t=1}^T \|\mu^{\pi, q^{t+1}} - \mu^{\pi, q^t}\|_{\infty, 1} \leq \sum_{t=1}^T \sum_{n=0}^{N-1} \sum_{x, a} \mu_n^{\pi, q^t}(x, a) \|q_{n+1}^{t+1}(\cdot|x, a) - q_{n+1}^t(\cdot|x, a)\|_1 \\
 839 \quad & \\
 840 \quad & \\
 841 \quad & \leq \sum_{t=1}^T \sum_{n=0}^{N-1} \sum_{x, a} \pi_n(a|x) \|q_{n+1}^{t+1}(\cdot|x, a) - q_{n+1}^t(\cdot|x, a)\|_1 \\
 842 \quad & \\
 843 \quad & \leq c \sum_{n=0}^{N-1} \sum_{x, a} \pi_n(a|x) \log(T) = c|\mathcal{X}|N \log(T). \\
 844 \quad & \\
 845 \quad & \\
 846 \quad & \\
 847 \quad & \\
 848 \quad & \square
 \end{aligned}$$

849 **Lemma A.6.** Consider a sequence of policies $(\pi^t)_{t \in [T]}$, and define a smoothed version of each policy $\tilde{\pi}^t$ for all $t \in [T]$ as $\tilde{\pi}^t := (1 - \alpha_t)\pi^t + \frac{\alpha_t}{|\mathcal{A}|}$, where $\alpha_t \in (0, 1)$. Let p and q be two probability transition kernels, denoted as $p := (p_n)_{n \in [N]}$ and $q := (q_n)_{n \in [N]}$, respectively. Therefore, for all $t \in [T]$,

$$850 \quad \|\mu^{\pi^t, p} - \mu^{\tilde{\pi}^t, q}\|_{\infty, 1} \leq \sum_{i=0}^{N-1} \sum_{x, a} \mu_i^{\pi^t, p}(x, a) \|p_{i+1}(\cdot|x, a) - q_{i+1}(\cdot|x, a)\|_1 + 2N\alpha_t.$$

851 *Proof.* See Lem. D.4 from (Moreno et al., 2024). □

860 A.2. Building a closed-form solution for each OMD iteration

861 In this subsection we argue that the MD optimization problem solved at each iteration in Lem. 2.1 has a closed-form solution. Define the convex function $G^t(\mu) := \tau \langle z^t, \mu \rangle + \Gamma(\mu, \tilde{\mu}^t)$, for $\tau > 0$.

862 Optimizing a convex objective function over policies is equivalent to optimizing it over state-action distributions in $\mathcal{M}_{\mu_0}^p$. Therefore, the optimization problem solved in Lem. 2.1 over the state-action distributions induced by q^{t+1} is equivalent to minimizing the same function over the space of policies:

$$863 \quad \underbrace{\min_{\mu \in \mathcal{M}_{\mu_0}^{q^{t+1}}} G^t(\mu)}_{(i): \text{ state-action problem}} \equiv \underbrace{\min_{\pi \in (\Delta_{\mathcal{A}})^{\mathcal{X} \times \mathcal{N}}} G^t(\mu^{\pi, q^{t+1}})}_{(ii): \text{ policy problem}}. \quad (19)$$

872 In Thm. 4.1 of (Moreno et al., 2024), it is shown that for each episode $t \in [T]$, an optimal policy for the problem

$$873 \quad \min_{\pi \in (\Delta_{\mathcal{A}})^{\mathcal{X} \times \mathcal{N}}} G^t(\mu^{\pi, q^{t+1}}) := \tau \langle z^t, \mu \rangle + \Gamma(\mu, \tilde{\mu}^t), \quad (20)$$

874 defined in Eq. (19), denoted by π^{t+1} , can be computed using an auxiliary sequence of functions $(\tilde{Q}_n^t)_{n \in [N]}$, where $\tilde{Q}_n^t : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$. The sequence starts with $Q_N^t(x, a) = -z_N^t(x, a)$, and for $n \in \{N, \dots, 1\}$, the following recursion is

used:

$$\begin{aligned}\pi_{n+1}^{t+1}(a|x) &= \frac{\tilde{\pi}_{n+1}^t(a|x) \exp(\tau \tilde{Q}_{n+1}^t(x, a))}{\sum_{a'} \tilde{\pi}_{n+1}^t(a'|x) \exp(\tau \tilde{Q}_{n+1}^t(x, a'))}, \\ \tilde{Q}_n^t(x, a) &= -z_n^t(x, a) + \sum_{x'} q_{n+1}^{t+1}(x'|x, a) \sum_{a'} \pi_{n+1}^{t+1}(a'|x') \left[-\frac{1}{\tau} \log \left(\frac{\pi_{n+1}^{t+1}(a'|x')}{\tilde{\pi}_{n+1}^t(a'|x')} \right) + \tilde{Q}_{n+1}^t(x', a') \right].\end{aligned}$$

The core idea of the proof is to show that, due to the specific divergence used (defined in Eq. (4)), Eq. (20) can be solved using dynamic programming. For further details, the reader is referred to Appendix B of (Moreno et al., 2024). A similar result was also obtained by (Cammardella et al., 2023), though they approached the optimization problem using Lagrangian multipliers instead of dynamic programming.

Problem (i) of Eq. (19) is convex, and the theoretical analysis are given in Lem. 2.1. Thanks to the equivalence between problems (i) and (ii) in Eq. (19), we can use the analysis of problem (i) to provide theoretical guarantees for the closed-form solution policy of problem (ii).

B. Missing results and proofs

Lemma B.1. For any strategy $\pi \in (\Delta_{\mathcal{A}})^{\mathcal{X} \times N}$, for any two probability kernels $p = (p_n)_{n \in [N]}$ and $q = (q_n)_{n \in [N]}$ such that $p_n, q_n : \mathcal{X} \times \mathcal{A} \times \mathcal{X} \rightarrow [0, 1]$, and $n \in [N]$,

$$\|\mu_n^{\pi, p} - \mu_n^{\pi, q}\|_1 \leq \sum_{i=0}^{n-1} \sum_{x, a} \mu_i^{\pi, p}(x, a) \|p_{i+1}(\cdot|x, a) - q_{i+1}(\cdot|x, a)\|_1.$$

Proof. From the definition of a state-action distribution sequence induced by a policy π in a probability kernel p in Eq. (2), we have that for all $(x, a) \in \mathcal{X} \times \mathcal{A}$ and $n \in [N]$,

$$\mu_n^{\pi, p}(x, a) = \sum_{x', a'} \mu_{n-1}^{\pi, p}(x', a') p_n(x|x', a') \pi_n(a|x).$$

Thus,

$$\begin{aligned}\|\mu_n^{\pi, p} - \mu_n^{\pi, q}\|_1 &= \sum_{x, a} |\mu_n^{\pi, p}(x, a) - \mu_n^{\pi, q}(x, a)| \\ &= \sum_{x, a} \sum_{x', a'} |\mu_{n-1}^{\pi, p}(x', a') p_n(x|x', a') - \mu_{n-1}^{\pi, q}(x', a') q_n(x|x', a')| \pi_n(a|x) \\ &= \sum_x \sum_{x', a'} |\mu_{n-1}^{\pi, p}(x', a') p_n(x|x', a') - \mu_{n-1}^{\pi, q}(x', a') q_n(x|x', a')| \\ &= \sum_x \sum_{x', a'} |\mu_{n-1}^{\pi, p}(x', a') p_n(x|x', a') - \mu_{n-1}^{\pi, p}(x', a') q_n(x|x', a') \\ &\quad + \mu_{n-1}^{\pi, p}(x', a') q_n(x|x', a') - \mu_{n-1}^{\pi, q}(x', a') q_n(x|x', a')| \\ &\leq \sum_{x', a'} \mu_{n-1}^{\pi, p}(x', a') \|p_n(\cdot|x', a') - q_n(\cdot|x', a')\|_1 + \sum_{x', a'} |\mu_{n-1}^{\pi, p}(x', a') - \mu_{n-1}^{\pi, q}(x', a')| \\ &= \sum_{x', a'} \mu_{n-1}^{\pi, p}(x', a') \|p_n(\cdot|x', a') - q_n(\cdot|x', a')\|_1 + \|\mu_{n-1}^{\pi, p} - \mu_{n-1}^{\pi, q}\|_1.\end{aligned}$$

Since for $n = 0$, $\|\mu_0^{\pi, p} - \mu_0^{\pi, q}\|_1 = 0$, by induction we get that

$$\|\mu_n^{\pi, p} - \mu_n^{\pi, q}\|_1 \leq \sum_{i=0}^{n-1} \sum_{x', a'} \mu_i^{\pi, p}(x', a') \|p_{i+1}(\cdot|x', a') - q_{i+1}(\cdot|x', a')\|_1.$$

□

B.1. Proof of Prop. 3.1

Proof. In the analysis, we explicitly write the term $n = 0$ separately from the other $n \in [N]$. We begin with the following decomposition:

$$\sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} \rangle + \sum_{t=1}^T \langle b_0^t, \mu_0 \rangle = \underbrace{\sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi^t, p} \rangle}_{(1)} + \underbrace{\sum_{t=1}^T \langle b^t, \mu^{\pi^t, p} \rangle + \sum_{t=1}^T \langle b_0^t, \mu_0 \rangle}_{(2)}.$$

Term (1) analysis. Using Holder's inequality, we have that

$$(1) \leq \sum_{t=1}^T \sum_{n=1}^N \|b_n^t\|_\infty \|\mu_n^{\pi^t, \hat{p}^t} - \mu_n^{\pi^t, p}\|_1.$$

From the definition of the bonus sequence, we have that for all $n \in [N]$, $\|b_n^t\|_\infty \leq L(N-n)C_\delta$. Hence,

$$\begin{aligned} (1) &\leq LC_\delta \sum_{t=1}^T \sum_{n=1}^N (N-n) \sum_{i=0}^{n-1} \sum_{x,a} \mu_i^{\pi^t, p}(x, a) \|p_{i+1}(\cdot|x, a) - \hat{p}_{i+1}^t(\cdot|x, a)\|_1 \\ &\leq LC_\delta |\mathcal{X}| N^3 \left[3\sqrt{2|\mathcal{A}|T \log \left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta} \right)} + 2\sqrt{2T \log \left(\frac{N}{\delta} \right)} \right] \end{aligned}$$

where the first inequality comes from Lem. B.1, and the second inequality is achieved for any $\delta \in (0, 1)$, with probability at least $1 - 2\delta$, using Lem. A.1.

Term (2) analysis. Using the definition of the bonus sequence in equation (11), and recalling that the initial state-action distribution μ_0 is always the same, we have that, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$\begin{aligned} (2) &= LC_\delta \sum_{t=1}^T \sum_{n=0}^N (N-n) \sum_{x,a} \frac{\mu_n^{\pi^t, p}(x, a)}{\sqrt{\max\{1, N_n^t(x, a)\}}} \\ &\leq LC_\delta N^2 \left[3\sqrt{|\mathcal{X}||\mathcal{A}|T} + |\mathcal{X}| \sqrt{2T \log \left(\frac{N}{\delta} \right)} \right], \end{aligned}$$

where the inequality comes from Lem. A.2.

Joining the upper bounds in term (1) and (2). Putting both upper bounds together we get that for any $\delta \in (0, 1)$, with probability at least $1 - 3\delta$, and from the definition of C_δ ,

$$\begin{aligned} \sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} \rangle + \sum_{t=1}^T \langle b_0^t, \mu_0 \rangle &\leq LC_\delta |\mathcal{X}| N^3 \left[3\sqrt{2|\mathcal{A}|T \log \left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta} \right)} + 2\sqrt{2T \log \left(\frac{N}{\delta} \right)} \right] \\ &\quad + LC_\delta N^2 \left[3\sqrt{|\mathcal{X}||\mathcal{A}|T} + |\mathcal{X}| \sqrt{2T \log \left(\frac{N}{\delta} \right)} \right] \\ &= O \left(LN^3 |\mathcal{X}|^{3/2} \sqrt{|\mathcal{A}|T \log \left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta} \right)} \right). \end{aligned}$$

□

B.2. Proof of Thm. 3.2 (Main result)

For proving the main result we join together all the pieces we presented in the main paper and the appendix.

Proof. We start by decomposing the regret and using the convexity of the objective function obtaining that

$$\begin{aligned} R_T(\pi) &= \sum_{t=1}^T F^t(\mu^{\pi^t, p}) - F^t(\mu^{\pi^t, \hat{p}^t}) + \sum_{t=1}^T F^t(\mu^{\pi^t, \hat{p}^t}) - F^t(\mu^{\pi, p}) \\ &\leq \underbrace{\sum_{t=1}^T \langle \nabla F^t(\mu^{\pi^t, p}), \mu^{\pi^t, p} - \mu^{\pi^t, \hat{p}^t} \rangle}_{R_T^{\text{MDP}}} + \underbrace{\sum_{t=1}^T \langle \nabla F^t(\mu^{\pi^t, \hat{p}^t}), \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, p} \rangle}_{R_T^{\text{policy}}}. \end{aligned}$$

We analyse each term separately:

Analysis of R_T^{MDP} . We begin by analyzing the term R_T^{MDP} , which represents the cost incurred due to not knowing the true probability kernel. First, we apply Hoeffding's inequality, and the fact that f_n^t is L -Lipschitz with respect to the norm $\|\cdot\|_1$. Following, we apply Lem. B.1, obtaining that

$$\begin{aligned} R_T^{\text{MDP}} &\leq \sum_{t=1}^T \sum_{n=1}^N \|\nabla f_n^t(\mu^{\pi^t, p})\|_\infty \|\mu_n^{\pi^t, p} - \mu_n^{\pi^t, \hat{p}^t}\|_1 \\ &\leq L \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x, a} \mu_i^{\pi^t, p}(x, a) \|p_{i+1}(\cdot|x, a) - \hat{p}_{i+1}^t(\cdot|x, a)\|_1. \end{aligned}$$

We can now apply Lem. A.1 to obtain that for any $\delta \in (0, 1)$, with probability at least $1 - 2\delta$,

$$R_T^{\text{MDP}} \leq L3|\mathcal{X}|N^2 \sqrt{2|\mathcal{A}|T \log\left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta}\right)} + L2|\mathcal{X}|N^2 \sqrt{2T \log\left(\frac{N}{\delta}\right)}. \quad (21)$$

Analysis of R_T^{policy} . To analyse R_T^{policy} we further decompose it as

$$R_T^{\text{policy}} = \underbrace{\sum_{t=1}^T \langle \nabla F^t(\mu^{\pi^t, \hat{p}^t}) - b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle}_{R_T^{\text{policy/MD}}} + \underbrace{\sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle + \sum_{t=1}^T \langle \nabla F^t(\mu^{\pi^t, \hat{p}^t}), \mu^{\pi, \hat{p}^t} - \mu^{\pi, p} \rangle}_{R_T^{\text{policy/bonus}}},$$

where recall that $b^t := (b_n^t)_{n \in [N]}$ is the bonus vector defined in Eq. (11).

Analysis of $R_T^{\text{policy/MD}}$. We begin by addressing the term that accounts for the regret incurred by using online Mirror Descent with changing constraint sets.

From Lemmas A.3 and A.4, we know that the probability sequence $(\hat{p}^t)_{t \in [T]}$ satisfies the condition that $\sum_{t=1}^T \|\hat{p}_n^{t+1} - \hat{p}_n^t\|_1 \leq c \log(T)$ for $c = e$. Additionally, at each time step t , since F^t is L_F -Lipschitz with respect to the norm $\|\cdot\|_{\infty, 1}$, we have $\|\nabla F^t(\mu)\|_{1, \infty} \leq L_F = LN$ for any state-action distribution μ . From the definition of the bonus vector, we also have that $\|b^t\|_{1, \infty} \leq LN^2 C_\delta$. Consequently, $\|\nabla F^t(\mu) - b^t\|_{1, \infty} \leq 2LN^2 C_\delta$. Therefore, as we compute μ^{t+1} by solving

$$\mu^{t+1} := \arg \min_{\mu \in \mathcal{M}_{\mu_0}^{\hat{p}^{t+1}}} \left\{ \tau \langle \nabla F^t(\mu^{\pi^t, \hat{p}^t}) - b^t, \mu \rangle + \Gamma(\mu, \tilde{\mu}^t) \right\},$$

by applying Lem. 2.1 with $\nu^t := \mu^{\pi, \hat{p}^t}$, $\zeta = 2LN^2 C_\delta$, and the sequence of probability transition kernels $(\hat{p}^t)_{t \in [T]}$, we obtain that for the optimal parameter $\tau = \sqrt{\frac{b}{\zeta^2 T}}$, where

$$b := N \left(\log(T) (e|\mathcal{X}||\mathcal{A}| + 4) + \log(|\mathcal{A}|) + 2Ne|\mathcal{X}| \log(T)^2 \log(|\mathcal{A}|) \right),$$

and recalling that $C_\delta := \sqrt{2|\mathcal{X}| \log(|\mathcal{X}||\mathcal{A}|NT/\delta)}$,

$$\begin{aligned}
 R_T^{\text{policy/MD}} &\leq 2\zeta\sqrt{bT} + \zeta N e |\mathcal{X}| \log(T) \\
 &\leq 2LN^2 \sqrt{2|\mathcal{X}| \log\left(\frac{|\mathcal{X}||\mathcal{A}|TN}{\delta}\right)} (2\sqrt{bT} + N \log(T) e |\mathcal{X}|) \\
 &= O\left(LN^2 |\mathcal{X}| \sqrt{T \log\left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta}\right)} \left(\sqrt{N|\mathcal{A}| \log(T)} + N\sqrt{\log(|\mathcal{A}|) \log(T)}\right)\right).
 \end{aligned} \tag{22}$$

Analysis of $R_T^{\text{policy/bonus}}$. We start by analysing the second term of the sum in $R_T^{\text{policy/bonus}}$. For any $\delta \in (0, 1)$, with probability at least $1 - \delta$, we have that

$$\begin{aligned}
 \sum_{t=1}^T \langle \nabla F^t(\mu^{\pi^t, \hat{p}^t}), \mu^{\pi, \hat{p}^t} - \mu^{\pi, p} \rangle &\leq \sum_{t=1}^T \sum_{n=1}^N \|\nabla f_n^t(\mu^{\pi, p})\|_\infty \|\mu_n^{\pi, p} - \mu_n^{\pi, \hat{p}^t}\|_1 \\
 &\leq L \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} \mu_i^{\pi, \hat{p}^t}(x, a) \|p_{i+1}(\cdot|x, a) - \hat{p}_{i+1}^t(\cdot|x, a)\|_1 \\
 &\leq L \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x,a} \mu_i^{\pi, \hat{p}^t}(x, a) \frac{C_\delta}{\sqrt{\max\{1, N_i^t(x, a)\}}} \\
 &= L \sum_{t=1}^T \sum_{n=0}^{N-1} (N-n) \sum_{x,a} \mu_n^{\pi, \hat{p}^t}(x, a) \frac{C_\delta}{\sqrt{\max\{1, N_n^t(x, a)\}}} \\
 &= \sum_{t=1}^T \langle b^t, \mu^{\pi, \hat{p}^t} \rangle + \sum_{t=1}^T \langle b_0^t, \mu_0 \rangle
 \end{aligned}$$

where the first inequality follows from Holder's inequality, the second from the fact that f_n^t is L -Lipschitz with respect to the norm $\|\cdot\|_1$ and Lem. B.1, the third from the concentration bound in Lem. 2.2 where we define $C_\delta := \sqrt{2|\mathcal{X}| \log(|\mathcal{X}||\mathcal{A}|NT/\delta)}$, and the last equality comes from the definition of the bonus vector in Eq. (11).

Replacing it at the $R_T^{\text{policy/bonus}}$ term we have that

$$\begin{aligned}
 R_T^{\text{policy/bonus}} &= \sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle + \sum_{t=1}^T \langle \nabla F^t(\mu^{\pi^t, \hat{p}^t}), \mu^{\pi, \hat{p}^t} - \mu^{\pi, p} \rangle \\
 &\leq \sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle + \sum_{t=1}^T \langle b^t, \mu^{\pi, \hat{p}^t} \rangle + \sum_{t=1}^T \langle b_0^t, \mu_0 \rangle \\
 &\leq \sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} \rangle + \sum_{t=1}^T \langle b_0^t, \mu_0 \rangle.
 \end{aligned}$$

Lastly, we apply Prop. 3.1 to achieve that, for any $\delta \in (0, 1)$, with probability at least $1 - 4\delta$,

$$R_T^{\text{policy/bonus}} \leq \sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} \rangle + \sum_{t=1}^T \langle b_0^t, \mu_0 \rangle = O\left(LN^3 |\mathcal{X}|^{3/2} \sqrt{|\mathcal{A}|T} \log\left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta}\right)\right). \tag{23}$$

Final upper bound on R_T^{policy} . Joining the upper bounds on $R_T^{\text{policy/MD}}$ and $R_T^{\text{policy/bonus}}$ from Eq.s (22) and (23) respectively, we achieve that for any $\delta \in (0, 1)$, with probability at least $1 - 4\delta$, ignoring logarithmic terms,

$$R_T^{\text{policy}} \leq \tilde{O}(LN^3 |\mathcal{X}|^{3/2} \sqrt{|\mathcal{A}|T}). \tag{24}$$

1100 **Joining the upper bounds on R_T^{MDP} and on R_T^{policy} .** Note that the terms in the upper bound on R_T^{policy} from Eq. (24)
 1101 dominate those in the upper bound on R_T^{MDP} from Eq. (21). Therefore, when combining both terms to complete the upper
 1102 bound on the regret, we obtain that, with probability $1 - 6\delta$,

$$1103 \quad R_T(\pi) \leq \tilde{O}(LN^3|\mathcal{X}|^{3/2}\sqrt{|\mathcal{A}|T}),$$

1104 concluding the proof.
 1105
 1106

□

1107 C. Proof of Lem. 2.1: Online Mirror Descent with varying constraint sets

1108 Before stating the proof we recall a few results from Bregman divergences and in particular the divergence Γ defined in
 1109 Eq. (4) that are used throughout the proof.

1110 To simplify notation, for any probability measure $\eta \in \Delta_E$, where E is any finite space, we define the neg-entropy function,
 1111 using the convention that $0 \log(0) = 0$, as $\phi(\eta) := \sum_{x \in E} \eta(x) \log \eta(x)$. For any $\mu := (\mu_n)_{n \in [N]} \in (\Delta_{\mathcal{X} \times \mathcal{A}})^N$, we define
 1112 $\rho_n(x) := \sum_{a \in \mathcal{A}} \mu_n(x, a)$ for all $n \in [N]$ and $x \in \mathcal{X}$, representing the marginal distribution over the state space. The
 1113 function inducing the divergence Γ , defined in Eq. (4), is given by

$$1114 \quad \psi(\mu) := \sum_{n=1}^N \phi(\mu_n) - \sum_{n=1}^N \phi(\rho_n). \quad (25)$$

1115 By definition of a Bregman divergence, for any two probability transition kernels p, q , for all $\mu \in \mathcal{M}_{\mu_0}^p$ and $\mu' \in \mathcal{M}_{\mu_0}^{q,*}$,
 1116 where $\mathcal{M}_{\mu_0}^{q,*}$ is the subset of $\mathcal{M}_{\mu_0}^q$ where the corresponding policies π satisfy $\pi_n(a|x) \neq 0$, we then have that

$$1117 \quad \Gamma(\mu, \mu') := \psi(\mu) - \psi(\mu') - \langle \nabla \psi(\mu'), \mu - \mu' \rangle. \quad (26)$$

1118 Additionally, for any probability transition kernel p , the function ψ is 1-strongly convex with respect to $\|\cdot\|_{\infty,1}$ within $\mathcal{M}_{\mu_0}^p$
 1119 (see Thm. 4.1 from (Moreno et al., 2024)). Consequently, a consequence from a known property of Bregman divergences
 1120 (Shalev-Shwartz, 2012) is that, for any $\mu \in \mathcal{M}_{\mu_0}^p$ and $\mu' \in \mathcal{M}_{\mu_0}^{p,*}$,

$$1121 \quad \Gamma(\mu, \mu') \geq \frac{1}{2} \|\mu - \mu'\|_{\infty,1}^2. \quad (27)$$

1122 *Lemma.* Let $(q^t)_{t \in [T]}$ be a sequence of probability transition kernels, and $(z^t)_{t \in [T]}$ a sequence of vectors in $\mathbb{R}^{N \times |\mathcal{X}| \times |\mathcal{A}|}$,
 1123 such that $\|z^t\|_{1,\infty} \leq \zeta$ for all $t \in [T]$. Initialize $\pi_n^1(a|x) := 1/|\mathcal{A}|$ as the uniform policy. For every $t \in [T]$, let
 1124 $\tilde{\pi}^t := (1 - \alpha_t)\pi^t + \alpha_t|\mathcal{A}|^{-1}$ be a smoothed version of the policy with $\alpha_t := 1/(t+1)$ and $\tilde{\mu}^t := \mu^{\tilde{\pi}^t, q^t}$. For each $t \in [T]$,
 1125 compute iteratively

$$1126 \quad \mu^{t+1} \in \arg \min_{\mu \in \mathcal{M}_{\mu_0}^{q^{t+1}}} \tau \langle z^t, \mu \rangle + \Gamma(\mu, \tilde{\mu}^t). \quad (28)$$

1127 Hence, there is a $\tau > 0$ such that, for any sequence $(\nu^t)_{t \in [T]}$, with $\nu^t := \nu^{\pi, q^t}$ for a common policy π ,

$$1128 \quad \sum_{t=1}^T \langle z^t, \mu^t - \nu^t \rangle \leq O(\zeta N \sqrt{V_T |\mathcal{X}| \log(|\mathcal{A}|) T \log(T)}),$$

1129 where $V_T \geq 1 + \max_{(n,x,a)} \sum_{t=1}^{T-1} \|q_n^t(\cdot|x, a) - q_n^{t+1}(\cdot|x, a)\|_1$.

