

Policy Optimization as Online Learning with Mediator Feedback

Alberto Maria Metelli*, Matteo Papini*, Pierluca D’Oro, Marcello Restelli

Dipartimento di Elettronica, Informazione e Bioingegneria, Politecnico di Milano

Piazza Leonardo da Vinci, 32, 20133, Milano, Italy

{albertomaria.metelli, matteo.papini, marcello.restelli}@polimi.it, pierluca.doro@mail.polimi.it

Abstract

Policy Optimization (PO) is a widely used approach to address continuous control tasks. In this paper, we introduce the notion of *mediator feedback* that frames PO as an online learning problem over the policy space. The additional available information, compared to the standard bandit feedback, allows reusing samples generated by one policy to estimate the performance of other policies. Based on this observation, we propose an algorithm, *RANDOMIZED-exploration policy Optimization via Multiple Importance Sampling with Truncation* (RANDOMIST), for regret minimization in PO, that employs a randomized exploration strategy, differently from the existing optimistic approaches. When the policy space is finite, we show that under certain circumstances, it is possible to achieve constant regret, while always enjoying logarithmic regret. We also derive problem-dependent regret lower bounds. Then, we extend RANDOMIST to compact policy spaces. Finally, we provide numerical simulations on finite and compact policy spaces, in comparison with PO and bandit baselines.

1 Introduction

Policy Optimization (PO, Deisenroth, Neumann, and Peters 2013) is a family of Reinforcement Learning (RL, Sutton and Barto 2018) algorithms based on the explicit optimization of the policy parameters. It represents the most promising approach for learning large-scale continuous control tasks and has already achieved marvelous results in video games (e.g., Vinyals et al. 2019) and robotics (e.g., Peng et al. 2020). These achievements, however, rely on massive amounts of simulation rollouts. The efficient use of experience data is essential both to reduce computational costs and to make learning online from real interaction possible. This is still largely an open problem and calls for better theoretical understanding. Any online-learning agent must face the exploration-exploitation dilemma: whether to leverage on its current knowledge to maximize performance or consider new alternatives. Fortunately, the Multi-Armed Bandit (MAB) literature (Bubeck and Cesa-Bianchi 2012; Latimore and Szepesvári 2018) provides a theoretical framework for the problem of efficient exploration under *bandit*

feedback, i.e., observing the effects of the chosen actions. The dilemma is addressed by minimizing the cumulative *regret* of the online performance w.r.t. the optimal one. The most popular exploration strategies are based on the Optimism in the Face of Uncertainty (OFU, Lai and Robbins 1985), of which UCB1 (Auer, Cesa-Bianchi, and Fischer 2002) is the prototypical algorithm, and on Thompson Sampling (TS, Thompson 1933). Both suffer only sublinear regret (Auer, Cesa-Bianchi, and Fischer 2002; Agrawal and Goyal 2012; Kaufmann, Korda, and Munos 2012). TS typically performs better in practice (Chapelle and Li 2011), but it is only computationally efficient in artificial settings (Kveton et al. 2019b). More recent randomized algorithms such as PHE (Perturbed History Exploration) (Kveton et al. 2019a) are able to match the theoretical and practical advantages of TS without the computational burden, and with no assumptions on the payoff distribution.

The OFU principle has been applied to RL (Jaksch, Ortner, and Auer 2010) and recently also to PO (Chowdhury and Gopalan 2019; Efroni et al. 2020), at the level of action selection. These methods are promising but limited to finite actions. A different perspective is proposed by Papini et al. (2019), where the decision problem is not defined over the agent’s actions but over the policy parameters. This change of viewpoint allows exploiting the special structure of the PO problem: for each policy, a sequence of states and actions performed by the agent is collected, constituting, alongside the rewards, a vastly richer signal than the simple bandit feedback. In this paper, we call it *mediator feedback* since this extra information acts as a mediator variable between the policy parameters and the return. OPTIMIST (Papini et al. 2019) is an OFU algorithm that uses Multiple Importance Sampling (MIS, Veach and Guibas 1995) to exploit the mediator feedback, so that the results of one policy provide information on all the others. This allows, in principle, to optimize over an infinite policy space with only finite samples and no regularity assumptions on the underlying process. There are two important limitations in Papini et al. (2019). First, the advantages of the mediator feedback over the bandit feedback are not clear from a theoretical perspective since the regret of OPTIMIST is comparable with that of UCB1 with finite policy space. Second, the policy selection of OPTIMIST requires maximizing a non-convex and non-differentiable index. In the continuous setting, this

*Equal contribution.

is addressed via discretization, with clear scalability issues.

