

NO-REGRET LEARNING WITH REVEALED TRANSITIONS IN ADVERSARIAL MARKOV DECISION PROCESSES

Anonymous authors

Paper under double-blind review

ABSTRACT

When learning in Adversarial Markov Decision Processes (MDPs), agents must deal with a sequence of arbitrarily chosen transition models and losses. In this paper, we consider the setting in which the transition model chosen by the adversary is revealed at the end of each episode. We propose the notion of *smoothed MDP* whose transition model aggregates with a generic function f_t the ones experienced so far. Coherently, we define the concept of *smoothed regret*, and we devise Smoothed Online Mirror Descent (SOMD), an enhanced version of OMD that leverages a novel regularization term to effectively learn in this setting. For specific choices of the aggregation function f_t defining the smoothed MDPs we retrieve, under full-feedback, a regret bound of order $\tilde{O}(L^{3/2}\sqrt{TL} + L\bar{C}_f^P)$ where T is the number of episodes, L is the horizon of the episode, and \bar{C}_f^P is a novel index of the degree of maliciousness of the adversarially chosen transitions. Under bandit feedback on the losses, we obtain a bound of order $\tilde{O}(L^{3/2}\sqrt{XAT} + L\bar{C}_f^P)$ using a simple importance weighted estimator on the losses.

1 INTRODUCTION

Reinforcement Learning (RL) studies sequential decision-making problems often modeled through the framework of the Markov Decision Processes (MDPs, Puterman, 2014). This framework allows the application of RL to a large variety of challenging problems and paved the way for the growing success of RL we witnessed in the last decade (e.g., Kormushev et al., 2013; Nagabandi et al., 2018; Dulac-Arnold et al., 2021). Nevertheless, MDPs are based on some grounding assumptions that may limit their modeling power in real-world scenarios. In particular, they assume that the environment dynamics P (i.e., transition model) and loss¹ ℓ are fixed throughout the whole interaction. However, in the real world, such elements may change due to external factors which might be the effect of *nature* (e.g., system anomaly, aging effects) or *strategic* actors (e.g., adversarial attacks). While the former case is usually captured by *non-stationary* MDPs (Lecarpentier & Rachelson, 2019; Cheung et al., 2020), the latter scenario is more challenging as it assumes the presence of another agent (i.e., an adversary) acting with an objective possibly conflicting with that of the agent.

Since the early work of Even-Dar et al. (2009), this class of problems has been addressed drawing inspiration from Online Learning (OL) literature (Orabona, 2019). Adversarial Markov Decision Processes (AMDPs, Even-Dar et al., 2009) have been designed to model the scenario in which the agent interacts for $T \in \mathbb{N}$ rounds facing an adversarially chosen MDP \mathcal{M}_t at every round $t \in \llbracket T \rrbracket$. Here, the performance of the agent’s policy π_t is evaluated in terms of the expected regret competing against a *fixed comparator* policy π° .²

When the transition model P is known or fixed (possibly stochastic) and the loss ℓ_t is adversarially chosen, several works (Zimin & Neu, 2013; Rosenberg & Mansour, 2019; Jin et al., 2020) have

¹We comply with the convention of the adversarial literature of using losses instead of rewards.

²This is a notion of *static regret* which is different from the notion of *dynamic regret* typically adopted in non-stationary MDPs in which the comparator policy π_t° is allowed to change over rounds. It is known that, even in the simpler bandit setting, when the environment is adversarial, the no-regret property is not achievable for the dynamic regret (Bubeck et al., 2012).

054 achieved compelling regret guarantees of order $\tilde{O}(\sqrt{T})$. However, when the adversary is allowed to
 055 select the transition model P_t at every round, the problem acquires new significant *computational*
 056 and *learning* challenges. In the *full-feedback* setting, where the transition model P_t is revealed
 057 to the learner at the end of round t , Abbasi Yadkori et al. (2013) showed a $\tilde{O}(\sqrt{T})$ regret bound
 058 with a computationally inefficient algorithm running (a variation of) Exp3 (Auer et al., 1995) on a
 059 covering of all stochastic policies. Indeed, the computational barrier has been formalized by (Liu
 060 et al., 2022) showing that no algorithm can simultaneously achieve the no-regret property and be
 061 computationally efficient. The scenario is even less encouraging under the *bandit feedback*, where the
 062 learner observes only the collected experience. Here, Tian et al. (2021) proved a regret lower bound
 063 of order $\Omega(\min\{T, 2^H\})$, being H the episode horizon. This formalizes a learning barrier showing
 064 that exponential dependence on the horizon is unavoidable.³

065 Accepting these impossibility results, the research effort has been directed towards the design of
 066 *computationally efficient* algorithms, allowing their regret to depend on the degree of adversary
 067 maliciousness in the choice of the transition model P_t . Inspired from the *corruption-robust* RL
 068 (Lykouris et al., 2021; Chen et al., 2021; Wei et al., 2022), the degree of maliciousness is formalized
 069 in a parameter C^P that quantifies the cumulative dissimilarity between the experienced transition
 070 models P_t and a nominal one P :

$$071 \quad C^P := \min_{P \in \mathcal{P}} \sum_{t \in \llbracket T \rrbracket} \sum_{k \in \llbracket 0, L-1 \rrbracket} \max_{(x,a) \in \mathcal{X}_k \times \mathcal{A}} \|P(\cdot|x, a) - P_t(\cdot|x, a)\|_1. \quad (1)$$

074 In the worst case, this C^P term leads to linear regret but, in general cases, may be significantly smaller.
 075 With these premises, Jin et al. (2023) proposes a state-of-the-art algorithm working under bandit
 076 feedback for both the losses and the transition models. With the knowledge of C^P , the algorithm
 077 achieves a regret dependence that gracefully degrades with the level of corruption $\tilde{O}(\sqrt{T} + C^P)$.
 078 However, when C^P is unknown, the resulting algorithm necessitates a complex inner subroutine
 079 (based on approaches from the *Corral* literature) making the final regret guarantee less explicit and,
 080 possibly, affected by large constants. Furthermore, the practicality of the algorithm in terms of
 081 computational complexity remains uncertain.

082 These results, therefore, leave several open questions. First, it is not clear whether the maliciousness
 083 parameter C^P inherited from the corruption-robust RL appropriately captures the challenges of the
 084 AMDP learning problem. Indeed, comparing P_t against a nominal MDP does not comply with the
 085 conventional behavior of an OL algorithm, i.e., adapt according to the experience observed so far.
 086 *Can different definitions of the degree of maliciousness highlight new properties that OL algorithms*
 087 *for AMDPs enjoy?* Second, the approach proposed (Jin et al., 2023) aspires to directly address the
 088 bandit feedback on both the losses and transitions models. This leaves open the question of whether
 089 with transition models revealed at the end of the episode (and possibly full or bandit feedback on the
 090 losses), tighter results can be achieved, especially without the knowledge of the corruption parameter
 091 C^P . *When the transition models are revealed at the end of the episode, can OL algorithms achieve*
 092 *better performances without the knowledge of C^P ?* This paper aims to address these open questions.

093 **Original Contributions.** The contributions of the paper can be summarized as follows.

- 094 • In Section 3, we introduce the novel concept of *smoothed* MDP and the related *smoothed* regret.
 095 Since OL algorithms make decisions based on past experience, the smoothed MDP is defined
 096 through a transition model $\bar{P}_t = f_t(P_1, \dots, P_t)$ that aggregates with a generic function f_t the
 097 previously experienced ones P_1, \dots, P_t . This allows introducing a novel quantification of the
 098 adversary maliciousness \bar{C}_f^P by evaluating the dissimilarity between the chosen transition model
 099 P_t and that of the smoothed MDP \bar{P}_t . This constant matches the corruption parameter C^P up
 100 to logarithmic terms when using an *averaging smoothing function* (i.e., \bar{P}_t is the average of
 101 P_1, \dots, P_t). Coherently, we define a smoothed regret measuring the learner’s performance against
 102 a comparator policy π° acting in such a smoothed MDP.
- 103 • In Section 4, we propose a novel regret-minimization algorithm *Smoothed Online Mirror Descent*
 104 (SOMD). Our approach is built upon a simple yet novel instance of Online Mirror Descent (OMD)
 105 with a well-calibrated entropic regularization. Interestingly, the algorithm does not require any
 106 knowledge about \bar{C}_f^P . We analyze SOMD for general smoothing functions f_t and we show that
 107

³This implies that the no-regret property is achievable only when $H = o(\log T)$ which is often unrealistic.

it achieves $\tilde{O}(L^{3/2}\sqrt{T})$ smoothed regret, which translates into $\tilde{O}(L^{3/2}\sqrt{T} + LC_f^P)$ for the choice of the averaging smoothing function. The regret analysis requires managing non-trivial aspects, which also represents a key component of our contribution.

- In Section 5, through an importance-weighted (IW) estimator of the loss function, we extend SOMD to the bandit feedback of losses, showing comparable performances on the regret. Still, we do not require the knowledge of the degree of corruption.

2 PROBLEM FORMULATION

Notation and Definitions. In the following, given $a, b \in \mathbb{N}$ with $a \leq b$, we denote with $\llbracket a \rrbracket := \{1, 2, \dots, a\}$ and with $\llbracket a, b \rrbracket := \{a, a + 1, \dots, b\}$. Given two vectors \mathbf{x} and \mathbf{y} , we denote with $\langle \mathbf{x}, \mathbf{y} \rangle$ the inner product $\mathbf{x}^\top \mathbf{y}$. Given a set \mathcal{A} , we denote with $|\mathcal{A}|$ the cardinality of the set, and with $\Delta(\mathcal{A})$ the probability simplex over the set. For two generic (discrete) distributions $q, q' \in \Delta(\mathcal{A})$ we define the *negative Shannon Entropy* as $\psi(q) = \sum_{a \in \mathcal{A}} q(a) \log q(a)$ and the *Bregman Divergence* $D_\psi(q, q')$ which, for the negative Shannon Entropy, corresponds to the generalized KL-divergence $D_\psi(q, q') = \sum_{a \in \mathcal{A}} (q(a) \log(q(a)/q'(a)) - q(a) + q'(a))$. Finally, we denote with $h(p) = -p \log p - (1-p) \log(1-p)$ the binary entropy function for $p \in [0, 1]$.

Setting. We consider the framework of episodic loop-free Markov Decision Processes (MDPs, Puterman, 2014). Specifically, we assume the agent is interacting with a sequence of $T \in \mathbb{N}$ MDPs $\{\mathcal{M}_t\}_{t \in \llbracket T \rrbracket}$ with $\mathcal{M}_t = (\mathcal{X}, \mathcal{A}, P_t, \ell_t)$. Here, \mathcal{X} is the state space, \mathcal{A} is the action space, P_t is the transition function $P_t : \mathcal{X} \times \mathcal{A} \rightarrow \Delta(\mathcal{X})$ such that $P_t(x'|x, a)$ is the probability of transitioning to state x' after taking action a in state x , and ℓ_t is the loss function $\ell_t : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$, defined such that $\ell_t(x, a)$ is the loss the agent incurs selecting a in x . We consider *finite* state and action sets with cardinality $|\mathcal{X}| < \infty$ and $|\mathcal{A}| = A < \infty$, respectively. As common in the literature (Jin et al., 2020; 2023) and w.l.o.g. the state space is assumed to be decomposed into $L + 1$ disjoint layers, namely $\mathcal{X} = \bigcup_{k \in \llbracket 0, L \rrbracket} \mathcal{X}_k$, and $\mathcal{X}_i \cap \mathcal{X}_j = \{\}, \forall i, j \in \llbracket 0, L \rrbracket, i \neq j$. Furthermore, the first and the last layers are assumed to be singletons: $\mathcal{X}_0 = \{x_0\}$ and $\mathcal{X}_L = \{x_L\}$. The layered structure also imposes $P_t(x'|x, a) > 0$ only if $x \in \mathcal{X}_k$ and $x' \in \mathcal{X}_{k+1}$, for some $k \in \llbracket 0, L - 1 \rrbracket$. Finally, to ease the exposition, we assume the cardinality of each layer to be $|\mathcal{X}_k| = X_k \leq X$.

Interaction Protocol. The interaction proceeds as follows. An adversary selects obliviously the sequence of transition models and losses $\{(P_t, \ell_t)\}_{t \in \llbracket T \rrbracket}$ before the interaction with the agent starts. Then, in each episode $t \in \llbracket T \rrbracket$, the agent sequentially decides which action to play following a *stochastic Markovian policy* $\pi_t : \mathcal{X} \rightarrow \Delta(\mathcal{A})$, where $\pi_t(a|x)$ denotes the probability of playing action a in state x . More in detail, starting from the fixed initial state $x_{t,0} = x_0$, for each layer $k \in \llbracket L - 1 \rrbracket$, the agent selects the action $a_{t,k} \sim \pi_t(\cdot|x_{t,k})$, the environment evolves to the next state $x_{t,k+1} \sim P_t(\cdot|x_{t,k}, a_{t,k})$, the agent observes the loss $\ell_t(x_{t,k}, a_{t,k})$, and the interaction proceeds until the terminal state $x_{t,L} = x_L$ is reached. At the end of each episode t , the full transition model P_t is revealed to the learner. Furthermore, in the *full-feedback model*, the full loss ℓ_t is revealed to the agent, while in the *bandit-feedback model* the loss is not revealed.

Occupancy Measures. As customary in the literature (Zimin & Neu, 2013; Jin et al., 2020), the problem will be treated in the space of *occupancy measures*. For a generic transition function P and a policy π the occupancy measure $q^{P,\pi}$ is defined as:

$$q^{P,\pi}(x, a, x') := \mathbb{P}[x_k = x, a_k = a, x_{k+1} = x' | P, \pi], \quad (2)$$

where $x \in \mathcal{X}_k, x' \in \mathcal{X}_{k+1}$, and $a \in \mathcal{A}$. This quantity represents the marginal probability of experiencing the transition (x, a, x') when deploying policy π in an MDP with transition model P . Similarly, we make use of $q^{P,\pi}(x, a) = \mathbb{P}[x_k = x, a_k = a | P, \pi] = \sum_{x' \in \mathcal{X}_{k+1}} q(x, a, x')$, and $q^{P,\pi}(x) = \mathbb{P}[x_k = x | P, \pi] = \sum_{a \in \mathcal{A}} q^{P,\pi}(x, a)$. Importantly, according to (Rosenberg & Mansour, 2019, Lemma 3.1), any valid occupancy measure is such that each layer of the MDP is visited exactly once, and thus, for every $k \in \llbracket 0, L - 1 \rrbracket$ it holds that $\sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k+1}} q(x, a, x') = 1$, and for every $k \in \llbracket L - 1 \rrbracket$ and every state $x \in \mathcal{X}_k$ it holds $\sum_{x' \in \mathcal{X}_{k-1}} \sum_{a \in \mathcal{A}} q(x', a, x) = \sum_{x' \in \mathcal{X}_{k+1}} \sum_{a \in \mathcal{A}} q(x, a, x')$. This implies that the probability of entering a state coming from the previous layer is equal to the probability of leaving that same state when going to the next layer.

Finally, denoting Δ as the set of occupancies q satisfying the previously defined properties, we have that every $q \in \Delta$ induces both a transition function P^q and a policy π^q , computed as:

$$P^q(x'|x, a) = \frac{q(x, a, x')}{\sum_{x' \in \mathcal{X}_{k(x)+1}} q(x, a, x')}, \quad \pi^q(a|x) = \frac{\sum_{x' \in \mathcal{X}_{k(x)+1}} q(x, a, x')}{\sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k(x)+1}} q(x, a, x')}, \quad (3)$$

where $k(x) \in \llbracket 0, L \rrbracket$ identifies the index of the layer to which the state x belongs. Throughout the analysis, when we need to refer to the occupancy at layer k , we use a superscript e.g., q^k in favor of compactness. For a transition function P , we denote as $\Delta(P) \subseteq \Delta$ the set of occupancies that induce exactly the transition P . Similarly, given a set of transition functions $\mathcal{P}' \subseteq \mathcal{P}$, being \mathcal{P} the set of all transition models, we denote as $\Delta(\mathcal{P}')$ the set of q 's such that $P^q \in (\mathcal{P}')$. Finally, we denote with Π the set of Markovian stochastic policies.

Learning Objectives. Let $V_t^\pi(x_0) = \mathbb{E}[\sum_{k \in \llbracket 0, L-1 \rrbracket} \ell_t(x_{t,k}, a_{t,k}) | P_t, \pi, x_0]$ be the expected cumulative loss suffered by the agent experiencing the trajectory $\{(s_{t,k}, a_{t,k})\}_{k \in \llbracket L-1 \rrbracket}$ generated under the state-action distribution induced by transition function P_t and policy π in \mathcal{M}_t . The *expected static regret* of the agent against any comparator policy $\pi^\circ \in \Pi$ is defined as:⁴

$$\mathcal{R}_T(\pi^\circ) := \mathbb{E} \left[\sum_{t \in \llbracket T \rrbracket} V_t^{\pi^\circ}(x_0) - \sum_{t \in \llbracket T \rrbracket} V_t^\pi(x_0) \right], \quad (4)$$

where the expectation taken w.r.t. the internal randomization of the algorithm and on the possible stochasticity of the environment. With occupancy measures, the expected cumulative loss can be conveniently rewritten as $V_t^\pi(x_0) = \langle q^{P_t, \pi}; \ell_t \rangle$. This allows to frame the task of the agent to select occupancy measures q_t instead of policies. Similarly, the expected regret can be expressed as $\mathcal{R}_T(\pi^\circ) = \mathbb{E}[\sum_{t \in \llbracket T \rrbracket} \langle q^{P_t, \pi_t} - q^{P_t, \pi^\circ}, \ell_t \rangle]$.

3 SMOOTHED MDPs AND SMOOTHED REGRET

As mentioned previously, our goal is to design computationally efficient algorithms for adversarial MDPs that achieve a regret scaling with a notion of the degree of maliciousness of the adversary in selecting the transition functions. However, unlike Jin et al. (2023), we aim for such a performance guarantee without introducing a notion of nominal transition function.

Smoothed MDP. With such an objective in mind, starting from the sequence of T MDPs $\{\mathcal{M}_t\}_{t \in \llbracket T \rrbracket}$, we define a *smoothed MDP* as $\{\overline{\mathcal{M}}_t\}_{t \in \llbracket T \rrbracket}$ with $\overline{\mathcal{M}}_t = (\mathcal{X}, \mathcal{A}, \overline{P}_t, \ell_t)$ in which the transition model P_t gets replaced with \overline{P}_t , named *smoothed transition model*. In general, \overline{P}_t can be any function of the history of transition functions. Formally, for every $t \in \llbracket T \rrbracket$, let $f_t : \mathcal{P}^t \rightarrow \mathcal{P}$ and $\overline{P}_t := f_t(P_1, \dots, P_t)$. This allows the introduction of a novel index for measuring the maliciousness of the adversary in selecting the transitions, that we call *smoothed transition error*, defined as follows:

$$\overline{C}_f^P := \sum_{t \in \llbracket T \rrbracket} \sum_{k \in \llbracket 0, L-1 \rrbracket} \max_{(x,a) \in \mathcal{X}_k \times \mathcal{A}} \|\overline{P}_t(\cdot|x, a) - P_t(\cdot|x, a)\|_1. \quad (5)$$

For adequately chosen smoothing functions, this constant interpolates between 0 when the adversary is absent, i.e., constant transitions, and $2LT$ in the worst case. As we shall see, when $\overline{C}_f^P = 0$, we will incur a regret of $\tilde{O}(\sqrt{T})$ even when the loss functions are completely arbitrary. To exemplify the opportunities of the smoothed transition error, we discuss the following examples.