1130 *Proof.* Throughout this proof, for all $t \in [T]$ we denote by π^t the policy inducing μ^t , meaning that $\mu^t := \mu^{\pi^t, q^t}$ and
 1131 $\tilde{\mu}^t := \mu^{\tilde{\pi}^t, q^t}$. We assume here that $\max_{(n,x,a)} \sum_{t=1}^{T-1} \|q_n^t(\cdot|x, a) - q_n^{t+1}(\cdot|x, a)\|_1 \leq c \log(T)$ for c a constant, as this is the
 1132 case for all the transition estimators we use to obtain the main results of the article.

1133 As $\mathcal{M}_{\mu_0}^{q^{t+1}}$ is a convex set (only linear constraints), the optimality conditions and the definition of a Bregman divergence in
 1134 Eq. (26) imply that for all $\nu^{t+1} \in \mathcal{M}_{\mu_0}^{q^{t+1}}$,

$$1135 \quad \langle \tau z^t + \nabla \psi(\mu^{t+1}) - \nabla \psi(\tilde{\mu}^t), \nu^{t+1} - \mu^{t+1} \rangle \geq 0.$$

Re-arranging the terms and using the three points inequality for Bregman divergences (Bubeck, 2015) we get that,

$$\tau \langle z^t, \mu^{t+1} - \nu^{t+1} \rangle \leq \langle \nabla \psi(\mu^{t+1}) - \nabla \psi(\tilde{\mu}^t), \nu^{t+1} - \mu^{t+1} \rangle = \Gamma(\nu^{t+1}, \tilde{\mu}^t) - \Gamma(\nu^{t+1}, \mu^{t+1}) - \Gamma(\mu^{t+1}, \tilde{\mu}^t).$$

Therefore, by adding and subtracting $\tau \langle z^t, \mu^t - \nu^t \rangle$ on the left-hand side,

$$\begin{aligned} \tau \langle z^t, \mu^{t+1} - \nu^{t+1} \rangle + \tau \langle z^t, \mu^t - \nu^t \rangle - \tau \langle z^t, \mu^t - \nu^t \rangle &\leq \Gamma(\nu^{t+1}, \tilde{\mu}^t) - \Gamma(\nu^{t+1}, \mu^{t+1}) - \Gamma(\mu^{t+1}, \tilde{\mu}^t) \\ \Rightarrow \tau \langle z^t, \mu^t - \nu^t \rangle &\leq \tau \langle z^t, \nu^{t+1} - \nu^t \rangle + \tau \langle z^t, \mu^t - \mu^{t+1} \rangle + \Gamma(\nu^{t+1}, \tilde{\mu}^t) - \Gamma(\nu^{t+1}, \mu^{t+1}) - \Gamma(\mu^{t+1}, \tilde{\mu}^t). \end{aligned}$$

Then, by summing over $t \in [T]$, we obtain that

$$\begin{aligned} \sum_{t=1}^T \langle z^t, \mu^t - \nu^t \rangle &\leq \underbrace{\frac{1}{\tau} \sum_{t=1}^T [\tau \langle z^t, \mu^t - \mu^{t+1} \rangle - \Gamma(\mu^{t+1}, \tilde{\mu}^t)]}_A + \underbrace{\frac{1}{\tau} \sum_{t=1}^T [\Gamma(\nu^{t+1}, \tilde{\mu}^t) - \Gamma(\nu^{t+1}, \mu^{t+1})]}_B \\ &\quad + \underbrace{\sum_{t=1}^T \langle z^t, \nu^{t+1} - \nu^t \rangle}_C. \end{aligned} \quad (29)$$

The term A arises due to our lack of knowledge of z^t at the beginning of episode t for all episodes (adversarial losses hypothesis). To address this, we employ Young's inequality and the strong convexity of Γ . For the term B , in the classic Online Mirror Descent proof (Shalev-Shwartz, 2012), where the set of constraints is fixed, the sum of the differences between the Bregman divergences telescopes (as would be the case with a fixed ν). However, because we are dealing with time-varying constraint sets, this telescoping effect does not occur in our situation. We will now proceed to derive an upper bound for each term, starting with term C that is straightforward.

Step 0: upper bound on C . Applying Holder's inequality, Lem. A.5 with a fixed policy π , and the hypothesis that $\|z^t\|_{1,\infty} \leq \zeta$,

$$C = \sum_{t=1}^T \langle z^t, \nu^{t+1} - \nu^t \rangle \leq \sum_{t=1}^T \|z^t\|_{1,\infty} \|\nu^{t+1} - \nu^t\|_{\infty,1} \leq \zeta c |\mathcal{X}| N \log(T). \quad (30)$$

Step 1: upper bound on B . We now analyse the second term of the sum in Eq. (29). To make the Bregman divergence terms telescope we add and subtract $\Gamma(\nu^t, \mu^t) - \Gamma(\nu^t, \tilde{\mu}^t)$, obtaining

$$\begin{aligned} \sum_{t=1}^T \Gamma(\nu^{t+1}, \tilde{\mu}^t) - \Gamma(\nu^{t+1}, \mu^{t+1}) &= \underbrace{\sum_{t=1}^T \Gamma(\nu^{t+1}, \tilde{\mu}^t) - \Gamma(\nu^t, \tilde{\mu}^t)}_{(i)} + \underbrace{\sum_{t=1}^T \Gamma(\nu^t, \tilde{\mu}^t) - \Gamma(\nu^t, \mu^t)}_{(ii)} \\ &\quad + \underbrace{\sum_{t=1}^T \Gamma(\nu^t, \mu^t) - \Gamma(\nu^{t+1}, \mu^{t+1})}_{(iii)}. \end{aligned} \quad (31)$$

We analyze each term. Using the definition of a Bregman divergence induced by ψ in Eq. (26) we get that

$$\begin{aligned} (i) &= \sum_{t=1}^T \psi(\nu^{t+1}) - \psi(\tilde{\mu}^t) - \langle \nabla \psi(\tilde{\mu}^t), \nu^{t+1} - \tilde{\mu}^t \rangle - \psi(\nu^t) + \psi(\tilde{\mu}^t) + \langle \nabla \psi(\tilde{\mu}^t), \nu^t - \tilde{\mu}^t \rangle \\ &= \sum_{t=1}^T \psi(\nu^{t+1}) - \psi(\nu^t) + \sum_{t=1}^T \langle \nabla \psi(\tilde{\mu}^t), \nu^t - \nu^{t+1} \rangle \\ &\leq -\psi(\nu^1) + \sum_{t=1}^T \|\nabla \psi(\tilde{\mu}^t)\|_{1,\infty} \|\nu^t - \nu^{t+1}\|_{\infty,1}, \end{aligned}$$

1210 where in the last inequality we use the telescoping nature of the first term and applied Hölder's inequality to the second
 1211 term. Recall that for $v := (v_n)_{n \in [N]}$ such that $v_n \in \mathbb{R}^{\mathcal{X} \times \mathcal{A}}$, we defined $\|v\|_{\infty,1} := \sup_{n \in [N]} \|v_n\|_1$. We now also define
 1212 $\|\omega\|_{1,\infty} := \sup_v \{|\langle \omega, v \rangle|, \|v\|_{\infty,1} \leq 1\} = \sum_{n=1}^N \sup_{x,a} |\omega_n(x,a)|$ as the respective dual norm.

1213
 1214 With our choice of Bregman divergence, and given that

$$1215 \quad \tilde{\pi}^t := (1 - \alpha_t)\pi^t + \alpha_t \frac{1}{|\mathcal{A}|},$$

1216
 1217 for each $n \in [N]$, $(x,a) \in \mathcal{X} \times \mathcal{A}$, $|\nabla \psi(\tilde{\mu}^t)(n, x, a)| = |\log(\tilde{\pi}_n^t(a|x))| \leq \log(|\mathcal{A}|/\alpha_t)$.

1218 From the Lemma hypothesis, there is a common policy π such that for all $t \in [T]$, $\nu^t := \nu^{\pi, q^t}$. Hence, from the result
 1219 above,

$$1220 \quad (i) \leq -\psi(\nu^1) + \sum_{t=1}^T N \log\left(\frac{|\mathcal{A}|}{\alpha_t}\right) \|\nu^{\pi, q^t} - \nu^{\pi, q^{t+1}}\|_{\infty,1}$$

$$1221 \quad \leq -\psi(\nu^1) + N \log\left(\frac{|\mathcal{A}|}{\min_{t \in [T]} \alpha_t}\right) c|\mathcal{X}|N \log(T),$$

1222 where the last inequality comes from Lem. A.5 with a fixed π .

1223 As for the second term, using our definition of Γ , we obtain that

$$1224 \quad (ii) = \sum_{t=1}^T \sum_{n,x,a} \nu_n^{\pi, q^t}(x,a) \log\left(\frac{\pi_n(a|x)}{\tilde{\pi}_n^t(a|x)}\right) - \sum_{n,x,a} \nu_n^{\pi, q^t}(x,a) \log\left(\frac{\pi_n(a|x)}{\pi_n^t(a|x)}\right)$$

$$1225 \quad = \sum_{t=1}^T \sum_{n,x,a} \nu_n^{\pi, q^t}(x,a) \log\left(\frac{\pi_n^t(a|x)}{\tilde{\pi}_n^t(a|x)}\right)$$

$$1226 \quad = \sum_{t=1}^T \sum_{n,x,a} \nu_n^{\pi, q^t}(x,a) \log\left(\frac{\pi_n^t(a|x)}{(1 - \alpha_t)\pi_n^t(a|x) + \alpha_t/|\mathcal{A}|}\right)$$

$$1227 \quad \leq N \sum_{t=1}^T (-\log(1 - \alpha_t)) \leq 2N \sum_{t=1}^T \alpha_t,$$

1228 where the last inequality is valid if $0 \leq \alpha_t \leq 0.5$.

1229 The third term telescopes, hence, since $-\Gamma(\nu^{T+1}, \mu^{T+1}) \leq 0$ because a Bregman divergence is always non-negative,

$$1230 \quad (iii) \leq \Gamma(\nu^1, \mu^1).$$

1231 Before adding back the three terms, note that, for $\pi_n^1(a|x) = 1/|\mathcal{A}|$, we have $\Gamma(\nu^1, \mu^1) - \psi(\nu^1) = -\psi(\mu^1)$. Furthermore,
 1232 $-\psi(\mu^1) \leq N \log(|\mathcal{A}|)$. Therefore,

$$1233 \quad \Gamma(\nu^1, \mu^1) - \psi(\nu^1) \leq N \log(|\mathcal{A}|). \quad (32)$$

1234 Summing over our bounds and using the Inequality (32), we get that B is upper bounded as

$$1235 \quad \frac{1}{\tau} \sum_{t=1}^T [\Gamma(\nu^{t+1}, \tilde{\mu}^t) - \Gamma(\nu^{t+1}, \mu^{t+1})] \leq \frac{1}{\tau} [(i) + (ii) + (iii)]$$

$$1236 \quad \leq \frac{N}{\tau} \log(|\mathcal{A}|) + \frac{N^2 c |\mathcal{X}|}{\tau} \log\left(\frac{|\mathcal{A}|}{\min_{t \in [T]} \alpha_t}\right) \log(T) + \frac{2N}{\tau} \sum_{t=1}^T \alpha_t. \quad (33)$$

1237
 1238 **Step 2: Upper bound on A .** It remains to upper bound term A from Eq. (29),

$$1239 \quad A = \frac{1}{\tau} \left[\sum_{t=1}^T \tau \langle z^t, \mu^t - \mu^{t+1} \rangle - \Gamma(\mu^{t+1}, \tilde{\mu}^t) \right], \quad (34)$$

1265 representing what we pay for not knowing the loss function in advance. For that we use Young's inequality (Beck &
1266 Teboulle, 2003): for any $\sigma > 0$ to be optimized later, and for each episode $t \in [T]$,

$$1267 \tau \langle z^t, \mu^t - \mu^{t+1} \rangle - \Gamma(\mu^{t+1}, \tilde{\mu}^t) \leq \frac{\tau^2 \|z^t\|_{1,\infty}^2}{2\sigma} + \frac{\sigma}{2} \|\mu^t - \mu^{t+1}\|_{\infty,1}^2 - \Gamma(\mu^{t+1}, \tilde{\mu}^t). \quad (35)$$

1270 From the definition of Γ in Eq. (4), we have that

$$1271 \Gamma(\mu^{t+1}, \tilde{\mu}^t) = \sum_{n=1}^N \sum_{(x,a)} \mu_n^{t+1}(x, a) \log \left(\frac{\pi_n^{t+1}(a|x)}{\tilde{\pi}^t(a|x)} \right) = \Gamma(\mu^{t+1}, \mu^{\tilde{\pi}^t, q^{t+1}}).$$

1272 From the strong convexity of ψ , as $\mu^{t+1} \in \mathcal{M}_{\mu_0}^{q^{t+1}}$ and $\mu^{\tilde{\pi}^t, q^{t+1}} \in \mathcal{M}_{\mu_0}^{q^{t+1},*}$, we then have from Eq. (27) that

$$1273 \Gamma(\mu^{t+1}, \tilde{\mu}^t) = \Gamma(\mu^{t+1}, \mu^{\tilde{\pi}^t, q^{t+1}}) \geq \frac{1}{2} \|\mu^{t+1} - \mu^{\tilde{\pi}^t, q^{t+1}}\|_{\infty,1}^2. \quad (36)$$

1274 Using the fact that for any vectors $a, b, c \in \mathbb{R}^d$ and for any norm $\|\cdot\|$, the inequality $\|a - b\|^2 \leq 2(\|a - c\|^2 + \|b - c\|^2)$
1275 holds, we then have by Eq. (36)

$$1276 \begin{aligned} \frac{1}{4} \|\mu^t - \mu^{t+1}\|_{\infty,1}^2 - \Gamma(\mu^{t+1}, \tilde{\mu}^t) &\leq \frac{1}{4} \|\mu^t - \mu^{t+1}\|_{\infty,1}^2 - \frac{1}{2} \|\mu^{\tilde{\pi}^t, q^{t+1}} - \mu^{t+1}\|_{\infty,1}^2 \\ &\leq \frac{1}{2} (\|\mu^t - \mu^{\tilde{\pi}^t, q^{t+1}}\|_{\infty,1}^2 + \|\mu^{\tilde{\pi}^t, q^{t+1}} - \mu^{t+1}\|_{\infty,1}^2) - \frac{1}{2} \|\mu^{\tilde{\pi}^t, q^{t+1}} - \mu^{t+1}\|_{\infty,1}^2 \\ &= \frac{1}{2} \|\mu^t - \mu^{\tilde{\pi}^t, q^{t+1}}\|_{\infty,1}^2. \end{aligned} \quad (37)$$

1277 For any $n \in [N]$, we have $\|\mu_n^t - \mu_n^{\tilde{\pi}^t, q^{t+1}}\|_1 \leq 2$. Using this result along with Lem. A.6 for $p = \hat{p}^t$ and $q = \hat{p}^{t+1}$, we derive
1278 the first inequality below. To obtain the second inequality, we apply Lem. A.5 with the sequence of policies $(\pi^t)_{t \in [T]}$.

$$1279 \begin{aligned} \sum_{t=1}^T \|\mu^t - \mu^{\tilde{\pi}^t, q^{t+1}}\|_{\infty,1}^2 &\leq 2 \sum_{t=1}^T \sup_{n \in [N]} \sum_{i=0}^{n-1} \sum_{x,a} \mu_i^t(x, a) \|q_{i+1}^t(\cdot|x, a) - q_{i+1}^{t+1}(\cdot|x, a)\|_1 + 4N \sum_{t=1}^T \alpha_t \\ &\leq 2c|\mathcal{X}||\mathcal{A}|N \log(T) + 4N \sum_{t=1}^T \alpha_t. \end{aligned} \quad (38)$$

1280 Therefore, summing Eq. (35) over $t \in [T]$ with $\sigma = 1/2$, and plugging the inequality above, yields

$$1281 \sum_{t=1}^T \tau \langle z^t, \mu^t - \mu^{t+1} \rangle - \Gamma(\mu^{t+1}, \tilde{\mu}^t) \leq \tau^2 \sum_{t=1}^T \|z^t\|_{1,\infty}^2 + c|\mathcal{X}||\mathcal{A}|N \log(T) + 2N \sum_{t=1}^T \alpha_t.$$

1282 Using that $\|z^t\|_{1,\infty} \leq \zeta$ and dividing by τ entails:

$$1283 A \leq \tau \zeta^2 T + \frac{N}{\tau} \left(c|\mathcal{X}||\mathcal{A}| \log(T) + 2 \sum_{t=1}^T \alpha_t \right). \quad (39)$$

1284 **Conclusion.** Finally, by replacing the final bounds of Eqs. (33), (39), and (30), we obtain

$$1285 \begin{aligned} \sum_{t=1}^T \langle z^t, \mu^t - \nu^t \rangle &\leq A + B + C \\ &\leq \tau T \zeta^2 + \frac{N}{\tau} \left(c|\mathcal{X}||\mathcal{A}| \log(T) + 2 \sum_{t=1}^T \alpha_t \right) + \frac{N}{\tau} \log(|\mathcal{A}|) \\ &\quad + \frac{N^2 c |\mathcal{X}|}{\tau} \log \left(\frac{|\mathcal{A}|}{\min_{t \in [T]} \alpha_t} \right) \log(T) + \frac{2N}{\tau} \sum_{t=1}^T \alpha_t + \zeta N c |\mathcal{X}| \log(T). \end{aligned}$$

In particular, for $\alpha_t = 1/(t+1)$,

$$\sum_{t=1}^T \langle z^t, \mu^t - \nu^t \rangle \leq \tau T \zeta^2 + \frac{1}{\tau} \underbrace{N \left[\log(T) (c|\mathcal{X}||\mathcal{A}| + 4) + \log(|\mathcal{A}|) + 2Nc|\mathcal{X}| \log(T)^2 \log(|\mathcal{A}|) \right]}_{=:b} + \zeta Nc|\mathcal{X}| \log(T).$$

Optimising over $\tau = \sqrt{\frac{b}{\zeta^2 T}}$,

$$\sum_{t=1}^T \langle z^t, \mu^t - \nu^t \rangle \leq 2\zeta \sqrt{bT} + \zeta Nc|\mathcal{X}| \log(T) = O(\zeta \sqrt{cN|\mathcal{X}||\mathcal{A}|T \log(T)} + \zeta N \sqrt{c|\mathcal{X}| \log(|\mathcal{A}|)T} \log(T) + \zeta Nc|\mathcal{X}| \log(T)),$$

concluding the proof. \square

D. Bandit feedback with bonus in RL

Notations. Throughout App. D, we define the trajectory observed by the learner in episode t as $o^t := (x_n^t, a_n^t, \ell_n^t(x_n^t, a_n^t))_{n \in [N]}$. Let \mathcal{F}^t denote the σ -algebra generated by the observations up to episode t , i.e., $\mathcal{F}^t := \sigma(o^1, \dots, o^{t-1})$. We use \mathbb{E}_t to represent the conditional expectation with respect to the observations up to episode t , i.e., $\mathbb{E}_t[\cdot] := \mathbb{E}[\cdot | \mathcal{F}^t]$.

Overview of existing approaches. To adapt Alg. 1 to the bandit case, we need to estimate the loss function for each MD update. A classic choice using importance sampling is:

$$\hat{\ell}_n^t(x, a) = \frac{\ell_n^t(x, a)}{\mu_n^{\pi^t, p}(x, a)} \mathbb{1}_{\{x_n^t = x, a_n^t = a\}}.$$

This update is unbiased, as $\mathbb{E}[\mathbb{1}_{\{x_n^t = x, a_n^t = a\}}] = \mathbb{E}[\mathbb{E}[\mathbb{1}_{\{x_n^t = x, a_n^t = a\}} | \mathcal{F}^t]] = \mu_n^{\pi^t, p}(x, a)$. However, since we do not know the true transition probability p , we cannot use this estimate directly. In (Rosenberg & Mansour, 2019a), they use μ^{π^t, \hat{p}^t} with UC-O-REPS and achieve a regret of $O(T^{3/4})$.

Consider the following confidence set, that is further detailed in Eq. (41),

$$\Omega^t := \{q \mid |q_n(x'|x, a) - \hat{p}_n^t(x'|x, a)| \leq \varepsilon_n^t(x'|x, a), \forall (x, a, x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}, n \in [N]\}.$$

In (Jin et al., 2020), the authors incorporate a parameter γ for implicit exploration, an idea from multi-armed bandits (Neu, 2015), and use the following estimate:

$$\hat{\ell}_n^t(x, a) = \frac{\ell_n^t(x, a)}{\bar{\mu}_n^t(x, a) + \gamma} \mathbb{1}_{\{x_n^t = x, a_n^t = a\}}, \quad (40)$$

where $\bar{\mu}_n^t(x, a) := \max_{q \in \Omega^t} \mu^{\pi^t, q}$. Although this is a biased estimate ($\mu_n^{\pi^t, p}(x, a) \leq \bar{\mu}_n^t(x, a)$), Ω^t is constructed to ensure that the bias introduced is reasonably small. They also argue that $\bar{\mu}$ can be computed efficiently through dynamic programming. They demonstrate that running UC-O-REPS from (Rosenberg & Mansour, 2019b) with this estimate achieves $\tilde{O}(\sqrt{T})$ regret, improving upon previous results.

In Alg. 2 we detail our method for solving the RL problem with adversarial losses, unknown probability transitions and bandit feedback. We proceed to the regret analysis.