In this work, we provide two major advancements. From the theoretical side, we provide regret lower bounds for the policy optimization problem with finite policy space, and we show that OPTIMIST actually enjoys *constant* regret under the assumptions made in (Papini et al. 2019). In fact, mediator feedback is so special that, under strong-enough assumptions, a greedy algorithm enjoys the same guarantees. We also devise a PHE-inspired randomized algorithm, called RANDOMIST (RANDOMized-exploration policy Optimization via Multiple Importance Sampling with Truncation), with similar regret guarantees as OPTIMIST. From the practical side, this allows replacing the unfeasible index maximization of OPTIMIST with a sampling procedure. Although our regret guarantees apply to the finite setting only, we propose a heuristic version of RANDOMIST for continuous problems, using a Markov Chain Monte Carlo (MCMC, Owen 2013). We show the advantages of this algorithm over continuous OPTIMIST in terms of computational complexity and performance.

The structure of the paper is as follows. We start in Section 2 with the basic background. In Section 3, we formalize the concept of mediator feedback in PO and derive two regret lower bounds. We illustrate, in Section 4, a possible way to exploit mediator feedback, based on importance sampling. Section 5 is devoted to the discussion of deterministic algorithms, providing the improved regret guarantees for OPTIMIST. In Section 6, we present RANDOMIST with its regret guarantees and the heuristic extension to the continuous case. In Section 7, we compare empirically RANDOMIST with relevant baselines on both illustrative examples and continuous-control problems. In Section 8, we discuss relationships with similar approaches from the bandit and RL literature. We conclude in Section 9, discussing the obtained results and proposing future research directions. The proofs of all the results can be found in Appendix D.

2 Preliminaries

In this section, we introduce some notation, the background on Markov decision processes and policy optimization.

Mathematical Background Let $(\mathcal{X}, \mathcal{F})$ be a measurable space, we denote with $\mathcal{P}(\mathcal{X})$ the set of probability measures over \mathcal{X} . Let $P, Q \in \mathcal{P}(\mathcal{X})$ such that $P \ll Q$,¹ for any $\alpha \in [0, \infty]$ the α -Rényi divergence (Rényi 1961) is defined as:²

$$D_\alpha(P\|Q) = \frac{1}{\alpha - 1} \log \int_{\mathcal{X}} \left(\frac{dP}{dQ} \right)^\alpha dQ.$$

We denote with $d_\alpha(P\|Q) = \exp[D_\alpha(P\|Q)]$ the exponentiated Rényi divergence (Cortes, Mansour, and Mohri 2010).

Markov Decision Processes and Policy Optimization A discrete-time Markov Decision Process (MDP, Puterman

¹ P is absolutely continuous w.r.t. Q , i.e., for every measurable set $\mathcal{Y} \subseteq \mathcal{X}$ we have $Q(\mathcal{Y}) = 0 \Rightarrow P(\mathcal{Y}) = 0$.

²In the limit, for $\alpha \rightarrow 1$ we have $D_1(P\|Q) = D_{\text{KL}}(P\|Q)$ and for $\alpha \rightarrow \infty$ we have $D_\infty(P\|Q) = \text{ess sup}_{\mathcal{X}} \frac{dP}{dQ}$.

1994) is a 6-tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \mu)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, \mathcal{P} is the transition model that for each $(s, a) \in \mathcal{S} \times \mathcal{A}$ provides the probability distribution of the next state $\mathcal{P}(\cdot|s, a) \in \mathcal{P}(\mathcal{S})$, $\mathcal{R}(s, a) \in \mathbb{R}$ is the reward function, $\gamma \in [0, 1]$ is the discount factor, and $\mu \in \mathcal{P}(\mathcal{S})$ is the initial-state distribution. In Policy Optimization (PO, Peters and Schaal 2008), we model the agent’s behavior by means of a policy $\pi_\theta(\cdot|s) \in \mathcal{P}(\mathcal{A})$ belonging to a space of parametric policies $\Pi_\Theta = \{\pi_\theta : \theta \in \Theta\}$. The interaction between an agent and an MDP generates a sequence of state-action pairs, named *trajectory*: $\tau = (s_0, a_0, s_1, a_1, \dots, s_{H-1}, a_{H-1})$ where $s_0 \sim \mu$, for all $h \in \{0, \dots, H-1\}$ we have $a_h \sim \pi_\theta(\cdot|s_h)$, $s_{h+1} \sim \mathcal{P}(\cdot|s_h, a_h)$ and $H \in \mathbb{N}$ is the trajectory length. Each parameter $\theta \in \Theta$ determines a policy $\pi_\theta \in \Pi_\Theta$ which, in turn, induces a probability measure $p_\theta \in \mathcal{P}(\mathcal{T})$ over the trajectory space \mathcal{T} . To every trajectory $\tau \in \mathcal{T}$, we associate an index of performance $\mathcal{R}(\tau) = \sum_{h=0}^{H-1} \gamma^h \mathcal{R}(s_h, a_h)$, called *return*. Without loss of generality we assume that $\mathcal{R}(\tau) \in [0, 1]$. Thus, we can evaluate the performance of a policy $\pi_\theta \in \Pi_\Theta$ by means of its *expected return*: $J(\theta) = \mathbb{E}_{\tau \sim p_\theta} [\mathcal{R}(\tau)]$. The goal of the agent consists in finding an optimal parameter, i.e., any θ^* maximizing $J(\theta)$.³