Example 3.1. *Let us consider the smoothing function f_t so that $\overline{P}_t = P_{t-1}$, i.e., the smoothed MDP $\overline{\mathcal{M}}_t$ has the last revealed transition model. In such a case, the smoothed transition error reduces to the model variation V_T common in the non-stationary literature as in Cheung et al. (2020):*

$$\overline{C}_f^P = \sum_{t \in \llbracket T \rrbracket} \sum_{k \in \llbracket 0, L-1 \rrbracket} \max_{(x,a) \in \mathcal{X}_k \times \mathcal{A}} \|P_t(\cdot|x, a) - P_{t-1}(\cdot|x, a)\|_1 = V_T. \quad (6)$$

Example 3.2 (Average Smoothing Function). *Let us consider the smoothing function f_t so that $\overline{P}_t = \frac{1}{t} \sum_{t' \in \llbracket t \rrbracket} P_{t'}$, i.e., the smoothed MDP $\overline{\mathcal{M}}_t$ has as transition function the average of the*

⁴Usually, the comparator is assumed to be the optimal policy in hindsight: $\pi^* \in \arg \min_{\pi \in \Pi} \sum_{t \in \llbracket T \rrbracket} V_t^\pi(x_0)$.

transition functions experienced so far. In this case, we can conveniently relate our smoothed transition error with the C^P of Jin et al. (2023):

$$\frac{1}{\log T + 2} \leq \frac{\overline{C}_f^P}{C^P} \leq \log T + 2. \quad (7)$$

This double-sided inequality is proved in Lemma A.1 and shows the equivalence of the two maliciousness measures, up to logarithmic terms.

Together with the smoothed transition error, we introduce an index of the variability of the smoothed MDP between consecutive episodes, named *smoothed transition variability*, defined as follows:

$$\overline{D}_f^P := \sum_{t \in \llbracket 2, T \rrbracket} \max_{x, a \in \mathcal{X} \times \mathcal{A}} \|\overline{P}_t(\cdot|x, a) - \overline{P}_{t-1}(\cdot|x, a)\|_1. \quad (8)$$

Smoothed Regret. The notion of smoothed MDP allows us to rewrite the static regret in (4) as

$$\mathcal{R}_T(\pi^\circ) = \mathbb{E} \left[\underbrace{\sum_{t \in \llbracket T \rrbracket} \langle q^{\overline{P}_t, \pi_t} - q^{\overline{P}_t, \pi^\circ}; \ell_t \rangle}_{\text{Smoothed Regret } \overline{\mathcal{R}}_T(\pi^\circ)} + \underbrace{\sum_{t \in \llbracket T \rrbracket} \langle q^{P_t, \pi_t} - q^{\overline{P}_t, \pi_t}; \ell_t \rangle}_{\text{(Policy) Proxy Regret}} + \underbrace{\sum_{t \in \llbracket T \rrbracket} \langle q^{\overline{P}_t, \pi^\circ} - q^{P_t, \pi^\circ}; \ell_t \rangle}_{\text{(Comparator) Proxy Regret}} \right].$$

The *smoothed regret* $\overline{\mathcal{R}}_T(\pi^\circ)$, accounts for the regret incurred by the agent when acting in the smoothed MDP $\overline{\mathcal{M}}$. On the other hand, the two *proxy regrets* account for the fact that both the agent’s policy and the comparator policy are employed in the true MDP. Thus, the latter depends on the adversary and captures its maliciousness, while the *smoothed regret* is fully dependent on the algorithm. $\overline{\mathcal{R}}_T(\pi^\circ)$ will be treated in the next sections while the following result shows how the smoothed transition error can be employed to bound the proxy regret terms.

Lemma 3.1 (Proxy Regret Upper Bound). *For any policy sequence $\{\pi_t\}_{t=1}^T$, and loss functions $\{\ell_t\}_{t=1}^T$ such that $\ell_t : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$ for any $t \in \{1, \dots, T\}$ it holds that:*

$$\text{Proxy Regret} = \sum_{t \in \llbracket T \rrbracket} \langle q^{\overline{P}_t, \pi_t} - q^{P_t, \pi_t}; \ell_t \rangle \leq L \overline{C}_f^P. \quad (9)$$

A reference to the proof can be found in Appendix A. Having highlighted the role of proxy regrets, in the following, we design and analyze our algorithm, namely *smoothed OMD*.

4 SMOOTHED OMD UNDER FULL-FEEDBACK ON LOSSES AND REVEALED TRANSITIONS

In this section, we first present a novel, smoothed, version of OMD, namely *Smoothed Online Mirror Descent* (SOMD) that exploits the intrinsic structure of smoothed MDPs, and, then, we provide the analysis of its no-smooth-regret property. We start considering generic smoothing functions f_t , then we focus on the average smoothing function (Example 3.2).

Algorithm Design. Coherently with the OMD algorithmic blueprint, at the end of each episode $t \in \llbracket T \rrbracket$, SOMD (Algorithm 1) computes an occupancy measure q_{t+1} that trades off between minimizing the loss of the round ℓ_t and not diverging excessively from a specific regularization reference \overline{q}_t . In mathematical terms, SOMD solves the constrained convex program:

$$q_{t+1} = \arg \min_{q \in \Delta(\overline{P}_t)} \langle q, \ell_t \rangle + \frac{1}{\eta} D_\psi(q, \overline{q}_t), \quad (10)$$

where $\eta > 0$ is a regularization hyperparameter. It is worth noting that the resulting occupancy q_{t+1} is constrained into $\Delta(\overline{P}_t)$, i.e., the set of occupancies realizable with the smoothed transition model \overline{P}_t . More importantly, the resulting policy $\pi_{t+1} := \pi^{q_{t+1}}$ will be evaluated in the environment paired with the smoothed transition model \overline{P}_{t+1} , leading to the occupancy $q^{\overline{P}_{t+1}, \pi_{t+1}}$. In general, we have that $q^{\overline{P}_{t+1}, \pi_{t+1}} \neq q_{t+1} = q^{\overline{P}_t, \pi_{t+1}}$. Intuitively, our program in Equation (10) “pretends” that π_{t+1} will be played in the current smoothed MDP \overline{P}_t , instead it will be played in \overline{P}_{t+1} . This

Algorithm 1: Smoothed Online Mirror Descent in Full-Feedback

Input : state space \mathcal{X} , action space \mathcal{A} , episode number T , learning rate $\eta > 0$, mixing parameter $\alpha \in [0, 1]$, smoothing functions $\{f_t\}_{t \in [T]}$.

Initialize : Set $\pi_1 = \pi^{q_1}$, $q_1(x, a, x') = u(x, a, x') \forall k \in [L - 1], (x, a, x') \in \mathcal{X}_k \times \mathcal{A} \times \mathcal{X}_{k+1}$.

```

1 for  $t = 1, \dots, T$  do
2   Execute policy  $\pi_t$  in  $\mathcal{M}_t$  and observe  $(P_t, \ell_t)$ .
3   Compute smoothed transition  $\bar{P}_t = f_t(P_1, \dots, P_t)$ .
4   Compute smoothed regularization point  $\bar{q}_t = (1 - \alpha)q_t + \alpha u$ 
5   Perform mirror descent step  $q_{t+1} = \arg \min_{q \in \Delta(\bar{P}_t)} \langle q, \ell_t \rangle + \frac{1}{\eta} D_\psi(q, \bar{q}_t)$ 
6   Update policy  $\pi_{t+1} = \pi^{q_{t+1}}$ 
7 end

```

represents a fundamental feature of SOMD that deviates from the classical OMD algorithms, which are often designed to deal with a fixed decision set $\Delta(P)$ (Jin et al., 2020) or a sequence of nested sets (Jin et al., 2023). SOMD manages the mismatch between such domains by: (i) leveraging the transition variability of smooth MDPs, encoded in term \bar{D}_f^P ; (ii) computing a *smoothed regularization reference*, \bar{q}_t , defined as a mixture between the SOMD decision at the previous step, q_t , and the uniform occupancy measure $u(x, a, x') = 1/(X_k A X_{k+1})$ for every $k \in [0, L - 1]$. Specifically, $\bar{q}_t = (1 - \alpha)q_t + \alpha u$, where $\alpha \in [0, 1]$ is a hyperparameter to be specified later that acts as a further source of regularization.⁵ The choice of $\alpha > 0$ will generate a bias term whose effect on the regret will be controlled through a proper selection of the value of α . Clearly, as supported by intuition, the solution delivered by the program in Equation (10) will be a good representative of the actual occupancy only when two consecutive smoothed MDPs are sufficiently similar $\bar{P}_t \approx \bar{P}_{t+1}$. This is encoded in the variability term \bar{D}_f^P that emerges in the regret analysis.

Smoothed Regret Analysis. The following provides the smoothed regret upper bound for SOMD.

Theorem 4.1 (Smoothed-Regret Bound for SOMD under full-feedback). *Let $\eta = \sqrt{(10L \log(2X^2 AT) \rho_T^f)}/T$ and $\alpha = 1/(1+T)$, then for any comparator policy $\pi^\circ \in \Pi$ Algorithm 1 suffers a smoothed regret of:*

$$\bar{\mathcal{R}}_T(\pi^\circ) \leq \mathcal{O} \left(L^2 \bar{D}_f^P + L^{3/2} \sqrt{T \log(X^2 AT) \rho_T^f} \right), \quad (11)$$

where $\rho_T^f := \log(T) + \bar{D}_f^P + \mathcal{H}_T(\bar{D}_f^P)$ and $\mathcal{H}_T(\bar{D}_f^P) = TLh\left(\frac{(L^2+L)\bar{D}_f^P}{2TL}\right)$.

As one could expect, the behavior of the constant $\mathcal{H}_T(\bar{D}_f^P)$ depends on the choice of the smoothing function to be used. For specific choices of smoothing functions, as it is for average smoothing functions, this term will be sub-linear in T . For the average choice of Example 3.2, we have that $\mathcal{H}_T(\bar{D}_f^P) = O(L^2 \log(T)^2)$; the interested reader can refer to Appendix C for in-detail analysis.

Now, we provide the reader with additional insight into our approach, by specializing the analysis and results to the specific class of smoothing functions that satisfy $\bar{P}_t := \frac{1}{t} \sum_{t' \in [t]} P_{t'}$ called “average smoothing function”. With such an additional structure, our algorithm is able to guarantee,

Corollary 4.2 (Smoothed-Regret Bound for SOMD under full-feedback and average smoothing). *Let $\eta = 3\sqrt{(2L \log(2X^2 AT) \log(T))}/T$, $\alpha = 1/(T + 1)$ and smoothing functions such that $\bar{P}_t := \frac{1}{t} \sum_{t' \in [t]} P_{t'}$, for any comparator policy $\pi^\circ \in \Pi$, Algorithm 1 suffers a smoothed regret of:*

$$\bar{\mathcal{R}}_T(\pi^\circ) \leq \mathcal{O} \left(L^2 \log(T) + L^{3/2} \sqrt{T \log(X^2 AT) \log(T)} \right). \quad (12)$$

While the full derivation can be found in Appendix B, here we outline the most relevant steps.

Proof Sketch. From now on we will overload the notation with $q_t^{\bar{P}_{t-1}, \pi_t} = q_t$ and $q_t^{\bar{P}_t, \pi^\circ} = q_t^\circ$ for compactness. SOMD computes q_t based on \bar{P}_{t-1} . Thus we isolate the part of the regret solely

⁵Parameter α is needed for technical reasons that will become clear in the proof of Theorem 4.1.

dependent on this algorithmic choice:

$$\bar{\mathcal{R}}_T(\pi^\circ) = \mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_t, \pi_t} - q^{\bar{P}_t, \pi^\circ}; \ell_t \rangle \right] = \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t - q_t^\circ; \ell_t \rangle \right]}_{\text{Algorithmic Regret } (\bar{\mathcal{R}}_T^A)} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_t, \pi_t} - q^{\bar{P}_{t-1}, \pi_t}; \ell_t \rangle \right]}_{\text{Update Regret } (\bar{\mathcal{R}}_T^U)}.$$

The *Update Regret* captures the mismatch in using \bar{P}_{t-1} to compute π^{q_t} and then using the same policy in \bar{P}_t . This term is affected by the magnitude of \bar{D}_f^P , thus by the “slowly”-changing behavior of the smoothed MDP, as analyzed in Lemma B.1. The *Algorithmic Regret* is the one specifically controlled by the algorithm, and it can be decomposed into the following terms:

$$\bar{\mathcal{R}}_T^A = \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \ell_t \rangle \right]}_{\text{Descent Regret } (\bar{\mathcal{R}}_T^D)} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t - \bar{q}_t; \ell_t \rangle \right]}_{\text{Regularization Regret } (\bar{\mathcal{R}}_T^R)}.$$

where the *Regularization Regret* describes the degree to which \bar{q}_t is different from q_t and it is controlled by the mixing coefficient α (Lemma B.3). The *Descent Regret* captures how performing an OMD step in the smoothed MDP affects the overall performance and satisfies:

$$\eta \bar{\mathcal{R}}_T^D \leq \underbrace{\mathbb{E} \left[\sum_{t=1}^T D_\psi(\bar{q}_t, \tilde{q}_{t+1}) \right]}_{\text{“Stability” term}} + \underbrace{\frac{\alpha}{(1-\alpha)} \sum_{t=1}^T D_\psi(q_t^\circ, u)}_{\text{“Residual” term}} + \underbrace{\sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \bar{q}_{t+1}))}_{\text{“Penalty” term}},$$

where \tilde{q}_{t+1} is the solution to the unconstrained OMD problem, as can be seen in Lemma B.4. The related “Stability” term can be bounded by standard OMD analysis as done in Jin et al. (2020). The “Residual” term catches the effect of mixing the output of the descent step with uniform distributions but it can be bounded pretty easily as in Lemma B.9. Finally and more interestingly, the “Penalty” term is due to the presence of a regularizer in optimizing the cumulative loss (see Lemma B.10). In standard OL analysis, a similar term can be easily bounded using telescoping arguments. However, the time-varying nature of smooth MDPs prevents us from using such arguments and bounding this term requires some machinery. Specifically, we first rewrite the penalty term as,

$$\sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \bar{q}_{t+1})) \leq D_\psi(q_1^\circ, \bar{q}_1) + \sum_{t=2}^T D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_{t-1}^\circ, \bar{q}_t) \quad (13)$$

While the first term can be easily bounded by the maximum range of the regulariser, as shown in Lemma B.11, the second term requires more machinery despite its non-telescoping behavior. First, we leverage the properties of the KL-divergence to obtain:

$$\sum_{t=2}^T D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_{t-1}^\circ, \bar{q}_t) \leq \sum_{t=2}^T |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| + \sum_{t=2}^T \|\log(\bar{q}_t)\|_\infty \|q_t^\circ - q_{t-1}^\circ\|_1.$$

The second summation can be easily bounded via Lemmas B.21 and B.20. For the first summation, we employ Theorem 3 by Sason (2013) to bound the absolute entropy differences: we first identify two different time regimes separated by $\bar{t} := \lceil L/1 - \frac{1}{X^2 A} \rceil$. Now, for $t < \bar{t}$ we simply apply the above-cited theorem. For $t \geq \bar{t}$ instead, the use of averaging as smoothing function allows us to further bound the total variation distance $d_{TV}(q_t^{\circ, k}, q_{t-1}^{\circ, k}) \leq L/t$ as for Lemma F.4. This allows us to bound the summation with sub-linear terms:

$$\sum_{t=\bar{t}}^T |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \leq \sum_{t=\bar{t}}^T \sum_{k=0}^{L-1} h(\epsilon_t) + \log(X^2 AT + X^2 A) \epsilon_t$$

Finally, optimizing for η in the *Algorithmic Regret* and combining all the single terms in the decomposition returns the final result. \square

Overall Regret Analysis. Now that we have proven that Algorithm 1 is no-smooth regret, what is simply left do is to combine this result with Lemma 3.1, leading to the following result.

Algorithm 2: Smoothed Online Mirror Descent in Bandit-Feedback

Input : state space \mathcal{X} , action space \mathcal{A} , episode number T , learning rate $\eta > 0$, mixing parameter $\alpha \in [0, 1]$, estimator parameter $\gamma > 0$, smoothing functions $\{f_t\}_{t \in [T]}$.

Initialize : Set $\pi_1 = \pi^{q_1}$, $q_1(x, a, x') = u(x, a, x') \forall k \in \llbracket 0, L-1 \rrbracket$, $(x, a, x') \in \mathcal{X}_k \times \mathcal{A} \times \mathcal{X}_{k+1}$.

```

1 for  $t = 1, \dots, T$  do
2   Execute policy  $\pi_t$  in  $\mathcal{M}_t$  and collect trajectory  $\{(x_{t,k}, a_{t,k}, \ell_t(x_{t,k}, a_{t,k}))\}_{t \in [T], k \in \llbracket 0, L-1 \rrbracket}$ 
3   Observe  $P_t$  and compute smoothed transition  $\bar{P}_t = f_t(P_1, \dots, P_t)$ .
4   Construct loss estimator:  $\hat{\ell}_t(x, a) = \frac{\ell_t(x, a)}{q^{P_t, \pi_t}(x, a) + \gamma} \mathbb{1}_t(x, a)$ ,  $\forall (x, a) \in \mathcal{X} \times \mathcal{A}$ 
5   Compute smoothed regularization point  $\bar{q}_t = (1 - \alpha) q_t + \alpha u$ 
6   Perform mirror descent step  $q_{t+1} = \arg \min_{q \in \Delta(\bar{P}_t)} \langle q, \hat{\ell}_t \rangle + \frac{1}{\eta} D_\psi(q, \bar{q}_t)$ 
7   Update policy  $\pi_{t+1} = \pi^{q_{t+1}}$ 
8 end

```

Corollary 4.3 (Regret Bound for SOMD under full-feedback and average smoothing). *For $\eta = 3\sqrt{(2L \log(2X^2AT) \log(T))}/T$, $\alpha = 1/(T+1)$, smoothing functions such that $\bar{P}_t := \frac{1}{t} \sum_{t' \in [t]} P_{t'}$ and any comparator policy $\pi^\circ \in \Pi$, Algorithm 1 suffers a regret of:*

$$\mathcal{R}_T(\pi^\circ) \leq \tilde{\mathcal{O}} \left(L^{3/2} \sqrt{T} + L \bar{C}_f^P \right). \quad (14)$$

5 SMOOTHED OMD UNDER BANDIT-FEEDBACK ON LOSSES AND REVEALED TRANSITIONS

In this section, we extend the previous results to the case in which losses are observed under bandit feedback. We show that SOMD can be adapted to this setting with limited adjustments.

Algorithmic Design. To face this challenging scenario, a common way to go is to construct loss estimators based on observations only. In particular, inverse importance-weighted estimators as of Jaksch et al. (2010) can be used to weight the estimation on the experienced trajectory. Thus, we will simply substitute the true feedback with an estimator, namely:

$$\hat{\ell}_t(x, a) := \frac{\ell_t(x, a)}{q^{P_t, \pi_t}(x, a) + \gamma} \mathbb{1}_t(x, a), \quad (15)$$

where $\gamma > 0$ is a parameter to be specified later that allows bounding the variance of the estimator, and $\mathbb{1}_t(x, a)$ is the indicator random variable for the event that the (x, a) is visited at round t . As it emerges from the analysis, the SOMD algorithm can be employed by just replacing ℓ_t with $\hat{\ell}_t$. The intrinsic properties of smoothing in smoothed MDPs will take care of most of the remaining complexity of the problem and the rest of the SOMD algorithm can be employed as is (Algorithm 2).

Smoothed Regret Analysis. We again proceed in bounding the regret in the smoothed MDP. In particular, we can state that using generic smoothing functions leads to the following:

Theorem 5.1 (Smoothed-Regret Bound for SOMD under bandit-feedback). *Let $\eta = \sqrt{(13L \log(2X^2AT) \rho_T^f)/(2XAT)}$, $\alpha = 1/(T+1)$, $\gamma = \eta$, generic smoothing functions and any comparator policy $\pi^\circ \in \Pi$, Algorithm 2 suffers a smoothed regret of:*

$$\bar{\mathcal{R}}_T(\pi^\circ) \leq \mathcal{O} \left(L^2 \bar{D}_f^P + L \bar{C}_f^P + L^{3/2} \sqrt{XAT \rho_T^f \log(X^2AT)} \right). \quad (16)$$

The corresponding result with the choice of the average smoothing is the following.

Corollary 5.2 (Smoothed-Regret Bound for SOMD under bandit-feedback and average smoothing). *Let $\eta = \sqrt{(21L \log(2X^2AT)(\log(T)))/(2XAT)}$, $\alpha = 1/(T+1)$, $\gamma = \eta$, smoothing functions such that $\bar{P}_t := \frac{1}{t} \sum_{t' \in [t]} P_{t'}$ and any comparator policy $\pi^\circ \in \Pi$, Algorithm 2 suffers a smoothed regret of:*

$$\bar{\mathcal{R}}_T(\pi^\circ) \leq \mathcal{O} \left(L \bar{C}_f^P + L^2 \log(T) + L^{3/2} \sqrt{XAT \log(X^2AT) \log(T)} \right) \quad (17)$$

The proof for such a result builds upon the same structure of the full feedback case and we refer the reader to Appendix D for a detailed explanation. Interestingly, and differently from the full-feedback setting, the smoothing transition error \bar{C}_f^P term appears even in the *smoothed regret*. This is inherited by the bias of the loss estimator $\hat{\ell}_t$ of Equation 15.