D.1. Auxiliary lemmas

Lemma D.1 (Lem. A.2 of (Luo et al., 2021), adapted from Lem. 1 of (Neu, 2015)). *Let $(z_n^t(x, a))_{t \in [T]}$ be a sequence of functions \mathcal{F}_t -measurable, such that $z_n^t(x, a) \in [0, R]$ for each $(x, a) \in \mathcal{X} \times \mathcal{A}$, and $n \in [N]$. Let $Z_n^t(x, a) \in [0, R]$ be a random variable such that $\mathbb{E}_t[Z_n^t(x, a)] = z_n^t(x, a)$. Then with probability $1 - \delta$,*

$$\sum_{t=1}^T \sum_{n=1}^N \sum_{x, a} \left(\frac{\mathbb{1}_{\{x_n^t = x, a_n^t = a\}} Z_n^t(x, a)}{\bar{\mu}_n^t(x, a) + \gamma} - \frac{\mu_n^{\pi^t, p}(x, a) z_n^t(x, a)}{\bar{\mu}_n^t(x, a)} \right) \leq \frac{RN}{2\gamma} \log \left(\frac{N}{\delta} \right).$$

1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429

Algorithm 2 MD-CURL with Additive Bonus for Bandit feedback RL

- 1: **Input:** number of episodes T , initial policy $\pi^1 \in \Pi$, initial state-action distribution μ_0 , initial state-action distribution sequence $\mu^1 = \tilde{\mu}^1 = \mu^{\pi^1, \hat{p}^1}$ with $\hat{p}_n^1(\cdot|x, a) = 1/|\mathcal{A}|$ for all (n, x, a) , learning rate $\tau > 0$, exploration parameter $\gamma = \tau$ (tuned in the proof of Thm. 4.1), sequence of parameters $(\alpha_t)_{t \in [T]}$ with $\alpha_t = 1/(t+1)$.
- 2: **Init.:** $\forall (n, x, a, x'), N_n^1(x, a) = M_n^1(x'|x, a) = 0$
- 3: **for** $t = 1, \dots, T$ **do**
- 4: Agent starts at $(x_0^t, a_0^t) \sim \mu_0(\cdot)$
- 5: **for** $n = 1, \dots, N$ **do**
- 6: Env. draws new state $x_n^t \sim p_n(\cdot|x_{n-1}^t, a_{n-1}^t)$
- 7: Update counts

$$N_{n-1}^{t+1}(x_{n-1}^t, a_{n-1}^t) \leftarrow N_{n-1}^t(x_{n-1}^t, a_{n-1}^t) + 1$$

$$M_{n-1}^{t+1}(x_n^t|x_{n-1}^t, a_{n-1}^t) \leftarrow M_{n-1}^t(x_n^t|x_{n-1}^t, a_{n-1}^t) + 1$$

- 8: Agent chooses an action $a_n^t \sim \pi_n^t(\cdot|x_n^t)$
- 9: Observe local loss $\ell_n^t(x_n^t, a_n^t)$
- 10: **end for**
- 11: Update transition estimate for all (n, x, a, x') : $\hat{p}_n^{t+1}(x'|x, a) := \frac{M_n^{t+1}(x'|x, a)}{\max\{1, N_n^t(x, a)\}}$
- 12: Compute bonus sequence for all (n, x, a) : $b_n^t(x, a) := \frac{(N-n)C_\delta}{\sqrt{\max\{1, N_n^{t+1}(x, a)\}}}$
- 13: Compute optimistic state-action distribution for all (n, x, a) : $\bar{\mu}_n^t(x, a) := \max_{q \in \Omega^t} \mu^{\pi^t, q}$, where Ω^t is defined as in Eq. (41)
- 14: Compute loss estimate for all (n, x, a) : $\hat{\ell}_n^t(x, a) = \frac{\ell_n^t(x, a)}{\bar{\mu}_n^t(x, a) + \gamma} \mathbb{1}_{\{x_n^t = x, a_n^t = a\}}$
- 15: Compute policy $\pi_n^{t+1}(x, a)$ by solving

$$\mu^{t+1} \in \arg \min_{\mu \in \mathcal{M}_{\mu_0}^{\hat{p}^{t+1}}} \{ \tau \langle \hat{\ell}^t - b^t, \mu \rangle + \Gamma(\mu, \tilde{\mu}^t) \},$$

which has a closed-form solution for π^{t+1} (see Sec. A.2)

- 16: Compute $\tilde{\pi}^{t+1}$, the smooth version of π^{t+1} :

$$\tilde{\pi}^{t+1} = (1 - \alpha_t)\pi^{t+1} + \alpha_t/|\mathcal{A}|$$

and the associated state-action distribution $\tilde{\mu}^{t+1} := \mu^{\tilde{\pi}^{t+1}, \hat{p}^{t+1}}$

- 17: **end for**
-

1430 We define the confidence interval used in Alg. 2 as

$$1431 \Omega^t := \{q \mid |q_n(x'|x, a) - \hat{p}_n^t(x'|x, a)| \leq \varepsilon_n^t(x'|x, a), \text{ for all } (x, a, x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}, n \in [N]\}, \quad (41)$$

1433 with

$$1434 \varepsilon_n^t(x'|x, a) := 2\sqrt{\frac{\hat{p}_n^t(x'|x, a) \log\left(\frac{TN|\mathcal{X}||\mathcal{A}|}{\delta}\right)}{\max\{1, N_{n-1}^t(x, a)\}}} + \frac{14 \log\left(\frac{TN|\mathcal{X}||\mathcal{A}|}{\delta}\right)}{3 \max\{1, N_{n-1}^t(x, a)\}}.$$

1438 We present two results regarding this confidence set. The first result, based on the empirical Bernstein inequality, shows
1439 that the true probability kernel p belongs to this confidence set with high probability. The second is a key lemma from (Jin
1440 et al., 2020), which explains how the confidence set shrinks over time. For the proofs, the reader is referred to the original
1441 references.

1442 **Lemma D.2** (Empirical Bernstein inequality, Thm. 4 (Maurer & Pontil, 2009)/Lem. 2 (Jin et al., 2020)). *With probability at
1443 least $1 - 4\delta$, we have that $p \in \Omega^t$ for all $t \in [T]$.*

1444 **Lemma D.3** (Lem. 4 (Jin et al., 2020)). *With probability at least $1 - 6\delta$, for any collection of transition functions $(p^{x,t})_{x \in \mathcal{X}}$
1445 such that $p^{x,t} \in \Omega^t$ for all x , we have*

$$1446 \sum_{t=1}^T \sum_{n=1}^N \sum_{x,a} |\mu_n^{\pi^t, p^{x,t}}(x, a) - \mu_n^{\pi^t, p}(x, a)| = O\left(N^2 |\mathcal{X}| \sqrt{\mathcal{A} T \log\left(\frac{TN|\mathcal{X}||\mathcal{A}|}{\delta}\right)}\right).$$

1451 D.2. Proof of Thm. 4.1

1452 *Proof.* We start by decomposing the static regret with respect to any policy $\pi \in (\Delta_{\mathcal{A}})^{\mathcal{X} \times N}$ as follows

$$1453 R_T(\pi) = \underbrace{\sum_{t=1}^T \langle \ell^t, \mu^{\pi^t, p} - \mu^{\pi^t, \hat{p}^t} \rangle}_{1, R_T^{\text{MDP}}} + \underbrace{\sum_{t=1}^T \langle \ell^t - b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle}_{2, R_T^{\text{policy/MD}}} \quad (42)$$

$$1454 + \underbrace{\sum_{t=1}^T \langle \ell^t, \mu^{\pi, \hat{p}^t} - \mu^{\pi, p} \rangle}_{3, R_T^{\text{policy/MD}}} - \underbrace{\sum_{t=1}^T \langle b^t, \mu^{\pi, \hat{p}^t} \rangle}_{4, \text{Bonus term}} + \underbrace{\sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} \rangle}_{4, \text{Bonus term}}.$$

1464 D.2.1. TERM 1: R_T^{MDP}

1465 The analysis of this term is already provided in App. B.2. Here, we can further leverage the fact that the objective function is
1466 linear and that, by definition, $\ell_n^t \in [0, 1]$. Therefore, with probability at least $1 - 2\delta$, we have:

$$1467 R_T^{\text{MDP}} \leq 3|\mathcal{X}|N^2 \sqrt{2|\mathcal{A}|T \log\left(\frac{|\mathcal{X}||\mathcal{A}|NT}{\delta}\right)} + 2|\mathcal{X}|N^2 \sqrt{2T \log\left(\frac{N}{\delta}\right)}. \quad (43)$$

1472 D.2.2. TERM 2: $R_T^{\text{policy/MD}}$

1473 In practice, the learner plays using the estimated loss function minus the bonus. Hence, $R_T^{\text{policy/MD}}$ accounts for both the bias
1474 introduced by the loss estimation and the standard mirror descent regret bound. We start with the following decomposition:

$$1475 R_T^{\text{policy/MD}} = \sum_{t=1}^T \langle \ell^t - b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle$$

$$1476 = \underbrace{\sum_{t=1}^T \langle \ell^t - \hat{\ell}^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle}_{\text{Bias terms}} + \underbrace{\sum_{t=1}^T \langle \hat{\ell}^t - b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle}_{\text{MD term}}$$

1485 **Mirror Descent term in $R_T^{\text{policy/MD}}$.** We begin by analyzing the error term from applying Mirror Descent with varying
 1486 constraint sets, which is similar to that of Lem. 2.1 with $z^t = \hat{\ell}^t - b^t$, and $(\hat{p}^t)_{t \in [T]}$ as the probability kernel sequence
 1487 defining the varying constraint sets. However, special attention is needed for the sup norm of the subgradient term since
 1488 it now involves an estimate of the loss function. Additionally, the optimal learning rate τ now also depends on both the
 1489 exploration parameter γ and the analysis of the bias terms, we provide a detailed explanation of how the entire analysis is
 1490 affected below.

1491 In proving Lem. 2.1 in App. C, the regret term for Mirror Descent is split into three terms: term A in Eq. (34), term B in
 1492 Eq. (31), and term C in Eq. (30). The analysis of term B remains unchanged since this term is independent of the chosen
 1493 loss function. We focus on what changes for terms A and C .

1495 As in the proof of Lem. 2.1, we use again the notation $\mu^t := \mu^{\pi^t, \hat{p}^t}$ for all $t \in [T]$. From Eq. (34) term A is defined as

$$1497 \quad A = \frac{1}{\tau} \sum_{t=1}^T [\tau \langle \hat{\ell}^t - b^t, \mu^t - \mu^{t+1} \rangle - \Gamma(\mu^{t+1}, \tilde{\mu}^t)].$$

1500 For a fixed t , from Young's inequality we have that

$$1502 \quad \sum_{n=1}^N \tau \langle \hat{\ell}_n^t - b_n^t, \mu_n^t - \mu_n^{t+1} \rangle \leq \sum_{n=1}^N \tau^2 \frac{\|\hat{\ell}_n^t - b_n^t\|_\infty^2}{2\sigma} + \frac{\sigma}{2} \|\mu_n^t - \mu_n^{t+1}\|_1^2.$$

1506 Following the analysis of term A in App. C, in special Eqs. (36), (37), and (38), we obtain that for $\sigma = 1/2$,

$$1508 \quad A \leq \sum_{t=1}^T \sum_{n=1}^N \tau \|\hat{\ell}_n^t - b_n^t\|_\infty^2 + \frac{eN|\mathcal{X}||\mathcal{A}| \log(T)}{\tau} + \frac{2N}{\tau} \sum_{t=1}^T \alpha_t. \quad (44)$$

1511 From Eq. (30), with the notation $\nu^t := \nu^{\pi^t, \hat{p}^t}$, term C is defined as $C = \frac{1}{\tau} \sum_{t=1}^T \tau \langle \hat{\ell}^t - b^t, \nu^{t+1} - \nu^t \rangle$. From Young's
 1512 inequality with $\sigma = 1/2$ we obtain that

$$1514 \quad C \leq \frac{1}{\tau} \sum_{t=1}^T \frac{\tau^2 \|\hat{\ell}^t - b^t\|_{1,\infty}^2}{2\sigma} + \frac{1}{\tau} \frac{\sigma}{2} \sum_{t=1}^T \|\nu^{t+1} - \nu^t\|_{\infty,1}^2$$

$$1517 \quad \leq \sum_{t=1}^T \sum_{n=1}^N \tau \|\hat{\ell}_n^t - b_n^t\|_\infty^2 + \frac{eN|\mathcal{X}| \log(T)}{2\tau}. \quad (45)$$

1521 **Bounding the sup norm of the estimated loss function.** Recall from the definition of the bonus function in Eq. (11),
 1522 with $L = 1$, that $\|b_n^t\|_\infty \leq NC_\delta =: b$ for all $n \in [N]$ and $t \in [T]$. As $\|\hat{\ell}_n^t - b_n^t\|_\infty^2 \leq \|\hat{\ell}_n^t\|_\infty^2 + \|b_n^t\|_\infty^2$, we can focus on the
 1523 term involving the sup norm of the estimated loss function.

1525 We apply Lem. D.1 with $Z_n^t(x, a) = \frac{\mathbb{1}_{\{x_n^t=x, a_n^t=a\}} \ell_n^t(x, a)^2}{\bar{\mu}_n^t(x, a) + \gamma}$ and $z_n^t(x, a) = \frac{\mu_n^{\pi^t, p}(x, a) \ell_n^t(x, a)^2}{\bar{\mu}_n^t(x, a) + \gamma}$. Note that $Z_n^t(x, a), z_n^t(x, a) \leq$
 1526 $\frac{1}{\gamma}$, that $z_n^t(x, a)$ is \mathcal{F}^t -measurable, and that $\mathbb{E}_t[Z_n^t(x, a)] = z_n^t(x, a)$. Therefore, with probability $1 - \delta$,

$$1528 \quad \tau \sum_{t=1}^T \sum_{n=1}^N \|\hat{\ell}_n^t\|_\infty^2 \leq \tau \sum_{t=1}^T \sum_{n=1}^N \sum_{x,a} \hat{\ell}_n^t(x, a)^2$$

$$1531 \quad = \tau \sum_{t=1}^T \sum_{n=1}^N \sum_{x,a} \frac{\mathbb{1}_{\{x_n^t=x, a_n^t=a\}} \ell_n^t(x, a)}{\bar{\mu}_n^t(x, a) + \gamma} \frac{\mathbb{1}_{\{x_n^t=x, a_n^t=a\}} \ell_n^t(x, a)}{\bar{\mu}_n^t(x, a) + \gamma}$$

$$1534 \quad = \tau \sum_{t=1}^T \sum_{n=1}^N \sum_{x,a} \frac{\mathbb{1}_{\{x_n^t=x, a_n^t=a\}} Z_n^t(x, a)}{\bar{\mu}_n^t(x, a) + \gamma}$$

$$1537 \quad \stackrel{\text{Lem. D.1}}{\leq} \tau \sum_{t=1}^T \sum_{n=1}^N \sum_{x,a} \frac{\mu_n^{\pi^t, p}(x, a)}{\bar{\mu}_n^t(x, a)} \frac{\mu_n^{\pi^t, p}(x, a) \ell_n^t(x, a)^2}{\bar{\mu}_n^t(x, a) + \gamma} + \tau \frac{1}{\gamma} \frac{N}{2\gamma} \log\left(\frac{N}{\delta}\right).$$

1540 Lem. D.2 states that the true probability transition $p \in \Omega^t$ with probability $1 - 4\delta$ for all $t \in [T]$. Hence, with probability
 1541 $1 - 4\delta$, $\bar{\mu}_n^t(x, a) \geq \mu_n^{\pi^t, p}(x, a)$. Consequently, by replacing it in the previous inequality, we obtain that with probability
 1542 $1 - 5\delta$,

$$\begin{aligned} 1543 & \tau \sum_{t=1}^T \sum_{n=1}^N \|\hat{\ell}_n^t\|_\infty^2 \leq \tau \sum_{t=1}^T \sum_{n=1}^N \sum_{x,a} \ell_n^t(x, a)^2 + \frac{N\tau}{2\gamma^2} \log\left(\frac{N}{\delta}\right) \\ 1544 & \leq \tau TN|\mathcal{X}||\mathcal{A}| + \frac{N\tau}{2\gamma^2} \log\left(\frac{N}{\delta}\right), \end{aligned}$$

1545 where for the last inequality we use that $\ell_n^t \in [0, 1]$.

1546 Therefore, by replacing it in the upper bound of the terms A and C in Eqs. (44) and (45) respectively, we obtain that

$$1547 \quad A + C \leq \tau 4bTN|\mathcal{X}||\mathcal{A}| + \frac{N\tau}{\gamma^2} \log\left(\frac{N}{\delta}\right) + \frac{3eN|\mathcal{X}||\mathcal{A}| \log(T)}{2\tau} + \frac{2N}{\tau} \sum_{t=1}^T \alpha_t.$$

1548 The upper bound on term B from Eq. (33) in App. C remains the same:

$$1549 \quad B \leq \frac{N}{\tau} \log(|\mathcal{A}|) + \frac{eN^2|\mathcal{X}|}{\tau} \log\left(\frac{|\mathcal{A}|}{\min_{t \in [T]} \alpha_t}\right) \log(T) + \frac{2N}{\tau} \sum_{t=1}^T \alpha_t.$$

1550 Thus, setting $\alpha_t = 1/(t+1)$, the final upper bound on the Mirror Descent term is, with high probability, given by

$$\begin{aligned} 1551 & \sum_{t=1}^T \langle \hat{\ell}^t - b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle \leq A + B + C \\ 1552 & \leq \tau 4bTN|\mathcal{X}||\mathcal{A}| + \frac{N\tau}{\gamma^2} \log\left(\frac{N}{\delta}\right) + 6 \frac{eN|\mathcal{X}||\mathcal{A}| \log(T)}{\tau} + \frac{N}{\tau} \log(|\mathcal{A}|) + \frac{eN^2|\mathcal{X}|}{\tau} \log(|\mathcal{A}|T) \log(T). \end{aligned} \quad (46)$$

1553 Before tuning the optimal parameter τ , we must first analyze the bias terms.

1554 **Bias terms.** We now proceed to analyze the bias terms. Our approach is similar to the one used in (Jin et al., 2020), with a
 1555 key difference: they utilize confidence sets in their Mirror Descent iterations, whereas we perform iterations over the set
 1556 induced by \hat{p}^t . We start by dividing the bias term in two:

$$\begin{aligned} 1557 & \sum_{t=1}^T \langle \ell^t - \hat{\ell}^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle = \underbrace{\sum_{t=1}^T \langle \ell^t - \hat{\ell}^t, \mu^{\pi^t, \hat{p}^t} \rangle}_{\text{Bias 1}} + \underbrace{\sum_{t=1}^T \langle \hat{\ell}^t - \ell^t, \mu^{\pi, \hat{p}^t} \rangle}_{\text{Bias 2}}. \end{aligned}$$

1558 **Bias 1.** Since μ^{π^t, \hat{p}^t} is \mathcal{F}^t -measurable, we have that $\mathbb{E}_t[\langle \ell^t - \hat{\ell}^t, \mu^{\pi^t, \hat{p}^t} \rangle] = \langle \mathbb{E}_t[\ell^t - \hat{\ell}^t], \mu^{\pi^t, \hat{p}^t} \rangle$. For any couple (x, a) ,
 1559 and for any time step $n \in [N]$,

$$1560 \quad \mathbb{E}_t[\ell_n^t(x, a) - \hat{\ell}_n^t(x, a)] = \ell_n^t(x, a) - \frac{\ell_n^t(x, a) \mu_n^{\pi^t, p}(x, a)}{\bar{\mu}_n^t(x, a) + \gamma} = \ell_n^t(x, a) \left(\frac{\bar{\mu}_n^t(x, a) + \gamma - \mu_n^{\pi^t, p}(x, a)}{\bar{\mu}_n^t(x, a) + \gamma} \right).$$

1561 Hence,

$$1562 \quad \mathbb{E}_t[\text{Bias 1}] = \sum_{t=1}^T \sum_{n=1}^N \sum_{x,a} \mu_n^{\pi^t, \hat{p}^t}(x, a) \ell_n^t(x, a) \left(\frac{\bar{\mu}_n^t(x, a) + \gamma - \mu_n^{\pi^t, p}(x, a)}{\bar{\mu}_n^t(x, a) + \gamma} \right).$$

1563 From Lem. D.2, and from the definition of $\bar{\mu}$, we have that with high probability, $\bar{\mu}_n^t(x, a) \geq \mu_n^{\pi^t, \hat{p}^t}(x, a)$, therefore,

$$1564 \quad \mathbb{E}_t[\text{Bias 1}] \leq \sum_{t=1}^T \sum_{n=1}^N \sum_{x,a} |\bar{\mu}_n^t(x, a) + \gamma - \mu_n^{\pi^t, p}(x, a)| \leq \sum_{t=1}^T \sum_{n=1}^N \sum_{x,a} |\bar{\mu}_n^t(x, a) - \mu_n^{\pi^t, p}(x, a)| + \gamma |\mathcal{X}||\mathcal{A}|NT.$$

1595 Note that $\bar{\mu}_n^t(x, a) = \max_{p \in \Omega^t} \mu_n^{\pi^t, p}(x, a) = \pi_n^t(a|x) \max_{p \in \Omega^t} \rho_n^{\pi^t, p}(x)$, where $\rho_n^{\pi^t, p}(x) := \sum_{a \in \mathcal{A}} \mu_n^{\pi^t, p}(x, a)$. Therefore,
 1596 for each $x \in \mathcal{X}$, there is a $p^{x,t} \in \Omega^t$ such that $\bar{\mu}_n^t(x, a) = \mu_n^{\pi^t, p^{x,t}}(x, a)$.

1597 From Lem. D.3 we obtain that

$$1598 \sum_{t=1}^T \sum_{n=1}^N \sum_{x,a} |\mu_n^{\pi^t, p^{x,t}}(x, a) - \mu_n^{\pi^t, p}(x, a)| = O\left(N^2 |\mathcal{X}| \sqrt{|\mathcal{A}| T \log\left(\frac{|\mathcal{X}| |\mathcal{A}| NT}{\delta}\right)}\right).$$

1602 Therefore,

$$1603 \mathbb{E}_t[\text{Bias 1}] = O\left(N^2 |\mathcal{X}| \sqrt{|\mathcal{A}| T \log\left(\frac{|\mathcal{X}| |\mathcal{A}| NT}{\delta}\right)}\right) + \gamma |\mathcal{X}| |\mathcal{A}| NT.$$

1607 As we have that $\text{Bias 1} = \mathbb{E}_t[\text{Bias 1}] + \sum_{t=1}^T \langle \mathbb{E}_t[\hat{\ell}^t] - \hat{\ell}^t, \mu^{\pi^t, \hat{p}^t} \rangle$, all that remains is to treat the second term of the sum.
 1608 With high probability, $\bar{\mu}_n^t(x, a) \geq \mu_n^{\pi^t, \hat{p}^t}(x, a)$, therefore

$$1610 \sum_{n=1}^N \sum_{x,a} \hat{\ell}_n^t(x, a) \mu_n^{\pi^t, \hat{p}^t}(x, a) \leq \sum_{n=1}^N \sum_{x,a} \ell_n^t(x, a) \mathbb{1}_{\{x_n^t = x, a_n^t = a\}} \leq N.$$

1613 Thus, Azuma's inequality gives us that

$$1615 \sum_{t=1}^T \langle \mathbb{E}_t[\hat{\ell}^t] - \hat{\ell}^t, \mu^{\pi^t, \hat{p}^t} \rangle \leq N \sqrt{2T \log\left(\frac{1}{\delta}\right)},$$

1618 which is of a smaller order than the terms previously appearing in the bias bound. Hence,

$$1620 \text{Bias 1} = O\left(N^2 |\mathcal{X}| \sqrt{|\mathcal{A}| T \log\left(\frac{|\mathcal{X}| |\mathcal{A}| NT}{\delta}\right)}\right) + \gamma |\mathcal{X}| |\mathcal{A}| NT.$$

1623 **Bias 2.** The result follows directly from Lem. 14 of (Jin et al., 2020) using μ^{π, \hat{p}^t} instead of $\mu^{\pi, p}$:

$$1625 \text{Bias 2} = \sum_{t=1}^T \langle \hat{\ell}^t - \ell^t, \mu^{\pi, \hat{p}^t} \rangle = O\left(\frac{N \log\left(\frac{|\mathcal{X}| |\mathcal{A}| N}{\delta}\right)}{\gamma}\right).$$

1628 **Optimizing the learning and exploration parameters τ and γ .** By joining the Mirror Descent term from Eq. (46)
 1629 along with the bounds on Bias 1 and Bias 2 terms, and setting $\gamma = \tau$, we obtain that

$$1631 R_T^{\text{policy/MD}} = O\left(\tau b N |\mathcal{X}| |\mathcal{A}| T + \frac{N}{\tau} \left[\log\left(\frac{N}{\delta}\right) + |\mathcal{X}| |\mathcal{A}| \log(T) + \log(|\mathcal{A}|) + N |\mathcal{X}| \log(|\mathcal{A}| T) \log(T) \right]\right. \\ 1634 \left. + N^2 |\mathcal{X}| \sqrt{|\mathcal{A}| T \log\left(\frac{|\mathcal{X}| |\mathcal{A}| NT}{\delta}\right)} + \tau |\mathcal{X}| |\mathcal{A}| NT + \frac{N \log\left(\frac{|\mathcal{X}| |\mathcal{A}| N}{\delta}\right)}{\tau}\right).$$

1637 Let $\varphi_1 := (b+1)N |\mathcal{X}| |\mathcal{A}|$, and

$$1638 \varphi_2 := N \left[\log\left(\frac{N}{\delta}\right) + |\mathcal{X}| |\mathcal{A}| \log(T) + \log(|\mathcal{A}|) + N |\mathcal{X}| \log(|\mathcal{A}| T) \log(T) + \log\left(\frac{|\mathcal{X}| |\mathcal{A}| N}{\delta}\right) \right].$$

1641 For $\tau \propto \sqrt{\varphi_2 / \varphi_1 T}$, recalling that $b := N C_\delta$, and that $C_\delta := \sqrt{2 |\mathcal{X}| \log(|\mathcal{X}| |\mathcal{A}| NT / \delta)}$, we obtain that, with high probability,

$$1643 R_T^{\text{policy/MD}} = 2\sqrt{\varphi_1 \varphi_2 T} = \tilde{O}\left(N^{3/2} |\mathcal{X}|^{5/4} |\mathcal{A}| \sqrt{T} + N^2 |\mathcal{X}|^{5/4} \sqrt{|\mathcal{A}| T}\right). \quad (47)$$

1645 D.2.3. TERM 3: $R_T^{\text{policy/MDP}}$

1646 The upper bound for this term directly follows from the analysis of adding the bonus term to compensate for insufficient
 1647 exploration, as discussed in Subsec. 3.2 of the main paper, and is detailed in App. B.2. Thus, we have that with high
 1648 probability, $R_T^{\text{policy/MDP}} \leq 0$.