3 Online Policy Optimization and Mediator Feedback

The online PO protocol works as follows. At each round $t \in [n]$, we evaluate a parameter vector $\theta_t \in \Theta$ by running policy π_{θ_t} , collecting one (or more) trajectory $\tau_t \in \mathcal{T}$ and observing the corresponding return $\mathcal{R}(\tau_t)$. Then, based on the history $\mathcal{H}_t = \{(\theta_i, \tau_i, \mathcal{R}(\tau_i))\}_{i=1}^t$, we update θ_t to get θ_{t+1} . From an *online learning* perspective, the goal of the agent consists in maximizing the sum of the expected returns over n rounds or, equivalently, minimizing the cumulative regret $R(n)$:

$$\max_{\theta_1, \dots, \theta_n \in \Theta} \sum_{t=1}^n J(\theta_t) \Leftrightarrow \min_{\theta_1, \dots, \theta_n \in \Theta} R(n) = \sum_{t=1}^n \Delta(\theta_t),$$

where $\Delta(\theta) = J^* - J(\theta)$ is the optimality gap of $\theta \in \Theta$ and $J^* = \sup_{\theta \in \Theta} J(\theta)$. Thus, whenever policy π_{θ_t} is executed the agent receives the trajectory-return pair $(\tau_t, \mathcal{R}(\tau_t))$, that we name *mediator feedback* (MF). The term “mediator” refers to the side information, the trajectory τ_t , that *mediates* between the parameter choice θ_t and the return $\mathcal{R}(\tau_t)$. By naïvely approaching PO as an online-learning problem over policy space, we would only consider *bandit feedback*, in which just the return $\mathcal{R}(\tau_t)$ is observable. In comparison, the MF allows to better exploit the *structure* underlying the PO problem (Figure 1).⁴ Indeed, while the return function \mathcal{R}

³To simplify the presentation, we frame our results for the usual *action-based* PO. Our findings directly extend to *parameter-based* exploration (Sehnke et al. 2008), in which policies are indirectly optimized by learning a hyperpolicy that outputs the policy parameters. Coherently with Papini et al. (2019), the empirical evaluation of Section 7 is carried out in the parameter-based framework.

⁴In this paper, we employ the wording “bandit feedback” with a different meaning compared to some provably efficient approaches to PO (e.g., Efroni et al. 2020). See also Section 8.

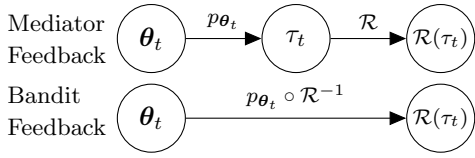


Figure 1: Graphical models comparing mediator and bandit feedbacks.

is unknown, the trajectory distribution p_θ is *partially* known:

$$p_\theta(\tau) = \mu(s_0) \prod_{h=0}^{H-1} \pi_\theta(a_h | s_h) P(s_{h+1} | s_h, a_h). \quad (1)$$

The policy factors π_θ , that depend on θ , are known to the agent, whereas the factors due to the environment (μ and P) are unknown but do not depend on θ . Intuitively, if two policies π_θ and $\pi_{\theta'}$ are sufficiently “similar”, given a trajectory τ from policy π_θ , the return $\mathcal{R}(\tau)$ provides information on the expected return of policy $\pi_{\theta'}$ too.

3.1 Regret Lower Bounds for Finite Policy Space

We focus on the intrinsic complexity of PO with finite policy space, deriving two lower bounds to the regret. The results are phrased, for simplicity, for the case of two policies, i.e., $|\Theta| = 2$, and the proof techniques are inspired to (Bubeck, Perchet, and Rigollet 2013). We start showing that, with enough structure between the policies, i.e., when the KL-divergence between the trajectory distributions is bounded, the best achievable regret is constant.

Theorem 3.1. *There exist an MDP and a parameter space $\Theta = \{\theta_1, \theta_2\}$ with $D_{\text{KL}}(p_{\theta_1} \| p_{\theta_2}) < \infty$, $D_{\text{KL}}(p_{\theta_2} \| p_{\theta_1}) < \infty$ and $J(\theta_1) - J(\theta_2) = \Delta$ such that, for sufficiently large n , all algorithms suffer regret $\mathbb{E} R(n) \geq \frac{1}{32\Delta}$.*

Instead, the presence of policies that are uninformative of one another, i.e., with infinite KL-divergence between the trajectory distributions, leads to a logarithmic regret.