Overall Regret Analysis. As for the full-feedback case, now that we have proven the smoothed regret bound for Algorithm 2 we combine this result with Lemma 3.1, leading to the following:

Corollary 5.3 (Regret Bound for SOMD under bandit-feedback and average smoothing). *Let $\eta = \sqrt{(21L \log(2X^2AT)(\log(T)))/(2XAT)}$, $\alpha = 1/(T + 1)$, $\gamma = \eta$, smoothing functions such that $\bar{P}_t := \frac{1}{t} \sum_{t' \in [t]} P_{t'}$ and any comparator policy π° , Algorithm 2 suffers a regret of:*

$$\mathcal{R}_T(\pi^\circ) = \tilde{O} \left(L\bar{C}_f^P + L^{3/2}\sqrt{XAT} \right). \quad (18)$$

6 RELATED WORKS

The presence of disturbances, adversarial attacks, and non-stationary behaviors in MDPs have been extensively treated through various lenses. We now highlight how this work is related to each subfield.

Adversarial MDPs. When only the losses are adversarially chosen and the transition functions are either known or fixed (i.e., $\bar{C}_f^P = 0$), many are cases of success in obtaining compelling regret guarantees (e.g. Zimin & Neu, 2013; Rosenberg & Mansour, 2019; Jaksch et al., 2010; Jin et al., 2020), where Zimin & Neu (2013); Jin et al. (2020) in particular take advantage of performing OMD steps in the occupancy measure space. Under adversarially chosen transitions as well, the only related work up to our knowledge is Jin et al. (2023), where bandit feedback for both the losses and the transition model is considered. However, when the degree of maliciousness is unknown, the resulting computational complexity remains uncertain and the regret guarantees are rather implicit. Furthermore, Jin et al. (2023) does not take into account the intermediate setting of revealed transitions, but applying their Algorithm 1 to such a setting would lead to a regret of $\mathcal{O}(L^2X\sqrt{AT\log(LXAT^2)} + L^5X^4A\log(LXAT^2) + C^PL^5X^4\log(LXAT^2))$ even with a known degree of maliciousness C^P , which is significantly worse than the performances of SOMD, constant-wise. On the other hand, we positively leverage the intermediate setting of revealed transitions, introducing ad-hoc degrees of maliciousness the notion of smoothed MDP formalism, and finally recovering comparable performances via computationally efficient OMD-based algorithms.

Corruption Robust RL. Works in this line (e.g. Lykouris et al., 2021; Chen et al., 2021) typically assume the presence of an adversary corrupting some of the rewards and/or transitions, compared to a nominal underlying MDP. These works then address a different notion of regret, namely the one defined coherently with respect to the loss incurred by the best policy in the nominal MDP and denoting the number of corrupted episodes by C , Wei et al. (2022) is the first to achieve a regret of $\tilde{O}(\min\{\frac{1}{\Delta}, \sqrt{T}\} + C)$ in a bandit feedback setting without requiring the knowledge of C , with Δ being the reward gap between the best and the second-best.

Robust MDPs & RL. Robust MDPs (e.g. Nilim & Ghaoui, 2005; Wiesemann et al., 2013) and Robust Reinforcement Learning (e.g. Morimoto & Doya, 2005; Lim et al., 2016) focus on computing policies that exhibit robustness in face of uncertainties over the transition and/or loss models so that to withstand potential mismatches between the models and the ground truth. Usually, though, minimax solutions against the worst-case scenario are sought, failing to adapt to easier instances and to smoothly interpolate performance based on the degree of mismatch.

Non-Stationary RL. Works in this line (e.g. Lecarpentier & Rachelson, 2019; Wei & Luo, 2021; Cheung et al., 2023) allow the MDP model to change arbitrarily over time and the performance metric employed is the dynamic regret, i.e. competing against a comparator policy varying over rounds. It is known that, even in the simpler bandit setting, when the environment is adversarial the no-regret property is not achievable for the dynamic regret (Bubeck et al., 2012), so the results presented in this work are not directly comparable with them.

7 DISCUSSION AND CONCLUSIONS

In this paper, we have addressed the OL problem in MDPs in which the transition functions are chosen by an adversary. Starting from the known computational and statistical limits of this challenging setting, we have the scenario in which the transition functions are revealed at the end of the episode. We have introduced the notion of smoothed MDP and, based on it, we designed suitable indexes \overline{C}_f^P and \overline{D}_f^P to assess the degree of maliciousness of the adversary that relate and generalize the existing ones. These indexes allowed us to design a computationally efficient algorithm, SOMD, that enjoys regret guarantees of order $\tilde{O}(\sqrt{T} + \overline{C}_f^P)$. These results are in line with the literature (Jin et al., 2023) but require no knowledge of the maliciousness index \overline{C}_f^P and are obtained with simple yet computationally efficient algorithms. Future works include the extension of the proposed approach to a *complete bandit feedback* setting where the transition functions are not revealed at the end of the episode. Furthermore, their generality makes our indexes \overline{C}_f^P and \overline{D}_f^P suitable to be employed beyond the specific averaging smooth of Example 3.2 and capture more sophisticated relations in the sequence of the transition models, such as *bounded variation* (Example 3.1).

REFERENCES

- Yasin Abbasi Yadkori, Peter L Bartlett, Varun Kanade, Yevgeny Seldin, and Csaba Szepesvári. Online learning in markov decision processes with adversarially chosen transition probability distributions. *Advances in neural information processing systems*, 26, 2013.
- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *Proceedings of IEEE annual foundations of computer science*, pp. 322–331. IEEE, 1995.
- Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- Gong Chen and Marc Teboulle. Convergence analysis of a proximal-like minimization algorithm using bregman functions. *SIAM Journal on Optimization*, 3(3):538–543, 1993.
- Yifang Chen, Simon S. Du, and Kevin Jamieson. Improved corruption robust algorithms for episodic reinforcement learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 139 of *Proceedings of Machine Learning Research*, pp. 1561–1570. PMLR, 2021.
- Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Reinforcement learning for non-stationary markov decision processes: The blessing of (more) optimism. In *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 1843–1854. PMLR, 2020.
- Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Nonstationary reinforcement learning: The blessing of (more) optimism. *Management Science*, 69(10):5722–5739, 2023.
- Gabriel Dulac-Arnold, Nir Levine, Daniel J Mankowitz, Jerry Li, Cosmin Paduraru, Sven Gowal, and Todd Hester. Challenges of real-world reinforcement learning: definitions, benchmarks and analysis. *Machine Learning*, 110(9):2419–2468, 2021.
- Eyal Even-Dar, Sham M. Kakade, and Yishay Mansour. Online markov decision processes. *Mathematics of Operations Research*, 34(3):726–736, 2009.
- Eyal Even-Dar, Sham M Kakade, and Yishay Mansour. Online markov decision processes. *Mathematics of Operations Research*, 34(3):726–736, 2009.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11:1563–1600, 2010.
- Chi Jin, Tiancheng Jin, Haipeng Luo, Suvrit Sra, and Tiancheng Yu. Learning adversarial markov decision processes with bandit feedback and unknown transition. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 119 of *Proceedings of Machine Learning Research*, pp. 4860–4869. PMLR, 2020.

- 540 Tiancheng Jin, Junyan Liu, Chloé Rouyer, William Chang, Chen-Yu Wei, and Haipeng Luo. No-
541 regret online reinforcement learning with adversarial losses and transitions. In *Advances in Neural*
542 *Information Processing Systems (NeurIPS)*, 2023.
- 543
544 Petar Kormushev, Sylvain Calinon, and Darwin G Caldwell. Reinforcement learning in robotics:
545 Applications and real-world challenges. *Robotics*, 2(3):122–148, 2013.
- 546 Erwan Lecarpentier and Emmanuel Rachelson. Non-stationary markov decision processes, a worst-
547 case approach using model-based reinforcement learning. *Advances in Neural Information Pro-*
548 *cessing Systems (NeurIPS)*, 32, 2019.
- 549
550 Shiau Hong Lim, Huan Xu, and Shie Mannor. Reinforcement learning in robust markov decision
551 processes. *Mathematics of Operations Research*, 41(4):1325–1353, 2016.
- 552 Qinghua Liu, Yuanhao Wang, and Chi Jin. Learning markov games with adversarial opponents:
553 Efficient algorithms and fundamental limits. In *Proceedings of the International Conference*
554 *on Machine Learning (ICML)*, volume 162 of *Proceedings of Machine Learning Research*, pp.
555 14036–14053. PMLR, 2022.
- 556
557 Thodoris Lykouris, Max Simchowitz, Alex Slivkins, and Wen Sun. Corruption-robust exploration in
558 episodic reinforcement learning. In *Proceedings of the Annual Conference on Learning Theory*
559 *(COLT)*, volume 134 of *Proceedings of Machine Learning Research*, pp. 3242–3245. PMLR, 2021.
- 560
561 Jun Morimoto and Kenji Doya. Robust reinforcement learning. *Neural Computation*, 17(2):335–359,
2005.
- 562
563 Anusha Nagabandi, Ignasi Clavera, Simin Liu, Ronald S Fearing, Pieter Abbeel, Sergey Levine, and
564 Chelsea Finn. Learning to adapt in dynamic, real-world environments through meta-reinforcement
565 learning. In *International Conference on Learning Representations (ICLR)*, 2018.
- 566
567 Arnab Nilim and Laurent El Ghaoui. Robust control of markov decision processes with uncertain
transition matrices. *Operations Research*, 53(5):780–798, 2005.
- 568
569 Francesco Orabona. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*,
2019.
- 570
571 Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John
572 Wiley & Sons, 2014.
- 573
574 Aviv Rosenberg and Yishay Mansour. Online stochastic shortest path with bandit feedback and
575 unknown transition function. In *Advances in Neural Information Processing Systems (NeurIPS)*,
pp. 2209–2218, 2019.
- 576
577 Igal Sason. Entropy bounds for discrete random variables via maximal coupling. *IEEE Transactions*
578 *on Information Theory*, 59(11):7118–7131, 2013.
- 579
580 Yi Tian, Yuanhao Wang, Tiancheng Yu, and Suvrit Sra. Online learning in unknown markov games.
581 In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 139 of
Proceedings of Machine Learning Research, pp. 10279–10288. PMLR, 2021.
- 582
583 Chen-Yu Wei and Haipeng Luo. Non-stationary reinforcement learning without prior knowledge:
584 an optimal black-box approach. In *Proceedings of the Annual Conference on Learning Theory*
585 *(COLT)*, volume 134 of *Proceedings of Machine Learning Research*, pp. 4300–4354. PMLR, 2021.
- 586
587 Chen-Yu Wei, Christoph Dann, and Julian Zimmert. A model selection approach for corruption
588 robust reinforcement learning. In *International Conference on Algorithmic Learning Theory (ALT)*,
volume 167 of *Proceedings of Machine Learning Research*, pp. 1043–1096. PMLR, 2022.
- 589
590 Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust markov decision processes. *Mathe-*
591 *matics of Operations Research*, 38(1):153–183, 2013.
- 592
593 Alexander Zimin and Gergely Neu. Online learning in episodic markovian decision processes by
relative entropy policy search. In *Advances in Neural Information Processing Systems (NeurIPS)*,
pp. 1583–1591, 2013.

A THEORETICAL ANALYSIS FOR SMOOTHED MDPs

Lemma A.1. Let $\bar{P}_t = \frac{1}{t} \sum_{t' \in [t]} P_{t'}$. Then, it holds that:

$$\frac{1}{\log T + 2} \leq \frac{\bar{C}_f^P}{C^P} \leq \log T + 2. \quad (19)$$

Proof. Let P be the nominal MDP with which C^P is defined. We have:

$$\begin{aligned} \bar{C}_f^P &= \sum_{t \in [T]} \sum_{k \in [0, L-1]} \max_{(x, a) \in \mathcal{X}_k \times \mathcal{A}} \|\bar{P}_t(\cdot | x, a) - P_t(\cdot | x, a)\|_1 \\ &\leq \sum_{t \in [T]} \sum_{k \in [0, L-1]} \max_{(x, a) \in \mathcal{X}_k \times \mathcal{A}} \|\bar{P}_t(\cdot | x, a) - P(\cdot | x, a)\|_1 \\ &\quad + \sum_{t \in [T]} \sum_{k \in [0, L-1]} \max_{(x, a) \in \mathcal{X}_k \times \mathcal{A}} \|P(\cdot | x, a) - P_t(\cdot | x, a)\|_1 \\ &= \sum_{t \in [T]} \sum_{k \in [0, L-1]} \max_{(x, a) \in \mathcal{X}_k \times \mathcal{A}} \left\| \frac{1}{t} \sum_{t' \in [t]} P_{t'}(\cdot | x, a) - P(\cdot | x, a) \right\|_1 + C^P \\ &\leq \sum_{t \in [T]} \frac{1}{t} \sum_{t' \in [t]} \sum_{k \in [0, L-1]} \max_{(x, a) \in \mathcal{X}_k \times \mathcal{A}} \|P_{t'}(\cdot | x, a) - P(\cdot | x, a)\|_1 + C^P \\ &\leq \sum_{t \in [T]} \frac{1}{t} \sum_{t' \in [T]} \sum_{k \in [0, L-1]} \max_{(x, a) \in \mathcal{X}_k \times \mathcal{A}} \|P_{t'}(\cdot | x, a) - P(\cdot | x, a)\|_1 + C^P \\ &\leq (\log T + 2)C^P, \end{aligned}$$

where the first inequality comes from triangle inequality. And the second from Jensen's. The lower bound comes from analogous derivation. \square

Lemma 3.1 (Proxy Regret Upper Bound). For any policy sequence $\{\pi_t\}_{t=1}^T$, and loss functions $\{\ell_t\}_{t=1}^T$ such that $\ell_t : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$ for any $t \in \{1, \dots, T\}$ it holds that:

$$\text{Proxy Regret} = \sum_{t \in [T]} \langle q^{\bar{P}_t, \pi_t} - q^{P_t, \pi_t}; \ell_t \rangle \leq L \bar{C}_f^P. \quad (9)$$

Proof. The proof follows the same steps as in the original derivation of Lemma F.7 in (Jin et al., 2023), we invite the reader to see the original work for a detailed proof. \square

B THEORETICAL ANALYSIS WITH FULL-FEEDBACK AND AVERAGE SMOOTHING FUNCTIONS

In this section, we present the proofs of the results discussed in Section 4 with the specific smoothing functions $f_t(P_1, \dots, P_t) = \frac{1}{t} \sum_{t'=1}^t P_{t'}$.

B.1 MAIN RESULTS

Corollary 4.2 (Smoothed-Regret Bound for SOMD under full-feedback and average smoothing). Let $\eta = 3\sqrt{(2L \log(2X^2 AT) \log(T))/T}$, $\alpha = 1/(T + 1)$ and smoothing functions such that $\bar{P}_t := \frac{1}{t} \sum_{t' \in [t]} P_{t'}$, for any comparator policy $\pi^\circ \in \Pi$, Algorithm 1 suffers a smoothed regret of:

$$\bar{\mathcal{R}}_T(\pi^\circ) \leq \mathcal{O} \left(L^2 \log(T) + L^{3/2} \sqrt{T \log(X^2 AT) \log(T)} \right). \quad (12)$$

Proof. We first define $q_t^{\bar{P}_{t-1}, \pi_t} = q_t$ and $q_t^{\bar{P}_t, \pi^\circ} = q_t^\circ$ for notational simplicity. It follows that the smoothed regret term can be decomposed as

$$\bar{\mathcal{R}}_T(\pi^\circ) = \mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_t, \pi_t} - q^{\bar{P}_t, \pi^\circ}; \ell_t \rangle \right] = \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t - q_t^\circ; \ell_t \rangle \right]}_{\substack{\text{Algorithmic Regret } (\bar{\mathcal{R}}_T^A) \\ \text{Lemma B.2}}} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_t, \pi_t} - q^{\bar{P}_{t-1}, \pi_t}; \ell_t \rangle \right]}_{\substack{\text{Update Regret } (\bar{\mathcal{R}}_T^U) \\ \text{Lemma B.1}}},$$

Then, the result follows directly from the combination of Lemma B.1 for the Update Regret and Lemma B.2 for the Algorithmic Regret, namely:

$$\begin{aligned} \bar{\mathcal{R}}_T(\pi^\circ) &\leq 2(L^2 + L)(1 + \log(T)) + 2L + 8L\sqrt{TL \log(2X^2AT) \log(T)} \\ &= 4L + 2L^2 + 2(L^2 + L)(1 + \log(T)) + 6L\sqrt{2TL \log(2X^2AT) \log(T)} \end{aligned}$$

which leads to the final result. \square

Corollary 4.3 (Regret Bound for SOMD under full-feedback and average smoothing). *For $\eta = 3\sqrt{(2L \log(2X^2AT) \log(T))/T}$, $\alpha = 1/(T+1)$, smoothing functions such that $\bar{P}_t := \frac{1}{t} \sum_{t' \in [t]} P_{t'}$ and any comparator policy $\pi^\circ \in \Pi$, Algorithm 1 suffers a regret of:*

$$\mathcal{R}_T(\pi^\circ) \leq \tilde{\mathcal{O}} \left(L^{3/2} \sqrt{T} + LC_f^P \right). \quad (14)$$

Proof. The result follows straightforwardly from combining the results from Corollary 4.2 and Lemma 3.1. \square

Lemma B.1 (Update Regret Bound). *For smoothing functions such that $\bar{P}_t := \frac{1}{t} \sum_{t' \in [t]} P_{t'}$ and $\ell_t : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$, for any policy sequence $\{\pi_t\}_{t \in [T]}$ we have that,*

$$\bar{\mathcal{R}}_T^U = \mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_t, \pi_t} - q^{\bar{P}_{t-1}, \pi_t}; \ell_t \rangle \right] \leq 2(L^2 + L)(1 + \log(T)).$$

Proof.

$$\begin{aligned} \bar{\mathcal{R}}_T^U &= \mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_t, \pi_t} - q^{\bar{P}_{t-1}, \pi_t}; \ell_t \rangle \right] \leq \mathbb{E} \left[\sum_{t=1}^T \left\| q^{\bar{P}_t, \pi_t} - q^{\bar{P}_{t-1}, \pi_t} \right\|_1 \right] \\ &\leq 2(L^2 + L) \sum_{t=1}^T \frac{1}{t} \\ &\leq 2(L^2 + L)(\log(T) + 1) \end{aligned}$$

The first inequality follows from Holder's and the fact that $\|\ell_t\|_\infty \leq 1$ by definition. The second inequality follows from Lemma F.10, and the last one comes from the bound on the harmonic series. \square

Lemma B.2 (Algorithmic Regret Bound). *Choosing $\eta = 3\sqrt{\frac{2L \log(2X^2AT) \log(T)}{T}}$, $\alpha = \frac{1}{1+T}$ Algorithmic Regret is bounded by*

$$\bar{\mathcal{R}}_T^A = \mathbb{E} \left[\sum_{t=1}^T \langle q_t - q_t^\circ; \ell_t \rangle \right] \leq 2L + 6L\sqrt{TL \log(2X^2AT) \log(T)}$$

Proof. The proof relies on the following decomposition:

$$\eta \bar{\mathcal{R}}_T^A = \eta \mathbb{E} \left[\sum_{t=1}^T \langle q_t - q_t^\circ; \ell_t \rangle \pm \eta \sum_{t=1}^T \langle \bar{q}_t; \ell_t \rangle \right] = \underbrace{\eta \mathbb{E} \left[\sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \ell_t \rangle \right]}_{\substack{\text{Descent Regret } (\bar{\mathcal{R}}_T^D) \\ \text{Lemma B.4}}} + \underbrace{\eta \mathbb{E} \left[\sum_{t=1}^T \langle q_t - \bar{q}_t; \ell_t \rangle \right]}_{\substack{\text{Regularization Regret } (\bar{\mathcal{R}}_T^R) \\ \text{Lemma B.3}}}.$$