1650 D.2.4. TERM 4: BONUS TERM

 1651 The analysis of this term follows directly from Prop. 3.1 from the main paper: for any $\delta \in (0, 1)$, with probability at least
 1652 $1 - 3\delta$, we have that

1653
$$\sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} \rangle = \tilde{O}(N^3 |\mathcal{X}|^{3/2} \sqrt{|\mathcal{A}|T}). \quad (48)$$

 1657 **D.3. Final bound**

 1658 Replacing the upper bound on all four terms from Eq.s (43), (47), (48), and that $R_T^{\text{policy/MDP}} \leq 0$ into the regret decomposition
 1659 in Eq. (42), we obtain that playing Alg. 2 for the RL problem with bandit feedback on adversarial loss functions has, with
 1660 high probability, a static regret of order

1662
$$R_T(\pi) = \tilde{O}(N^3 |\mathcal{X}|^{3/2} \sqrt{|\mathcal{A}|T} + N^{3/2} |\mathcal{X}|^{5/4} |\mathcal{A}| \sqrt{T}).$$

 1664 \square

 1666 **E. Curl with bandit feedback**

 1668 **Notation** Let S and A be two positive integers. For convenience and brevity, we suppose in what follows that $\mathcal{X} = [S]$
 1669 and $\mathcal{A} = [A]$. Accordingly, we will often use S and A in place of $|\mathcal{X}|$ and $|\mathcal{A}|$ respectively. Let $\dot{A} \in \{A, A-1\}$. For
 1670 a vector $\xi \in \mathbb{R}^{NS\dot{A}}$ and $(n, x, a) \in [N] \times [S] \times [\dot{A}]$, we use $\xi_n(x, a)$ as a shorthand for $\xi((n-1)S\dot{A} + (x-1)\dot{A} + a)$.
 1671 Similarly, let $\dot{A} \in \{A, A-1\}$. Then, for an $NS\dot{A} \times NS\dot{A}$ matrix M , we use $M(n, x, a, n', x', a')$ to denote the item in row
 1672 $(n-1)S\dot{A} + (x-1)\dot{A} + a$ and column $(n'-1)S\dot{A} + (x'-1)\dot{A} + a'$ of M . For any $d \in \mathbb{Z}_+$, let $\mathbb{1}_d \in \mathbb{R}^d$ be the vector
 1673 with all entries equal to one and \mathbf{I}_d the $d \times d$ identity matrix.

 1675 **E.1. An alternative representation for the decision sets**

 1677 In the following, we will fix an arbitrary transition kernel $p := (p_n)_{n \in [N]}$. We recall the notation that for $\zeta \in \mathcal{M}_{\mu_0}^p$,
 1678 $\rho_n^\zeta(x) := \sum_{a \in \mathcal{A}} \zeta_n(x, a)$ for $(n, x) \in [N] \times \mathcal{X}$, which satisfies $\rho_n^\zeta(x) = \sum_{x', a' \in \mathcal{X} \times \mathcal{A}} \zeta_{n-1}(x', a') p_n(x|x', a')$ for $n \geq 2$.
 1679 At the first step, we define $\rho_1^p(x) := \sum_{x', a' \in \mathcal{X} \times \mathcal{A}} \mu_0(x', a') p_1(x|x', a')$, which satisfies $\rho_1^p(x) = \rho_1^\zeta(x)$ for every $\zeta \in \mathcal{M}_{\mu_0}^p$
 1680 and $x \in \mathcal{X}$ since the initial state distribution is the same for all occupancy measures in $\mathcal{M}_{\mu_0}^p$.

 1682 We describe here the mapping alluded to in Sec. 4.2.1 of $\mathcal{M}_{\mu_0}^p$ to a lower-dimensional space where it could have a
 1683 non-empty interior. This is analogous to how one can define a bijective map between the simplex Δ_d and the set
 1684 $\{x \in \mathbb{R}^{d-1} : \mathbb{1}_{d-1}^\top x \leq 1 \text{ and } x_i \geq 0 \forall i \in [d-1]\}$, which is the intersection of the positive orthant of \mathbb{R}^{d-1} with the L_1 unit
 1685 ball, see (Jézéquel et al., 2022, Section 2). This can be done since any coordinate x_{i^*} of a vector $x \in \Delta_d$ can be recovered
 1686 from the rest of the coordinates: $x_{i^*} = 1 - \sum_{i \neq i^*} x_i$. In our case, denoting by a^* the last action in \mathcal{A} (i.e., $a^* = A$, recalling
 1687 that $\mathcal{A} = [A]$), we will represent the occupancy measures as vectors in $\mathbb{R}^{NS(A-1)}$ by omitting all coordinates that correspond
 1688 to this action. We can afford to do so, since for any $\mu \in \mathcal{M}_{\mu_0}^p$, we have that $\mu_n(x, a^*) = \rho_n^\mu(x) - \sum_{a \neq a^*} \mu_n(x, a)$ where
 1689 ρ_n^μ is recoverable from μ_{n-1} and given in the first step by the initial state distribution $\rho_1^p(x)$, which does not depend on μ . In
 1690 the following, we use this idea to define the sought mapping.

 1691 Define the $A \times (A-1)$ matrix

1693
$$G := \begin{bmatrix} \mathbf{I}_{A-1} \\ -\mathbb{1}_{A-1}^\top \end{bmatrix}$$

 1696 and let H be the $NSA \times NS(A-1)$ matrix obtained via taking the direct sum of NS copies of G : $H := \bigoplus_{i=1}^{NS} G$. Define
 1697 $w^{p,1} \in \mathbb{R}^{NSA}$ such that

1698
$$w_n^{p,1}(x, a) := \rho_1^p(x) \mathbb{I}\{n=1, a=a^*\}.$$

 1699 Next, for every $2 \leq m \leq N$, we define $W^{p,m}$ as the $NSA \times NSA$ matrix where

1701
$$W^{p,m}(n, x, a, n', x', a') := \mathbb{I}\{n=m, n'=m-1, a=a^*\} p_m(x|x', a').$$

 1703 ¹For an $n \times m$ matrix M and an $n' \times m'$ matrix M' , $M \oplus M'$ is the $(n+n') \times (m+m')$ block matrix $\begin{bmatrix} M & \mathbf{0} \\ \mathbf{0} & M' \end{bmatrix}$.

Then, we define the $NSA \times NS(A-1)$ matrix

$$B^p := (\mathbf{I}_{NSA} + W^{p,N}) \dots (\mathbf{I}_{NSA} + W^{p,3})(\mathbf{I}_{NSA} + W^{p,2})H$$

and the vector

$$\beta^p := (\mathbf{I}_{NSA} + W^{p,N}) \dots (\mathbf{I}_{NSA} + W^{p,3})(\mathbf{I}_{NSA} + W^{p,2})\mathbf{w}^{p,1}.$$

Finally, define the function $\Xi_p: \mathbb{R}^{NS(A-1)} \rightarrow \mathbb{R}^{NSA}$ where

$$\Xi_p(\xi) := B^p \xi + \beta^p$$

for $\xi \in \mathbb{R}^{NS(A-1)}$.

To explain the semantics of Ξ_p , let $\mu \in \mathcal{M}_{\mu_0}^p$ and $\tilde{\mu} \in \mathbb{R}^{NS(A-1)}$ be such that $\tilde{\mu}_n(x, a) := \mu_n(x, a)$ for all $(n, x, a) \in [N] \times \mathcal{X} \times \mathcal{A} \setminus a^*$. It then holds that $\Xi_p(\tilde{\mu}) = \mu$. To see this, note that $H\tilde{\mu}$ expands $\tilde{\mu}$ setting $(H\tilde{\mu})_n(x, a^*) = -\sum_{a \neq a^*} \mu_n(x, a)$. To fully recover $\mu_n(x, a^*)$, what remains is to add $\rho_n^\mu(x)$. This is achieved at $n = 1$ by adding $\mathbf{w}^{p,1}$ to $H\tilde{\mu}$ since $\mathbf{w}_n^{p,1}(x, a^*) = \rho_1^\mu(x) = \rho_1^\mu(x)$ and $\mathbf{w}_n^{p,1}(x, a) = 0$ for $a \neq a^*$. Next, at $n = 2$, the matrix $W^{p,2}$ extracts the values $\rho_2^\mu(x)$ when operated on $H\tilde{\mu} + \mathbf{w}^{p,1}$ such that $\mu_2(x, a^*)$ is recovered at coordinate $(2, x, a^*)$ of $(\mathbf{I}_{NSA} + W^{p,2})(H\tilde{\mu} + \mathbf{w}^{p,1})$. Iterating this procedure until step N allows us to fully recover μ from $\tilde{\mu}$. While for a generic $\xi \in \mathbb{R}^{NS(A-1)}$, $\Xi_p(\xi)$ is the unique vector in \mathbb{R}^{NSA} satisfying $(\Xi_p(\xi))_n(x, a) = \xi_n(x, a)$ for all n, x , and $a \neq a^*$; $(\Xi_p(\xi))_1(x, a^*) = \rho_1^\mu(x) - \sum_{a \neq a^*} (\Xi_p(\xi))_1(x, a)$ for all x ; and $(\Xi_p(\xi))_n(x, a^*) = \sum_{x', a'} (\Xi_p(\xi))_{n-1}(x', a') p_n(x|x', a') - \sum_{a \neq a^*} (\Xi_p(\xi))_n(x, a)$ for all x and $n \geq 2$.

Note that B^p has full column rank since for any $\xi \in \mathbb{R}^{NS(A-1)}$, $B^p \xi$ is only an expansion of ξ ; hence, we can define its left pseudo-inverse $(B^p)^+ := ((B^p)^\top B^p)^{-1} (B^p)^\top$, which satisfies $(B^p)^+ B^p = \mathbf{I}_{NS(A-1)}$. On the other hand, the matrix $B^p (B^p)^+$ projects vectors in \mathbb{R}^{NSA} onto the column space of B^p , which is given by

$$\left\{ \mu \in \mathbb{R}^{NSA} : \sum_a \mu_n(x, a) = \sum_{x', a'} \mu_{n-1}(x', a') p_n(x|x', a') \forall x \in \mathcal{X}, 2 \leq n \leq N \text{ and } \sum_a \mu_1(x, a) = 0 \forall x \in \mathcal{X} \right\}. \quad (49)$$

It is easy to verify that for any $\mu, \mu' \in \mathcal{M}_{\mu_0}^p$, $\mu - \mu'$ lies in the column space of B^p (recall that $\sum_a \mu_1(x, a) = \sum_a \mu'_1(x, a) = \rho_1^\mu(x)$). Moreover, $\beta^p \in \mathcal{M}_{\mu_0}^p$ as it corresponds to a policy π where $\pi_n(a^*|x) = 1$ for all n and x . Therefore, for any $\mu \in \mathcal{M}_{\mu_0}^p$, $\mu - \beta^p$ belongs to the column space of B^p , and we consequently have that

$$\Xi_p((B^p)^+(\mu - \beta^p)) = B^p (B^p)^+(\mu - \beta^p) + \beta^p = \mu.$$

Hence, by the definition of Ξ_p , $(B^p)^+(\mu - \beta^p)$ coincides with μ on all coordinates $(n, x, a) \in [N] \times \mathcal{X} \times \mathcal{A} \setminus \{a^*\}$ (since the map Ξ_p only expands the input vector adding the coordinates corresponding to action a^*), and is then the only point in $\mathbb{R}^{NS(A-1)}$ that Ξ_p maps to μ . In light of this, we define

$$(\mathcal{M}_{\mu_0}^p)^- := \{ \xi \in \mathbb{R}^{NS(A-1)} \mid \Xi_p(\xi) \in \mathcal{M}_{\mu_0}^p \},$$

the pre-image of $\mathcal{M}_{\mu_0}^p$ under Ξ_p . Accordingly, we define $\Lambda_p: (\mathcal{M}_{\mu_0}^p)^- \rightarrow \mathcal{M}_{\mu_0}^p$ as the restriction of Ξ_p to $(\mathcal{M}_{\mu_0}^p)^-$; that is,

$$\Lambda_p := \Xi_p|_{(\mathcal{M}_{\mu_0}^p)^-}.$$

This then is a bijective function, with $\Lambda_p^{-1}(\mu) = (B^p)^+(\mu - \beta^p)$.

Still, $(\mathcal{M}_{\mu_0}^p)^-$ is not guaranteed to have a non-empty interior. Suppose that some state x^* is not reachable at a certain step n^* ; that is, for every state x and action a , $p_{n^*}(x^*|x, a) = 0$ if $n^* \geq 2$, or just that $\rho_1^\mu(x^*) = 0$ if $n^* = 1$. Then, for any $\mu \in \mathcal{M}_{\mu_0}^p$, $\mu_{n^*}(x^*, a) = 0$ for every action a . This implies that for every $\xi \in (\mathcal{M}_{\mu_0}^p)^-$, $\xi_{n^*}(x^*, a) = 0$ for all $a \neq a^*$ (since these coordinates are preserved under Λ_p), and hence, $(\mathcal{M}_{\mu_0}^p)^-$ has an empty interior. To remedy this, we rely on Asm. 4.4, which is equivalent to imposing that for every state x , $\rho_1^\mu(x) > 0$ and there exists for every step n a state-action pair (x', a') such that $p_{n+1}(x|x', a') > 0$. We show next that this condition is sufficient for $(\mathcal{M}_{\mu_0}^p)^-$ to have a non-empty interior. We first present an alternative characterization of $(\mathcal{M}_{\mu_0}^p)^-$.

Lemma E.1. *It holds that*

$$(\mathcal{M}_{\mu_0}^p)^- = \{ \xi \in \mathbb{R}^{NS(A-1)} : B^p \xi \geq -\beta^p \}.$$

1760 Hence, $(\mathcal{M}_{\mu_0}^p)^-$ is a polyhedral set formed by NSA constraints; namely that for $n, x, a \in [N] \times \mathcal{X} \times \mathcal{A}$,
 1761 $B^p(n, x, a, \cdot, \cdot, \cdot)^\top \xi + \beta_n^p(x, a) \geq 0$.

1762
 1763 *Proof.* Any $\xi \in (\mathcal{M}_{\mu_0}^p)^-$ clearly satisfies $B^p \xi \geq -\beta^p$ since $B^p \xi + \beta^p = \Lambda_p(\xi) \in \mathcal{M}_{\mu_0}^p$, whose coordinates are non-
 1764 negative. Conversely, assume that $B^p \xi \geq -\beta^p$ for some $\xi \in \mathbb{R}^{NS(A-1)}$, and let $\mu := \Xi_p(\xi) = B^p \xi + \beta^p$. Showing that
 1765 $\xi \in (\mathcal{M}_{\mu_0}^p)^-$ is equivalent, by definition, to showing that $\Xi_p(\xi) \in \mathcal{M}_{\mu_0}^p$. Since $\beta^p \in \mathcal{M}_{\mu_0}^p$ and $B^p \xi$ belongs to the column
 1766 space of B^p specified in (49), it holds that

$$1767 \quad \sum_a \mu_n(x, a) = \sum_{x', a'} \mu_{n-1}(x', a') p_n(x|x', a')$$

1771 for every $n \geq 2$, and that $\sum_a \mu_1(x, a) = \sum_a \beta_1^p(x, a) = \rho_1^{\beta^p}(x) = \rho_1^p(x)$. Then, to show that $\mu \in \mathcal{M}_{\mu_0}^p$, it remains to
 1772 show that $\mu \in (\Delta_{\mathcal{X} \times \mathcal{A}})^N$. By assumption, μ only has non-negative coordinates; therefore, we only have to show that
 1773 $\sum_{x,a} \mu_n(x, a) = 1$ at every n . This easily done via induction: $\sum_{x,a} \mu_1(x, a) = \sum_x \rho_1^p(x) = 1$, and for $n \geq 2$,

$$1774 \quad \sum_{x,a} \mu_n(x, a) = \sum_x \sum_{x', a'} \mu_{n-1}(x', a') p_n(x|x', a') = \sum_{x', a'} \mu_{n-1}(x', a').$$

1777 □

1779 **Lemma E.2.** $(\mathcal{M}_{\mu_0}^p)^-$ has a non-empty interior if and only if Asm. 4.4 holds.

1781 *Proof.* Necessity is immediate as argued before. We prove sufficiency utilizing an argument from the proof of Proposition
 1782 2.3 in (Wolsey & Nemhauser, 1999). For every step-state-action triple (n, x, a) , it is easy to verify that Asm. 4.4 implies the
 1783 existence of some $\mu \in \mathcal{M}_{\mu_0}^p$ such that $\mu_n(x, a) > 0$. Taking a convex combination with full support of one such occupancy
 1784 measure for every (n, x, a) results, via the convexity of $\mathcal{M}_{\mu_0}^p$, in an occupancy measure $\mu^* \in \mathcal{M}_{\mu_0}^p$ whose entries are all
 1785 strictly positive. Hence, $\xi^* := \Lambda_p^{-1}(\mu^*)$ is an interior point of the polyhedral set $(\mathcal{M}_{\mu_0}^p)^-$ as it satisfies with strict inequality
 1786 all the constraints defining it. □

1788 E.2. Entropic Regularization Approach

1790 E.2.1. FITTING A EUCLIDEAN BALL IN THE CONSTRAINT SET

1791 For the following, fix $\varepsilon \in (0, 1/S)$. From Sec. 4.2.1, recall the definition $\kappa := \varepsilon/(A-1 + \sqrt{A-1})$. We now show that
 1792 $\kappa \mathbb{1}_{NS(A-1)} + \kappa \mathbf{v} \in (\mathcal{M}_{\mu_0}^p)^-$ for any $\mathbf{v} \in \mathbb{B}^{NS(A-1)}$ assuming the transition kernel $p := (p_n)_{n \in [N]}$ satisfies the condition
 1793 of Asm. 4.2; that is, $p_n(x'|x, a) \geq \varepsilon$ for all $(n, x, x', a) \in [N] \times \mathcal{X}^2 \times \mathcal{A}$. Take $\zeta^{v,p} := \Xi_p(\kappa \mathbb{1}_{NS(A-1)} + \kappa \mathbf{v})$. Note that
 1794 showing that $\kappa \mathbb{1}_{NS(A-1)} + \kappa \mathbf{v} \in (\mathcal{M}_{\mu_0}^p)^-$ is equivalent to showing that $\zeta^{v,p} \in \mathcal{M}_{\mu_0}^p$. In the following, we proceed with the
 1795 latter.

1797 Note that via Lem. E.1, it suffices to show that $\zeta^{v,p}$ is non-negative. We use induction in the following to show more
 1798 particularly that $\zeta^{v,p} \in (\Delta_{\mathcal{X} \times \mathcal{A}})^N$. By the definition of $\zeta^{v,p}$, we have that for $(n, x) \in [N] \times \mathcal{X}$,

$$1800 \quad \zeta_n^{v,p}(x, a) = \frac{\varepsilon}{A-1 + \sqrt{A-1}} (1 + \mathbf{v}_n(x, a)) \quad \forall a \in \mathcal{A} \setminus a^* \text{ and } \zeta_n^{v,p}(x, a^*) = \rho_n^{\zeta^{v,p}}(x) - \sum_{a \neq a^*} \zeta_n^{v,p}(x, a), \quad (50)$$

1803 where $\rho_n^{\zeta^{v,p}}(x) = \sum_{a', x' \in \mathcal{A} \times \mathcal{X}} \zeta_{n-1}^{v,p}(x', a') p_n(x|x', a')$ for $n \geq 2$ and $\rho_1^{\zeta^{v,p}}(x) = \sum_{a', x' \in \mathcal{A} \times \mathcal{X}} \mu_0(x', a') p_1(x|x', a') =$
 1804 $\rho_1^p(x)$. For $a \neq a^*$, clearly $\zeta_n^{v,p}(x, a) \geq 0$ as $\mathbf{v}_n(x, a) \geq -1$. Note that at any step n and state x , the Cauchy-Schwarz
 1805 inequality and the fact that $\mathbf{v} \in \mathbb{B}^{NS(A-1)}$ yield that

$$1807 \quad \sum_{a \neq a^*} \mathbf{v}_n(x, a) \leq \sqrt{A-1} \sqrt{\sum_{a \neq a^*} |\mathbf{v}_n(x, a)|^2} \leq \sqrt{A-1}.$$

1810 Hence,

$$1812 \quad \sum_{a \neq a^*} \zeta_n^{v,p}(x, a) = \sum_{a \neq a^*} \frac{\varepsilon}{A-1 + \sqrt{A-1}} (1 + \mathbf{v}_n(x, a)) \leq \varepsilon.$$

1815 On the other hand, Asm. 4.2 implies that $\rho_1^{\zeta^{v,p}}(x) \geq \varepsilon$ for every x . Hence, (50) gives that $\zeta_1^{v,p}(x, a^*) \geq 0$ at every x .
 1816 Moreover, (50) also implies that $\sum_{x,a} \zeta_1^{v,p}(x, a) = \sum_x \rho_1^p(x) = 1$, yielding that $\zeta_1^{v,p} \in \Delta_{\mathcal{X} \times \mathcal{A}}$. For $n \geq 2$, assuming
 1817 that $\zeta_{n-1}^{v,p} \in \Delta_{\mathcal{X} \times \mathcal{A}}$, Asm. 4.2 implies again that $\rho_n^{\zeta^{v,p}}(x) \geq \varepsilon$ for every x . We then get via (50) that $\zeta_n^{v,p}(x, a^*) \geq 0$
 1818 and that $\zeta_n^{v,p} \in \Delta_{\mathcal{X} \times \mathcal{A}}$ since $\sum_{x,a} \zeta_n^{v,p}(x, a) = \sum_x \rho_n^{\zeta^{v,p}}(x) = 1$, which holds again via (50) and the assumption that
 1819 $\zeta_{n-1}^{v,p} \in \Delta_{\mathcal{X} \times \mathcal{A}}$. Induction then establishes that $\zeta^{v,p} \in (\Delta_{\mathcal{X} \times \mathcal{A}})^N$ as sought. As mentioned above, this implies via Lem. E.1
 1820 that $\zeta^{v,p} \in \mathcal{M}_{\mu_0}^p$, or equivalently, that $\kappa \mathbb{1}_{NS(A-1)} + \kappa \mathbf{v} \in (\mathcal{M}_{\mu_0}^p)^-$ and $\zeta^{v,p} = \Lambda_p(\kappa \mathbb{1}_{NS(A-1)} + \kappa \mathbf{v})$.
 1821

1822 E.2.2. ESTIMATING THE TRANSITION KERNEL

1824 In this section, we define and analyze an alternative transition kernel estimator to the one given in Eq. (7). What we seek in
 1825 this new estimator is (1) that it estimates well the true transition kernel, with a guarantee similar to that of Lem. 2.2; (2) that
 1826 it drifts across rounds in a controlled manner, satisfying the bound of Lem. A.3 up to a constant; and (3) that, at the same
 1827 time, it satisfies the condition of Asm. 4.2 almost surely, supposing, naturally, that it is satisfied by the true kernel.