Theorem 3.2. *There exist an MDP and a parameter space $\Theta = \{\theta_1, \theta_2\}$ with $D_{\text{KL}}(p_{\theta_1} \| p_{\theta_2}) = \infty$ or $D_{\text{KL}}(p_{\theta_2} \| p_{\theta_1}) = \infty$, and $J(\theta_1) - J(\theta_2) = \Delta$ such that, for any $n \geq 1$, all algorithms suffer regret $\mathbb{E} R(n) \geq \frac{1}{8\Delta} \log(\Delta^2 n)$.*

4 Exploiting Mediator Feedback with Importance Sampling

In this section, we illustrate how Importance Sampling techniques (IS, Cochran 1977; Owen 2013) can be employed to effectively exploit the mediator feedback in PO.⁵

Monte Carlo Estimation With the bandit feedback at each round $t \in [n]$, the agent has access to the history of parameter-return pairs $\mathcal{H}_t = \{(\theta_i, \mathcal{R}(\tau_i))_{i=1}^{t-1}$. Let $T_t(\theta) = \sum_{i=1}^{t-1} \mathbb{1}\{\theta_i = \theta\}$ be the number of trajectories collected with policy $\pi_\theta \in \Pi_\Theta$ up to round $t - 1$. To estimate the

⁵We stress that IS is just *one* method, and not necessarily the best one, to exploit the structure of the PO problem.

expected return $J(\theta)$, if no additional structure is available, we can only use the samples collected when executing π_θ , leading to the Monte Carlo (MC) estimator:

$$\hat{J}_t^{\text{MC}}(\theta) = \frac{1}{T_t(\theta)} \sum_{i=1}^{t-1} \mathcal{R}(\tau_i) \mathbb{1}\{\theta_i = \theta\}. \quad (2)$$

\hat{J}_t^{MC} is unbiased for $J(\theta)$ and its variance scales with $\text{Var}[\hat{J}_t^{\text{MC}}(\theta)] \leq 1/T_t(\theta)$. Clearly, $\hat{J}_t^{\text{MC}}(\theta)$ can be computed only for the policies that have been executed at least once.

Multiple Importance Sampling Estimation With the mediator feedback, at each round $t \in [n]$ we have access to additional information, i.e., the history of parameter-trajectory-return triples $\mathcal{H}_t = \{(\theta_i, \tau_i, \mathcal{R}(\tau_i))_{i=1}^{t-1}$. Thanks to the factorization in Equation (1), we can compute the trajectory distribution ratios without knowing P and μ :

$$\frac{p_\theta(\tau)}{p_{\theta'}(\tau)} = \prod_{h=0}^{H-1} \frac{\pi_\theta(a_h | s_h)}{\pi_{\theta'}(a_h | s_h)}.$$

Thus, we can use *all* the samples to estimate the expected return of *any* policy. Let $\Phi_t = \sum_{j=1}^{t-1} \frac{1}{t-1} p_{\theta_j}$ be the mixture induced by the policies executed up to time $t - 1$: if $p_\theta \ll \Phi_t$, we can employ a Multiple Importance Sampling (MIS, Veach and Guibas 1995) estimator (with balance heuristic):⁶

$$\hat{J}_t(\theta) = \frac{1}{t-1} \sum_{i=1}^{t-1} \omega_{\theta,t}(\tau_i) \mathcal{R}(\tau_i), \quad (3)$$

where $\omega_{\theta,t}(\tau_i) = p_\theta(\tau_i)/\Phi_t(\tau_i)$ is the *importance weight*. Thus, for estimating the expected return $J(\theta)$ of policy π_θ we do not need to execute π_θ , but just require the absolute continuity $p_\theta \ll \Phi_t$ (surely fulfilled if $T_t(\theta) \geq 1$). The statistical properties of the MIS estimator can be phrased in terms of the Rényi divergence. We can prove that $0 \leq \hat{J}_t(\theta) \leq d_\infty(p_\theta \| \Phi_t)$ and the variance can be bounded as $\text{Var}[\hat{J}_t(\theta)] \leq d_2(p_\theta \| \Phi_t)/(t-1)$ (Metelli et al. 2018; Papini et al. 2019; Metelli et al. 2020). Since the variance of $\hat{J}_t(\theta)$ scales with $d_2(p_\theta \| \Phi_t)/(t-1)$ instead of $1/T_t(\theta)$, as for $\hat{J}_t^{\text{MC}}(\theta)$, we refer to $\eta_t(\theta) := (t-1)/d_2(p_\theta \| \Phi_t)$ as the *effective number of trajectories*. It is worth noting that $\eta_t(\theta) \geq T_t(\theta)$ (Lemma C.4); thus, thanks to the structure introduced by the mediator feedback, the MIS estimator variance is always smaller than the MC estimator variance.⁷

Truncated Multiple Importance Sampling Estimation

The main limitation of the MIS estimator is that the importance weight $\omega_{\theta,t}$ displays a *heavy-tail* behavior, preventing exponential concentration, unless $d_\infty(p_\theta \| \Phi_t)$ is finite (Metelli et al. 2018). A common solution consists in

⁶For an extensive discussion of importance sampling and heuristics (e.g., balance heuristic) refer to (Owen 2013).