By applying Lemma B.3 for the Regularization Regret and Lemma B.4 for the Descent Regret we can rewrite the Algorithmic Regret as

$$\eta \bar{\mathcal{R}}_T^A \leq 2\eta L + \eta^2 T L + 18L^2 \log(2X^2 AT) \log(T)$$

We first chose the optimal η as:

$$\eta = 3\sqrt{\frac{2L \log(2X^2 AT) \log(T)}{T}}$$

and then we upper bound the Algorithmic Regret after substituting the optimal η , namely:

$$\begin{aligned} \bar{\mathcal{R}}_T^A &\leq 2L + \eta T L + \frac{18L^2 \log(2X^2 AT) \log(T)}{\eta} \\ &\leq 2L + 3L\sqrt{2TL \log(2X^2 AT) \log(T)} + 3L\sqrt{2TL \log(2X^2 AT) \log(T)} \end{aligned}$$

which leads to the final result. \square

Lemma B.3 (Regularisation Regret Bound). For $\bar{q}_t = (1 - \alpha)q_t + \alpha u$, $u(x, a, x') = \frac{1}{X_k \cdot A \cdot X_{k+1}} \forall k \in \llbracket L - 1 \rrbracket$ and $\alpha = \frac{1}{T+1}$, it holds that

$$\eta \bar{\mathcal{R}}_T^R = \eta \mathbb{E} \left[\sum_{t=1}^T \langle q_t - \bar{q}_t; \ell_t \rangle \right] \leq 2\eta L$$

Proof.

$$\begin{aligned} \eta \bar{\mathcal{R}}_T^R &= \eta \sum_{t=1}^T \langle q_t - \bar{q}_t; \ell_t \rangle = \alpha \eta \sum_{t=1}^T \langle q_t - u; \ell_t \rangle \\ &\leq \alpha \eta \sum_{t=1}^T \|q_t - u\|_1 \\ &\leq \alpha \eta \sum_{t=1}^T \sum_{k=0}^{L-1} \|q_t^k - u\|_1 \\ &\leq 2\alpha \eta L T = \frac{2}{1+T} \eta L T \\ &\leq 2\eta L \end{aligned}$$

Where the first equality comes from the definition of \bar{q}_t , the first inequality follows from Holder's and the fact that $\|\ell_t\|_\infty \leq 1$. Finally, the last step derives from bounding the 1-norm between distributions by 2 and the definition of $\alpha = \frac{1}{T+1}$. \square

Lemma B.4 (Descent Regret Bound). For $\alpha = \frac{1}{T+1}$, $\eta > 0$ the following fact holds,

$$\eta \bar{\mathcal{R}}_T^D = \eta \sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \ell_t \rangle \leq \eta^2 T L + 18L^2 \log(2X^2 AT) \log(T)$$

Proof. We start by decomposing the regret at a generic step t into

$$\eta \langle \bar{q}_t - q_t^\circ; \ell_t \rangle = \underbrace{\eta \langle \bar{q}_t - \tilde{q}_{t+1}; \ell_t \rangle}_{\text{Lemma B.5}} + \underbrace{\eta \langle \tilde{q}_{t+1} - q_t^\circ; \ell_t \rangle}_{\text{Lemma B.6}}$$

Where the \tilde{q}_{t+1} represents the unconstrained solution to the OMD optimization problem.

$$q_{t+1} = \arg \min_{q \in \mathcal{P}_t} D_\psi(q, \tilde{q}_{t+1}), \quad \tilde{q}_{t+1} = \arg \min_q \langle q; \ell_t \rangle + \frac{1}{\eta} D_\psi(q, \bar{q}_t)$$

We apply Lemma B.5 on $\eta\langle\bar{q}_t - \tilde{q}_{t+1}; \ell_t\rangle$ and Lemma B.6 on $\eta\langle\tilde{q}_{t+1} - q_t^\circ; \ell_t\rangle$ obtaining:

$$\begin{aligned}\eta\langle\bar{q}_t - q_t^\circ; \ell_t\rangle &= D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(\tilde{q}_{t+1}, \bar{q}_t) + D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(\tilde{q}_{t+1}, \bar{q}_t) - D_\psi(q_t^\circ, \tilde{q}_{t+1}) \\ &= D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \tilde{q}_{t+1})\end{aligned}$$

Since only the optimality condition for the unconstrained problem and the three-point lemma have been employed so far, this result holds for any q_t°, \bar{q}_t . Now, considering the result of the projection step of OMD, q_{t+1} , the generalized Pythagorean Theorem allow us to state that $D_\psi(q_t^\circ, \tilde{q}_{t+1}) \geq D_\psi(q_t^\circ, q_{t+1})$. Such a result, together with the application of Lemma B.7 to $D_\psi(q_t^\circ, q_{t+1})$ provide:

$$\begin{aligned}\eta\langle\bar{q}_t - q_t^\circ; \ell_t\rangle &= D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \tilde{q}_{t+1}) \\ &\leq D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, q_{t+1}) \\ &\leq D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(q_t^\circ, \bar{q}_t) - \underbrace{D_\psi(q_t^\circ, \bar{q}_{t+1}) + \frac{\alpha}{(1-\alpha)}D_\psi(q_t^\circ, u)}_{\text{Lemma B.7}}\end{aligned}$$

Now, summing over $t \in \llbracket T \rrbracket$ we get

$$\begin{aligned}\eta\bar{\mathcal{R}}_T^D &= \eta \sum_{t=1}^T \langle\bar{q}_t - q_t^\circ; \ell_t\rangle \leq \underbrace{\sum_{t=1}^T D_\psi(\bar{q}_t, \tilde{q}_{t+1})}_{\substack{\text{“Stability” term} \\ \text{Lemma B.8}}} + \underbrace{\frac{\alpha}{(1-\alpha)} \sum_{t=1}^T D_\psi(q_t^\circ, u)}_{\substack{\text{“Residual” term} \\ \text{Lemma B.9}}} \\ &\quad + \underbrace{\sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \bar{q}_{t+1}))}_{\substack{\text{“Penalty” term} \\ \text{Lemma B.10}}}.\end{aligned}$$

Now we combine the results from Lemma B.8 for the “Stability” term, Lemma B.9 for the “Residual” term and Lemma B.10 for the “Penalty” term, obtaining:

$$\eta\bar{\mathcal{R}}_T^D \leq \eta^2 TL + L \log(X^2 A) + 17L^2 \log(2X^2 AT) \log(T)$$

□

B.2 AUXILIARY LEMMAS

Lemma B.5. For $\bar{q}_t \in \Delta$ and $\tilde{q}_{t+1} = \arg \min_q \langle q; \ell_t \rangle + \frac{1}{\eta} D_\psi(q, \bar{q}_t)$, we have that,

$$\eta\langle\bar{q}_t - \tilde{q}_{t+1}; \ell_t\rangle = D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(\tilde{q}_{t+1}, \bar{q}_t)$$

Proof. Applying the first order optimality conditions for unconstrained optimisation problems,

$$\langle\bar{q}_t - \tilde{q}_{t+1}; \ell_t + \frac{1}{\eta}(\nabla\psi(\tilde{q}_{t+1}) - \nabla\psi(\bar{q}_t))\rangle = 0$$

Rearranging and applying the three-point inequality (Chen & Teboulle, 1993),

$$\begin{aligned}\eta\langle\tilde{q}_{t+1} - \bar{q}_t; \ell_t\rangle &= \langle\bar{q}_t - \tilde{q}_{t+1}; \nabla\psi(\tilde{q}_{t+1}) - \nabla\psi(\bar{q}_t)\rangle \\ &= D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(\tilde{q}_{t+1}, \bar{q}_t).\end{aligned}$$

□

Lemma B.6. For $q_t^\circ \in \Delta$ and $\tilde{q}_{t+1} = \arg \min_q \langle q; \ell_t \rangle + \frac{1}{\eta} D_\psi(q, \bar{q}_t)$ we have that,

$$\eta\langle\tilde{q}_{t+1} - q_t^\circ; \ell_t\rangle = D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(\tilde{q}_{t+1}, \bar{q}_t) - D_\psi(q_t^\circ, \tilde{q}_{t+1})$$

Proof. We apply first-order optimality conditions

$$\langle q_t^\circ - \tilde{q}_{t+1}; \ell_t + \frac{1}{\eta}(\nabla\psi(\tilde{q}_{t+1}) - \nabla\psi(\bar{q}_t))\rangle = 0$$

Rearranging and applying the three-point inequality (Chen & Teboulle, 1993),

$$\begin{aligned} \eta \langle \tilde{q}_{t+1} - q_t^\circ; \ell_t \rangle &= \langle q_t^\circ - \tilde{q}_{t+1}; \nabla \psi(\tilde{q}_{t+1}) - \nabla \psi(\bar{q}_t) \rangle \\ &= D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(\tilde{q}_{t+1}, \bar{q}_t) - D_\psi(q_t^\circ, \tilde{q}_{t+1}) \end{aligned}$$

□

Lemma B.7. For $\psi: \Delta \rightarrow \mathbb{R}$ as the negative shannon entropy and for $q_t^\circ, q_{t+1} \in \Delta$ and $\bar{q}_{t+1} = (1 - \alpha)q_{t+1} + \alpha u$ we have:

$$D_\psi(q_t^\circ, q_{t+1}) \geq D_\psi(q_t^\circ, \bar{q}_{t+1}) - \frac{\alpha}{(1 - \alpha)} D_\psi(q_t^\circ, u).$$

Proof. Due to the convexity of the regularizer, namely the KL divergence, we observe that

$$\begin{aligned} D_\psi(q_t^\circ, \bar{q}_{t+1}) &= D_\psi(q_t^\circ, (1 - \alpha)q_{t+1} + \alpha u) \leq (1 - \alpha)D_\psi(q_t^\circ, q_{t+1}) + \alpha D_\psi(q_t^\circ, u) \\ \frac{1}{(1 - \alpha)} D_\psi(q_t^\circ, \bar{q}_{t+1}) - \frac{\alpha}{(1 - \alpha)} D_\psi(q_t^\circ, u) &\leq D_\psi(q_t^\circ, q_{t+1}) \end{aligned}$$

We retrieve that:

$$D_\psi(q_t^\circ, q_{t+1}) \geq D_\psi(q_t^\circ, \bar{q}_{t+1}) - \frac{\alpha}{(1 - \alpha)} D_\psi(q_t^\circ, u)$$

Where the last step comes from the fact that $1 - \alpha \leq 1$. □

Lemma B.8 (“Stability” Term Bound). For $\eta > 0$, $\ell_t: \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$, $\bar{q}_t \in \Delta$ and $\tilde{q}_{t+1} = \arg \min_q \langle q; \ell_t \rangle + \frac{1}{\eta} D_\psi(q, \bar{q}_t)$ we have that:

$$(STAB) \quad \sum_{t=1}^T D_\psi(\bar{q}_t, \tilde{q}_{t+1}) \leq \eta^2 TL.$$

Proof. The term can be easily bounded following standard analysis of OMD:

$$\begin{aligned} &\sum_{t=1}^T D_\psi(\bar{q}_t, \tilde{q}_{t+1}) \\ &= \sum_{t=1}^T \left[\sum_{x, a, x'} \bar{q}_t(x, a, x') \log \left(\frac{\bar{q}_t(x, a, x')}{\tilde{q}_{t+1}(x, a, x')} \right) - \sum_{x, a, x'} (\bar{q}_t(x, a, x') - \tilde{q}_{t+1}(x, a, x')) \right] \\ &= \sum_{t=1}^T \left[\sum_{x, a, x'} \eta \ell_t(x, a) \bar{q}_t(x, a, x') + \bar{q}_t(x, a, x') \exp(-\eta \ell_t(x, a)) - \bar{q}_t(x, a, x') \right] \\ &\leq \sum_{t=1}^T \left[\sum_{(x, a, x')} \eta^2 \bar{q}_t(x, a, x') \ell_t^2(x, a) \right] \\ &\leq \eta^2 TL \end{aligned}$$

Where the first equality comes from the definition of generalized KL divergence, the second by applying the definition of the solution of the unconstrained optimization problem, namely:

$$\tilde{q}_{t+1}(x, a, x') = \bar{q}_t(x, a, x') \exp(-\eta \ell_t(x, a)), \forall (x, a, x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}$$

and further simplifications. Finally, we use the standard bound of the exponential function $e^{-x} \leq 1 - x + x^2, \forall x \geq 0$, the fact that losses are bounded, $\ell_t(x, a) \in [0, 1]$ and as a last step we use the fundamental property of occupancies $\sum_{(x, a, x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}} \bar{q}_t(x, a, x') = L$. □

Lemma B.9 (“Residual” Term Bound). For $\alpha = \frac{1}{T+1}$, $u(x, a, x') = \frac{1}{X_k A X_{k+1}}$ and $q_t^\circ \in \Delta$ it holds that

$$(RES) \quad \frac{\alpha}{1 - \alpha} \sum_{t=1}^T D_\psi(q_t^\circ, u) \leq L \log(X^2 A)$$

864 *Proof.* We first notice that

$$\begin{aligned}
865 D_\psi(q_t^\circ, u) &= \sum_{k=0}^{L-1} \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k+1}} q_t^\circ(x, a, x') \log \left(\frac{q_t^\circ(x, a, x')}{\frac{1}{X_k A X_{k+1}}} \right) \\
866 &\leq \sum_{k=0}^{L-1} \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k+1}} q_t^\circ(x, a, x') \log(X^2 A) \\
867 &\leq L \log(X^2 A)
\end{aligned}$$

873 The first inequality comes from neglecting negative terms and the last using the fundamental property
874 of occupancy measures.

875 Thus, considering the sum over t :

$$877 \frac{\alpha}{1-\alpha} \sum_{t=1}^T D_\psi(q_t^\circ, u) \leq \frac{\alpha}{1-\alpha} T L \log(X^2 A)$$

880 And choosing α so that $\frac{\alpha}{1-\alpha} T = 1$, namely $\alpha = \frac{1}{T+1}$ leads to the desired bound. \square

882 **Lemma B.10** (“Penalty” Term Bound). *For $\bar{q}_t = (1 - \alpha)q_t + \alpha u$ and for $\alpha = \frac{1}{T+1}$*

$$884 (PEN) \quad \sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \bar{q}_{t+1})) \leq 17L^2 \log(2X^2 AT) \log(T)$$

887 *Proof.* First, we unravel the summation:

$$\begin{aligned}
888 \sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \bar{q}_{t+1})) &= D_\psi(q_1^\circ, \bar{q}_1) - D_\psi(q_1^\circ, \bar{q}_2) + D_\psi(q_T^\circ, \bar{q}_T) - D_\psi(q_T^\circ, \bar{q}_{T+1}) \\
889 &\leq \underbrace{D_\psi(q_1^\circ, \bar{q}_1)}_{\substack{\text{Range of } \psi \\ \text{Lemma B.11}}} + \underbrace{\sum_{t=2}^T D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_{t-1}^\circ, \bar{q}_t)}_{\substack{\text{Non-telescoping term} \\ \text{Lemma B.12}}}
\end{aligned}$$

892 Where the inequality comes from neglecting negative terms. What is left to do is to combine Lemma
893 B.11, that expresses the range of the negative shannon entropy, and Lemma B.12 providing a bound
894 for the non-telescoping term.

$$900 \sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \bar{q}_{t+1})) \leq L \log(X^2 A) + 16L^2 \log(2X^2 AT) \log(T)$$

903 \square

904 **Lemma B.11** (Range of ψ Bound). *For $q_t^\circ \in \Delta$ and $\bar{q}_1 = (1 - \alpha)q_1 + \alpha u$ we have that,*

$$907 D_\psi(q_1^\circ, \bar{q}_1) \leq L \log(X^2 A)$$

908 *Proof.* The steps are similar to the ones of Lemma B.9:

$$\begin{aligned}
909 D_\psi(q_1^\circ, \bar{q}_1) &= D_\psi(q_1^\circ, u) = \sum_{k=0}^{L-1} \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k+1}} q_1^\circ(x, a, x') \log \left(\frac{q_1^\circ(x, a, x')}{\frac{1}{X_k A X_{k+1}}} \right) \\
910 &\leq \sum_{k=0}^{L-1} \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k+1}} q_1^\circ(x, a, x') \log(X^2 A) \\
911 &\leq L \log(X^2 A).
\end{aligned}$$

912 \square

Lemma B.12 (Non-telescoping term Bound). For $q_t^\circ \in \Delta$ and $\bar{q}_t = (1 - \alpha)q_t + \alpha u$ we have that:

$$\sum_{t=2}^T D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_{t-1}^\circ, \bar{q}_t) \leq 16L^2 \log(2X^2AT) \log(T).$$

Proof. We first focus on each time-step t . From the definition of the Bregman Divergence we have:

$$\begin{aligned} D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_{t-1}^\circ, \bar{q}_t) &= \psi(q_t^\circ) - \psi(q_{t-1}^\circ) + \langle \nabla \psi(\bar{q}_t); q_{t-1}^\circ - q_t^\circ \rangle \\ &= \psi(q_t^\circ) - \psi(q_{t-1}^\circ) + \langle \log(\bar{q}_t); q_{t-1}^\circ - q_t^\circ \rangle + \langle \mathbf{1}; q_{t-1}^\circ - q_t^\circ \rangle \\ &\leq |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| + \|\log(\bar{q}_t)\|_\infty \|q_{t-1}^\circ - q_t^\circ\|_1 \end{aligned}$$

Where the second equality comes from the computation of the derivative of the Negative Shannon Entropy:

$$\nabla \psi(f) = \left(\frac{\partial \psi(f)}{\partial f} \right)^T = (\log(f) + 1, \dots)^T$$

Finally, the last step comes from Holder's inequality and from taking the absolute value of the difference of the entropies.