1828 To recall the notation, for each round $t \in [T]$, o^t denotes a random trajectory obtained by executing the policy π^t in the
 1829 environment; that is, $o^t := (x_1^t, a_1^t, \dots, x_N^t, a_N^t)$ where $a_n^t \sim \pi^t(\cdot | x_n^t)$ and $x_n^t \sim p_n(\cdot | x_{n-1}^t, a_{n-1}^t)$.² We also recall the
 1830 definitions

$$1831 \quad N_n^t(x, a) := \sum_{s=1}^{t-1} \mathbb{1}_{\{x_n^s = x, a_n^s = a\}} \quad \text{and} \quad M_n^t(x' | x, a) := \sum_{s=1}^{t-1} \mathbb{1}_{\{x_{n+1}^s = x', x_n^s = x, a_n^s = a\}}.$$

1835 Fix $n, x, a \in [N] \times \mathcal{X} \times \mathcal{A}$. As an intermediate step, we compute at the beginning of each round t the Laplace (add-one)
 1836 estimator for $p_n^t(\cdot | x, a)$; that is, for $x' \in \mathcal{X}$,

$$1837 \quad \tilde{p}_n^t(x' | x, a) := \frac{M_{n-1}^t(x' | x, a) + 1}{N_{n-1}^t(x, a) + S}. \quad (51)$$

1841 To obtain a guarantee on the accuracy of this estimator, we firstly describe a slightly different setting. Let $(\hat{x}_n^s)_{s=1}^T$ be an
 1842 i.i.d. sequence of states such that $\hat{x}_n^s \sim p_n(\cdot | x, a)$. Then, for $k \in [T]$, we define the Laplace estimator

$$1843 \quad \hat{p}_n^k(x' | x, a) := \frac{1 + \sum_{s=1}^k \mathbb{1}_{\{\hat{x}_n^s = x'\}}}{k + S}.$$

1847 Notice that in our setting, the distribution $\tilde{p}_n^t(\cdot | x, a)$ is equivalent to $\hat{p}_n^{N_{n-1}^t(x, a)}(\cdot | x, a)$, keeping in mind that the number of
 1848 samples $N_{n-1}^t(x, a)$ is random and dependent on the observed samples. Let $D_{\text{KL}}(p \| q)$ denote the KL-divergence between
 1849 distributions (probability mass functions) p and q . We derive the following result concerning the divergence between \tilde{p}^t and
 1850 p using known properties of the Laplace estimator and a union bound argument.

1851 **Lemma E.3.** For fixed $t, n, x, a \in [T] \times [N] \times \mathcal{X} \times \mathcal{A}$, it holds with probability at least $1 - \delta$ that

$$1852 \quad D_{\text{KL}}(p_n(\cdot | x, a) \| \tilde{p}_n^t(\cdot | x, a)) \leq \frac{161S + 6\sqrt{S} \log^{5/2} \frac{ST}{4\delta} + 310}{\max\{1, N_{n-1}^t(x, a)\}}.$$

1857 *Proof.* For a fixed $k \in [T]$, Thm. 2 in (Canonne et al., 2023) and Prop. 1 in (Mourtada & Gaïffas, 2022) imply that

$$1858 \quad P\left(D_{\text{KL}}(p_n(\cdot | x, a) \| \hat{p}_n^k(\cdot | x, a)) > \frac{161S + 6\sqrt{S} \log^{5/2} \frac{S}{4\delta} + 310}{k}\right) \leq \delta.$$

1861 Via a union bound, we obtain that³

$$1862 \quad P\left(D_{\text{KL}}(p_n(\cdot | x, a) \| \tilde{p}_n^t(\cdot | x, a)) > \frac{161S + 6\sqrt{S} \log^{5/2} \frac{S}{4\delta} + 310}{\max\{1, N_{n-1}^t(x, a)\}}\right)$$

1866 ²Recall that $(x_0^t, a_0^t) \sim \mu_0(\cdot, \cdot)$.

1867 ³Note that if $N_{n-1}^t(x, a) = 0$, then $\tilde{p}_n^t(\cdot | x, a)$ is the uniform distribution and $D_{\text{KL}}(p_n(\cdot | x, a) \| \tilde{p}_n^t(\cdot | x, a)) \leq \log S$; hence, the
 1868 bound trivially holds.

$$\leq P\left(\exists k \in [T]: D_{\text{KL}}(p_n(\cdot|x, a) \parallel \hat{p}_n^k(\cdot|x, a)) > \frac{161S + 6\sqrt{S} \log^{5/2} \frac{S}{4\delta} + 310}{k}\right) \leq \delta T.$$

The lemma then follows after rescaling δ . \square

Note that the distribution $\tilde{p}_n^t(\cdot|x, a)$ does not necessarily satisfy the conditions of Asm. 4.2 uniformly. Next, we define for a given $\varepsilon \in [0, 1/S]$ the set

$$\Delta_{\mathcal{X}}^\varepsilon := \{x \in \mathbb{R}^d : \mathbb{1}_d^\top x = 1 \text{ and } x_i \geq \varepsilon \forall i \in [d]\} \subseteq \Delta_{\mathcal{X}},$$

which is the set of state distribution assigning probability at least ε to every state. We then define $\hat{p}_n^t(\cdot|x, a)$ as the information projection of $\tilde{p}_n^t(\cdot|x, a)$ onto $\Delta_{\mathcal{X}}^\varepsilon$; that is,

$$\hat{p}_n^t(\cdot|x, a) := \arg \min_{q \in \Delta_{\mathcal{X}}^\varepsilon} D_{\text{KL}}(q \parallel \tilde{p}_n^t(\cdot|x, a)), \quad (52)$$

which exists and is unique since $\Delta_{\mathcal{X}}^\varepsilon$ is compact and $D_{\text{KL}}(\cdot \parallel \tilde{p}_n^t(\cdot|x, a))$ is continuous and strictly convex where it is finite (note that $\tilde{p}_n^t(\cdot|x, a)$ never assigns zero probability to any state; hence, $D_{\text{KL}}(q \parallel \tilde{p}_n^t(\cdot|x, a))$ is finite for any $q \in \Delta_{\mathcal{X}}^\varepsilon$). If $\tilde{p}_n^t(\cdot|x, a)$ is not already in $\Delta_{\mathcal{X}}^\varepsilon$, this projection can only bring us closer to $p_n(\cdot|x, a)$ in the KL -divergence sense as the following inequality (Cover & Thomas, 2006, Thm. 11.6.1) states:

$$D_{\text{KL}}(p_n(\cdot|x, a) \parallel \hat{p}_n^t(\cdot|x, a)) \leq D_{\text{KL}}(p_n(\cdot|x, a) \parallel \tilde{p}_n^t(\cdot|x, a)) - D_{\text{KL}}(\hat{p}_n^t(\cdot|x, a) \parallel \tilde{p}_n^t(\cdot|x, a)). \quad (53)$$

With this fact in mind, we can arrive at the following result, a parallel of Lem. 2.2.

Lemma E.4. *With probability at least $1 - \delta$, it holds for all $t, n, x, a \in [T] \times [N] \times \mathcal{X} \times \mathcal{A}$ simultaneously that*

$$\|p_n(\cdot|x, a) - \hat{p}_n^t(\cdot|x, a)\|_1 \leq \sqrt{\frac{322S + 12\sqrt{S} \log^{5/2} \frac{S^2 AN T^2}{4\delta} + 620}{\max\{1, N_{n-1}^t(x, a)\}}}.$$

Proof. The statement is a consequence of (53), Lem. E.3, and Pinsker's inequality; followed by an application of a union bound over all rounds, steps, and state-action pairs. \square

What remains now is to show that there exists a constant $c > 0$ such that

$$\|\hat{p}_{n+1}^{t+1}(\cdot|x, a) - \hat{p}_{n+1}^t(\cdot|x, a)\|_1 \leq c \frac{\mathbb{1}_{\{x_n^t = x, a_n^t = a\}}}{\max\{1, N_n^{t+1}(x, a)\}}.$$

This can be easily shown to hold for \tilde{p}^t , i.e., before the projection step, as states the following lemma.

Lemma E.5. *For all $n \in [N - 1]$, $(x, a, x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}$, and $t \in [T]$; $\tilde{p}_{n+1}^t(x'|x, a)$ as defined in (51) satisfies*

$$\|\tilde{p}_{n+1}^{t+1}(\cdot|x, a) - \tilde{p}_{n+1}^t(\cdot|x, a)\|_1 \leq \frac{2\mathbb{1}_{\{x_n^t = x, a_n^t = a\}}}{N_n^{t+1}(x, a) + S}.$$

Proof. The derivation follows along the same lines as the proof of Lem. A.3. We have that

$$\begin{aligned} \tilde{p}_{n+1}^{t+1}(x'|x, a) &= \frac{\mathbb{1}_{\{x_{n+1}^t = x', x_n^t = x, a_n^t = a\}} + M_n^t(x'|x, a) + 1}{N_n^{t+1}(x, a) + S} \\ &= \frac{\mathbb{1}_{\{x_{n+1}^t = x', x_n^t = x, a_n^t = a\}}}{N_n^{t+1}(x, a) + S} + \frac{N_n^t(x, a) + S}{N_n^{t+1}(x, a) + S} \tilde{p}_{n+1}^t(x'|x, a). \end{aligned}$$

Hence,

$$\begin{aligned} \tilde{p}_{n+1}^{t+1}(x'|x, a) - \tilde{p}_{n+1}^t(x'|x, a) &= \frac{\mathbb{1}_{\{x_{n+1}^t = x', x_n^t = x, a_n^t = a\}}}{N_n^{t+1}(x, a) + S} + \tilde{p}_{n+1}^t(x'|x, a) \frac{N_n^t(x, a) - N_n^{t+1}(x, a)}{N_n^{t+1}(x, a) + S} \\ &= \frac{\mathbb{1}_{\{x_n^t = x, a_n^t = a\}}}{N_n^{t+1}(x, a) + S} (\mathbb{1}_{\{x_{n+1}^t = x'\}} - \tilde{p}_{n+1}^t(x'|x, a)). \end{aligned}$$

1925 Finally, we conclude that

$$\begin{aligned}
 1926 \quad & \|\tilde{p}_{n+1}^{t+1}(\cdot|x, a) - \tilde{p}_{n+1}^t(\cdot|x, a)\|_1 = 2 \sum_{x': \tilde{p}_{n+1}^{t+1}(x'|x, a) \geq \tilde{p}_{n+1}^t(x'|x, a)} (\tilde{p}_{n+1}^{t+1}(x'|x, a) - \tilde{p}_{n+1}^t(x'|x, a)) \\
 1927 \quad & = \frac{2\mathbb{1}_{\{x_n^t=x, a_n^t=a\}}}{N_n^{t+1}(x, a) + S} (1 - \tilde{p}_{n+1}^t(x_{n+1}^t|x, a)) \leq \frac{2\mathbb{1}_{\{x_n^t=x, a_n^t=a\}}}{N_n^{t+1}(x, a) + S}. \\
 1928 \quad & \\
 1929 \quad & \\
 1930 \quad & \\
 1931 \quad & \\
 1932 \quad & \square
 \end{aligned}$$

1934 To derive a similar bound for the projected estimator \hat{p}^t , we firstly derive a more explicit characterization of the information
 1935 projection onto $\Delta_{\mathcal{X}}^\varepsilon$. For a fixed $\varepsilon \in (0, 1/S)$, define the function $g_\varepsilon: \mathbb{R} \times \Delta_{\mathcal{X}} \rightarrow \mathbb{R}$ as

$$1936 \quad g_\varepsilon(r; p) := \sum_{x \in \mathcal{X}} \max\{r\varepsilon, p(x)\}.$$

1937 **Lemma E.6.** *For any given $p \in \Delta_{\mathcal{X}}$, the map $r \mapsto g_\varepsilon(r; p)$ is εS -Lipschitz and has a unique fixed point. Moreover, denoting
 1938 this fixed point by r^* , it holds that $r^* \in [1, \max_x p(x)/\varepsilon]$, and that $g_\varepsilon(r; p) > r$ for $r < r^*$ and $g_\varepsilon(r; p) < r$ for $r > r^*$.*

1939 *Proof.* We firstly note that $g_\varepsilon(\cdot; p)$ can be easily verified to be convex. For any $r \in \mathbb{R}$ and any subgradient h of $g_\varepsilon(\cdot; p)$ at r ,
 1940 it holds that $|h| \leq \varepsilon S$. Hence, the convexity of $g_\varepsilon(\cdot; p)$ implies that $|g_\varepsilon(r; p) - g_\varepsilon(r'; p)| \leq \varepsilon S|r - r'|$ for any $r, r' \in \mathbb{R}$, or
 1941 that $g_\varepsilon(\cdot; p)$ is εS -Lipschitz. This implies, since $\varepsilon S < 1$ by assumption, that $g_\varepsilon(\cdot; p)$ is a contraction mapping; hence, via
 1942 Banach's fixed point theorem, it admits a unique fixed point $r^* \in \mathbb{R}$. For $r < 1$, it holds that $g_\varepsilon(r; p) \geq \sum_x p(x) = 1 > r$.
 1943 While for $r \geq \max_x p(x)/\varepsilon$, $g_\varepsilon(r; p) = r\varepsilon S < r$. Therefore, $r^* \in [1, \max_x p(x)/\varepsilon]$. Moreover, for any $r < r^*$ ($r > r^*$),
 1944 it must hold that $g_\varepsilon(r; p) > r$ ($g_\varepsilon(r; p) < r$); as otherwise, the intermediate value theorem, applied to $g_\varepsilon(r; p) - r$, would
 1945 imply the existence of another fixed point, a contradiction. \square

1946 Next, we define $r_\varepsilon: \Delta_{\mathcal{X}} \rightarrow \mathbb{R}$ as the function that maps a distribution $p \in \Delta_{\mathcal{X}}$ to the fixed point of $g_\varepsilon(\cdot; p)$. This function is
 1947 well-defined as implied by Lem. E.6. We now show that the solution of the information projection problem onto $\Delta_{\mathcal{X}}^\varepsilon$ can be
 1948 expressed in terms of the function r_ε . For $p \in \Delta_{\mathcal{X}}$, we define $p_\varepsilon \in \Delta_{\mathcal{X}}^\varepsilon$ as

$$1949 \quad p_\varepsilon(x) := \frac{\max\{r_\varepsilon(p)\varepsilon, p(x)\}}{\sum_{x' \in \mathcal{X}} \max\{r_\varepsilon(p)\varepsilon, p(x')\}} = \max\{\varepsilon, p(x)/r_\varepsilon(p)\}.$$

1950 **Lemma E.7.** *For $p \in \Delta_{\mathcal{X}}$, it holds that $p_\varepsilon = \arg \min_{q \in \Delta_{\mathcal{X}}^\varepsilon} D_{\text{KL}}(q \| p)$.*

1951 *Proof.* We assume without loss of generality that $p(x) > 0$ for all $x \in \mathcal{X}$; as otherwise, we can cast the problem into a
 1952 lower dimensional one considering only the elements $x \in \mathcal{X}$ for which $p(x) > 0$. Since the constraint set is compact and
 1953 the objective is continuous and strictly convex, this minimization problem admits a unique optimal solution. We start by
 1954 rewriting the problem as

$$\begin{aligned}
 1955 \quad & \min_{q \in \mathbb{R}^S} \sum_{x \in \mathcal{X}} q(x) \log \frac{q(x)}{p(x)} \\
 1956 \quad & \text{subject to} \quad \varepsilon - q(x) \leq 0 \quad \forall x \in \mathcal{X} \\
 1957 \quad & \sum_{x \in \mathcal{X}} q(x) - 1 = 0
 \end{aligned}$$

1958 Define the Lagrangian

$$1959 \quad L(q, u, v) := \sum_{x \in \mathcal{X}} q(x) \log \frac{q(x)}{p(x)} + \sum_{x \in \mathcal{X}} u(x)(\varepsilon - q(x)) + v \left(\sum_{x \in \mathcal{X}} q(x) - 1 \right)$$

1960 for $v \in \mathbb{R}$ and $u \in \mathbb{R}_{\geq 0}^S$. We have that

$$1961 \quad \frac{\partial L}{\partial q(x)}(q, u, v) = \log \frac{q(x)}{p(x)} + 1 - u(x) + v.$$

We now show that we can satisfy the KKT conditions by choosing a solution pair q^* and u^*, v^* where

$$q^*(x) := p_\varepsilon(x) = \max\{\varepsilon, p(x)/r_\varepsilon(p)\}, \quad u^*(x) := \log \frac{\max\{\varepsilon, p(x)/r_\varepsilon(p)\}}{p(x)/r_\varepsilon(p)}, \quad \text{and} \quad v^* := -1 + \log r_\varepsilon(p).$$

Firstly, q^* indeed belongs to $\Delta_\lambda^\varepsilon$ by the definition of p_ε , and u^* is non-negative. Moreover, whenever $q^*(x) > \varepsilon$, we get that $u^*(x) = 0$; hence, complementary slackness holds. Finally,

$$\frac{\partial L}{\partial q(x)}(q^*, u^*, v^*) = \log \frac{\max\{\varepsilon, p(x)/r_\varepsilon(p)\}}{p(x)} + 1 - \log \frac{\max\{\varepsilon, p(x)/r_\varepsilon(p)\}}{p(x)/r_\varepsilon(p)} - 1 + \log r_\varepsilon(p) = 0.$$

Therefore, we conclude that p_ε is the optimal solution. \square

Computing p_ε , or the information projection of p onto $\Delta_\lambda^\varepsilon$, can be performed efficiently. In particular, the following characterization implies that $r_\varepsilon(p)$ can be computed exactly in a finite number of steps by iterating over the set of states.

Lemma E.8. *Let $\mathcal{X}_p^+ := \{x \in \mathcal{X} : g_\varepsilon(p(x)/\varepsilon; p) < p(x)/\varepsilon\}$ and $\mathcal{X}_p^- := \mathcal{X} \setminus \mathcal{X}_p^+$. Then,*

$$r_\varepsilon(p) = \frac{\sum_{x \in \mathcal{X}^+} p(x)}{1 - \varepsilon |\mathcal{X}_p^-|}.$$

Proof. As stated in the proof of Lem. E.6, for $r \geq \max_{x \in \mathcal{X}} p(x)/\varepsilon$, $g_\varepsilon(r; p) = r\varepsilon S < r$; hence \mathcal{X}_p^+ is non-empty as it at least includes $\arg \max_{x \in \mathcal{X}} p(x)$. Moreover, from the same lemma, we have that $r_\varepsilon(p) < \min_{x \in \mathcal{X}_p^+} p(x)/\varepsilon$ and $r_\varepsilon(p) \geq \max_{x \in \mathcal{X}_p^-} p(x)/\varepsilon$ (if \mathcal{X}_p^- is non-empty). Therefore,

$$r_\varepsilon(p) = g_\varepsilon(r_\varepsilon(p); p) = \sum_{x \in \mathcal{X}} \max\{r_\varepsilon(p)\varepsilon, p(x)\} = r_\varepsilon(p)\varepsilon |\mathcal{X}_p^-| + \sum_{x \in \mathcal{X}^+} p(x).$$

The previous lemma also implies that $r_\varepsilon(p) \leq (1 - \varepsilon S)^{-1}$. Returning back to our original objective, we show next that $\|p_\varepsilon - q_\varepsilon\|_1$ is no larger than a constant multiple of $\|p - q\|_1$ for any two distributions p and q . Towards that end, we first show that r_ε is Lipschitz continuous.

Lemma E.9. *For $\varepsilon \leq \frac{1}{2S}$, the function r_ε is 1-Lipschitz with respect to the $\|\cdot\|_1$ norm; that is,*

$$|r_\varepsilon(p) - r_\varepsilon(q)| \leq \|p - q\|_1$$

for any $p, q \in \Delta_\mathcal{X}$.

Proof. Note that, for any fixed $r \in \mathbb{R}$, $g_\varepsilon(r; \cdot)$ is convex; and that for any $p \in \Delta_\mathcal{X}$ and subgradient k of $g_\varepsilon(r; \cdot)$ at p , it holds that k is non-negative and satisfies $\|k\|_\infty \leq 1$. Hence, for any $p, q \in \Delta_\mathcal{X}$,

$$|g_\varepsilon(r; p) - g_\varepsilon(r; q)| \leq \sum_{x: p(x) > q(x)} (p(x) - q(x)) = \sum_{x: q(x) > p(x)} (q(x) - p(x)) = \frac{1}{2} \|p - q\|_1. \quad (54)$$

Then, we obtain that

$$\begin{aligned} |r_\varepsilon(p) - r_\varepsilon(q)| &= |g_\varepsilon(r_\varepsilon(p); p) - g_\varepsilon(r_\varepsilon(q); q)| \\ &\leq |g_\varepsilon(r_\varepsilon(p); p) - g_\varepsilon(r_\varepsilon(q); p)| + |g_\varepsilon(r_\varepsilon(q); p) - g_\varepsilon(r_\varepsilon(q); q)| \\ &\leq \varepsilon S |r_\varepsilon(p) - r_\varepsilon(q)| + \frac{1}{2} \|p - q\|_1 \leq \frac{1}{2} |r_\varepsilon(p) - r_\varepsilon(q)| + \frac{1}{2} \|p - q\|_1, \end{aligned}$$

where the second inequality follows from (54) and Lem. E.6, and the last inequality holds since $\varepsilon \leq \frac{1}{2S}$. The lemma then follows after rearranging the last result. \square

Lemma E.10. *Assuming $\varepsilon \leq \frac{1}{2S}$, it holds for any $p, q \in \Delta_\mathcal{X}$ that $\|p_\varepsilon - q_\varepsilon\|_1 \leq \frac{5}{2} \|p - q\|_1$.*

2035 *Proof.* We have that

$$\begin{aligned}
 2036 \quad q_\varepsilon(x) &= \frac{\max\{r_\varepsilon(q)\varepsilon, q(x)\}}{r_\varepsilon(q)} \\
 2037 \quad &= \frac{\max\{r_\varepsilon(q)\varepsilon, q(x)\} - \max\{r_\varepsilon(p)\varepsilon, p(x)\} + \max\{r_\varepsilon(p)\varepsilon, p(x)\}}{r_\varepsilon(q)} \\
 2038 \quad &= \frac{\max\{r_\varepsilon(q)\varepsilon, q(x)\} - \max\{r_\varepsilon(p)\varepsilon, p(x)\}}{r_\varepsilon(q)} + \frac{r_\varepsilon(p)}{r_\varepsilon(q)} p_\varepsilon(x).
 \end{aligned}$$

2044 Then,

$$2045 \quad q_\varepsilon(x) - p_\varepsilon(x) = \frac{1}{r_\varepsilon(q)} (\max\{r_\varepsilon(q)\varepsilon, q(x)\} - \max\{r_\varepsilon(p)\varepsilon, p(x)\} + (r_\varepsilon(p) - r_\varepsilon(q))p_\varepsilon(x)).$$

2048 Using Lem. E.9 and the fact that

$$\begin{aligned}
 2050 \quad |\max\{r_\varepsilon(q)\varepsilon, q(x)\} - \max\{r_\varepsilon(p)\varepsilon, p(x)\}| &\leq \max\{\varepsilon|r_\varepsilon(q) - r_\varepsilon(p)|, |q(x) - p(x)|\} \\
 2051 \quad &\leq \varepsilon|r_\varepsilon(q) - r_\varepsilon(p)| + |q(x) - p(x)|,
 \end{aligned}$$

2053 we obtain that

$$\begin{aligned}
 2054 \quad \|p_\varepsilon - q_\varepsilon\|_1 &= \sum_x |q_\varepsilon(x) - p_\varepsilon(x)| \\
 2055 \quad &\leq \frac{1}{r_\varepsilon(q)} \sum_x (\varepsilon|r_\varepsilon(q) - r_\varepsilon(p)| + |q(x) - p(x)| + p_\varepsilon(x)|r_\varepsilon(p) - r_\varepsilon(q)|) \\
 2056 \quad &\leq \frac{\|p - q\|_1}{r_\varepsilon(q)} (2 + \varepsilon S) \leq \frac{5}{2} \|p - q\|_1,
 \end{aligned}$$

2062 where the last step uses that $\varepsilon S \leq 1/2$ and $r_\varepsilon(q) \geq 1$. □

2064 Finally, we arrive at the sought result, a parallel of Lem. A.3.

2065 **Lemma E.11.** For all $n \in [N - 1]$, $(x, a, x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}$, and $t \in [T]$; $\hat{p}_{n+1}^t(x'|x, a)$ as defined in (52) with $\varepsilon \leq \frac{1}{2S}$ satisfies

$$2068 \quad \|\hat{p}_{n+1}^{t+1}(\cdot|x, a) - \hat{p}_{n+1}^t(\cdot|x, a)\|_1 \leq \frac{5 \mathbb{1}_{\{x_n^t = x, a_n^t = a\}}}{N_n^{t+1}(x, a) + S}.$$

2071 *Proof.* This is a direct consequence of the definition in (52) and Lems. E.5, E.7 and E.10. □

2073 E.2.3. THE ALGORITHM

2074 For $\delta \in (0, 1)$, define

$$2076 \quad C'_\delta := \sqrt{322S + 12\sqrt{S} \log^{5/2} \frac{S^2 ANT^2}{4\delta} + 620}, \quad (55)$$

2078 which is the leading factor in the confidence bound of Lem. E.4. For the purpose of exploration, much like the full information case, we will utilize at each round t a bonus reward vector $b^t \in \mathbb{R}^{NSA}$ to be subtracted from the estimated gradient, where

$$2082 \quad b_n^t(x, a) := L(N - n) \frac{C'_{1/T}}{\sqrt{\max\{1, N_n^t(x, a)\}}} \quad (56)$$

2085 for $(t, n, x, a) \in [T] \times (\{0\} \cup [N]) \times \mathcal{X} \times \mathcal{A}$.