⁷The effective number of trajectories $\eta_t(\theta)$ is, in fact, the *effective sample size* of $\hat{J}_t(\theta)$ (Martino, Elvira, and Louzada 2017).

truncating the estimator (Ionides 2008) at the cost of introducing a negative bias. Given a (time-variant and policy-dependent) truncation threshold $M_t(\boldsymbol{\theta}) < \infty$, the Truncated MIS (TMIS) was introduced by Papini et al. (2019):

$$\check{J}_t(\boldsymbol{\theta}) = \frac{1}{t-1} \sum_{i=1}^{t-1} \check{\omega}_{\boldsymbol{\theta},t}(\tau_i) \mathcal{R}(\tau_i), \quad (4)$$

where $\check{\omega}_{\boldsymbol{\theta},t}(\tau_i) = \min\{M_t(\boldsymbol{\theta}), \omega_{\boldsymbol{\theta},t}(\tau_i)\}$. TMIS enjoys more desirable theoretical properties than plain MIS. While its variance scales similarly to $\hat{J}_t(\boldsymbol{\theta})$ since $\text{Var}[\check{J}_t(\boldsymbol{\theta})] \leq d_2(p_{\boldsymbol{\theta}}\|\Phi_t)/(t-1)$, the range can be bounded as $0 \leq \check{J}_t(\boldsymbol{\theta}) \leq M_t(\boldsymbol{\theta})$. Thus, the range is controlled by $M_t(\boldsymbol{\theta})$ and no longer by the divergence $d_{\infty}(p_{\boldsymbol{\theta}}\|\Phi_t)$, which may be infinite. Similarly, the bias can be bounded as $J(\boldsymbol{\theta}) - \mathbb{E}_{\tau_i \sim p_{\boldsymbol{\theta}_i}}[\check{J}_t(\boldsymbol{\theta})] \leq d_2(p_{\boldsymbol{\theta}}\|\Phi_t)/M_t(\boldsymbol{\theta})$ (see Papini et al. (2019) and Lemma C.1 for details). If we are interested in minimizing the joint contribution of bias and variance, this suggests to increase $M_t(\boldsymbol{\theta})$ progressively over the rounds.

5 Deterministic Algorithms

In this section, we consider finite policy spaces ($|\Theta| < \infty$) and discuss algorithms for PO that select policies deterministically, i.e., $\boldsymbol{\theta}_t$ is a deterministic function of history \mathcal{H}_{t-1} .

Follow The Leader The simplest algorithm accounting for the mediator feedback is Follow The Leader (FTL). It maintains a TMIS estimator $\check{J}_t(\boldsymbol{\theta})$ and selects the policy with the highest estimated expected return, i.e., $\boldsymbol{\theta}_t \in \arg \max_{\boldsymbol{\theta} \in \Theta} \check{J}_t(\boldsymbol{\theta})$. This is a pure-exploitation algorithm, unsuited for bandit feedback. Surprisingly, under a strong form of mediator feedback, FTL enjoys *constant* regret.

Theorem 5.1. *Let $\Theta = [K]$, $v(\boldsymbol{\theta}) = \max_{\boldsymbol{\theta}' \in \Theta} d_2(p_{\boldsymbol{\theta}}\|p_{\boldsymbol{\theta}'})$ for all $\boldsymbol{\theta} \in \Theta$ and $v^*(\boldsymbol{\theta}) = \max\{v(\boldsymbol{\theta}), v(\boldsymbol{\theta}^*)\}$, where $\pi_{\boldsymbol{\theta}^*}$ is an optimal policy. If $v := \max_{\boldsymbol{\theta} \in \Theta} v(\boldsymbol{\theta}) < \infty$, then, for any $\alpha > 1$, the expected regret of FTL using TMIS with truncation $M_t(\boldsymbol{\theta}) = \sqrt{\frac{td_2(p_{\boldsymbol{\theta}}\|\Phi_t)}{\alpha \log t}}$ is bounded as:*

$$\mathbb{E} R(n) \leq \sum_{\boldsymbol{\theta} \in \Theta: \Delta(\boldsymbol{\theta}) > 0} \frac{48\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})} \log \frac{24\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} + \Delta(\boldsymbol{\theta}_1) + \frac{2K}{\alpha-1} \min\left\{1, \sqrt{2 \log v}\right\}. \quad (5)$$

We refer to the condition when all pairwise Rényi divergences are finite (i.e., $v < \infty$) as *perfect mediator feedback*. In such case, we have the remarkable property that running *any* policy in Π_{Θ} provides information for *all* the others. Indeed, the effective number of trajectories satisfies $\eta_t(\boldsymbol{\theta}) \geq (t-1)/v$ (Lemma C.4). Unfortunately, when $v = \infty$, FTL degenerates to *linear* regret (Fact D.1).