We combine the results from Lemma B.13 and Lemma B.14 to obtain over the whole time horizon:

$$\begin{aligned} \sum_{t=2}^T D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_{t-1}^\circ, \bar{q}_t) &\leq \underbrace{\sum_{t=2}^T |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)|}_{\text{Lemma B.13}} + \underbrace{\sum_{t=2}^T \|\log(\bar{q}_t)\|_\infty \|q_{t-1}^\circ - q_t^\circ\|_1}_{\text{Lemma B.14}} \\ &\leq 8L^2 \log(2X^2AT) \log(T) + 8L^2 \log(X^2AT + X^2A) \log(T) \end{aligned}$$

□

Lemma B.13. *it holds that*

$$\sum_{t=2}^T |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \leq 8L^2 \log(2X^2AT) \log(T)$$

Proof. The proof comes from applying Lemma B.15 and Lemma B.16 to the following decomposition:

$$\begin{aligned} \sum_{t=2}^T |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| &= \underbrace{\sum_{t=2}^{\underline{t}} |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)|}_{\text{Lemma B.15}} + \underbrace{\sum_{t=\bar{t}}^T |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)|}_{\text{Lemma B.16}} \\ &\leq 2L^2 \log(X^2A) + 6L^2 \log(2X^2AT) \log(T) \end{aligned}$$

Simplifying the expression with another upper bound completes the proof. □

Lemma B.14. Defined $\bar{t} := \left\lceil \frac{L}{1 - \frac{1}{X^2A}} \right\rceil$ and $\underline{t} := \left\lfloor \frac{L}{1 - \frac{1}{X^2A}} \right\rfloor$, for $\bar{q}_t = (1 - \alpha)q_t + \alpha u$ and $\alpha = \frac{1}{T+1}$ it holds that,

$$\sum_{t=2}^T \|\log(\bar{q}_t)\|_\infty \|q_{t-1}^\circ - q_t^\circ\|_1 \leq 8L^2 \log(X^2AT + X^2A) \log(T)$$

Proof. The proof comes from applying Lemma B.20 and Lemma B.21 to the following decomposition,

$$\sum_{t=2}^T \|\log(\bar{q}_t)\|_\infty \|q_{t-1}^\circ - q_t^\circ\|_1 = \underbrace{\sum_{t=2}^{\underline{t}} \|\log(\bar{q}_t)\|_\infty \|q_{t-1}^\circ - q_t^\circ\|_1}_{\text{Lemma B.20}} + \underbrace{\sum_{t=\bar{t}}^T \|\log(\bar{q}_t)\|_\infty \|q_{t-1}^\circ - q_t^\circ\|_1}_{\text{Lemma B.21}}.$$

□

972 **Lemma B.15.** For $X^2A \geq 2$ and $\underline{t} := \left\lfloor \frac{L}{1 - \frac{1}{X^2A}} \right\rfloor$

$$973 \sum_{t=2}^{\underline{t}} |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \leq 2L^2 \log(X^2A)$$

974
975
976
977 *Proof.*

$$978 \sum_{t=2}^{\underline{t}} |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \leq \sum_{t=2}^{\underline{t}} \sum_{k=0}^{L-1} |\psi(q_t^{\circ,k}) - \psi(q_{t-1}^{\circ,k})|$$

$$979 \leq \sum_{t=2}^{\underline{t}} L \log(X^2A)$$

$$980 \leq L \log(X^2A) \underline{t}$$

$$981 \leq 2L^2 \log(X^2A)$$

982
983
984
985
986
987
988 Where the second inequality comes from Lemma B.17, while the last step comes from observing that,
989 for $X^2A \geq 2$,

$$990 \underline{t} \leq \frac{L}{1 - \frac{1}{X^2A}}$$

$$991 \leq L \frac{X^2A}{X^2A - 1}$$

$$992 \leq 2L$$

993
994
995
996
997 \square

998 **Lemma B.16.** For $\bar{t} := \left\lfloor \frac{L}{1 - \frac{1}{X^2A}} \right\rfloor$ and $\epsilon_t = \frac{L}{t}$ it holds that,

$$1000 \sum_{t=\bar{t}}^T |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \leq 6L^2 \log(2X^2AT) \log(T)$$

1001
1002
1003 *Proof.*

$$1004 \sum_{t=\bar{t}}^T |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \leq \sum_{t=\bar{t}}^T \sum_{k=0}^{L-1} h(\epsilon_t) + \log(X^2AT + X^2A) \epsilon_t$$

$$1005 \leq 2L^2 \log(T) + L^2 \log^2(T) + L^2$$

$$1006 + L^2 \log(X^2AT + X^2A) (\log(T) + 1)$$

1007
1008
1009
1010
1011
1012 Where the first inequality comes from Lemma B.18, the second inequality comes from Lemma B.19
1013 and from the bound on the harmonic series. Further upperbounding to simplify the expression finishes
1014 the proof. \square

1015 **Lemma B.17** (Per-step entropy difference bound). For any $t \leq \underline{t} := \left\lfloor \frac{L}{1 - \frac{1}{X^2A}} \right\rfloor$ it holds that,

$$1016 |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \leq L \log(X^2A)$$

1017
1018
1019 *Proof.*

$$1020 |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \leq \sum_{k=0}^{L-1} |\psi(q_t^{\circ,k}) - \psi(q_{t-1}^{\circ,k})|$$

$$1021 \leq \sum_{k=0}^{L-1} \log(X_k A X_{k+1})$$

$$\begin{aligned}
1026 & \\
1027 & \leq \sum_{k=0}^{L-1} \log(X^2 A) \\
1028 & \\
1029 & \leq L \log(X^2 A)
\end{aligned}$$

1030 Again, the first inequality comes from triangle inequality, the rest is a direct consequence of (Theorem 3, Sason, 2013). \square

1033 **Lemma B.18** (Per-step entropy difference bound). *For any $t \geq \bar{t} := \left\lceil \frac{L}{1 - \frac{1}{X^2 A}} \right\rceil$ and $\epsilon_t = \frac{L}{t}$ it holds that,*

$$1036 \quad |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \leq \sum_{k=0}^{L-1} h(\epsilon_t) + \log(X^2 AT + X^2 A) \epsilon_t$$

1039 *Proof.* Denoting with $\epsilon_{t,k} = d_{TV}(q_t^{\circ,k}, q_{t-1}^{\circ,k})$ the total variation distance between the occupancies at layer k , we notice that for any $t \in \llbracket T \rrbracket$.

$$\begin{aligned}
1042 & \\
1043 & |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \leq \sum_{k=0}^{L-1} |\psi(q_t^{\circ,k}) - \psi(q_{t-1}^{\circ,k})| \\
1044 & \\
1045 & \leq \sum_{k=0}^{L-1} h(\epsilon_{t,k}) + \epsilon_{t,k} \log(X^2 A - 1) \\
1046 & \\
1047 & \leq \sum_{k=0}^{L-1} h(\epsilon_{t,k}) + \epsilon_{t,k} \log(X^2 AT + X^2 A) \\
1048 & \\
1049 & \leq \sum_{k=0}^{L-1} h(\epsilon_t) + \epsilon_t \log(X^2 AT + X^2 A) \\
1050 & \\
1051 & \leq \sum_{k=0}^{L-1} h(\epsilon_t) + \epsilon_t \log(X^2 AT + X^2 A) \\
1052 & \\
1053 &
\end{aligned}$$

1054 The first inequality comes from triangle inequality. The rest comes from a straightforward application of (Theorem 3, Sason, 2013). \square

1056 **Lemma B.19** (Binary Entropy Bound). *For $\bar{t} := \left\lceil \frac{L}{1 - \frac{1}{X^2 A}} \right\rceil$ and $\epsilon_t = \frac{L}{t}$ it holds that,*

$$1058 \quad \sum_{t=\bar{t}}^T \sum_{k=0}^{L-1} h(\epsilon_t) \leq 2L^2 \log(T) + L^2 \log^2(T) + L^2$$

1062 *Proof.* For each $t \in [\bar{t}, T]$ we have that,

$$\begin{aligned}
1063 & \\
1064 & h(\epsilon_t) = \frac{L}{t} \log\left(\frac{t}{L}\right) + \left(1 - \frac{L}{t}\right) \log\left(\frac{t}{t-L}\right) \\
1065 & \\
1066 & \leq \frac{L}{t} \log\left(\frac{t}{L}\right) + \log\left(\frac{t}{t-L}\right) \\
1067 & \\
1068 & \leq L \log(T) \frac{1}{t} + \log\left(\frac{t}{t-L}\right) \\
1069 & \\
1070 & \leq L \log(T) \frac{1}{t} + \frac{L}{t'}.
\end{aligned}$$

1072 The last inequality makes use of

$$1073 \quad \log\left(\frac{t}{t-L}\right) = \log\left(\frac{t'+L}{t'}\right) = \log\left(1 + \frac{L}{t'}\right) \leq \frac{L}{t'}, \quad (t = t' + L)$$

1076 Now considering the summations over k and t ,

$$1077 \quad \sum_{t=\bar{t}}^T \sum_{k=0}^{L-1} h(\epsilon_t) \leq L \left(\sum_{t=\bar{t}}^T L \log(T) \frac{1}{t} + \frac{L}{t'} \right)$$

$$\begin{aligned}
&\leq L^2 \log(T) \sum_{t=1}^T \frac{1}{t} + L^2 \sum_{t'=1}^{T-L} \frac{1}{t'} \\
&\leq 2L^2 \log(T) + L^2 \log^2(T) + L^2
\end{aligned}$$

where the last inequality comes from the fact that $\sum_{t'=1}^{T-L} \frac{1}{t'} \leq 1 + \log(T-L) \leq 1 + \log(T)$. \square

Lemma B.20. For $t \in \llbracket 2, \bar{t} \rrbracket$, where $\bar{t} := \left\lceil \frac{L-1}{1-X^2A} \right\rceil$, it holds that,

$$\sum_{t=2}^{\bar{t}} \|\log(\bar{q}_t)\|_{\infty} \|q_t^{\circ} - q_{t-1}^{\circ}\|_1 \leq 4L^2 \log(X^2AT + X^2A) \log(T)$$

Proof.

$$\begin{aligned}
\sum_{t=2}^{\bar{t}} \|\log(\bar{q}_t)\|_{\infty} \|q_t^{\circ} - q_{t-1}^{\circ}\|_1 &\leq 2 \sum_{t=2}^{\bar{t}} \sum_{k=0}^{L-1} \log(X^2AT + X^2A) \epsilon_{t,k} \\
&\leq 2 \sum_{t=2}^{\bar{t}} \sum_{k=0}^{L-1} \log(X^2AT + X^2A) \epsilon_t \\
&\leq 2 \sum_{t=2}^{\bar{t}} \sum_{k=0}^{L-1} \log(X^2AT + X^2A) \frac{L}{t} \\
&= 2L^2 \log(X^2AT + X^2A) \sum_{t=2}^{\bar{t}} \frac{1}{t} \\
&\leq 2L^2 \log(X^2AT + X^2A) (\log(\bar{t}) + 1) \\
&\leq 2L^2 \log(X^2AT + X^2A) (\log(T) + 1)
\end{aligned}$$

Where the first inequality comes from Lemma B.22, the second from the bound Lemma F.4 and the last step comes from the monotonicity of the logarithm. \square

Lemma B.21. For $\bar{t} := \left\lceil \frac{L-1}{1-X^2A} \right\rceil$, $\epsilon_{t,k} := d_{TV}(q_t^{\circ,k}, q_{t-1}^{\circ,k})$ and $\epsilon_t = \frac{L}{t}$ it holds that,

$$\sum_{t=\bar{t}}^T \|\log(\bar{q}_t)\|_{\infty} \|q_t^{\circ} - q_{t-1}^{\circ}\|_1 \leq 4L^2 \log(X^2AT + X^2A) \log(T)$$

Proof.

$$\begin{aligned}
\sum_{t=\bar{t}}^T \|\log(\bar{q}_t)\|_{\infty} \|q_t^{\circ} - q_{t-1}^{\circ}\|_1 &\leq 2 \sum_{t=\bar{t}}^T \sum_{k=0}^{L-1} \log(X^2AT + X^2A) \epsilon_{t,k} \\
&\leq 2 \sum_{t=\bar{t}}^T \sum_{k=0}^{L-1} \log(X^2AT + X^2A) \epsilon_t \\
&\leq 2L^2 \log(X^2AT + X^2A) (\log(T) + 1)
\end{aligned}$$

Where the first inequality comes from Lemma B.22, the second from the bound Lemma F.4. The last step comes from the bound on the harmonic series. \square

Lemma B.22. For any $t \in \llbracket T \rrbracket$ and $\epsilon_{t,k} = d_{TV}(q_t^{\circ,k}, q_{t-1}^{\circ,k})$ it holds that:

$$\|\log(\bar{q}_t)\|_{\infty} \|q_t^{\circ} - q_{t-1}^{\circ}\|_1 \leq 2 \sum_{k=0}^{L-1} \log(X^2AT + X^2A) \epsilon_{t,k}$$

1134 *Proof.* First, we notice that for any $t \in \llbracket T \rrbracket$:

$$\begin{aligned}
1135 & \\
1136 & \|\log(\bar{q}_t)\|_\infty = \left\| \left(\log(\bar{q}_t(x, a, x')), \dots \right)^\top \right\|_\infty \\
1137 & \leq \left\| \left(\log(\alpha u(x, a, x')), \dots \right)^\top \right\|_\infty \\
1138 & = \left\| \left(\log \left(\alpha \frac{1}{X^2 A X} \right), \dots \right)^\top \right\|_\infty \\
1139 & \\
1140 & = \log \left(\frac{X^2 A}{\alpha} \right) \\
1141 & \\
1142 & = \log (X^2 A T + X^2 A)
\end{aligned}$$

1146 Where the first inequality comes from the definition $\bar{q}_t(x, a, x') := (1 - \alpha)q_t(x, a, x') +$
1147 $\alpha u(x, a, x'), \forall (x, a, x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}$ and the monotonicity of the logarithm.

1148 Following this result, it holds that for any $t \in \llbracket T \rrbracket$:

$$\begin{aligned}
1150 & \\
1151 & \|\log(\bar{q}_t)\|_\infty \|q_t^\circ - q_{t-1}^\circ\|_1 = \sum_{k=0}^{L-1} \log (X^2 A T + X^2 A) \|q_t^{\circ,k} - q_{t-1}^{\circ,k}\|_1 \\
1152 & \\
1153 & = 2 \sum_{k=0}^{L-1} \log (X^2 A T + X^2 A) d_{TV} (q_t^{\circ,k}, q_{t-1}^{\circ,k}) \\
1154 & \\
1155 & = 2 \sum_{k=0}^{L-1} \log (X^2 A T + X^2 A) \epsilon_{t,k} \\
1156 & \\
1157 & \\
1158 &
\end{aligned}$$

□

1161 C THEORETICAL ANALYSIS WITH FULL-FEEDBACK AND GENERIC 1162 SMOOTHING FUNCTIONS 1163

1164 In this section, we present the proofs of the results discussed in Section 4 related to the regret analysis
1165 agnostic of the smoothing function to be used. For convenience, we will restate the Theorems and
1166 Lemmas before providing a detailed analysis of each and report just the Lemmas that differentiate
1167 from Appendix B.

1169 C.1 MAIN RESULTS 1170

1171 For the proposed Algorithm, we can provide the following:

1172 **Theorem 4.1** (Smoothed-Regret Bound for SOMD under full-feedback). *Let $\eta =$*
1173 *$\sqrt{(10L \log(2X^2 AT) \rho_T^f)/T}$ and $\alpha = 1/(1+T)$, then for any comparator policy $\pi^\circ \in \Pi$ Algorithm 1*
1174 *suffers a smoothed regret of:*

$$1175 \\
1176 \bar{\mathcal{R}}_T(\pi^\circ) \leq \mathcal{O} \left(L^2 \bar{D}_f^P + L^{3/2} \sqrt{T \log(X^2 AT) \rho_T^f} \right), \quad (11)$$

1177 where $\rho_T^f := \log(T) + \bar{D}_f^P + \mathcal{H}_T(\bar{D}_f^P)$ and $\mathcal{H}_T(\bar{D}_f^P) = T L h \left(\frac{(L^2+L)\bar{D}_f^P}{2TL} \right)$.

1181 *Proof.* We first define $q_t^{\bar{P}_{t-1}, \pi_t} = q_t$ and $q_t^{\bar{P}_t, \pi^\circ} = q_t^\circ$ for notational simplicity. It follows that the
1182 smoothed regret term can be decomposed as

$$\begin{aligned}
1184 & \bar{\mathcal{R}}_T(\pi^\circ) = \mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_t, \pi_t} - q^{\bar{P}_t, \pi^\circ}; \ell_t \rangle \right] = \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t - q_t^\circ; \ell_t \rangle \right]}_{\text{Algorithmic Regret } (\bar{\mathcal{R}}_T^A)} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_t, \pi_t} - q^{\bar{P}_{t-1}, \pi_t}; \ell_t \rangle \right]}_{\text{Update Regret } (\bar{\mathcal{R}}_T^U)}, \\
1185 & \\
1186 & \\
1187 &
\end{aligned}$$

Then, the result follows directly from the combination of Lemma C.2 for the Update Regret and Lemma C.3 for the Algorithmic Regret, namely:

$$\bar{\mathcal{R}}_T(\pi^\circ) \leq (L^2 + L)\bar{D}_f^P + 2L + 2L\sqrt{10LT \log(2X^2 AT)\rho_T^f}$$

which leads to the final result. \square

Lemma C.1 (Trend of $\mathcal{H}_T(\bar{D}_f^P)$ for average smoothing functions). *Let $\bar{P}_t := \frac{1}{t} \sum_{t' \in [t]} P_{t'}$, then $\bar{D}_f^P \approx L \log T$ and*

$$\begin{aligned} \mathcal{H}_T(\bar{D}_f^P) &= TLh\left(\frac{(L^2 + L)\bar{D}_f^P}{2TL}\right) \\ &= TL \frac{(L^2 + L) \log T}{2T} \log \frac{2T}{(L^2 + L) \log T} + \\ &\quad + TL \left(1 - \frac{(L^2 + L) \log T}{2T}\right) \log \left(\frac{2T}{2T - (L^2 + L) \log T}\right) \\ &\approx L^2 (\log T)^2. \end{aligned}$$

Lemma C.2 (Update Regret Bound). *For $\ell_t : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$, for any policy sequence $\{\pi_t\}_{t \in [T]}$ and defined $\bar{D}_f^P = \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x,a) - \bar{P}_{t-1}(\cdot|x,a)\|_1$, then for a generic smoothing function f , we have that,*

$$\bar{\mathcal{R}}_T^U = \mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_t, \pi_t} - q^{\bar{P}_{t-1}, \pi_t}; \ell_t \rangle \right] \leq (L^2 + L)\bar{D}_f^P$$

Proof.

$$\begin{aligned} \bar{\mathcal{R}}_T^U &= \mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_t, \pi_t} - q^{\bar{P}_{t-1}, \pi_t}; \ell_t \rangle \right] \leq \mathbb{E} \left[\sum_{t=1}^T \left\| q^{\bar{P}_t, \pi_t} - q^{\bar{P}_{t-1}, \pi_t} \right\|_1 \right] \\ &\leq (L^2 + L) \sum_{t=1}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x,a) - \bar{P}_{t-1}(\cdot|x,a)\|_1 \\ &\leq (L^2 + L)\bar{D}_f^P \end{aligned}$$

The first inequality follows from Holder's and the fact that $\|\ell_t\|_\infty \leq 1$. The second inequality comes from Lemma F.10 and the remaining step follow from the definition of \bar{D}_f^P . \square

Lemma C.3 (Algorithmic Regret Bound). *Choosing $\eta = \sqrt{\frac{10L \log(2X^2 AT)\rho_T^f}{T}}$, $\alpha = \frac{1}{1+T}$, $\epsilon_{t,k} = d_{TV}(q_t^{\circ,k}, q_{t-1}^{\circ,k})$, $\mathcal{H}_T(\bar{D}_f^P) = TLh\left(\frac{(L^2+L)\bar{D}_f^P}{2TL}\right)$, $\bar{D}_f^P = \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x,a) - \bar{P}_{t-1}(\cdot|x,a)\|_1$ and for $\rho_T^f = \log(T) + \bar{D}_f^P + \mathcal{H}_T(\bar{D}_f^P)$, the Algorithmic Regret is bounded by*

$$\bar{\mathcal{R}}_T^A = \mathbb{E} \left[\sum_{t=1}^T \langle q_t - q_t^\circ; \ell_t \rangle \right] \leq 2L + 2L\sqrt{10LT \log(2X^2 AT)\rho_T^f} \quad (20)$$

Proof. The proof relies on the following decomposition:

$$\eta \bar{\mathcal{R}}_T^A = \eta \mathbb{E} \left[\sum_{t=1}^T \langle q_t - q_t^\circ; \ell_t \rangle \pm \eta \sum_{t=1}^T \langle \bar{q}_t; \ell_t \rangle \right] = \underbrace{\eta \mathbb{E} \left[\sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \ell_t \rangle \right]}_{\text{Descent Regret } (\bar{\mathcal{R}}_T^D) \text{ Lemma C.4}} + \underbrace{\eta \mathbb{E} \left[\sum_{t=1}^T \langle q_t - \bar{q}_t; \ell_t \rangle \right]}_{\text{Regularization Regret } (\bar{\mathcal{R}}_T^R) \text{ Lemma B.3}}.$$

By applying Lemma B.3 for the Regularization Regret and Lemma C.4 for the Descent Regret we can rewrite the Algorithmic Regret as

$$\begin{aligned}\eta \bar{\mathcal{R}}_T^A &\leq 2\eta L + \eta^2 TL + \mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2 AT) \bar{D}_f^P + 10L^2 \log(2X^2 AT) \log(T) \\ &\leq 2\eta L + \eta^2 TL + 10(L^2 \log(2X^2 AT)) (\log(T) + \bar{D}_f^P + \mathcal{H}_T(\bar{D}_f^P)) \\ &= 2\eta L + \eta^2 TL + 10L^2 \log(2X^2 AT) \rho_T^f\end{aligned}$$

where $\rho_T^f = \log(T) + \bar{D}_f^P + \mathcal{H}_T(\bar{D}_f^P)$. Choosing the optimal η :

$$\eta = \sqrt{\frac{10L \log(2X^2 AT) \rho_T^f}{T}}$$

and then we upper-bound the Algorithmic Regret after substituting the optimal η , namely:

$$\begin{aligned}\bar{\mathcal{R}}_T^A &\leq 2L + \eta TL + \frac{10L^2 \log(2X^2 AT) \rho_T^f}{\eta} \\ &\leq 2L + L\sqrt{10LT \log(2X^2 AT) \rho_T^f} + L\sqrt{10LT \log(2X^2 AT) \rho_T^f} \\ &\leq 2L + 2L\sqrt{10LT \log(2X^2 AT) \rho_T^f}\end{aligned}$$

which leads to the final result. \square

Lemma C.4 (Descent Regret Bound). For $\alpha = \frac{1}{T+1}$, $\epsilon_{t,k} = d_{TV}(q_t^{\circ,k}, q_{t-1}^{\circ,k})$, $\mathcal{H}_T(\bar{D}_f^P) = TLh\left(\frac{(L^2+L)\bar{D}_f^P}{2TL}\right)$ and $\bar{D}_f^P = \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\|_1$ the following fact holds

$$\eta \bar{\mathcal{R}}_T^D = \eta \sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \ell_t \rangle \leq \eta^2 TL + \mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2 AT) \bar{D}_f^P + 10L^2 \log(2X^2 AT) \log(T)$$