2086 Finally, with all its components detailed, we present Alg. 3, our first approach for CURL with bandit feedback. As mentioned in Sec. 4.2.1, the main changes compared to Alg. 1 are the use of spherical estimation to obtain a surrogate for the gradient and the use of a suitably altered transition kernel estimator.

Algorithm 3 Bonus O-MD-CURL (bandit feedback)

input: learning rate $\tau > 0$, perturbation rate $\delta \in (0, 1]$, sequence of exploration parameters $(\alpha_t)_{t \in [T]} \in (0, 1)^T$

initialization: $\hat{p}_n^1(x'|x, a) \leftarrow 1/S \forall (n, x, x', a)$, $\mu^1 \in \arg \min_{\mu \in \mathcal{M}_{\mu_0}^{\hat{p}_0^1}} \psi(\mu)$

for $t = 1, \dots, T$ **do**

draw $\mathbf{u}^t \in \mathbb{S}^{NS(A-1)}$ uniformly at random

$\zeta^t \leftarrow \zeta^{\mathbf{u}^t, \hat{p}^t} = \Lambda_{\hat{p}^t}(\kappa \mathbb{1}_{NS(A-1)} + \kappa \mathbf{u}^t)$

$\hat{\mu}^t \leftarrow (1 - \delta)\mu^t + \delta\zeta^t$

$\pi_n^t(a|x) \leftarrow \hat{\mu}^t(x, a) / \sum_{a \in \mathcal{A}} \hat{\mu}^t(x, a)$

execute π^t and observe $F^t(\mu^{\pi^t, p})$ and a sampled trajectory $o^t := (x_1^t, a_1^t, \dots, x_N^t, a_N^t)$

$g^t \leftarrow \frac{1-\delta}{\delta\kappa} NS(A-1) F^t(\mu^{\pi^t, p}) \mathbf{u}^t$

construct $\hat{g}^t \in \mathbb{R}^{NSA}$ as $\hat{g}_n^t(x, a) \leftarrow g_n^t(x, a)$ for $a \neq a^*$ and $\hat{g}_n^t(x, a^*) \leftarrow 0$

construct bonus vector b^t as in (56)

$\tilde{\pi}_n^t(a|x) \leftarrow (1 - \alpha_t)\mu^t(x, a) / \sum_{a \in \mathcal{A}} \mu^t(x, a) + \alpha_t/A$

construct the new estimated kernel \hat{p}^{t+1} via (51) and (52)

set $\mu^{t+1} \in \arg \min_{\mu \in \mathcal{M}_{\mu_0}^{\hat{p}^{t+1}}} \tau \langle \hat{g}^t - b^t, \mu \rangle + \Gamma(\mu, \mu^{\tilde{\pi}^t, \hat{p}^t})$

end for

E.2.4. AUXILIARY LEMMAS

Lemma E.12. For $0 < \delta < 1$ and \hat{p}^t as defined in (52), it holds that

$$\sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x, a} \mu_i^{\pi^t, p}(x, a) \|p_{i+1}(\cdot|x, a) - \hat{p}_{i+1}^t(\cdot|x, a)\|_1 \leq 3C'_\delta N^2 \sqrt{SAT} + 2SN^2 \sqrt{2T \log\left(\frac{N}{\delta}\right)}$$

with probability at least $1 - 2\delta$.

Proof. This lemma can be proved in the same manner as its full information version Lem. A.1 with only two small changes; we use the bound of Lem. E.4 instead of Lem. B.1 and we modify the definition of the filtration to be $\mathcal{F}_t := \sigma(\mathbf{u}^1, o^1, \dots, \mathbf{u}^{t-1}, o^{t-1}, \mathbf{u}^t)$. \square

Lemma E.13. For any $0 < \delta < 1$,

$$\sum_{t=1}^T \sum_{n=0}^N (N-n) \sum_{x, a} \frac{\mu_n^{\pi^t, p}(x, a)}{\sqrt{\max\{1, N_n^t(x, a)\}}} \leq 3N^2 \sqrt{SAT} + SN^2 \sqrt{2T \log\left(\frac{N}{\delta}\right)},$$

holds with probability at least $1 - \delta$.

Proof. The proof is the same as for Lem. A.2 (the version proved in the full information case), except that, again, the filtration used in the proof would be defined as $\mathcal{F}_t := \sigma(\mathbf{u}^1, o^1, \dots, \mathbf{u}^{t-1}, o^{t-1}, \mathbf{u}^t)$. \square

Proposition E.14. Let b^t and \hat{p}^t be as defined in (56) and (52) respectively. Then, for any $\delta \in (0, 1)$, with probability at least $1 - 3\delta$,

$$\begin{aligned} \sum_{t=1}^T \langle \mu^{\pi^t, \hat{p}^t}, b^t \rangle + \sum_{t=1}^T \langle \mu_0, b_0^t \rangle &\leq LC'_{1/T} N^3 \left(3C'_\delta \sqrt{SAT} + 2S \sqrt{2T \log\left(\frac{N}{\delta}\right)} \right) \\ &\quad + LC'_{1/T} N^2 \left(3\sqrt{SAT} + S \sqrt{2T \log\left(\frac{N}{\delta}\right)} \right). \end{aligned}$$

Proof. The proof is the same as that of Prop. 3.1 except that we would rely on Lems. E.12 and E.13 in place of Lems. A.1 and A.2, and use the definition of b^t in (56) instead of (11). \square

2145 **Lemma E.15.** Let X be a random variable taking values in \mathbb{R} , $z_1, z_2 \geq 0$ be two constants, and $\delta' \in (0, 1)$. If $X \leq z_2$
 2146 uniformly and $P(X > z_1) \leq \delta'$, then $\mathbb{E}[X] \leq z_1 + \delta' z_2$.

2147
 2148 *Proof.* Simply, $\mathbb{E}[X] = \mathbb{E}[\mathbb{I}\{X \leq z_1\}X] + \mathbb{E}[\mathbb{I}\{X > z_1\}X] \leq z_1 + z_2 P(x > z_1) \leq z_1 + \delta' z_2$. \square
 2149

2150 E.2.5. REGRET ANALYSIS

2151 The following theorem, a restatement of Thm. 4.3, provides a regret bound for Alg. 3. Recall that we have adopted in this
 2152 section the shorthand notation $S = |\mathcal{X}|$ and $A = |\mathcal{A}|$.

2154 **Theorem E.16.** Under Asm. 4.2, Alg. 3 with a suitable tuning of τ , δ , and $(\alpha_t)_{t \in [T]}$ satisfies for any policy $\pi \in (\Delta_{\mathcal{A}})^{\mathcal{X} \times N}$
 2155 that

$$2156 \mathbb{E}[R_T(\pi)] \lesssim \sqrt{\frac{L(L+1)}{\varepsilon}} S^{5/4} A^{5/4} N^3 T^{3/4} + \frac{L+1}{\varepsilon} S^2 A^{5/2} N^4 \sqrt{T},$$

2157 where \lesssim signifies that the inequality holds up to factors logarithmic in T , N , S , and A .
 2158

2160 *Proof.* Fixing $\pi \in (\Delta_{\mathcal{A}})^{\mathcal{X} \times N}$, we have that

$$2161 \mathbb{E}[R_T(\pi)] = \mathbb{E} \sum_{t=1}^T (F^t(\mu^{\pi^t, p}) - F^t(\mu^{\pi, p}))$$

$$2162 = \mathbb{E} \underbrace{\sum_{t=1}^T (F^t(\mu^{\pi^t, p}) - F^t(\mu^{\pi^t, \hat{p}^t}))}_{\textcircled{1}} + \mathbb{E} \underbrace{\sum_{t=1}^T (F^t(\mu^{\pi^t, \hat{p}^t}) - F^t(\mu^{\pi, \hat{p}^t}))}_{\textcircled{2}} + \mathbb{E} \underbrace{\sum_{t=1}^T (F^t(\mu^{\pi, \hat{p}^t}) - F^t(\mu^{\pi, p}))}_{\textcircled{3}}.$$

2170 It holds with probability at least $1 - \frac{2}{T}$ that

$$2171 \textcircled{1} \leq L \sum_{t=1}^T \|\mu^{\pi^t, p} - \mu^{\pi^t, \hat{p}^t}\|_1 = L \sum_{t=1}^T \sum_{n=1}^N \|\mu_n^{\pi^t, p} - \mu_n^{\pi^t, \hat{p}^t}\|_1$$

$$2172 \leq L \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x, a} \mu_i^{\pi^t, p}(x, a) \|p_{i+1}(\cdot | x, a) - \hat{p}_{i+1}^t(\cdot | x, a)\|_1$$

$$2173 \leq 3LN^2 \sqrt{SATC'_{1/T}} + 2LSN^2 \sqrt{2T \log(NT)},$$

2174 where the first inequality uses the Lipschitz continuity of F^t , the second inequality follows from Lem. B.1, and the last
 2175 inequality follows from Lem. E.12. Hence, since $L \sum_{t=1}^T \|\mu^{\pi^t, p} - \mu^{\pi^t, \hat{p}^t}\|_1 \leq 2NLT$, it holds via Lem. E.15 (with $\delta' = \frac{2}{T}$)
 2176 that

$$2177 \mathbb{E}[\textcircled{1}] \leq 3LN^2 \sqrt{SATC'_{1/T}} + 2LSN^2 \sqrt{2T \log(NT)} + 4LN. \quad (57)$$

2180 For the third sum, we use again the Lipschitz continuity of F^t , Lem. B.1, and Lem. E.4 to get that with probability at least
 2181 $1 - \frac{1}{T}$,

$$2182 \textcircled{3} \leq L \sum_{t=1}^T \|\mu^{\pi, \hat{p}^t} - \mu^{\pi, p}\|_1 \leq L \sum_{t=1}^T \sum_{n=1}^N \|\mu_n^{\pi, \hat{p}^t} - \mu_n^{\pi, p}\|_1$$

$$2183 \leq L \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x, a} \mu_i^{\pi, \hat{p}^t}(x, a) \|p_{i+1}(\cdot | x, a) - \hat{p}_{i+1}^t(\cdot | x, a)\|_1$$

$$2184 \leq L \sum_{t=1}^T \sum_{n=1}^N \sum_{i=0}^{n-1} \sum_{x, a} \mu_i^{\pi, \hat{p}^t}(x, a) \frac{C'_{1/T}}{\sqrt{\max\{1, N_i^t(x, a)\}}}$$

$$2185 = L \sum_{t=1}^T \sum_{n=0}^N (N-n) \sum_{x, a} \mu_n^{\pi, \hat{p}^t}(x, a) \frac{C'_{1/T}}{\sqrt{\max\{1, N_n^t(x, a)\}}}$$

$$\begin{aligned}
 &= \sum_{t=1}^T \langle \mu^{\pi, \hat{p}^t}, b^t \rangle + \sum_{t=1}^T \langle \mu_0, b_0^t \rangle \\
 &= \sum_{t=1}^T \langle \mu^{\pi^t, \hat{p}^t}, b^t \rangle + \sum_{t=1}^T \langle \mu_0, b_0^t \rangle + \sum_{t=1}^T \langle \mu^{\pi, \hat{p}^t} - \mu^{\pi^t, \hat{p}^t}, b^t \rangle.
 \end{aligned}$$

Via Prop. E.14, it holds with probability $1 - \frac{3}{T}$ that

$$\begin{aligned}
 \sum_{t=1}^T \langle \mu^{\pi^t, \hat{p}^t}, b^t \rangle + \sum_{t=1}^T \langle \mu_0, b_0^t \rangle &\leq LC'_{1/T} N^3 [3C'_{1/T} \sqrt{SAT} + 2S \sqrt{2T \log(NT)}] \\
 &\quad + LC'_{1/T} N^2 [3\sqrt{SAT} + S \sqrt{2T \log(NT)}].
 \end{aligned}$$

Hence, chaining these last two results and using a union bound, we get via Lem. E.15 (with $\delta' = \frac{4}{T}$) that

$$\begin{aligned}
 \mathbb{E} \left[\textcircled{3} - \sum_{t=1}^T \langle \mu^{\pi, \hat{p}^t} - \mu^{\pi^t, \hat{p}^t}, b^t \rangle \right] &\leq LC'_{1/T} N^3 [3C'_{1/T} \sqrt{SAT} + 2S \sqrt{2T \log(NT)}] \\
 &\quad + LC'_{1/T} N^2 [3\sqrt{SAT} + S \sqrt{2T \log(NT)}] + 8LN(1 + C'_{1/T} N), \quad (58)
 \end{aligned}$$

where we have used that

$$\textcircled{3} - \sum_{t=1}^T \langle \mu^{\pi, \hat{p}^t} - \mu^{\pi^t, \hat{p}^t}, b^t \rangle \leq L \sum_{t=1}^T \|\mu^{\pi, \hat{p}^t} - \mu^{\pi^t, \hat{p}^t}\|_1 + \sum_{t=1}^T \|\delta^t\|_\infty \|\mu^{\pi, \hat{p}^t} - \mu^{\pi^t, \hat{p}^t}\|_1 \leq 2LNT(1 + C'_{1/T} N).$$

Define $\hat{F}^t: (\Delta_{\mathcal{X} \times \mathcal{A}})^N \rightarrow \mathbb{R}$ as

$$\hat{F}^t(\mu) = \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[F^t((1-\delta)\mu + \delta\zeta^{\mathbf{v}, \hat{p}^t}) \right].$$

As \hat{p}^t satisfies the condition of Asm. 4.2 by design, $\zeta^{\mathbf{v}, \hat{p}^t} \in \mathcal{M}_{\mu_0}^{\hat{p}^t} \subset (\Delta_{\mathcal{X} \times \mathcal{A}})^N$ as argued in App. E.2.1; thus, \hat{F}^t is well-defined. Similarly, since $\mathbf{u}^t \in \mathbb{S}^{NS(A-1)} \subset \mathbb{B}^{NS(A-1)}$ and $\zeta^t = \zeta^{\mathbf{u}^t, \hat{p}^t}$, it holds that $\zeta^t \in \mathcal{M}_{\mu_0}^{\hat{p}^t}$. Via the convexity of $\mathcal{M}_{\mu_0}^{\hat{p}^t}$, the fact that $\mu^t \in \mathcal{M}_{\mu_0}^{\hat{p}^t}$, and the definition of $\hat{\mu}^t$; it holds that $\hat{\mu}^t \in \mathcal{M}_{\mu_0}^{\hat{p}^t}$. This yields that $\hat{\mu}^t = \mu^{\pi^t, \hat{p}^t}$, recalling the definition of π^t in Alg. 3. Using the Lipschitz smoothness of F^t , we have that

$$\begin{aligned}
 F^t(\mu^{\pi^t, \hat{p}^t}) - \hat{F}^t(\mu^t) &= F^t(\hat{\mu}^t) - \hat{F}^t(\mu^t) = F^t((1-\delta)\mu^t + \delta\zeta^t) - \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[F^t((1-\delta)\mu^t + \delta\zeta^{\mathbf{v}, \hat{p}^t}) \right] \\
 &\leq \delta L \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \|\zeta^t - \zeta^{\mathbf{v}, \hat{p}^t}\|_1 \leq 2\delta LN
 \end{aligned}$$

and that

$$\begin{aligned}
 \hat{F}^t(\mu^{\pi, \hat{p}^t}) - F^t(\mu^{\pi, \hat{p}^t}) &= \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[F^t((1-\delta)\mu^{\pi, \hat{p}^t} + \delta\zeta^{\mathbf{v}, \hat{p}^t}) \right] - F^t(\mu^{\pi, \hat{p}^t}) \\
 &\leq \delta L \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \|\zeta^{\mathbf{v}, \hat{p}^t} - \mu^{\pi, \hat{p}^t}\|_1 \leq 2\delta LN.
 \end{aligned}$$

Hence,

$$\begin{aligned}
 \textcircled{2} &= \sum_{t=1}^T (F^t(\mu^{\pi^t, \hat{p}^t}) - \hat{F}^t(\mu^t) + \hat{F}^t(\mu^t) - \hat{F}^t(\mu^{\pi, \hat{p}^t}) + \hat{F}^t(\mu^{\pi, \hat{p}^t}) - F^t(\mu^{\pi, \hat{p}^t})) \\
 &\leq \sum_{t=1}^T \langle \nabla \hat{F}^t(\mu^t), \mu^t - \mu^{\pi, \hat{p}^t} \rangle + 4\delta LNT \\
 &= \sum_{t=1}^T \langle \nabla \hat{F}^t(\mu^t) - b^t, \mu^t - \mu^{\pi, \hat{p}^t} \rangle + 4\delta LNT + \sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi, \hat{p}^t} \rangle + \sum_{t=1}^T \langle b^t, \mu^t - \mu^{\pi^t, \hat{p}^t} \rangle.
 \end{aligned}$$

2255 The last term is easily bounded as follows:

$$2256 \sum_{t=1}^T \langle b^t, \mu^t - \mu^{\pi^t, \hat{p}^t} \rangle = \sum_{t=1}^T \langle b^t, \mu^t - \hat{\mu}^t \rangle = \delta \sum_{t=1}^T \langle b^t, \mu^t - \zeta^t \rangle \leq \delta \sum_{t=1}^T \|b^t\|_\infty \|\mu^t - \zeta^t\|_1 \leq 2\delta C'_{1/T} L N^2 T.$$

2260 We then conclude that

$$2261 \mathbb{E} \left[\textcircled{2} + \sum_{t=1}^T \langle b^t, \mu^{\pi^t, \hat{p}^t} - \mu^{\pi^t, \hat{p}^t} \rangle \right] \leq \mathbb{E} \sum_{t=1}^T \langle \nabla \hat{F}^t(\mu^t) - b^t, \mu^t - \mu^{\pi^t, \hat{p}^t} \rangle + 4\delta L N T + 2\delta C'_{1/T} L N^2 T. \quad (59)$$

2265 Then, combining (57), (58), and (59) yields that

$$2266 \mathbb{E} [R_T(\pi)] \leq \mathbb{E} \sum_{t=1}^T \langle \nabla \hat{F}^t(\mu^t) - b^t, \mu^t - \mu^{\pi^t, \hat{p}^t} \rangle + 2\delta L(2 + C'_{1/T} N) N T \\ 2267 + 3LC'_{1/T} N^3 [3C'_{1/T} \sqrt{SAT} + 2S\sqrt{2T \log(NT)}] + 4LN(3 + 2C'_{1/T} N). \quad (60)$$

2272 Define $\mathbb{F}^t, \hat{\mathbb{F}}^t : (\mathcal{M}_{\mu_0}^{\hat{p}^t})^- \rightarrow \mathbb{R}$ as $\mathbb{F}^t(\xi) := F^t(\Lambda_{\hat{p}^t}(\xi))$ and $\hat{\mathbb{F}}^t(\xi) := \hat{F}^t(\Lambda_{\hat{p}^t}(\xi))$. Then, recalling that $\kappa := \frac{\varepsilon}{A-1+\sqrt{A-1}}$,

$$2273 \hat{\mathbb{F}}^t(\Lambda_{\hat{p}^t}^{-1}(\mu^t)) = \hat{F}^t(\mu^t) = \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[F^t((1-\delta)\mu^t + \delta\zeta^{\mathbf{v}, \hat{p}^t}) \right] \\ 2274 = \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[\mathbb{F}^t(\Lambda_{\hat{p}^t}^{-1}((1-\delta)\mu^t + \delta\zeta^{\mathbf{v}, \hat{p}^t})) \right] \\ 2275 = \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[\mathbb{F}^t((B^{\hat{p}^t})^+((1-\delta)\mu^t + \delta\zeta^{\mathbf{v}, \hat{p}^t} - \beta^{\hat{p}^t})) \right] \\ 2276 = \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[\mathbb{F}^t((1-\delta)(B^{\hat{p}^t})^+(\mu^t - \beta^{\hat{p}^t}) + \delta(B^{\hat{p}^t})^+(\zeta^{\mathbf{v}, \hat{p}^t} - \beta^{\hat{p}^t})) \right] \\ 2277 = \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[\mathbb{F}^t((1-\delta)\Lambda_{\hat{p}^t}^{-1}(\mu^t) + \delta\Lambda_{\hat{p}^t}^{-1}(\zeta^{\mathbf{v}, \hat{p}^t})) \right] \\ 2278 = \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[\mathbb{F}^t((1-\delta)\Lambda_{\hat{p}^t}^{-1}(\mu^t) + \delta\kappa \mathbf{1}_{NS(A-1)} + \delta\kappa \mathbf{v}) \right],$$

2286 where the fourth equality follows from the fact that $\Lambda_{\hat{p}^t}^{-1}(\mu) = (B^{\hat{p}^t})^+(\mu - \beta^{\hat{p}^t})$, and the last equality follows since $\zeta^{\mathbf{v}, \hat{p}^t} = \Lambda_{\hat{p}^t}(\kappa \mathbf{1}_{NS(A-1)} + \kappa \mathbf{v})$. Lem. 1 in (Flaxman et al., 2005) and the chain rule imply that

$$2287 \nabla \hat{\mathbb{F}}^t(\Lambda_{\hat{p}^t}^{-1}(\mu^t)) = \frac{1-\delta}{\delta\kappa} NS(A-1) \mathbb{E}_{\mathbf{u} \in \mathbb{S}^{NS(A-1)}} \left[\mathbb{F}^t((1-\delta)\Lambda_{\hat{p}^t}^{-1}(\mu^t) + \delta\kappa \mathbf{1}_{NS(A-1)} + \delta\kappa \mathbf{u}) \mathbf{u} \right] \\ 2288 = \frac{1-\delta}{\delta\kappa} NS(A-1) \mathbb{E}_{\mathbf{u} \in \mathbb{S}^{NS(A-1)}} \left[\mathbb{F}^t((1-\delta)\Lambda_{\hat{p}^t}^{-1}(\mu^t) + \delta\Lambda_{\hat{p}^t}^{-1}(\zeta^{\mathbf{u}, \hat{p}^t})) \mathbf{u} \right] \\ 2289 = \frac{1-\delta}{\delta\kappa} NS(A-1) \mathbb{E}_{\mathbf{u} \in \mathbb{S}^{NS(A-1)}} \left[F^t((1-\delta)\mu^t + \delta\zeta^{\mathbf{u}, \hat{p}^t}) \mathbf{u} \right] \\ 2290 = \frac{1-\delta}{\delta\kappa} NS(A-1) \mathbb{E}_{\mathbf{u}^t \in \mathbb{S}^{NS(A-1)}} \left[F^t((1-\delta)\mu^t + \delta\zeta^t) \mathbf{u}^t \right],$$