UCB1 We can always apply an algorithm for standard bandit feedback, like UCB1 (Lai and Robbins 1985; Auer, Cesa-Bianchi, and Fischer 2002), to PO with finite policy space, ignoring the mediator feedback. UCB1 maintains the sample mean $\hat{J}_t^{\text{MC}}(\boldsymbol{\theta})$ of the observed returns for

Algorithm 1 OPTIMIST

Input: initial parameter $\boldsymbol{\theta}_1, \alpha > 1$

Execute $\pi_{\boldsymbol{\theta}_1}$, observe $\tau_1 \sim p_{\boldsymbol{\theta}_1}$ and $\mathcal{R}(\tau_1)$

for $t = 2, \dots, n$ **do**

 Compute expected return estimate $\check{J}_t(\boldsymbol{\theta})$

 Compute index:

$$B_t(\boldsymbol{\theta}) = \check{J}_t(\boldsymbol{\theta}) + (1 + \sqrt{2}) \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta})}}$$

 Select $\boldsymbol{\theta}_t \in \arg \max_{\boldsymbol{\theta} \in \Theta} B_t(\boldsymbol{\theta})$

 Execute $\pi_{\boldsymbol{\theta}_t}$, observe $\tau_t \sim p_{\boldsymbol{\theta}_t}$ and $\mathcal{R}(\tau_t)$

end for

each $\boldsymbol{\theta} \in \Theta$ and selects the one that maximizes $\hat{J}_t^{\text{MC}}(\boldsymbol{\theta}) + \sqrt{(\alpha \log t)/T_t(\boldsymbol{\theta})}$. The optimistic bonus favors policies that have been selected less often, in accordance with the OFU principle. Being designed for bandit feedback, UCB1 guarantees $\mathcal{O}(\Delta^{-1} \log n)$ regret (Auer, Cesa-Bianchi, and Fischer 2002) even if $v = \infty$, but it cannot exploit mediator feedback when actually present.

In principle, we could employ FTL or UCB1 based on whether v is finite or infinite. There are two reasons why this approach might be inappropriate. First, we would disregard the possibility to share information among pairs of policies with finite divergence, losing possible practical benefits (not captured by the current regret analysis). Second, even when $v < \infty$, the regret of FTL is $\mathcal{O}(v\Delta^{-1} \log(v\Delta^{-2}))$ that, at finite time, might be worse than $\mathcal{O}(\Delta^{-1} \log n)$, especially for large v . Note that deriving the conditions on v so that the regret of UCB1 is smaller than that of FTL is not practical since it would require the knowledge of the gap Δ .

OPTIMIST The difficulty in combining the advantages of FTL and UCB1 is overcome by OPTIMIST (Algorithm 1), an OFU-based algorithm introduced by Papini et al. (2019).⁸ It selects policies as to maximize an *optimistic* TMIS expected return estimate that favors policies with a lower effective number of trajectories. In the original paper (Papini et al. 2019), OPTIMIST is only shown to enjoy sublinear regret in high probability under perfect mediator feedback ($v < \infty$). We show here that OPTIMIST actually enjoys constant regret under perfect mediator feedback (like FTL) without ever degenerating into linear regret (like UCB1).

Theorem 5.2. *Let $\Theta = [K]$ and $v(\boldsymbol{\theta}) = \max_{\boldsymbol{\theta}' \in \Theta} d_2(p_{\boldsymbol{\theta}}\|p_{\boldsymbol{\theta}'})$ for all $\boldsymbol{\theta} \in \Theta$ ($v(\boldsymbol{\theta})$ can be infinite). For any $\alpha > 1$, the expected regret of OPTIMIST with truncation $M_t(\boldsymbol{\theta}) = \sqrt{\frac{td_2(p_{\boldsymbol{\theta}}\|\Phi_t)}{\alpha \log t}}$ is bounded as:*

(a) if $v := \max_{\boldsymbol{\theta} \in \Theta} v(\boldsymbol{\theta}) < \infty$:

$$\mathbb{E} R(n) \leq \sum_{\boldsymbol{\theta} \in \Theta: \Delta(\boldsymbol{\theta}) > 0} \frac{48\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})} \log \frac{24\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2}$$

⁸We consider here a slight variant of OPTIMIST with an explicit exploration parameter α in place of the original confidence parameter δ from (Papini et al. 2019), since we focus on expected regret rather than high-probability regret.

Algorithm 2 RANDOMIST

Input: initial parameter θ_1 , scale $a \geq 0$, translation $b \geq 0$, $\alpha > 1$

Execute π_{θ_1} , observe $\tau_1 \sim p_{\theta_1}$ and $\mathcal{R}(\tau_1)$

for $t = 2, \dots, n$ **do**

 Compute expected return estimate $\check{J}_t(\theta)$

 Generate perturbation:

$$U_t(\theta) = \frac{1}{\eta_t(\theta)} \sum_{l=1}^{a\eta_t(\theta)} \tau_l + b, \text{ with } \tau_l \sim \text{Ber}(1/2)$$

 Select $\theta_t \in \arg \max_{\theta \in \Theta} \check{J}_t(\theta) + U_t(\theta)$

 Execute π_{θ_t} , observe $\tau_t \sim p_{\theta_t}$ and $\mathcal{R}(\tau_t)$

end for

$$+ \Delta(\theta_1) + \frac{2K}{\alpha - 1} \min \left\{ 1, \sqrt{2 \log v} \right\};$$

(b) in any case:

$$\mathbb{E} R(n) \leq \sum_{\theta \in \Theta: \Delta(\theta) > 0} \frac{24\alpha}{\Delta(\theta)} \log n + \frac{\alpha + 1}{\alpha - 1} K,$$

with an instance-independent expected regret of $\mathbb{E} R(n) \leq 4\sqrt{6\alpha K n \log n} + (\alpha + 1)K/(\alpha - 1)$.