Proof. We start by decomposing the regret at a generic step t into

$$\eta \langle \bar{q}_t - q_t^\circ; \ell_t \rangle = \eta \langle \bar{q}_t - \tilde{q}_{t+1}; \ell_t \rangle + \eta \langle \tilde{q}_{t+1} - q_t^\circ; \ell_t \rangle$$

Where the \tilde{q}_{t+1} represents the unconstrained solution to the OMD optimization problem.

$$q_{t+1} = \arg \min_{q \in \Delta(\bar{P}_t)} D_\psi(q, \tilde{q}_{t+1}), \quad \tilde{q}_{t+1} = \arg \min_q \langle q; \ell_t \rangle + \frac{1}{\eta} D_\psi(q, \bar{q}_t)$$

We apply Lemma B.5 on $\eta \langle \bar{q}_t - \tilde{q}_{t+1}; \ell_t \rangle$ and Lemma B.6 on $\eta \langle \tilde{q}_{t+1} - q_t^\circ; \ell_t \rangle$ obtaining:

$$\begin{aligned}\eta \langle \bar{q}_t - q_t^\circ; \ell_t \rangle &= D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(\tilde{q}_{t+1}, \bar{q}_t) + D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(\tilde{q}_{t+1}, \bar{q}_t) - D_\psi(q_t^\circ, \tilde{q}_{t+1}) \\ &= D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \tilde{q}_{t+1})\end{aligned}$$

Since only the optimality condition for the unconstrained problem and the three-point lemma have been employed so far, this result holds for any q_t°, \bar{q}_t . Now, considering the result of the projection step of OMD q_{t+1} , the generalized Pythagorean Theorem allow us to state that $D_\psi(q_t^\circ, \tilde{q}_{t+1}) \geq D_\psi(q_t^\circ, q_{t+1})$. Such a result, together with the application of Lemma B.7 to $D_\psi(q_t^\circ, q_{t+1})$ provide:

$$\begin{aligned}\eta \langle \bar{q}_t - q_t^\circ; \ell_t \rangle &= D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \tilde{q}_{t+1}) \\ &\leq D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, q_{t+1}) \\ &\leq D_\psi(\bar{q}_t, \tilde{q}_{t+1}) + D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \bar{q}_{t+1}) + \frac{\alpha}{(1-\alpha)} D_\psi(q_t^\circ, u)\end{aligned}$$

Now, summing over $t \in [T]$ we get

$$\eta \bar{\mathcal{R}}_T^D = \eta \sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \ell_t \rangle \leq \underbrace{\sum_{t=1}^T D_\psi(\bar{q}_t, \tilde{q}_{t+1})}_{\text{“Stability” term Lemma B.8}} + \frac{\alpha}{(1-\alpha)} \underbrace{\sum_{t=1}^T D_\psi(q_t^\circ, u)}_{\text{“Residual” term Lemma B.9}}$$

1296
1297
1298
1299
1300

$$+ \underbrace{\sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \bar{q}_{t+1}))}_{\substack{\text{"Penalty" term} \\ \text{Lemma C.5}}}.$$

1301 Now we combine the results from Lemma B.8 for the “Stability” term, Lemma B.9 for the “Residual”
1302 term and Lemma C.5 for the “Penalty” term, obtaining:

1303
1304
1305
1306

$$\eta \bar{\mathcal{R}}_T^D \leq \eta^2 TL + L \log(X^2 A) + \mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2 AT) \bar{D}_f^P + 9L^2 \log(2X^2 AT) \log(T)$$

□

1307
1308

C.2 FURTHER LEMMAS

1309
1310
1311
1312

Lemma C.5 (“Penalty” Bound). For $\bar{q}_t = (1 - \alpha)q_t + \alpha u$, $\alpha = \frac{1}{T+1}$, $\epsilon_{t,k} = d_{TV}(q_t^{\circ,k}, q_{t-1}^{\circ,k})$,
 $\mathcal{H}_T(\bar{D}_f^P) = TLh\left(\frac{(L^2+L)\bar{D}_f^P}{2TL}\right)$ and $\bar{D}_f^P = \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\|_1$ it
holds that

1313
1314
1315
1316

$$\sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \bar{q}_{t+1})) \leq \mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2 AT) \bar{D}_f^P + 9L^2 \log(2X^2 AT) \log(T)$$

1317
1318

Proof. First, we unravel the summation:

1319
1320
1321
1322
1323
1324
1325
1326

$$\begin{aligned} \sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \bar{q}_{t+1})) &= D_\psi(q_1^\circ, \bar{q}_1) - D_\psi(q_1^\circ, \bar{q}_2) + D_\psi(q_2^\circ, \bar{q}_2) - D_\psi(q_2^\circ, \bar{q}_3) + \dots \\ &\leq \underbrace{D_\psi(q_1^\circ, \bar{q}_1)}_{\substack{\text{Range of } \psi \\ \text{Lemma B.11}}} + \underbrace{\sum_{t=2}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_{t-1}^\circ, \bar{q}_t))}_{\substack{\text{non-telescoping term} \\ \text{Lemma C.6}}} \end{aligned}$$

1327
1328
1329

Where the inequality comes from neglecting negative terms. What is left to do is to combine Lemma B.11, which expresses the range of the negative Shannon Entropy, and Lemma C.6 providing a bound for the non-telescoping term.

1330
1331
1332

$$\begin{aligned} \sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \bar{q}_{t+1})) &\leq L \log(X^2 A) + \mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2 AT) \bar{D}_f^P \\ &\quad + 8L^2 \log(2X^2 AT) \log(T) \end{aligned}$$

1333
1334
1335
1336

□

1337
1338
1339
1340

Lemma C.6 (Bound on the non-telescoping term). For $q_t^\circ \in \Delta$, $\bar{q}_t = (1 - \alpha)q_t + \alpha u$, $\epsilon_{t,k} = d_{TV}(q_t^{\circ,k}, q_{t-1}^{\circ,k})$, $\mathcal{H}_T(\bar{D}_f^P) = TLh\left(\frac{(L^2+L)\bar{D}_f^P}{2TL}\right)$ and $\bar{D}_f^P = \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\|_1$ we have that:

1341
1342
1343
1344

$$\sum_{t=2}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_{t-1}^\circ, \bar{q}_t)) \leq \mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2 AT) \bar{D}_f^P + 8L^2 \log(2X^2 AT) \log(T)$$

1345
1346

Proof. We first focus on each time-step t . From the definition of the Bregman Divergence we have:

1347
1348
1349

$$\begin{aligned} D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_{t-1}^\circ, \bar{q}_t) &= \psi(q_t^\circ) - \psi(q_{t-1}^\circ) + \langle \nabla \psi(\bar{q}_t); q_{t-1}^\circ - q_t^\circ \rangle \\ &= \psi(q_t^\circ) - \psi(q_{t-1}^\circ) + \langle \log(\bar{q}_t); q_{t-1}^\circ - q_t^\circ \rangle + \langle \mathbf{1}; q_{t-1}^\circ - q_t^\circ \rangle \\ &\leq |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| + \|\log(\bar{q}_t)\|_\infty \|q_{t-1}^\circ - q_t^\circ\|_1 \end{aligned}$$

Where the second equality comes from the computation of the derivative of the Negative Shannon Entropy:

$$\nabla \psi(f) = \left(\frac{\partial \psi(f)}{\partial f} \right)^T = (\log(f) + 1, \dots)^T$$

Finally, the last step comes from Holder's inequality and from taking the absolute value of the difference of the entropies.

We combine the results from Lemma C.7 and Lemma B.14 to obtain over the whole time horizon:

$$\begin{aligned} \sum_{t=2}^T D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_{t-1}^\circ, \bar{q}_t) &\leq \underbrace{\sum_{t=2}^T |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)|}_{\text{Lemma C.7}} + \underbrace{\sum_{t=2}^T \|\log(\bar{q}_t)\|_\infty \|q_t^\circ - q_{t-1}^\circ\|_1}_{\text{Lemma B.14}} \\ &\leq \mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2 AT) \bar{D}_f^P \\ &\quad + 8L^2 \log(X^2 AT + X^2 A) \log(T) \end{aligned}$$

□

Lemma C.7. For $\mathcal{H}_T(\bar{D}_f^P) = TLh\left(\frac{(L^2+L)\bar{D}_f^P}{2TL}\right)$ and $\bar{D}_f^P = \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x,a) - \bar{P}_{t-1}(\cdot|x,a)\|_1$ we have that:

$$\sum_{t=2}^T |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \leq \mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2 AT) \bar{D}_f^P.$$

Proof. Once defined $\epsilon_{t,k} = d_{\text{TV}}(q_t^{\circ,k}, q_{t-1}^{\circ,k})$, $\mathcal{H}_T(\bar{D}_f^P) = TLh\left(\frac{(L^2+L)\bar{D}_f^P}{2TL}\right)$ and $\bar{D}_f^P = \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x,a) - \bar{P}_{t-1}(\cdot|x,a)\|_1$ the proof comes from a direct application of (Theorem 3, Sason, 2013):

$$\begin{aligned} &\sum_{t=2}^T |\psi(q_t^\circ) - \psi(q_{t-1}^\circ)| \\ &\leq \sum_{t=2}^T \sum_{k=0}^{L-1} |\psi(q_t^{\circ,k}) - \psi(q_{t-1}^{\circ,k})| \\ &\leq \sum_{t=2}^T \sum_{k=0}^{L-1} h(\epsilon_{t,k}) + \epsilon_{t,k} \log(X^2 A - 1) \\ &\leq TL \sum_{t=2}^T \sum_{k=0}^{L-1} \frac{1}{TL} h(\epsilon_{t,k}) + \log(X^2 AT + X^2 A) \frac{1}{2} \sum_{t=2}^T \sum_{k=0}^{L-1} \|q_t^{\circ,k} - q_{t-1}^{\circ,k}\|_1 \\ &\leq TLh\left(\sum_{t=2}^T \sum_{k=0}^{L-1} \frac{1}{TL} \epsilon_{t,k}\right) \\ &\quad + \frac{1}{2}(L^2 + L) \log(X^2 AT + X^2 A) \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x,a) - \bar{P}_{t-1}(\cdot|x,a)\|_1 \\ &\leq TLh\left(\sum_{t=2}^T \sum_{k=0}^{L-1} \frac{\|q_t^{\circ,k} - q_{t-1}^{\circ,k}\|_1}{2TL}\right) + \frac{1}{2}(L^2 + L) \log(X^2 AT + X^2 A) \bar{D}_f^P \\ &\leq TLh\left(\sum_{t=2}^T \sum_{k=0}^{L-1} \frac{(L^2 + L) \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x,a) - \bar{P}_{t-1}(\cdot|x,a)\|_1}{2TL}\right) \end{aligned}$$

$$\begin{aligned}
&\leq TLh\left(\frac{(L^2 + L)\bar{D}_f^P}{2TL}\right) + \frac{1}{2}(L^2 + L)\log(X^2AT + X^2A)\bar{D}_f^P \\
&\leq \mathcal{H}_T(\bar{D}_f^P) + L^2\log(2X^2AT)\bar{D}_f^P.
\end{aligned}$$

where the fourth inequality follows from Jensens' over the binary entropy. Assuming that $\frac{(L^2+L)\bar{D}_f^P}{2TL} \leq 1/2$ and by applying Corollary F.10 we get the fifth inequality. Finally, we use the definition of $\mathcal{H}_T(\bar{D}_f^P)$ and further simplify the bound to achieve the final result. \square

D THEORETICAL ANALYSIS WITH BANDIT-FEEDBACK AND AVERAGE SMOOTHING FUNCTIONS

In this section, we present the proofs of the results discussed in Section 5 with the specific smoothing functions $f_t(P_1, \dots, P_t) = \frac{1}{t} \sum_{t'=1}^t P_{t'}$. For convenience, we will restate the Theorems and Lemmas before providing a detailed analysis of each and report just the Lemmas that differentiate from Appendix B.

D.1 MAIN RESULTS

Corollary 5.2 (Smoothed-Regret Bound for SOMD under bandit-feedback and average smoothing). *Let $\eta = \sqrt{(21L \log(2X^2AT)(\log(T)))/(2XAT)}$, $\alpha = 1/(T+1)$, $\gamma = \eta$, smoothing functions such that $\bar{P}_t := \frac{1}{t} \sum_{t' \in [t]} P_{t'}$ and any comparator policy $\pi^\circ \in \Pi$, Algorithm 2 suffers a smoothed regret of:*

$$\bar{\mathcal{R}}_T(\pi^\circ) \leq \mathcal{O}\left(L\bar{C}_f^P + L^2\log(T) + L^{3/2}\sqrt{XAT\log(X^2AT)\log(T)}\right) \quad (17)$$

Proof. We first define $q_t = q_t^{\bar{P}_{t-1}, \pi_t}$ and $q_t^\circ = q_t^{\bar{P}_t, \pi^\circ}$. It follows that the smoothed regret term can be decomposed as

$$\bar{\mathcal{R}}_T(\pi^\circ) \leq \sum_{t=1}^T \langle q_t^{\bar{P}_t, \pi_t} - q_t^{\bar{P}_t, \pi^\circ}; \ell_t \rangle = \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t - q_t^\circ; \ell_t \rangle \right]}_{\text{Algorithmic Regret } (\bar{\mathcal{R}}_T^A) \text{ Lemma D.1}} + \underbrace{\sum_{t=1}^T \langle q_t^{\bar{P}_t, \pi_t} - q_t^{\bar{P}_{t-1}, \pi_t}; \ell_t \rangle}_{\text{Update Regret } (\bar{\mathcal{R}}_T^U)}$$

Where the expectation is with respect to the internal randomisation of the agent. Then, the result follows directly from the combination of Lemma B.1 for the Update Regret and Lemma D.1 for the Algorithmic Regret, namely:

$$\begin{aligned}
\bar{\mathcal{R}}_T(\pi^\circ) &\leq 2(L^2 + L)(1 + \log(T)) + 2L\bar{C}_f^P + 4(L^2 + L)(1 + \log(T)) \\
&\quad + 14L\sqrt{LXAT\log(2X^2AT)\log(T)} \\
&\leq 2L\bar{C}_f^P + 12L^2 + 12L^2\log(T) + 14L\sqrt{LXAT\log(2X^2AT)\log(T)}
\end{aligned}$$

which leads to the final result. \square

Corollary 5.3 (Regret Bound for SOMD under bandit-feedback and average smoothing). *Let $\eta = \sqrt{(21L \log(2X^2AT)(\log(T)))/(2XAT)}$, $\alpha = 1/(T+1)$, $\gamma = \eta$, smoothing functions such that $\bar{P}_t := \frac{1}{t} \sum_{t' \in [t]} P_{t'}$ and any comparator policy π° , Algorithm 2 suffers a regret of:*

$$\mathcal{R}_T(\pi^\circ) = \tilde{\mathcal{O}}\left(L\bar{C}_f^P + L^{3/2}\sqrt{XAT}\right). \quad (18)$$

Proof. The result follows straightforwardly from combining the results from Theorem 5.2 and Lemma 3.1. \square

Lemma D.1 (Algorithmic Regret Bound). *Choosing $\eta = \sqrt{\frac{21L \log(2X^2AT)(\log(T))}{2XAT}}$, $\alpha = \frac{1}{T+1}$ and $\gamma = \eta$ the Algorithmic Regret is bounded by*

$$\bar{\mathcal{R}}_T^A = \mathbb{E} \left[\sum_{t=1}^T \langle q_t - q_t^\circ; \ell_t \rangle \right] \leq 2L\bar{C}_f^P + 4(L^2 + L)(\log(T) + 1) + 14L\sqrt{LXAT \log(2X^2AT) \log(T)}$$

Proof. The proof relies on the following decomposition:

$$\bar{\mathcal{R}}_T^A = \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \hat{\ell}_t \rangle \right]}_{\text{Descent Regret } (\bar{\mathcal{R}}_T^D)} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t - \bar{q}_t; \hat{\ell}_t \rangle \right]}_{\text{Regularization Regret } (\bar{\mathcal{R}}_T^R)} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t; \ell_t - \hat{\ell}_t \rangle \right]}_{\text{Bias 1 (B1)}} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t^\circ; \hat{\ell}_t - \ell_t \rangle \right]}_{\text{Bias 2 (B2)}}$$

By applying Lemma D.2 for the Regularization Regret, Lemma D.3 for the Descent Regret and Lemmas D.5, D.6 for the bias terms respectively, we can rewrite the Algorithmic Regret as

$$\begin{aligned} \bar{\mathcal{R}}_T^A &\leq \frac{L}{\eta} + \eta LXAT + 2(L^2 + L)(\log(T) + 1) + L\bar{C}_f^P + \frac{1}{\eta} 18L^2 \log(2X^2AT) \log(T) \\ &\quad + \frac{2L}{\eta} + \eta LXAT + L\bar{C}_f^P + 2(L^2 + L)(\log(T) + 1) \\ &\leq 2L\bar{C}_f^P + 4(L^2 + L)(\log(T) + 1) + \eta 2LXAT + \frac{1}{\eta} 21L^2 \log(2X^2AT) \log(T) \end{aligned}$$

We first chose the optimal η as:

$$\eta = \sqrt{\frac{21L \log(2X^2AT)(\log(T))}{2XAT}}$$

and then we upperbound the Algorithmic Regret after substituting the optimal η , namely:

$$\begin{aligned} \bar{\mathcal{R}}_T^A &\leq 2L\bar{C}_f^P + 4(L^2 + L)(\log(T) + 1) + \eta 2LXAT + \frac{1}{\eta} 21L^2 \log(2X^2AT) \log(T) \\ &\leq 2L\bar{C}_f^P + 4(L^2 + L)(\log(T) + 1) + 2L\sqrt{42LXAT \log(2X^2AT) \log(T)} \\ &\leq 2L\bar{C}_f^P + 4(L^2 + L)(\log(T) + 1) + 14L\sqrt{LXAT \log(2X^2AT) \log(T)} \end{aligned}$$

which leads to the final result. \square

Lemma D.2 (Regularisation Regret Bound). *For $\bar{q}_t = (1 - \alpha)q_t + \alpha u$, $u(x, a, x') = \frac{1}{X_k \cdot A \cdot X_{k+1}} \forall k \in \llbracket L - 1 \rrbracket$, $\alpha = \frac{1}{T+1}$ and $\gamma = \eta$, it holds that*

$$\bar{\mathcal{R}}_T^R = \mathbb{E} \left[\sum_{t=1}^T \langle q_t - \bar{q}_t; \ell_t \rangle \right] \leq 2\frac{L}{\gamma} = 2\frac{L}{\eta}$$

Proof.

$$\begin{aligned} \bar{\mathcal{R}}_T^R &= \sum_{t=1}^T \langle q_t - \bar{q}_t; \ell_t \rangle = \alpha \sum_{t=1}^T \langle q_t - u; \ell_t \rangle \\ &\leq \frac{\alpha}{\gamma} \sum_{t=1}^T \|q_t - u\|_1 \\ &\leq \frac{\alpha}{\gamma} \sum_{t=1}^T \sum_{k=0}^{L-1} \|q_t^k - u\|_1 \\ &\leq 2\frac{\alpha}{\gamma} LT = \frac{1}{\gamma} \frac{2}{1+T} LT \leq 2\frac{L}{\gamma} \end{aligned}$$

Where the first equality comes from the definition of \bar{q}_t , the first inequality follows from Holder's and the fact that $\|\hat{\ell}_t\|_\infty \leq \frac{1}{\gamma}$. The last from the trivial bound on the 1-norm between distributions and the definition of $\alpha = \frac{1}{T+1}$. \square