2299 where the last equality uses that $\zeta^t = \zeta^{\mathbf{u}^t, \hat{p}^t}$ and that both μ^t and \hat{p}^t are independent with respect to \mathbf{u}^t . And since $\nabla \hat{\mathbb{F}}^t(\Lambda_{\hat{p}^t}^{-1}(\mu^t)) = (B^{\hat{p}^t})^\top \nabla \hat{F}^t(\mu^t)$, we obtain that

$$2300 (B^{\hat{p}^t} (B^{\hat{p}^t})^+)^\top \nabla \hat{F}^t(\mu^t) = \frac{1-\delta}{\delta\kappa} NS(A-1) \mathbb{E}_{\mathbf{u}^t \in \mathbb{S}^{NS(A-1)}} \left[F^t((1-\delta)\mu^t + \delta\zeta^t) ((B^{\hat{p}^t})^+)^\top \mathbf{u}^t \right] \\ 2301 = \mathbb{E}_{\mathbf{u}^t \in \mathbb{S}^{NS(A-1)}} \left[((B^{\hat{p}^t})^+)^\top \hat{\mathbf{g}}^t \right], \quad (61)$$

2306 where

$$2307 \hat{\mathbf{g}}^t := \frac{1-\delta}{\delta\kappa} NS(A-1) F^t((1-\delta)\mu^t + \delta\zeta^t) \mathbf{u}^t = \frac{1-\delta}{\delta\kappa} NS(A-1) F^t(\hat{\mu}^t) \mathbf{u}^t.$$

2310 The vector \hat{g}^t differs from g^t (which is employed in Alg. 3) in that it is defined using $F^t(\hat{\mu}^t)$ instead of $F^t(\mu^{\pi^t, p})$. For
 2311 round $t \in [T]$, let $\mathcal{F}_t := \sigma(\mathbf{u}^1, o^1, \dots, \mathbf{u}^t, o^t)$ denote the σ -algebra generated by the random events up to the end of round
 2312 t ; and let $\mathbb{E}_t[\cdot] := \mathbb{E}[\cdot \mid \mathcal{F}_{t-1}]$ with \mathcal{F}_0 being the trivial σ -algebra. We then have that

$$\begin{aligned} 2313 \mathbb{E} \sum_{t=1}^T \langle \nabla \hat{F}^t(\mu^t), \mu^t - \mu^{\pi, \hat{p}^t} \rangle &= \mathbb{E} \sum_{t=1}^T \langle (B^{\hat{p}^t} (B^{\hat{p}^t})^+)^{\top} \nabla \hat{F}^t(\mu^t), \mu^t - \mu^{\pi, \hat{p}^t} \rangle \\ 2314 &= \mathbb{E} \sum_{t=1}^T \langle \mathbb{E}_t[\langle (B^{\hat{p}^t})^+ \rangle^{\top} \hat{g}^t], \mu^t - \mu^{\pi, \hat{p}^t} \rangle \\ 2315 &= \mathbb{E} \sum_{t=1}^T \langle \langle (B^{\hat{p}^t})^+ \rangle^{\top} \hat{g}^t, \mu^t - \mu^{\pi, \hat{p}^t} \rangle, \end{aligned}$$

2322 where the first equality holds via the fact that $B^{\hat{p}^t} (B^{\hat{p}^t})^+ (\mu^t - \mu^{\pi, \hat{p}^t}) = \mu^t - \mu^{\pi, \hat{p}^t}$ since $\mu^t - \mu^{\pi, \hat{p}^t}$ belongs to the
 2323 column space of $B^{\hat{p}^t}$ (see App. E.1), the second equality uses (61) and the fact that conditioned on \mathcal{F}_{t-1} , the only source
 2324 of randomness in $\langle (B^{\hat{p}^t})^+ \rangle^{\top} \hat{g}^t$ is \mathbf{u}^t , which is sampled independently in each round; and the last equality uses the
 2325 tower rule, linearity of expectation, and the fact that $\mu^t - \mu^{\pi, \hat{p}^t}$ is measurable with respect to \mathcal{F}_{t-1} . Since $\mu^t - \mu^{\pi, \hat{p}^t} =$
 2326 $B^{\hat{p}^t} (\Lambda_{\hat{p}^t}^{-1}(\mu^t) - \Lambda_{\hat{p}^t}^{-1}(\mu^{\pi, \hat{p}^t}))$, we have that

$$2327 (\hat{g}^t)^{\top} (B^{\hat{p}^t})^+ (\mu^t - \mu^{\pi, \hat{p}^t}) = (\hat{g}^t)^{\top} (B^{\hat{p}^t})^+ B^{\hat{p}^t} (\Lambda_{\hat{p}^t}^{-1}(\mu^t) - \Lambda_{\hat{p}^t}^{-1}(\mu^{\pi, \hat{p}^t})) = (\hat{g}^t)^{\top} (\Lambda_{\hat{p}^t}^{-1}(\mu^t) - \Lambda_{\hat{p}^t}^{-1}(\mu^{\pi, \hat{p}^t}))$$

2331 since $(B^{\hat{p}^t})^+ B^{\hat{p}^t} = \mathbf{I}_{NS(A-1)}$, see App. E.1. Therefore,

$$\begin{aligned} 2332 \mathbb{E} \sum_{t=1}^T \langle \nabla \hat{F}^t(\mu^t), \mu^t - \mu^{\pi, \hat{p}^t} \rangle &= \mathbb{E} \sum_{t=1}^T \langle \hat{g}^t, \Lambda_{\hat{p}^t}^{-1}(\mu^t) - \Lambda_{\hat{p}^t}^{-1}(\mu^{\pi, \hat{p}^t}) \rangle \\ 2333 &= \mathbb{E} \sum_{t=1}^T \langle g^t, \Lambda_{\hat{p}^t}^{-1}(\mu^t) - \Lambda_{\hat{p}^t}^{-1}(\mu^{\pi, \hat{p}^t}) \rangle + \mathbb{E} \sum_{t=1}^T \langle \hat{g}^t - g^t, \Lambda_{\hat{p}^t}^{-1}(\mu^t) - \Lambda_{\hat{p}^t}^{-1}(\mu^{\pi, \hat{p}^t}) \rangle \\ 2334 &= \mathbb{E} \sum_{t=1}^T \langle \hat{g}^t, \mu^t - \mu^{\pi, \hat{p}^t} \rangle + \mathbb{E} \sum_{t=1}^T \langle \hat{g}^t - g^t, \Lambda_{\hat{p}^t}^{-1}(\mu^t) - \Lambda_{\hat{p}^t}^{-1}(\mu^{\pi, \hat{p}^t}) \rangle, \end{aligned}$$

2341 where the last equality follows from the definition of \hat{g}^t (see Alg. 3) and the fact that μ^t and μ^{π, \hat{p}^t} are expansions of $\Lambda_{\hat{p}^t}^{-1}(\mu^t)$
 2342 and $\Lambda_{\hat{p}^t}^{-1}(\mu^{\pi, \hat{p}^t})$ respectively, augmented with the entries corresponding to action a^* . Focusing on the second sum, we have
 2343 that

$$\begin{aligned} 2344 \sum_{t=1}^T \langle \hat{g}^t - g^t, \Lambda_{\hat{p}^t}^{-1}(\mu^t) - \Lambda_{\hat{p}^t}^{-1}(\mu^{\pi, \hat{p}^t}) \rangle &\leq \sum_{t=1}^T \|\hat{g}^t - g^t\|_{\infty} \|\Lambda_{\hat{p}^t}^{-1}(\mu^t) - \Lambda_{\hat{p}^t}^{-1}(\mu^{\pi, \hat{p}^t})\|_1 \\ 2345 &\leq \sum_{t=1}^T \|\hat{g}^t - g^t\|_{\infty} \|\mu^t - \mu^{\pi, \hat{p}^t}\|_1 \\ 2346 &\leq 2N \sum_{t=1}^T \|\hat{g}^t - g^t\|_{\infty} \\ 2347 &= 2 \frac{1-\delta}{\delta\kappa} N^2 S(A-1) \sum_{t=1}^T \|\mathbf{u}^t\|_{\infty} |F^t(\hat{\mu}^t) - F^t(\mu^{\pi^t, p})| \\ 2348 &\leq \frac{4}{\varepsilon\delta} N^2 S A^2 \sum_{t=1}^T |F^t(\hat{\mu}^t) - F^t(\mu^{\pi^t, p})| \\ 2349 &= \frac{4}{\varepsilon\delta} N^2 S A^2 \sum_{t=1}^T |F^t(\mu^{\pi^t, \hat{p}^t}) - F^t(\mu^{\pi^t, p})| \\ 2350 &\leq \frac{4}{\varepsilon\delta} L N^2 S A^2 \sum_{t=1}^T \|\mu^{\pi^t, \hat{p}^t} - \mu^{\pi^t, p}\|_1, \end{aligned}$$

2365 where the fourth inequality uses that $\kappa \geq \frac{\varepsilon}{2A}$ and that $\mathbf{u}^t \in \mathbb{S}^{NS(A-1)}$, and the last inequality uses the Lipschitz continuity
 2366 of F^t . As shown in (57), we have that

$$2367 \mathbb{E} \sum_{t=1}^T \|\mu^{\pi^t, \hat{p}^t} - \mu^{\pi^t, p}\|_1 \leq 3N^2 \sqrt{SAT} C'_{1/T} + 2SN^2 \sqrt{2T \log(NT)} + 4N.$$

2371 Hence,

$$2372 \mathbb{E} \sum_{t=1}^T \langle \nabla \hat{F}^t(\mu^t), \mu^t - \mu^{\pi, \hat{p}^t} \rangle \leq \mathbb{E} \sum_{t=1}^T \langle \hat{g}^t, \mu^t - \mu^{\pi, \hat{p}^t} \rangle + \frac{4}{\varepsilon \delta} LN^3 SA^2 (3N \sqrt{SAT} C'_{1/T} + 2SN \sqrt{2T \log(NT)} + 4).$$

2376 Combining this result with (60) yields that

$$\begin{aligned} 2377 \mathbb{E}[R_T(\pi)] &\leq \mathbb{E} \sum_{t=1}^T \langle \hat{g}^t - b^t, \mu^t - \mu^{\pi, \hat{p}^t} \rangle + \frac{4}{\varepsilon \delta} LN^3 SA^2 (3N \sqrt{SAT} C'_{1/T} + 2SN \sqrt{2T \log(NT)} + 4) \\ 2378 &\quad + 2\delta L(2 + C'_{1/T} N) NT + 3LC'_{1/T} N^3 [3C'_{1/T} \sqrt{SAT} + 2S \sqrt{2T \log(NT)}] + 4LN(3 + 2C'_{1/T} N) \\ 2379 &\leq \mathbb{E} \sum_{t=1}^T \langle \hat{g}^t - b^t, \mu^t - \mu^{\pi, \hat{p}^t} \rangle + \underbrace{\delta 2L(2 + C'_{1/T} N) NT}_{=:\Xi_1} + 4LN(3 + 2C'_{1/T} N) \\ 2380 &\quad + \underbrace{\frac{1}{\delta} \frac{7}{\varepsilon} LC'_{1/T} N^3 SA^2 (3N \sqrt{AT} C'_{1/T} + 2N \sqrt{2ST \log(NT)} + 4)}_{=:\Xi_2}, \end{aligned} \quad (62)$$

2389 where the last inequality uses that $C'_{1/T} \geq \sqrt{S}$. Note that

$$\begin{aligned} 2390 \|\hat{g}^t\|_{1, \infty} &= \sum_{n=1}^N \|g_n^t\|_{\infty} \\ 2391 &= \frac{1 - \delta}{\varepsilon \delta} NS(A-1)(A-1 + \sqrt{A-1}) F^t(\mu^{\pi^t, p}) \sum_{n=1}^N \|\mathbf{u}_n^t\|_{\infty} \\ 2392 &\leq \frac{2}{\varepsilon \delta} N^2 SA^2 \sum_{n=1}^N \|\mathbf{u}_n^t\|_{\infty} \\ 2393 &\leq \frac{2}{\varepsilon \delta} N^2 SA^2 \sqrt{N} \sqrt{\sum_{n=1}^N \|\mathbf{u}_n^t\|_{\infty}^2} \leq \frac{2}{\varepsilon \delta} N^{5/2} SA^2 \sqrt{\sum_{n=1}^N \sum_{x,a} |\mathbf{u}_n^t(x, a)|^2} \leq \frac{2}{\varepsilon \delta} N^{5/2} SA^2, \end{aligned}$$

2404 where the second inequality uses Cauchy-Schwarz and the last inequality uses that $\mathbf{u}^t \in \mathbb{S}^{NS(A-1)}$. Moreover, we have that
 2405 $\|b^t\|_{1, \infty} = \sum_{n=1}^N \|b_n^t\|_{\infty} \leq \sum_{n=1}^N L(N-n) C'_{1/T} \leq LN^2 C'_{1/T}$. Hence, using that $C'_{1/T} \geq \sqrt{S}$,

$$2406 \|\hat{g}^t - b^t\|_{1, \infty} \leq \frac{2}{\varepsilon \delta} N^{5/2} SA^2 + LN^2 C'_{1/T} \leq \frac{2}{\varepsilon \delta} (L+1) C'_{1/T} \sqrt{S} A^2 N^{5/2}.$$

2410 Via Lems. A.4 and E.11,⁴ we can invoke Lem. 2.1 with $c = 5\varepsilon$, $\zeta = \frac{2}{\varepsilon \delta} (L+1) C'_{1/T} \sqrt{S} A^2 N^{5/2}$, and $\alpha_t = 1/(t+1)$ to get
 2411 that (from the proof of Lem. 2.1)

$$\begin{aligned} 2412 \sum_{t=1}^T \langle \hat{g}^t - b^t, \mu^t - \mu^{\pi, \hat{p}^t} \rangle &\leq \tau \left(\frac{2}{\varepsilon \delta} (L+1) C'_{1/T} \sqrt{S} A^2 N^{5/2} \right)^2 T + \frac{20e^2 SN \log(AT)^2 (N+A)}{\tau} \\ 2413 &\quad + \frac{10e^2}{\varepsilon \delta} (L+1) C'_{1/T} S^{3/2} A^2 N^{7/2} \log(T). \end{aligned}$$

2418 ⁴To invoke Lem. E.11, we assume without loss of generality that the constant ε specified in Asm. 4.2 satisfies $\varepsilon \leq \frac{1}{2S}$.

2420 Tuning τ optimally yields that

$$\begin{aligned}
 2421 & \\
 2422 & \sum_{t=1}^T \langle \hat{g}^t - b^t, \mu^t - \mu^{\pi, \hat{p}^t} \rangle \leq \frac{4}{\varepsilon \delta} (L+1) C'_{1/T} S A^2 N^{5/2} \sqrt{20e^2 N(N+A)T \log(AT)^2} \\
 2423 & \\
 2424 & \\
 2425 & \qquad \qquad \qquad + \frac{10e^2}{\varepsilon \delta} (L+1) C'_{1/T} S^{3/2} A^2 N^{7/2} \log(T) \\
 2426 & \\
 2427 & \leq \frac{1}{\delta} \frac{10e^2}{\varepsilon} (L+1) C'_{1/T} S A^2 N^3 \log(AT) \underbrace{(\sqrt{(N+A)T} + \sqrt{SN})}_{=:\Xi_3}. \\
 2428 & \\
 2429 &
 \end{aligned}$$

2430 Hence, plugging back into (62) yields that

$$2431 \mathbb{E} [R_T(\pi)] \leq \delta \Xi_1 + \frac{1}{\delta} (\Xi_2 + \Xi_3) + 4LN(3 + 2C'_{1/T}N).$$

2432 Setting $\delta := \min \left\{ 1, \sqrt{\frac{\Xi_2 + \Xi_3}{\Xi_1}} \right\}$, we get that

$$2433 \mathbb{E} [R_T(\pi)] \leq \max \left\{ 2\sqrt{\Xi_1(\Xi_2 + \Xi_3)}, 2(\Xi_2 + \Xi_3) \right\} + 4LN(3 + 2C'_{1/T}N).$$

2434 Consequently, the theorem follows after using the definition of $C'_{1/T}$ from Eq. (55) and ignoring log factors. \square

2441 E.3. Self-Concordant Regularization Approach

2442 We have used the set $(\mathcal{M}_{\mu_0}^p)^-$, the preimage of $\mathcal{M}_{\mu_0}^p$ under the map Ξ_p (or Λ_p), to represent in $\mathbb{R}^{NS(A-1)}$ the set of valid occupancy measures. A more concise characterization, given by Lem. E.1, is that

$$2443 (\mathcal{M}_{\mu_0}^p)^- = \{ \xi \in \mathbb{R}^{NS(A-1)} : B^p \xi \geq -\beta^p \};$$

2444 in other words, $(\mathcal{M}_{\mu_0}^p)^-$ is a convex polytope formed by the constraints $B^p(n, x, a, \cdot, \cdot, \cdot)^\top \xi + \beta_n^p(x, a) \geq 0$ for $n, x, a \in [N] \times \mathcal{X} \times \mathcal{A}$. Moreover, Lem. E.2 asserts that $\text{int}((\mathcal{M}_{\mu_0}^p)^-)$, the interior of $(\mathcal{M}_{\mu_0}^p)^-$, is not empty under Asm. 4.4.

2445 We consider then the function $\psi_{\text{lb}} : \text{int}((\mathcal{M}_{\mu_0}^p)^-) \rightarrow \mathbb{R}$ defined as

$$2446 \psi_{\text{lb}}(\xi) := - \sum_{n,x,a} \log(B(n, x, a, \cdot, \cdot, \cdot)^\top \xi + \beta_n(x, a)).$$

2447 As mentioned in Sec. 4.2.2, Corollary 3.1.1 in (Nemirovski, 2004) yields that ψ_{lb} is a ϑ -self-concordant barrier (see Definition 3.1.1 in Nemirovski, 2004) for $(\mathcal{M}_{\mu_0}^p)^-$ with $\vartheta = N \cdot S \cdot A$. The approach we analyze here is to perform OMD directly on the set $(\mathcal{M}_{\mu_0}^p)^-$ with ψ_{lb} as the regularizer.

2448 For $\xi \in \text{int}((\mathcal{M}_{\mu_0}^p)^-)$ and $y \in \mathbb{R}^{NS(A-1)}$, define the local norm $\|y\|_\xi := \sqrt{y^\top \nabla^2 \psi_{\text{lb}}(\xi) y}$. This is indeed a norm since the fact that $(\mathcal{M}_{\mu_0}^p)^-$ is bounded implies via Property II in (Nemirovski, 2004, Section 2.2) that the Hessian of ψ_{lb} is non-singular everywhere. Its dual norm is denoted as $\|y\|_{\xi,*} := \sqrt{y^\top (\nabla^2 \psi_{\text{lb}}(\xi))^{-1} y}$. The Dikin ellipsoid of radius r at $\xi \in \text{int}((\mathcal{M}_{\mu_0}^p)^-)$ is given by

$$2449 \mathcal{E}_r(\xi) := \{y \in \mathbb{R}^{NS(A-1)} : \|y - \xi\|_\xi \leq r\} = \xi + r(\nabla^2 \psi_{\text{lb}}(\xi))^{-1/2} \mathbb{B}^{NS(A-1)}.$$

2450 Via Property I in (Nemirovski, 2004, Section 2.2), $\mathcal{E}_1(\xi) \subseteq (\mathcal{M}_{\mu_0}^p)^-$ for any $\xi \in \text{int}((\mathcal{M}_{\mu_0}^p)^-)$.

2451 For $\xi, y \in \text{int}((\mathcal{M}_{\mu_0}^p)^-)$, we denote by $D_{\psi_{\text{lb}}}(y, \xi) := \psi_{\text{lb}}(y) - \psi_{\text{lb}}(\xi) - \langle y - \xi, \nabla \psi_{\text{lb}}(\xi) \rangle$ the Bregman divergence between y and ξ with respect to ψ_{lb} . From the proof of Thm. E.16, we recall the definition $\mathbb{F}^t := F^t \circ \Lambda_p$. As alluded to above, our OMD updates will take the form

$$2452 \xi^{t+1} \leftarrow \arg \min_{\xi \in (\mathcal{M}_{\mu_0}^p)^-} \tau \langle g^t, \xi \rangle + D_{\psi_{\text{lb}}}(\xi, \xi^t),$$

2453 where g^t will be chosen as a surrogate for $\nabla \mathbb{F}^t(\xi^t)$. Differently from the proof of Thm. E.16, we redefine the smoothed approximation $\hat{\mathbb{F}}^t : (\mathcal{M}_{\mu_0}^p)^- \rightarrow \mathbb{R}$ such that

$$2454 \hat{\mathbb{F}}^t(\xi) := \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} [\mathbb{F}^t((1-\delta)\xi + \delta(\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v}))].$$

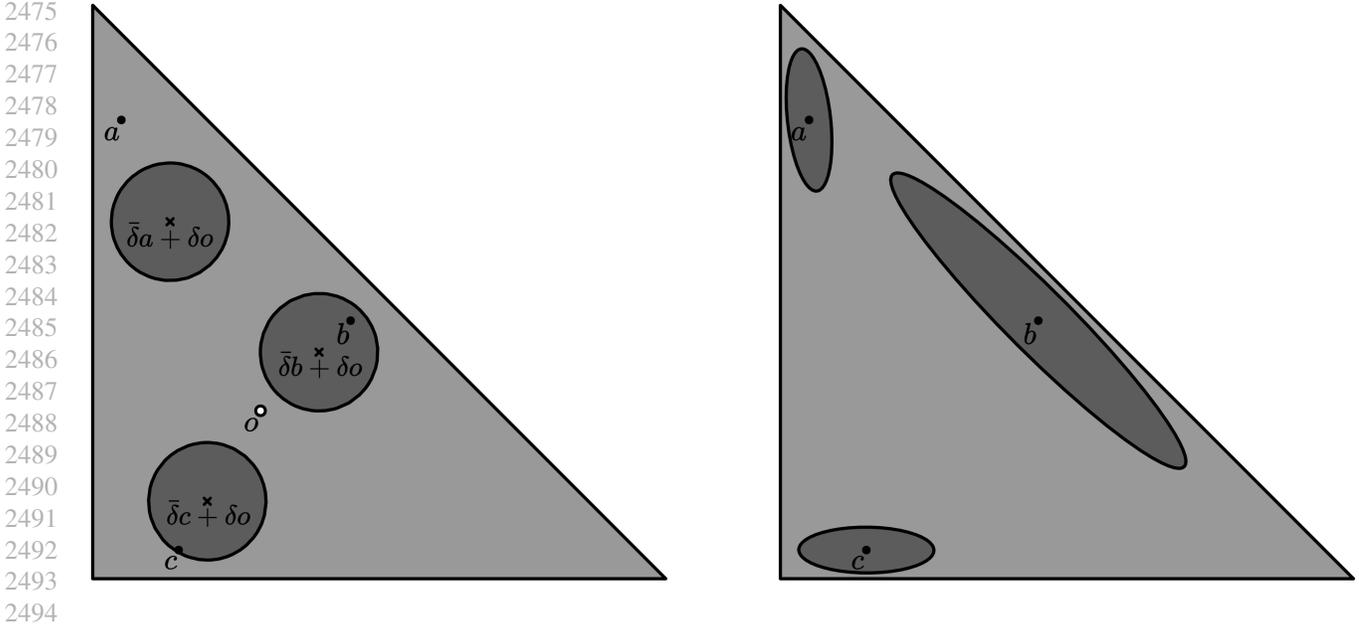


Figure 4. This figure provides a graphical comparison between the sampling approach used in Alg. 3, represented on the left, and that used in Alg. 4, represented on the right. The simplified domain here is $\{x \in [0, 1]^2 : \|x\|_1 \leq 1\}$. Both approaches are illustrated at three points: a , b , and c . In the first approach, with some $\delta \in (0, 1)$ and $\bar{\delta} := 1 - \delta$, we sample from a circle of radius $\delta/(2 + \sqrt{2})$ centered at a convex combination between the point of interest and $o := (1/(2 + \sqrt{2}), 1/(2 + \sqrt{2}))$. In the second approach, we consider the barrier $-\log(1 - x_1 - x_2) - \sum_{i=1,2} \log(x_i)$ and sample from the Dikin ellipsoid (of a certain common radius) induced by this function at each point.