Note also that the regret correctly goes to zero with the divergence (when $v = 1$, all the policies are equivalent). It is an interesting open problem whether better regret guarantees can be provided for the intermediate case, i.e., when some (but not all) the Rényi divergences are finite.

6 Randomized Algorithms

In this section, we propose a novel algorithm for regret minimization in PO that selects the policies with a randomized strategy. RANDOMIST (RANDOMized-exploration policy Optimization via Multiple Importance Sampling with Truncation, Algorithm 2) is based on PHE (Kveton et al. 2019a) and employs additional samples to *perturb* the TMIS expected return estimate $\check{J}_t(\theta)$, enforcing exploration.⁹ Clearly, RANDOMIST shares the randomized nature of exploration with the Bayesian approaches for bandits (e.g., Thompson Sampling (Thompson 1933)) although no prior-posterior mechanism is explicitly implemented and no assumption (apart for boundedness) on the return distribution is needed. At each round $t = 2, \dots, n$, we update the TMIS expected return estimate for each policy $\check{J}_t(\theta)$ and we generate the perturbation $U_t(\theta)$ that is obtained through $a\eta_t(\theta)$ *pseudo-rewards* sampled from a Bernoulli distribution $\text{Ber}(1/2)$. Then, we play the policy maximizing the *perturbed estimated expected return*, i.e., the sum of the estimated expected return $\check{J}_t(\theta)$ and the perturbation $U_t(\theta)$. The two hyperparameters are the *perturbation scale* $a > 0$ and the *perturbation translation* $b > 0$. Informally, a and b are responsible for the amount of exploration: a governs the variance of the perturbation, while b (which is absent in PHE) accounts for the negative bias introduced by the TMIS estimator. We now present the properties of RANDOMIST with finite parameter space and propose an extension to deal with compact parameter spaces.

⁹In this sense, RANDOMIST, as well as PHE, resembles the Follow the Perturbed Leader (Hannan 1957) strategy.

Finite Parameter Space If the policy space is finite, we can show that RANDOMIST enjoys guarantees similar to those of OPTIMIST on the expected regret.

Theorem 6.1. *Let $\Theta = [K]$, $v(\theta) = \max_{x' \in \Theta} d_2(p_\theta \| p_{\theta'})$ for all $\theta \in \Theta$ ($v(\theta)$ can be infinite) and $v^*(\theta) = \max\{v(\theta), v(\theta^*)\}$ where π_{θ^*} is an optimal policy. For any $\alpha > 1$, the expected regret of RANDOMIST with truncation $M_t(\theta) = \sqrt{\frac{td_2(p_\theta \| \Phi_t)}{\alpha \log t}}$ is bounded as follows:*

(a) if $v := \max_{\theta \in \Theta} v(\theta) < \infty$, $b \leq \sqrt{(\alpha \log t)/\eta_t(\theta)}$ and $a \geq 0$:

$$\begin{aligned} \mathbb{E} R(n) \leq & \sum_{\theta \in \Theta: \Delta(\theta) > 0} \frac{(188 + 32a)\alpha v^*(\theta)}{\Delta(\theta)} \log \frac{(94 + 16a)\alpha v^*(\theta)}{\Delta(\theta)^2} \\ & + \Delta(\theta_1) + \frac{\alpha + 3}{\alpha - 1} \min \left\{ 1, \sqrt{2 \log v} \right\} K; \end{aligned}$$

(b) no matter the value of v , if $a > 8$ and $J(\theta) - \mathbb{E}[\check{J}_t(\theta)] \leq b \leq \sqrt{(\alpha \log t)/\eta_t(\theta)}$:

$$\mathbb{E} R(n) \leq \sum_{\theta \in \Theta: \Delta(\theta) > 0} \frac{(52 + 110a)c\alpha}{\Delta(\theta)} \log n + 2 \frac{\alpha + 1}{\alpha - 1} K,$$

where $c = 2 + \frac{e^2 \sqrt{a}}{\sqrt{2\pi}} \exp\left[\frac{16}{a-8}\right] \left(1 + \sqrt{\frac{\pi a}{a-8}}\right)$, with an instance-independent regret of $\mathbb{E} R(n) \leq 2\sqrt{(52 + 110a)c\alpha K n \log n} + 2 \frac{\alpha + 1}{\alpha - 1} K$.