Lemma D.3 (Descent Regret Bound). *For $\alpha = \frac{1}{T+1}$ the following fact holds*

$$\eta \bar{\mathcal{R}}_T^D = \eta \mathbb{E} \left[\sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \hat{\ell}_t \rangle \right] \leq \frac{\eta^2}{\gamma^2} L + \eta^2 L A X T + \frac{\eta^2}{\gamma} 2(L^2 + L)(\log(T) + 1) + \frac{\eta^2}{\gamma} L \bar{C}_f^P \\ + 18L^2 \log(2X^2 AT) \log(T)$$

Proof. The proof follows the same steps as those in Lemma B.4 replacing ℓ_t with $\hat{\ell}_t$ to obtain,

$$\eta \bar{\mathcal{R}}_T^D = \eta \mathbb{E} \left[\sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \hat{\ell}_t \rangle \right] \\ \leq \mathbb{E} \left[\underbrace{\sum_{t=1}^T D_\psi(\bar{q}_t, \tilde{q}_{t+1})}_{\text{“Stability” term}} + \underbrace{\frac{\alpha}{(1-\alpha)} \sum_{t=1}^T D_\psi(q_t^\circ, u)}_{\text{“Residual” term}} + \underbrace{\sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \tilde{q}_{t+1}))}_{\text{“Penalty” term}} \right].$$

Now we combine the results from Lemma D.4 for the “Stability” term, Lemma B.9 for the “Residual” term and Lemma B.10 for the “Penalty” term, obtaining:

$$\eta \bar{\mathcal{R}}_T^D \leq \frac{\eta^2}{\gamma^2} L + \eta^2 L A X T + \frac{\eta^2}{\gamma} 2(L^2 + L)(\log(T) + 1) + \frac{\eta^2}{\gamma} L \bar{C}_f^P + L \log(X^2 A) \\ + 17L^2 \log(2X^2 AT) \log(T)$$

□

D.2 AUXILIARY LEMMAS FOR THE BANDIT FEEDBACK

Lemma D.4 (“Stability” Term Bound). *Choosing $\alpha = \frac{1}{T+1}$ and $\gamma > 0$,*

$$\mathbb{E} \left[\sum_{t=1}^T D_\psi(\bar{q}_t, \tilde{q}_{t+1}) \right] \leq \frac{\eta^2}{\gamma^2} L + \eta^2 L A X T + \frac{\eta^2}{\gamma} 2(L^2 + L)(\log(T) + 1) + \frac{\eta^2}{\gamma} L \bar{C}_f^P$$

Proof. The term can be bounded as follows:

$$\mathbb{E} \left[\sum_{t=1}^T D_\psi(\bar{q}_t, \tilde{q}_{t+1}) \right] = \mathbb{E} \left[\sum_{t=1}^T \left[\sum_{x,a,x'} \bar{q}_t(x,a,x') \log \left(\frac{\bar{q}_t(x,a,x')}{\tilde{q}_{t+1}(x,a,x')} \right) - \sum_{x,a,x'} (\bar{q}_t(x,a,x') - \tilde{q}_{t+1}(x,a,x')) \right] \right] \\ = \mathbb{E} \left[\sum_{t=1}^T \left[\sum_{x,a,x'} \eta \hat{\ell}_t(x,a) \bar{q}_t(x,a,x') + \bar{q}_t(x,a,x') \exp(-\eta \hat{\ell}_t(x,a)) - \bar{q}_t(x,a,x') \right] \right] \\ \leq \mathbb{E} \left[\sum_{t=1}^T \left[\sum_{x,a,x'} \eta^2 \bar{q}_t(x,a,x') \mathbb{E}_t \left[\hat{\ell}_t^2(x,a) \right] \right] \right] \\ \leq (1-\alpha) \eta^2 \sum_{t=1}^T \sum_{x,a} q^{\bar{P}_{t-1}, \pi_t}(x,a) \frac{q^{P_t, \pi_t}(x,a)}{q^{P_t, \pi_t}(x,a) (q^{P_t, \pi_t}(x,a) + \gamma)} \\ + \alpha \eta^2 \sum_{t=1}^T \sum_{x,a} \frac{u(x,a) q^{P_t, \pi_t}(x,a)}{(q^{P_t, \pi_t}(x,a) + \gamma)^2} \\ \leq \frac{\eta^2}{\gamma^2} \alpha T L + \eta^2 \sum_{t=1}^T \sum_{x,a} \frac{q^{\bar{P}_{t-1}, \pi_t}(x,a)}{(q^{P_t, \pi_t}(x,a) + \gamma)} \\ \leq \frac{\eta^2}{\gamma^2} L + \eta^2 \sum_{t=1}^T \sum_{x,a} \frac{|q^{\bar{P}_{t-1}, \pi_t}(x,a) \pm q^{\bar{P}_t, \pi_t}(x,a) \pm q^{P_t, \pi_t}(x,a)|}{(q^{P_t, \pi_t}(x,a) + \gamma)}$$

$$\begin{aligned}
&\leq \frac{\eta^2}{\gamma^2} L + \eta^2 \sum_{t=1}^T \sum_{x,a} \frac{\left| q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{\bar{P}_t, \pi_t}(x, a) \right|}{(q^{P_t, \pi_t}(x, a) + \gamma)} \\
&\quad + \eta^2 \sum_{t=1}^T \sum_{x,a} \frac{\left| q^{\bar{P}_t, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \right|}{(q^{P_t, \pi_t}(x, a) + \gamma)} + \eta^2 \sum_{t=1}^T \sum_{x,a} \frac{q^{P_t, \pi_t}(x, a)}{(q^{P_t, \pi_t}(x, a) + \gamma)} \\
&\leq \frac{\eta^2}{\gamma^2} L + \eta^2 L X A T + \frac{\eta^2}{\gamma} \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_t, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \right| \\
&\quad + \frac{\eta^2}{\gamma} \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{\bar{P}_t, \pi_t}(x, a) \right| \\
&\leq \frac{\eta^2}{\gamma^2} L + \eta^2 L X A T + \frac{\eta^2}{\gamma} \sum_{t=1}^T \sum_x \left| q^{\bar{P}_t, \pi_t}(x) - q^{P_t, \pi_t}(x) \right| \\
&\quad + \frac{\eta^2}{\gamma} \sum_{t=1}^T \left\| q^{\bar{P}_{t-1}, \pi_t} - q^{\bar{P}_t, \pi_t} \right\|_1 \\
&\leq \frac{\eta^2}{\gamma^2} L + \eta^2 L X A T + \frac{\eta^2}{\gamma} L \bar{C}_f^P + \frac{\eta^2}{\gamma} 2(L^2 + L)(\log(T) + 1)
\end{aligned}$$

Where the first equality comes from the definition of generalized KL divergence, the second by applying the definition of the solution of the unconstrained optimization problem, namely:

$$\tilde{q}_{t+1}(x, a, x') = \bar{q}_t(x, a, x') \exp\left(-\eta \hat{\ell}_t(x, a)\right), \forall (x, a, x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}.$$

The first inequality comes from the standard bound of the exponential function,

$$e^{-\eta \hat{\ell}_t(x, a)} \leq 1 - \eta \hat{\ell}_t(x, a) + \left(\eta \hat{\ell}_t(x, a)\right)^2, \forall \eta \hat{\ell}_t(x, a) \geq 0$$

which is satisfied $\forall \gamma > 0$. The fifth inequality comes from setting $\alpha = \frac{1}{T+1}$ and from triangle inequality. Finally, the last step comes from Corollary F.6 and Lemma B.1. \square

Lemma D.5 (Bias 1 Bound). For $\hat{\ell}_t$ as in Algorithm 2 it holds that,

$$(BIAS 1) \quad \mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_{t-1}, \pi_t}; \ell_t - \hat{\ell}_t \rangle \right] \leq \gamma L X A T + L \bar{C}_f^P + 2(L^2 + L)(\log(T) + 1)$$

Proof.

$$\begin{aligned}
\mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_{t-1}, \pi_t}; \ell_t - \hat{\ell}_t \rangle \right] &= \mathbb{E} \left[\sum_{t=1}^T \sum_{x,a} q^{\bar{P}_{t-1}, \pi_t}(x, a) (\ell_t(x, a) - \mathbb{E}_t[\hat{\ell}_t(x, a)]) \right] \\
&= \sum_{t=1}^T \sum_{x,a} q^{\bar{P}_{t-1}, \pi_t}(x, a) \ell_t(x, a) \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\
&\quad \pm \sum_{t=1}^T \sum_{x,a} q^{P_t, \pi_t}(x, a) \ell_t(x, a) \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\
&\leq \sum_{t=1}^T \sum_{x,a} q^{P_t, \pi_t}(x, a) \ell_t(x, a) \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\
&\quad + \sum_{t=1}^T \sum_{x,a} \left(q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \right) \ell_t(x, a) \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right)
\end{aligned}$$

$$\begin{aligned}
&\leq \sum_{t=1}^T \sum_{x,a} q^{P_t, \pi_t}(x, a) \ell_t(x, a) \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\
&\quad + \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \right| \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\
&\leq \gamma L X A T + \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \right| \\
&\leq \gamma L X A T + \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \pm q^{\bar{P}_t, \pi_t}(x, a) \right| \\
&\leq \gamma L X A T + \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{\bar{P}_t, \pi_t}(x, a) \right| \\
&\quad + \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_t, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \right| \\
&\leq \gamma L X A T + \sum_{t=1}^T \left\| q^{\bar{P}_{t-1}, \pi_t} - q^{\bar{P}_t, \pi_t} \right\|_1 \\
&\quad + \sum_{t=1}^T \sum_x \left| q^{\bar{P}_t, \pi_t}(x) - q^{P_t, \pi_t}(x) \right| \\
&\leq \gamma L X A T + 2(L^2 + L)(\log(T) + 1) + L \bar{C}_f^P
\end{aligned}$$

Where the last inequality comes from applying Lemma B.1 and Corollary F.6. \square

Lemma D.6 (Bias 2 Bound). *For $\hat{\ell}_t$ as in Algorithm 2 it holds that,*

$$(BIAS 2) \quad \mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_t, \pi^\circ}; \hat{\ell}_t - \ell_t \rangle \right] \leq 0$$

Proof. It is sufficient to recall that,

$$\ell_t(x, a) - \mathbb{E}_t \left[\hat{\ell}_t(x, a) \right] \in \left[0, \frac{\gamma \ell_t(x, a)}{q^{P_t, \pi_t}(x, a)} \right]$$

Namely, that we are underestimating the true loss. \square

E THEORETICAL ANALYSIS WITH BANDIT-FEEDBACK AND GENERIC SMOOTHING FUNCTION

In this section, we present the proofs of the results discussed in Section 5 related to the regret analysis agnostic of the smoothing function to be used. For convenience, we will restate the Theorems and Lemmas before providing a detailed analysis of each and report just the Lemmas that differentiate from Appendix B and Appendix D.

E.1 MAIN RESULTS

Theorem 5.1 (Smoothed-Regret Bound for SOMD under bandit-feedback). *Let $\eta = \sqrt{(13L \log(2X^2 AT) \rho_T^f) / (2X AT)}$, $\alpha = 1/(T + 1)$, $\gamma = \eta$, generic smoothing functions and any comparator policy $\pi^\circ \in \Pi$, Algorithm 2 suffers a smoothed regret of:*

$$\bar{\mathcal{R}}_T(\pi^\circ) \leq \mathcal{O} \left(L^2 \bar{D}_f^P + L \bar{C}_f^P + L^{3/2} \sqrt{X AT \rho_T^f \log(X^2 AT)} \right). \quad (16)$$

1674 *Proof.* We first define $q_t = q_t^{\bar{P}^{t-1}, \pi_t}$ and $q_t^\circ = q_t^{\bar{P}^t, \pi^\circ}$. It follows that the smoothed regret term can
 1675 be decomposed as
 1676

$$1677 \bar{\mathcal{R}}_T(\pi^\circ) \leq \sum_{t=1}^T \langle q_t^{\bar{P}^t, \pi_t} - q_t^{\bar{P}^t, \pi^\circ}; \ell_t \rangle = \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t - q_t^\circ; \ell_t \rangle \right]}_{\substack{\text{Algorithmic Regret } (\bar{\mathcal{R}}_T^A) \\ \text{Lemma E.1}}} + \underbrace{\sum_{t=1}^T \langle q_t^{\bar{P}^t, \pi_t} - q_t^{\bar{P}^{t-1}, \pi_t}; \ell_t \rangle}_{\text{Update Regret } (\bar{\mathcal{R}}_T^U)},$$

1680 Where the expectation is with respect to the internal randomisation of the agent. Then, the result
 1681 follows directly from the combination of Lemma C.2 for the Update Regret and Lemma E.1 for the
 1682 Algorithmic Regret, namely:
 1683

$$1684 \bar{\mathcal{R}}_T(\pi^\circ) \leq (L^2 + L)\bar{D}_f^P + 2(L^2 + L)\bar{D}_f^P + 2L\bar{C}_f^P + 12L\sqrt{LXAT \log(2X^2AT)}\rho_T^f$$

1685 which leads to the final result. \square

1686
 1687 **Lemma E.1** (Algorithmic Regret Bound). *Choosing $\eta = \sqrt{\frac{13L \log(2X^2AT)\rho_T^f}{2XAT}}$, $\alpha = \frac{1}{T+1}$, $\gamma = \eta$,
 1688 $\mathcal{H}_T(\bar{D}_f^P) = TLh\left(\frac{(L^2+L)\bar{D}_f^P}{2TL}\right)$, $\bar{D}_f^P = \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\|_1$ and for
 1689 $\rho_T^f = \log(T) + \bar{D}_f^P + \mathcal{H}_f$, the Algorithmic Regret is bounded by*

$$1690 \bar{\mathcal{R}}_T^A = \mathbb{E} \left[\sum_{t=1}^T \langle q_t - q_t^\circ; \ell_t \rangle \right] \leq 2(L^2 + L)\bar{D}_f^P + 2L\bar{C}_f^P + 12L\sqrt{LXAT \log(2X^2AT)}\rho_T^f$$

1691
 1692 *Proof.* The proof relies on the following decomposition:
 1693

$$1694 \bar{\mathcal{R}}_T^A = \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \hat{\ell}_t \rangle \right]}_{\text{Descent Regret } (\bar{\mathcal{R}}_T^D)} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t - \bar{q}_t; \hat{\ell}_t \rangle \right]}_{\text{Regularization Regret } (\bar{\mathcal{R}}_T^R)} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t; \ell_t - \hat{\ell}_t \rangle \right]}_{\text{Bias 1 (B1)}} + \underbrace{\mathbb{E} \left[\sum_{t=1}^T \langle q_t^\circ; \hat{\ell}_t - \ell_t \rangle \right]}_{\text{Bias 2 (B2)}}$$

1700 By applying Lemma D.2 for the Regularization Regret, Lemma E.2 for the Descent Regret and
 1701 Lemmas E.4, D.6 for the bias terms, we can rewrite the Algorithmic Regret as

$$1702 \begin{aligned} \bar{\mathcal{R}}_T^A &\leq 2\frac{L}{\eta} + \frac{\eta}{\gamma^2}L + \eta LXAT + \frac{\eta}{\gamma}L\bar{C}_f^P + \frac{\eta}{\gamma}(L^2 + L)\bar{D}_f^P + \gamma LXAT + L\bar{C}_f^P + (L^2 + L)\bar{D}_f^P \\ &\quad + \frac{1}{\eta} \left(\mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2AT)\bar{D}_f^P + 10L^2 \log(2X^2AT) \log(T) \right) \\ &\leq 2\frac{L}{\eta} + \frac{1}{\eta}L + \eta LXAT + L\bar{C}_f^P + (L^2 + L)\bar{D}_f^P + \eta LXAT + L\bar{C}_f^P + (L^2 + L)\bar{D}_f^P \\ &\quad + \frac{1}{\eta} \left(\mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2AT)\bar{D}_f^P + 10L^2 \log(2X^2AT) \log(T) \right) \\ &\leq \eta 2LXAT + 2L\bar{C}_f^P + 2(L^2 + L)\bar{D}_f^P \\ &\quad + \frac{1}{\eta} \left(\mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2AT)\bar{D}_f^P + 13L^2 \log(2X^2AT) \log(T) \right) \\ &\leq \eta 2LXAT + 2L\bar{C}_f^P + 2(L^2 + L)\bar{D}_f^P \\ &\quad + \frac{1}{\eta} 13L^2 \log(2X^2AT) \left(\mathcal{H}_T(\bar{D}_f^P) + \bar{D}_f^P + \log(T) \right) \\ &\leq \eta 2LXAT + 2(L^2 + L)\bar{D}_f^P + 2L\bar{C}_f^P + \frac{13}{\eta} L^2 \log(2X^2AT) \rho_T^f \end{aligned}$$

1723 We first chose the optimal η as:

$$1724 \eta = \sqrt{\frac{13L \log(2X^2AT)\rho_T^f}{2XAT}}$$

and then we upperbound the Algorithmic Regret after substituting the optimal η , namely:

$$\begin{aligned}\bar{\mathcal{R}}_T^A &\leq 2(L^2 + L)\bar{D}_f^P + 2L\bar{C}_f^P + 2L\sqrt{26LXAT \log(2X^2AT)}\rho_T^f \\ &\leq 2(L^2 + L)\bar{D}_f^P + 2L\bar{C}_f^P + 12L\sqrt{LXAT \log(2X^2AT)}\rho_T^f\end{aligned}$$

which leads to the final result. \square

Lemma E.2 (Descent Regret Bound). For $\alpha = \frac{1}{T+1}$ and defining $\mathcal{H}_T(\bar{D}_f^P) = TLh\left(\frac{(L^2+L)\bar{D}_f^P}{2TL}\right)$,

$\bar{D}_f^P = \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\|_1$, it holds that:

$$\begin{aligned}\eta\bar{\mathcal{R}}_T^D &= \eta\mathbb{E}\left[\sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \hat{\ell}_t \rangle\right] \leq \frac{\eta^2}{\gamma^2}L + \eta^2LXAT + \frac{\eta^2}{\gamma}L\bar{C}_f^P + \frac{\eta^2}{\gamma}(L^2 + L)\bar{D}_f^P \\ &\quad + \mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2AT)\bar{D}_f^P + 10L^2 \log(2X^2AT) \log(T)\end{aligned}$$

Proof. The proof follows the same steps as those in Lemma C.4 replacing ℓ_t with $\hat{\ell}_t$ to obtain,

$$\begin{aligned}\eta\bar{\mathcal{R}}_T^D &= \eta\mathbb{E}\left[\sum_{t=1}^T \langle \bar{q}_t - q_t^\circ; \hat{\ell}_t \rangle\right] \\ &\leq \mathbb{E}\left[\underbrace{\sum_{t=1}^T D_\psi(\bar{q}_t, \tilde{q}_{t+1})}_{\text{Stability term}} + \underbrace{\frac{\alpha}{(1-\alpha)} \sum_{t=1}^T D_\psi(q_t^\circ, u)}_{\text{Residual term}} + \underbrace{\sum_{t=1}^T (D_\psi(q_t^\circ, \bar{q}_t) - D_\psi(q_t^\circ, \tilde{q}_{t+1}))}_{\text{Penalty term}}\right].\end{aligned}$$

Now we combine the results from Lemma E.3 for the ‘‘Stability’’ term, Lemma B.9 for the ‘‘Residual’’ term and Lemma C.5 for the ‘‘Penalty’’ term, obtaining:

$$\begin{aligned}\eta\bar{\mathcal{R}}_T^D &\leq \frac{\eta^2}{\gamma^2}L + \eta^2LXAT + \frac{\eta^2}{\gamma}L\bar{C}_f^P + \frac{\eta^2}{\gamma}(L^2 + L)\bar{D}_f^P \\ &\quad + \mathcal{H}_T(\bar{D}_f^P) + L^2 \log(2X^2AT)\bar{D}_f^P + 10L^2 \log(2X^2AT) \log(T)\end{aligned}$$

\square

E.2 AUXILIARY LEMMAS FOR THE BANDIT FEEDBACK

Lemma E.3 (Bound of ‘‘Stability’’ term). Chosing $\alpha = \frac{1}{T+1}$, $\gamma > 0$ and defining $\mathcal{H}_T(\bar{D}_f^P) = TLh\left(\frac{(L^2+L)\bar{D}_f^P}{2TL}\right)$, $\bar{D}_f^P = \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\|_1$, it holds that:

$$\mathbb{E}\left[\sum_{t=1}^T D_\psi(\bar{q}_t, \tilde{q}_{t+1})\right] \leq \frac{\eta^2}{\gamma^2}L + \eta^2LXAT + \frac{\eta^2}{\gamma}L\bar{C}_f^P + \frac{\eta^2}{\gamma}(L^2 + L)\bar{D}_f^P$$