This is well-defined since we are evaluating \mathbb{F}^t on a convex combination of the argument ξ and a point inside the ellipsoid $\mathcal{E}_1(\xi^t)$, which is a subset of $(\mathcal{M}_{\mu_0}^p)^-$ as cited before. Via Corollary 6.8 in (Hazan, 2021) and the chain rule, we have that

$$\nabla \widehat{\mathbb{F}}^t(\xi) = \frac{(1 - \delta)}{\delta} NS(A - 1) \mathbb{E}_{\mathbf{u} \in \mathbb{S}^{NS(A-1)}} [\mathbb{F}^t((1 - \delta)\xi + \delta(\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{u})) (\nabla^2 \psi_{\text{lb}}(\xi^t))^{1/2} \mathbf{u}]. \quad (63)$$

Hence, with \mathbf{u}^t sampled uniformly from $\mathbb{S}^{NS(A-1)}$, we pick (as mentioned in Sec. 4.2.2)

$$g^t := \frac{(1 - \delta)}{\delta} NS(A - 1) \mathbb{F}^t(\xi^t + \delta(\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{u}^t) (\nabla^2 \psi_{\text{lb}}(\xi^t))^{1/2} \mathbf{u}^t \quad (64)$$

such that $\mathbb{E}_{\mathbf{u}^t}[g^t] = \nabla \widehat{\mathbb{F}}^t(\xi^t)$, see also (Saha & Tewari, 2011) for a similar estimator in another BCO setting. We summarize this approach in Alg. 4, and provide in Fig. 4 a graphical comparison with the sampling approach of Alg. 3 on a simple decision set. Before proving the regret bound of Thm. 4.5, we collect a few standard properties and auxiliary results concerning self-concordant barriers and their use as regularizers.

E.3.1. AUXILIARY LEMMAS

For $x, y \in \text{int}(\mathcal{M}_{\mu_0}^p)^-$, it holds via Property I in (Nemirovski, 2004, Section 2.2) that

$$(1 - \|y - x\|_x)^2 \nabla^2 \psi_{\text{lb}}(x) \preceq \nabla^2 \psi_{\text{lb}}(y) \preceq \frac{1}{(1 - \|y - x\|_x)^2} \nabla^2 \psi_{\text{lb}}(x) \quad (65)$$

whenever $\|y - x\|_x < 1$. We state the following auxiliary lemma, which will be used to assert the proximity between ξ^t and ξ^{t+1} for our algorithm. Establishing this ‘stability’ is a crucial step in the local norm analysis.

Lemma E.17. *Let $x \in \text{int}(\mathcal{M}_{\mu_0}^p)^-$ and $\ell \in \mathbb{R}^{NS(A-1)}$ be such that $\|\ell\|_{x,*} \leq \frac{1}{16}$, and define*

$$y := \arg \min_{\xi \in \text{int}(\mathcal{M}_{\mu_0}^p)^-} \langle \ell, \xi \rangle + D_{\psi_{\text{lb}}}(\xi, x).$$

Then, $y \in \mathcal{E}_{1/2}(x)$.

2530 **Algorithm 4** Bandit O-MD-CURL with logarithmic barrier regularization

2531 **input:** domain $(\mathcal{M}_{\mu_0}^p)^-$ with non-empty interior, learning rate $\tau > 0$, exploration parameter $\delta \in (0, 1]$
 2532 **initialization:** $\xi^1 \leftarrow \arg \min_{\xi \in \text{int}(\mathcal{M}_{\mu_0}^p)^-} \psi_{\text{lb}}(\xi)$
 2533 **for** $t = 1, \dots, T$ **do**
 2534 draw $\mathbf{u}^t \in \mathbb{S}^{NS(A-1)}$ uniformly at random
 2535 $\hat{\xi}^t \leftarrow \xi^t + \delta(\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{u}^t$
 2536 $\hat{\mu}^t \leftarrow \Lambda_p(\hat{\xi}^t)$
 2537 $\pi_n^t(a|x) \leftarrow \hat{\mu}^t(x, a) / \sum_{a \in \mathcal{A}} \hat{\mu}^t(x, a)$
 2538 output π^t and observe $F^t(\hat{\mu}^t)$
 2539 $g^t \leftarrow \frac{(1-\delta)}{\delta} NS(A-1) F^t(\hat{\mu}^t) (\nabla^2 \psi_{\text{lb}}(\xi^t))^{1/2} \mathbf{u}^t$
 2540 $\xi^{t+1} \leftarrow \arg \min_{\xi \in \text{int}(\mathcal{M}_{\mu_0}^p)^-} \tau \langle g^t, \xi \rangle + D_{\psi_{\text{lb}}}(\xi, \xi^t)$
 2541 **end for**

2542 *Proof.* For $\xi \in \text{int}(\mathcal{M}_{\mu_0}^p)^-$, let

$$2543 \quad g(\xi) := \langle \ell, \xi \rangle + D_{\psi_{\text{lb}}}(\xi, x) = \langle \ell, \xi \rangle + \psi_{\text{lb}}(\xi) - \psi_{\text{lb}}(x) - \langle \xi - x, \nabla \psi_{\text{lb}}(x) \rangle.$$

2544 Note that g is a self-concordant function on $\text{int}(\mathcal{M}_{\mu_0}^p)^-$ (Item (ii) in Nemirovski, 2004, Proposition 2.1.1), whose Hessian
 2545 (hence, local norms and Dikin ellipsoids) coincides with that of ψ_{lb} everywhere. Moreover, g is below bounded thanks to
 2546 $(\mathcal{M}_{\mu_0}^p)^-$ being a bounded set, which implies that g attains its minimum on $\text{int}(\mathcal{M}_{\mu_0}^p)^-$ (Property VI in Nemirovski, 2004,
 2547 Section 2.2). This minimum is also unique via strict convexity. Hence, y is well-defined.

2548 The rest of the proof is similar to the proof of Lem. 13 in (Wei & Luo, 2018) and Lem. 9 in (Van der Hoeven et al.,
 2549 2023). Thanks to the strict convexity of g , to show that $y \in \mathcal{E}_{1/2}(x)$ it suffices to show that for any ξ on the boundary of
 2550 $\mathcal{E}_{1/2}(x)$, $g(x) \leq g(\xi)$; this is because $x \in \mathcal{E}_{1/2}(x)$ and $y = \arg \min_{\xi \in (\mathcal{M}_{\mu_0}^p)^-} g(\xi)$. For any such ξ on the boundary of $\mathcal{E}_{1/2}(x)$,
 2551 Taylor's theorem implies that there exists some z on the line segment between x and ξ such that

$$2552 \quad \begin{aligned} g(\xi) - g(x) &= \langle \xi - x, \nabla g(x) \rangle + \frac{1}{2} (\xi - x)^\top \nabla^2 g(z) (\xi - x) \\ &= \langle \xi - x, \ell \rangle + \frac{1}{2} (\xi - x)^\top \nabla^2 \psi_{\text{lb}}(z) (\xi - x) \\ &\geq \langle \xi - x, \ell \rangle + \frac{1}{8} (\xi - x)^\top \nabla^2 \psi_{\text{lb}}(x) (\xi - x) \\ &= \langle \xi - x, \ell \rangle + \frac{1}{8} \|\xi - x\|_x^2 \\ &\geq -\|\xi - x\|_x \|\ell\|_{x,*} + \frac{1}{8} \|\xi - x\|_x^2 \\ &= -\frac{1}{2} \|\ell\|_{x,*} + \frac{1}{32} \geq 0, \end{aligned}$$

2553 where the second equality holds since $\nabla^2 g = \nabla^2 \psi_{\text{lb}}$ and $\nabla g(x) = \ell + \nabla \psi_{\text{lb}}(x) - \nabla \psi_{\text{lb}}(x) = \ell$, the first inequality holds
 2554 via (65) and the fact that $z \in \mathcal{E}_{1/2}(x)$, the second inequality holds via the definition of a dual norm, the last equality holds
 2555 since ξ is on the boundary of $\mathcal{E}_{1/2}(x)$, and the last inequality holds via the assumption that $\|\ell\|_{x,*} \leq \frac{1}{16}$. \square

2556 For $x \in \text{int}(\mathcal{M}_{\mu_0}^p)^-$, the Minkowski function of $(\mathcal{M}_{\mu_0}^p)^-$ with the pole at x is defined as (Nemirovski, 2004, Section 3.2)

$$2557 \quad \pi_x(y) := \inf \{ t > 0 : x + t^{-1}(y - x) \in (\mathcal{M}_{\mu_0}^p)^- \}$$

2558 for $y \in (\mathcal{M}_{\mu_0}^p)^-$. The following lemma readily follows from the properties of the Minkowski function. It is used in
 2559 the analysis to handle the bias term of the standard OMD regret guarantee, which is slightly more involved in this case
 2560 considering that the comparator need not belong to the interior of $(\mathcal{M}_{\mu_0}^p)^-$, where ψ_{lb} is defined (and finite).

2585 **Lemma E.18.** Let $x \in \text{int}(\mathcal{M}_{\mu_0}^p)^-$, $y \in (\mathcal{M}_{\mu_0}^p)^-$, $\delta \in (0, 1)$, and $z := (1 - \delta)y + \delta x$. Then,

$$2587 \quad \psi_{\text{lb}}(z) \leq \psi_{\text{lb}}(x) + NSA \log \delta^{-1}.$$

2588 Further, let $\dot{x} := \arg \min_{x \in \text{int}(\mathcal{M}_{\mu_0}^p)^-} \psi_{\text{lb}}(x)$ and $\dot{z} := (1 - \delta)y + \delta \dot{x}$. Then,

$$2591 \quad D_{\psi_{\text{lb}}}(\dot{z}, \dot{x}) \leq NSA \log \delta^{-1}.$$

2593 *Proof.* Since ψ_{lb} is an NSA -self-concordant barrier for $\mathcal{M}_{\mu_0}^p$, Property II in (Nemirovski, 2004, Section 3.2) implies that

$$2594 \quad \psi_{\text{lb}}(z) \leq \psi_{\text{lb}}(x) + NSA \log \frac{1}{1 - \pi_x(z)}.$$

2597 On the other hand,

$$2599 \quad x + (1 - \delta)^{-1}(z - x) = x + (1 - \delta)^{-1}((1 - \delta)y + \delta x - x) = x + y - x = y \in (\mathcal{M}_{\mu_0}^p)^-,$$

2601 implying that $\pi_x(z) \leq 1 - \delta$. Hence, $\psi_{\text{lb}}(z) \leq \psi_{\text{lb}}(x) + NSA \log \delta^{-1}$.

2603 Next, we note that the optimality of \dot{x} implies that

$$2604 \quad D_{\psi_{\text{lb}}}(\dot{z}, \dot{x}) = \psi_{\text{lb}}(\dot{z}) - \psi_{\text{lb}}(\dot{x}) - \langle \dot{z} - \dot{x}, \nabla \psi_{\text{lb}}(\dot{x}) \rangle \leq \psi_{\text{lb}}(\dot{z}) - \psi_{\text{lb}}(\dot{x}),$$

2606 which concludes the proof when combined with the first part. \square

2608 E.3.2. REGRET ANALYSIS

2610 We are now ready to prove the regret bound of Thm. 4.5, which is stated more explicitly in the following theorem.

2612 **Theorem E.19.** Under Asm. 4.4, Alg. 4 with $\tau = \frac{\delta}{16} \sqrt{\frac{\log T}{N^3 S A T}}$ and $\delta = \min \left\{ \sqrt{\frac{17}{4L} \frac{N^{3/4} S^{3/4} A^{3/4} (\log T)^{1/4}}{T^{1/4}}}, 1 \right\}$ satisfies for
2613 any policy $\pi \in \Pi$ that

$$2615 \quad \mathbb{E}[R_T(\pi)] \leq \max \left\{ 4\sqrt{17LN}^{7/4} (SAT)^{3/4} (\log T)^{1/4}, 34\sqrt{N^5 S^3 A^3 T \log T} \right\} + 2LN.$$

2618 *Proof.* Firstly, we assert that the iterates ξ^t are well defined; similar to what was argued in the proof of Lem. E.17, the
2619 functions $\psi_{\text{lb}}(\cdot)$ and $\tau \langle g^t, \cdot \rangle + D_{\psi_{\text{lb}}}(\cdot, \xi^t)$ are self-concordant on $\text{int}(\mathcal{M}_{\mu_0}^p)^-$ (Item (ii) in Nemirovski, 2004, Proposition
2620 2.1.1) and bounded from below thanks to $(\mathcal{M}_{\mu_0}^p)^-$ being a bounded set, implying via Property VI in (Nemirovski, 2004,
2621 Section 2.2) that each of these functions attains its minimum on $\text{int}(\mathcal{M}_{\mu_0}^p)^-$, which is also unique via strict convexity. Also
2622 note that indeed $\hat{\mu}^t \in \mathcal{M}_{\mu_0}^p$ since $\hat{\xi}^t \in (\mathcal{M}_{\mu_0}^p)^-$ as we argued before presenting the algorithm.

2624 Let $\mu^* \in \arg \min_{\mu \in \mathcal{M}_{\mu_0}^p} \sum_{t=1}^T F^t(\mu)$ and $\bar{R}_T := \mathbb{E} \sum_{t=1}^T (F^t(\mu^{\pi^t, p}) - F^t(\mu^*))$, which satisfies $\bar{R}_T = \max_{\pi \in \Pi} \mathbb{E}[R_T(\pi)]$.

2625 Define $\xi^* := (1 - \delta)\Lambda_p^{-1}(\mu^*) + \delta \xi^1$, where $\delta \in (0, 1)$ is a constant to be specified later. To start with, we have that

$$2627 \quad \bar{R}_T = \mathbb{E} \sum_{t=1}^T (F^t(\mu^{\pi^t, p}) - F^t(\mu^*)) = \mathbb{E} \sum_{t=1}^T (F^t(\hat{\mu}^t) - F^t(\mu^*)) = \mathbb{E} \sum_{t=1}^T (\mathbb{F}^t(\hat{\xi}^t) - \mathbb{F}^t(\Lambda_p^{-1}(\mu^*))).$$

2630 Next, we derive that

$$\begin{aligned} 2632 & \mathbb{F}^t(\hat{\xi}^t) - \hat{\mathbb{F}}^t(\xi^t) \\ 2633 &= \mathbb{F}^t(\xi^t + \delta(\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{u}^t) - \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[\mathbb{F}^t(\xi^t + \delta(\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v}) \right] \\ 2634 &\leq L \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left\| \Lambda_p(\xi^t + \delta(\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{u}^t) - \Lambda_p(\xi^t + \delta(\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v}) \right\|_1 \\ 2635 &= \delta L \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left\| \Lambda_p((\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{u}^t) - \Lambda_p((\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v}) \right\|_1 \\ 2636 &= \delta L \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left\| \Lambda_p(\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{u}^t) - \Lambda_p(\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v}) \right\|_1 \end{aligned}$$

$$\begin{aligned}
 &\leq \delta L \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[\|\Lambda_p(\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{u}^t)\|_1 + \|\Lambda_p(\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v})\|_1 \right] \\
 &\leq 2\delta LN,
 \end{aligned}$$

where the first inequality uses the Lipschitz smoothness of F^t and the fact that $\mathbb{F}^t = F^t \circ \Lambda_p$; the second and third equalities use the fact that Λ_p is an affine function; and the last inequality holds since both $\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{u}^t$, $\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v} \in \mathcal{E}_1(\xi^t) \subset (\mathcal{M}_{\mu_0}^p)^-$, and that for any $\xi \in (\mathcal{M}_{\mu_0}^p)^-$, $\Lambda_p(\xi) \in \mathcal{M}_{\mu_0}^p$ and therefore satisfies $\|\Lambda_p(\xi)\|_1 \leq N$. We similarly derive that

$$\begin{aligned}
 &\widehat{\mathbb{F}}^t(\xi^*) - \mathbb{F}^t(\Lambda_p^{-1}(\mu^*)) \\
 &= \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[\mathbb{F}^t((1-\delta)\xi^* + \delta(\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v})) \right] - \mathbb{F}^t(\Lambda_p^{-1}(\mu^*)) \\
 &= \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[\mathbb{F}^t((1-\delta)(1-\delta)\Lambda_p^{-1}(\mu^*) + (1-\delta)\delta\xi^1 + \delta(\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v})) \right] \\
 &\quad - \mathbb{F}^t(\Lambda_p^{-1}(\mu^*)) \\
 &\leq L \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left\| (1-\delta)(1-\delta)\mu^* + (1-\delta)\delta\Lambda_p(\xi^1) + \delta\Lambda_p(\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v}) - \mu^* \right\|_1 \\
 &\leq L \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[|\delta\delta - \delta - \delta| \|\mu^*\|_1 + (1-\delta)\delta \|\Lambda_p(\xi^1)\|_1 + \delta \|\Lambda_p(\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v})\|_1 \right] \\
 &\leq L \mathbb{E}_{\mathbf{v} \in \mathbb{B}^{NS(A-1)}} \left[(\delta + \delta) \|\mu^*\|_1 + \delta \|\Lambda_p(\xi^1)\|_1 + \delta \|\Lambda_p(\xi^t + (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1/2} \mathbf{v})\|_1 \right] \\
 &\leq 2\delta LN + 2\delta LN.
 \end{aligned}$$

Hence, using also the convexity of $\widehat{\mathbb{F}}^t$, we obtain that

$$\begin{aligned}
 \bar{R}_T &\leq \mathbb{E} \sum_{t=1}^T (\widehat{\mathbb{F}}^t(\xi^t) - \widehat{\mathbb{F}}^t(\xi^*)) + 4\delta LNT + 2\delta LNT \\
 &\leq \mathbb{E} \sum_{t=1}^T \langle \nabla \widehat{\mathbb{F}}^t(\xi^t), \xi^t - \xi^* \rangle + 4\delta LNT + 2\delta LNT.
 \end{aligned} \tag{66}$$

In this proof, let $\mathcal{F}_t := \sigma(u^1, \dots, u^t)$ denote the σ -algebra generated by u^1, \dots, u^t ; and let $\mathbb{E}_t[\cdot] := \mathbb{E}[\cdot | \mathcal{F}_{t-1}]$ with \mathcal{F}_0 being the trivial σ -algebra. We then have that

$$\widehat{\mathbb{F}}^t(\xi^t) = \mathbb{E}_{\mathbf{u}^t}[g^t] = \mathbb{E}_t[g^t],$$

where the first equality follows from (63) and the the second equality holds since conditioned on \mathcal{F}_{t-1} , \mathbf{u}^t is the only source of randomness in g^t and is sampled identically and independently in every round. Using that $\xi^t - \xi^*$ is measurable with respect to \mathcal{F}_{t-1} , we then obtain that

$$\bar{R}_T \leq \mathbb{E} \sum_{t=1}^T \langle \mathbb{E}_t g^t, \xi^t - \xi^* \rangle + 4\delta LNT + 2\delta LNT = \mathbb{E} \sum_{t=1}^T \langle g^t, \xi^t - \xi^* \rangle + 4\delta LNT + 2\delta LNT.$$

Via the definition of ξ^t and the fact that $\xi^* \in \text{int}(\mathcal{M}_{\mu_0}^p)^-$, Lem. 6.16 in (Orabona, 2023) implies that

$$\sum_{t=1}^T \langle g^t, \xi^t - \xi^* \rangle \leq \frac{D_{\psi_{\text{lb}}}(\xi^*, \xi^1)}{\tau} + \sum_{t=1}^T \frac{\tau}{2} \|g^t\|_{\zeta^t, *},$$

where ζ^t lies on the line segment between ξ^t and ξ^{t+1} . We firstly observe that

$$\begin{aligned}
 \|g^t\|_{\zeta^t, *}^2 &= \left(\frac{(1-\delta)}{\delta} NS(A-1) F^t(\hat{\mu}^t) \right)^2 \cdot (\mathbf{u}^t)^\top (\nabla^2 \psi_{\text{lb}}(\xi^t))^{1/2} (\nabla^2 \psi_{\text{lb}}(\xi^t))^{-1} (\nabla^2 \psi_{\text{lb}}(\xi^t))^{1/2} \mathbf{u}^t \\
 &= \left(\frac{(1-\delta)}{\delta} NS(A-1) F^t(\hat{\mu}^t) \right)^2 \leq \frac{1}{\delta^2} N^4 S^2 A^2,
 \end{aligned}$$

where we have used that $F^t(\hat{\mu}^t) \leq N$. Hence, if

$$\tau \leq \frac{\delta}{16N^2SA}, \tag{67}$$

2695 then $\tau \|g_t\|_{\xi^t, *}$ $\leq 1/16$. Consequently, Lem. E.17 (with $x = \xi^t$, $y = \xi^{t+1}$, and $\ell = \tau g_t$) would assert that $\xi^{t+1} \in \mathcal{E}_{1/2}(\xi^t)$;
 2696 and hence, $\zeta^t \in \mathcal{E}_{1/2}(\xi^t)$. It would then hold via (65) that

$$2697 \quad \|g_t\|_{\zeta^t, *}^2 \leq \frac{1}{(1 - \|\zeta^t - \xi^t\|_{\xi^t})^2} \|g_t\|_{\xi^t, *}^2 \leq 4 \|g_t\|_{\xi^t, *}^2,$$

2700 On the other hand, via Lem. E.18 and the definitions of ξ^1 and ξ^* , we have that

$$2702 \quad D_{\psi_b}(\xi^*, \xi^1) \leq NSA \log \dot{\delta}^{-1}.$$

2704 Hence, conditioned on (67), we obtain the following regret bound

$$2706 \quad \bar{R}_T \leq \frac{NSA \log \dot{\delta}^{-1}}{\tau} + \frac{2\tau}{\delta^2} N^4 S^2 A^2 T + 4\delta LNT + 2\dot{\delta} LNT. \quad (68)$$

2709 Setting

$$2711 \quad \dot{\delta} = \frac{1}{T}, \quad \tau = \frac{\delta}{16} \sqrt{\frac{\log T}{N^3 SAT}}, \quad \text{and} \quad \delta = \min \left\{ \sqrt{\frac{17}{4L} \frac{N^{3/4} S^{3/4} A^{3/4} (\log T)^{1/4}}{T^{1/4}}}, 1 \right\};$$

2714 we obtain that

$$\begin{aligned} 2715 \quad \bar{R}_T &\leq \frac{NSA \log T}{\tau} + \frac{2\tau}{\delta^2} N^4 S^2 A^2 T + 4\delta LNT + 2LN \\ 2716 &\leq \frac{17}{\delta} \sqrt{N^5 S^3 A^3 T \log T} + 4\delta LNT + 2LN \\ 2717 &\leq \max \left\{ 4\sqrt{17LN^{7/4}} (SAT)^{3/4} (\log T)^{1/4}, 34\sqrt{N^5 S^3 A^3 T \log T} \right\} + 2LN. \end{aligned} \quad (69)$$

2722 If $T \geq NSA \log(T)$, then our choice of τ indeed satisfies (67):

$$2724 \quad \tau = \frac{\delta}{16} \sqrt{\frac{\log T}{N^3 SAT}} \leq \frac{\delta}{16N^2 SA}.$$

2726 Otherwise, we can fall back to the trivial regret bound

$$2728 \quad \bar{R}_T \leq NT \leq N^2 SA \log(T),$$

2730 which is dominated by the bound in (69); hence, the theorem follows. \square