Proof. The term can be bounded as follows:

$$\begin{aligned}\mathbb{E}\left[\sum_{t=1}^T D_\psi(\bar{q}_t, \tilde{q}_{t+1})\right] &= \mathbb{E}\left[\sum_{t=1}^T \left[\sum_{x,a,x'} \bar{q}_t(x, a, x') \log\left(\frac{\bar{q}_t(x, a, x')}{\tilde{q}_{t+1}(x, a, x')}\right) - \sum_{x,a,x'} (\bar{q}_t(x, a, x') - \tilde{q}_{t+1}(x, a, x'))\right]\right] \\ &= \mathbb{E}\left[\sum_{t=1}^T \left[\sum_{x,a,x'} \eta\hat{\ell}_t(x, a)\bar{q}_t(x, a, x') + \bar{q}_t(x, a, x') \exp(-\eta\hat{\ell}_t(x, a)) - \bar{q}_t(x, a, x')\right]\right] \\ &\leq \mathbb{E}\left[\sum_{t=1}^T \left[\sum_{x,a,x'} \eta^2\bar{q}_t(x, a, x')\mathbb{E}_t[\hat{\ell}_t^2(x, a)]\right]\right]\end{aligned}$$

$$\begin{aligned}
&\leq (1 - \alpha)\eta^2 \sum_{t=1}^T \sum_{x,a} q^{\bar{P}_{t-1}, \pi_t}(x, a) \frac{q^{P_t, \pi_t}(x, a)}{q^{P_t, \pi_t}(x, a)(q^{P_t, \pi_t}(x, a) + \gamma)} \\
&\quad + \alpha\eta^2 \sum_{t=1}^T \sum_{x,a} \frac{u(x, a)q^{P_t, \pi_t}(x, a)}{(q^{P_t, \pi_t}(x, a) + \gamma)^2} \\
&\leq \frac{\eta^2}{\gamma^2} \alpha TL + \eta^2 \sum_{t=1}^T \sum_{x,a} \frac{q^{\bar{P}_{t-1}, \pi_t}(x, a)}{(q^{P_t, \pi_t}(x, a) + \gamma)} \\
&\leq \frac{\eta^2}{\gamma^2} L + \eta^2 \sum_{t=1}^T \sum_{x,a} \frac{|q^{\bar{P}_{t-1}, \pi_t}(x, a) \pm q^{\bar{P}_t, \pi_t}(x, a) \pm q^{P_t, \pi_t}(x, a)|}{(q^{P_t, \pi_t}(x, a) + \gamma)} \\
&\leq \frac{\eta^2}{\gamma^2} L + \eta^2 \sum_{t=1}^T \sum_{x,a} \frac{|q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{\bar{P}_t, \pi_t}(x, a)|}{(q^{P_t, \pi_t}(x, a) + \gamma)} \\
&\quad + \eta^2 \sum_{t=1}^T \sum_{x,a} \frac{|q^{\bar{P}_t, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a)|}{(q^{P_t, \pi_t}(x, a) + \gamma)} + \eta^2 \sum_{t=1}^T \sum_{x,a} \frac{q^{P_t, \pi_t}(x, a)}{(q^{P_t, \pi_t}(x, a) + \gamma)} \\
&\leq \frac{\eta^2}{\gamma^2} L + \eta^2 LXAT + \frac{\eta^2}{\gamma} \sum_{t=1}^T \sum_{x,a} |q^{\bar{P}_t, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a)| \\
&\quad + \frac{\eta^2}{\gamma} \sum_{t=1}^T \sum_{x,a} |q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{\bar{P}_t, \pi_t}(x, a)| \\
&\leq \frac{\eta^2}{\gamma^2} L + \eta^2 LXAT + \frac{\eta^2}{\gamma} \sum_{t=1}^T \sum_x |q^{\bar{P}_t, \pi_t}(x) - q^{P_t, \pi_t}(x)| \\
&\quad + \frac{\eta^2}{\gamma} \sum_{t=1}^T \|q^{\bar{P}_{t-1}, \pi_t} - q^{\bar{P}_t, \pi_t}\|_1 \\
&\leq \frac{\eta^2}{\gamma^2} L + \eta^2 LXAT + \frac{\eta^2}{\gamma} L\bar{C}_f^P + \frac{\eta^2}{\gamma} (L^2 + L)\bar{D}_f^P
\end{aligned}$$

Where the first equality comes from the definition of generalized KL divergence, the second by applying the definition of the solution of the unconstrained optimization problem, namely:

$$\bar{q}_{t+1}(x, a, x') = \bar{q}_t(x, a, x') \exp\left(-\eta \hat{\ell}_t(x, a)\right), \forall (x, a, x') \in \mathcal{X} \times \mathcal{A} \times \mathcal{X}$$

and further simplifications. The first inequality comes from the standard bound of the exponential function,

$$e^{-\eta \hat{\ell}_t(x, a)} \leq 1 - \eta \hat{\ell}_t(x, a) + \left(\eta \hat{\ell}_t(x, a)\right)^2, \forall \eta \hat{\ell}_t(x, a) \geq 0$$

which is satisfied $\forall \gamma > 0$. The fourth inequality comes from setting $\alpha = \frac{1}{T+1}$ and from triangle inequality. Finally, the last inequality comes from Corollary F.6 and Lemma C.2. \square

Lemma E.4 (Bound on Bias 1). For $\gamma > 0$ and defining $\bar{D}_f^P = \sum_{t=2}^T \max_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\|_1$, it holds that:

$$\mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_{t-1}, \pi_t}; \ell_t - \hat{\ell}_t \rangle \right] \leq \gamma LXAT + L\bar{C}_f^P + (L^2 + L)\bar{D}_f^P$$

Proof.

$$\mathbb{E} \left[\sum_{t=1}^T \langle q^{\bar{P}_{t-1}, \pi_t}; \ell_t - \hat{\ell}_t \rangle \right] = \mathbb{E} \left[\sum_{t=1}^T \sum_{x,a} q^{\bar{P}_{t-1}, \pi_t}(x, a) (\ell_t(x, a) - \mathbb{E}_t[\hat{\ell}_t(x, a)]) \right]$$

$$\begin{aligned}
&= \sum_{t=1}^T \sum_{x,a} q^{\bar{P}_{t-1}, \pi_t}(x, a) \ell_t(x, a) \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\
&\quad \pm \sum_{t=1}^T \sum_{x,a} q^{P_t, \pi_t}(x, a) \ell_t(x, a) \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\
&\leq \sum_{t=1}^T \sum_{x,a} q^{P_t, \pi_t}(x, a) \ell_t(x, a) \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\
&\quad + \sum_{t=1}^T \sum_{x,a} \left(q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \right) \ell_t(x, a) \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\
&\leq \sum_{t=1}^T \sum_{x,a} q^{P_t, \pi_t}(x, a) \ell_t(x, a) \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\
&\quad + \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \right| \left(\frac{\gamma}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\
&\leq \gamma L X A T + \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \right| \\
&\leq \gamma L X A T + \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \pm q^{\bar{P}_{t-1}, \pi_t}(x, a) \right| \\
&\leq \gamma L X A T + \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{\bar{P}_{t-1}, \pi_t}(x, a) \right| \\
&\quad + \sum_{t=1}^T \sum_{x,a} \left| q^{\bar{P}_{t-1}, \pi_t}(x, a) - q^{P_t, \pi_t}(x, a) \right| \\
&\leq \gamma L X A T + \sum_{t=1}^T \left\| q^{\bar{P}_{t-1}, \pi_t} - q^{P_t, \pi_t} \right\|_1 + \sum_{t=1}^T \sum_x \left| q^{\bar{P}_{t-1}, \pi_t}(x) - q^{P_t, \pi_t}(x) \right| \\
&\leq \gamma L X A T + (L^2 + L) \bar{D}_f^P + L \bar{C}_f^P
\end{aligned}$$

Where the last inequality comes from applying Lemma C.2 and Corollary F.6. \square

F INSTRUMENTAL LEMMAS

Here we report a few additional instrumental lemmas used throughout the proofs together with known lemmas from some references.

F.1 STATISTICAL PROPERTIES OF $\hat{\ell}_t$

Lemma F.1 (Bias of $\hat{\ell}_t$). *Given the estimator used by Algorithm 2, and for $\gamma > 0$ we have that,*

$$\ell_t(x, a) - \mathbb{E}_t[\ell_t(x, a)] \leq \frac{\gamma \ell_t(x, a)}{q^{P_t, \pi_t}(x, a)}$$

Proof.

$$\begin{aligned}
\ell_t(x, a) - \mathbb{E}_t[\ell_t(x, a)] &= \ell_t(x, a) - \frac{\ell_t(x, a) \mathbb{E}_t[\mathbb{1}_t(x, a)]}{q^{P_t, \pi_t}(x, a) + \gamma} \\
&= \ell_t(x, a) \left(1 - \frac{q^{P_t, \pi_t}(x, a)}{q^{P_t, \pi_t}(x, a) + \gamma} \right)
\end{aligned}$$

1890
1891
1892
1893
1894
1895
1896

$$\begin{aligned} &= \ell_t(x, a) \left(1 - \frac{q^{P_t, \pi_t}(x, a)}{q^{P_t, \pi_t}(x, a) + \gamma} \right) \\ &\leq \frac{\gamma \ell_t(x, a)}{q^{P_t, \pi_t}(x, a)} \end{aligned}$$

□

1897 **Lemma F.2** (Second-order moment of $\hat{\ell}_t$). *Given the estimator used by Algorithm 2, and for $\gamma > 0$*
1898 *we have that,*

1899
1900
1901

$$\mathbb{E}_t[\hat{\ell}_t^2(x, a)] \leq \frac{\ell_t^2(x, a)}{(q^{P_t, \pi_t}(x, a) + \gamma)}$$

1902 *Proof.*

1903
1904
1905
1906
1907
1908
1909
1910
1911
1912

$$\begin{aligned} \mathbb{E}_t[\hat{\ell}_t^2(x, a)] &\leq \left(\frac{\ell_t^2(x, a)}{(q^{P_t, \pi_t}(x, a) + \gamma)^2} \mathbb{E}_t[\mathbf{1}_t(x, a)] \right) \\ &= \frac{\ell_t^2(x, a) q^{P_t, \pi_t}(x, a)}{(q^{P_t, \pi_t}(x, a) + \gamma)^2} \\ &\leq \frac{\ell_t^2(x, a)}{(q^{P_t, \pi_t}(x, a) + \gamma)} \end{aligned}$$

□

1913
1914

F.2 PROPERTIES OF AVERAGE SMOOTHING

1915
1916
1917

Lemma F.3 (Bound on Smoothed Transitions). *Let $\bar{P}_t(\cdot|x, a) = \frac{1}{t} \sum_{t'=1}^t P_{t'}(\cdot|x, a)$, then*
1918 $\|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\|_1 \leq \frac{2}{t}, \forall t \in [T], \forall x, a \in \mathcal{X} \times \mathcal{A}$.

1919 *Proof.*

1920
1921
1922
1923
1924
1925
1926
1927
1928

$$\begin{aligned} \|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\|_1 &= \left\| \frac{1}{t} \sum_{t'=1}^t P_{t'}(\cdot|x, a) - \frac{1}{t-1} \sum_{t'=1}^{t-1} P_{t'}(\cdot|x, a) \right\|_1 \\ &\leq \left\| \frac{t-1-t}{t(t-1)} \sum_{t'=1}^{t-1} P_{t'}(\cdot|x, a) \right\|_1 + \left\| \frac{1}{t} P_t(\cdot|x, a) \right\|_1 \\ &\leq \frac{1}{t(t-1)} \sum_{t'=1}^{t-1} 1 + \frac{1}{t} = \frac{2}{t} \end{aligned}$$

1929
1930

Where the first inequality comes from triangle inequality and the second from the fact we are dealing with elements of simplexes. □

1931
1932

Lemma F.4 (Bound on $\epsilon_{t,k}$). *For $\epsilon_{t,k} = \frac{1}{2} \|q_t^{\circ,k} - q_{t-1}^{\circ,k}\|_1$ and $\bar{P}_t = \frac{1}{t} \sum_{t'=1}^t P_{t'}$ we have that,*

1933
1934
1935

$$\epsilon_{t,k} \leq \epsilon_t = \frac{L}{t}$$

1936 *Proof.*

1937
1938
1939
1940
1941
1942
1943

$$\begin{aligned} \epsilon_{t,k} &= \frac{1}{2} \|q_t^{\circ,k} - q_{t-1}^{\circ,k}\|_1 \\ &= \frac{1}{2} \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k+1}} \|q_t^{\circ,k}(x, a, x') - q_{t-1}^{\circ,k}(x, a, x')\|_1 \\ &= \frac{1}{2} \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k+1}} \left| q_t^{\circ,k}(x, a) \bar{P}_t(x'|x, a) - q_{t-1}^{\circ,k}(x, a) \bar{P}_{t-1}(x'|x, a) \right| \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k+1}} \left| q_t^{\circ, k}(x, a) \bar{P}_t(x'|x, a) - q_{t-1}^{\circ, k}(x, a) \bar{P}_{t-1}(x'|x, a) \pm q_{t-1}^{\circ, k}(x, a) \bar{P}_t(x'|x, a) \right| \\
&\leq \frac{1}{2} \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k+1}} \left| q_t^{\circ, k}(x, a) \bar{P}_t(x'|x, a) - q_{t-1}^{\circ, k}(x, a) \bar{P}_t(x'|x, a) \right| \\
&\quad + \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k+1}} \left| q_{t-1}^{\circ, k}(x, a) \bar{P}_t(x'|x, a) - q_{t-1}^{\circ, k}(x, a) \bar{P}_{t-1}(x'|x, a) \right| \\
&\leq \frac{1}{2} \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \left| q_t^{\circ, k}(x, a) - q_{t-1}^{\circ, k}(x, a) \right| \sum_{x' \in \mathcal{X}_{k+1}} \bar{P}_t(x'|x, a) \\
&\quad + \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} q_{t-1}^{\circ, k}(x, a) \sum_{x' \in \mathcal{X}_{k+1}} \left| \bar{P}_t(x'|x, a) - \bar{P}_{t-1}(x'|x, a) \right| \\
&\stackrel{\text{a}}{\leq} \frac{1}{2} \sum_{x \in \mathcal{X}_k} \sum_{a \in \mathcal{A}} \left| q_t^{\circ, k}(x, a) - q_{t-1}^{\circ, k}(x, a) \right| \sum_{x' \in \mathcal{X}_{k+1}} \bar{P}_t(x'|x, a) + \frac{1}{t} \\
&\stackrel{\text{b}}{\leq} \frac{1}{2} \sum_{x \in \mathcal{X}_k} \left| q_t^{\circ, k}(x) - q_{t-1}^{\circ, k}(x) \right| + \frac{1}{t} \\
&\stackrel{\text{c}}{\leq} \frac{k}{t} + \frac{1}{t} = \frac{k+1}{t} \leq \frac{L}{t}
\end{aligned}$$

Where “a” comes from Lemma F.3 and “b” from Lemma F.8 and finally “c” from Lemma F.9. \square

F.3 KNOWN LEMMAS AND COROLLARIES

Lemma F.5 (Lemma D.3.3, Jin et al. 2023). *Denote the set of tuples $\mathcal{X}_k \times \mathcal{A} \times \mathcal{X}_{k+1}$ by W_k . For any transition functions \bar{P}_t, P_t , and any policy π ,*

$$\begin{aligned}
q^{\bar{P}_t, \pi}(x) - q^{P_t, \pi}(x) &= \sum_{k=0}^{k(x)-1} \sum_{(u, v, w) \in W_k} q^{\bar{P}_t, \pi}(u, v) (\bar{P}_t(w|u, v) - P_t(w|u, v)) q^{P_t, \pi}(x|w) \\
&= \sum_{k=0}^{k(x)-1} \sum_{(u, v, w) \in W_k} q^{P_t, \pi}(u, v) (\bar{P}_t(w|u, v) - P_t(w|u, v)) q^{\bar{P}_t, \pi}(x|w)
\end{aligned}$$

where $q^{P', \pi}(x|w)$ is the probability of visiting x under π executed in P .

According to Lemma F.5 we can estimate the occupancy measure difference caused by the error on the transition function at episode t with the following corollary.

Corollary F.6 (Corollary D.3.6, Jin et al. 2023). *For any episode t and any policy π we have:*

$$\left| q^{\bar{P}_t, \pi}(x) - q^{P_t, \pi}(x) \right| \leq C^{\bar{P}_t}, \quad \forall x \neq x_L, \quad \text{and} \quad \sum_{x \neq x_L} \left| q^{\bar{P}_t, \pi}(x) - q^{P_t, \pi}(x) \right| \leq LC^{\bar{P}_t}$$

Corollary F.7 (Corollary D.3.7, Jin et al. 2023). *For any policy sequence $\{\pi_t\}_{t=1}^T$, and loss functions $\{\ell_t\}_{t=1}^T$ such that $\ell_t : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$ for any $t \in \{1, \dots, T\}$ it holds that*

$$\sum_{t=1}^T \langle q^{\bar{P}_t, \pi_t} - q^{P_t, \pi_t}; \ell_t \rangle \leq LC^{\bar{P}}$$

Lemma F.8 (Lemma D.2, Rosenberg & Mansour 2019). *Let π be a policy and let \bar{P}_t, \bar{P}_{t-1} be transition functions such that $\|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\| \leq \nu, \forall x, a \in \mathcal{X} \times \mathcal{A}$ then the following equations hold,*

$$\begin{aligned}
&\sum_{x \in \mathcal{X}} \left| q^{\bar{P}_t, \pi}(x) - q^{\bar{P}_{t-1}, \pi}(x) \right| = \sum_{x \in \mathcal{X}} \sum_{a \in \mathcal{A}} \left| q^{\bar{P}_t, \pi}(x, a) - q^{\bar{P}_{t-1}, \pi}(x, a) \right| \\
&\sum_{x \in \mathcal{X}} \sum_{a \in \mathcal{A}} \left| q^{\bar{P}_t, \pi}(x, a) - q^{\bar{P}_{t-1}, \pi}(x, a) \right| \leq \sum_{x \in \mathcal{X}} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k(x)+1}} \left| q^{\bar{P}_t, \pi}(x, a, x') - q^{\bar{P}_{t-1}, \pi}(x, a, x') \right|
\end{aligned}$$

$$\sum_{x \in \mathcal{X}} \sum_{a \in \mathcal{A}} \sum_{x' \in \mathcal{X}_{k(x)+1}} \left| q^{\bar{P}_t, \pi}(x, a, x') - q^{\bar{P}_{t-1}, \pi}(x, a, x') \right| = \sum_{x \in \mathcal{X}} \sum_{a \in \mathcal{A}} \left| q^{\bar{P}_t, \pi}(x, a) - q^{\bar{P}_{t-1}, \pi}(x, a) \right| + L\nu$$

Lemma F.9 (Lemma E.2, Rosenberg & Mansour 2019). *Let π be a policy and let \bar{P}_t, \bar{P}_{t-1} be transition functions such that $\|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\| \leq \nu, \forall x, a \in \mathcal{X} \times \mathcal{A}$. Then $\forall k \in \llbracket 0, L-1 \rrbracket$*

$$\sum_{x_k \in \mathcal{X}_k} \left| q^{\bar{P}_t, \pi}(x_k) - q^{\bar{P}_{t-1}, \pi}(x_k) \right| \leq k\nu$$

Corollary F.10 (Corollary E.2, Rosenberg & Mansour 2019). *Let π be a policy and let \bar{P}_t, \bar{P}_{t-1} be transition functions such that $\|\bar{P}_t(\cdot|x, a) - \bar{P}_{t-1}(\cdot|x, a)\|_1 \leq \nu, \forall x, a \in \mathcal{X} \times \mathcal{A}$. Then $\forall k \in \llbracket 0, L-1 \rrbracket$*

$$\left\| q^{\bar{P}_t, \pi} - q^{\bar{P}_{t-1}, \pi} \right\|_1 \leq L^2\nu + L\nu = \mathcal{O}(L^2\nu)$$