Under perfect mediator feedback RANDOMIST enjoys constant regret, like OPTIMIST, although with a dependence on $v^*(\theta)$, which involves the divergence w.r.t. an optimal policy. Moreover, in such case, since exploration is not needed, we could even set $a = b = 0$ reducing RANDOMIST to FTL. Similarly to OPTIMIST, when we allow $v = \infty$, the regret becomes logarithmic and the hyperparameters a and b must be carefully set to enforce exploration.

Compact Parameter Space When the parameter space is a compact set, i.e., $\Theta = [-M, M]^d$, the $\arg \max$ in Algorithm 2 cannot be explicitly computed. However, the random variable $\theta \in \arg \max_{\theta' \in \Theta} \check{J}_t(\theta) + U_t(\theta)$ can be seen as sampled from the distribution for θ of being the parameter in Θ with the largest perturbed estimated expected return, whose p.d.f. is given by (D'Eramo et al. 2017):

$$\begin{aligned} \mathbf{g}_t^*(\theta) &= g\left(\check{J}_t(\theta) + U_t(\theta) = \sup_{\theta' \in \Theta} \check{J}_t(\theta') + U_t(\theta') \mid \mathcal{H}_{t-1}\right) \\ &= \int_{\mathbb{R}} \frac{g_\theta(y)}{G_\theta(y)} \prod_{\Theta} G_{\theta'}(y) d\theta' dy, \end{aligned} \quad (6)$$

where $\prod_{\Theta} G_\theta(y) d\theta = \exp\left(\int_{\Theta} \log G_\theta(y) d\theta\right)$ is the *product integral* (Davis and Chatfield 1970), g_θ and G_θ are the p.d.f. and the c.d.f. of the random variable $\check{J}_t(\theta) + U_t(\theta)$ conditioned to the history \mathcal{H}_{t-1} . The *computation* of \mathbf{g}_t^* (even up to a constant) is challenging as the product integral requires a numerical integration over the parameter space Θ . Provided that an approximation (up to a constant) \mathbf{g}_t^\dagger of \mathbf{g}_t^* is available, we can use a Monte Carlo Markov Chain method (Owen

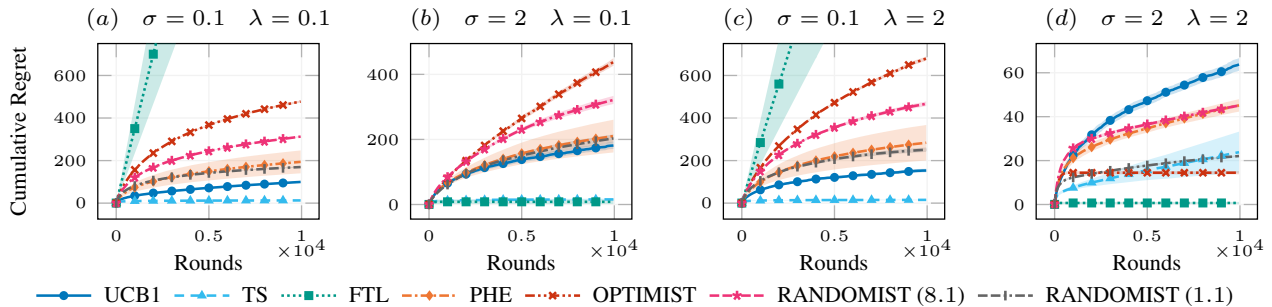


Figure 2: Cumulative regret on the illustrative PO for four values of σ and λ . 20 runs, 95% c.i.

2013) to generate a sample $\theta_t \sim \mathfrak{g}_t^\dagger$. As a practical approximation, we consider the p.d.f. for θ of having a perturbed estimated expected return larger than that of the previously executed policies:¹⁰ $\mathfrak{g}_t^\dagger(\theta) \propto \int_{\mathbb{R}} g_\theta(y) \prod_{i=1}^{t-1} G_{\theta_i}(y) dy$.

Since $\mathcal{O}(d)$ iterations of MCMC are sufficient to generate a sample (Beskos and Stuart 2009), where d is the dimensionality of Θ , and one evaluation of \mathfrak{g}_t^\dagger can be performed in time $\mathcal{O}(t^3)$, the per-round complexity of RANDOMIST is $\mathcal{O}(dt^3)$. This can be further reduced to $\mathcal{O}(dt^2)$ via clever caching (see Appendix F). OPTIMIST (Papini et al. 2019) can also be applied to continuous parameter spaces, with an $\tilde{\mathcal{O}}(\sqrt{vdn})$ high-probability regret bound. However, it is not clear how to perform the maximization step of OPTIMIST efficiently in this setting, since the optimistic index is non-differentiable and non-convex in the parameter variable. Discretization is adopted in (Papini et al. 2019), leading to $\mathcal{O}(t^{1+d/2})$ time complexity, that is exponential in d . The RANDOMIST variant proposed here, although heuristic, has only polynomial dependence on d , thus scaling more favorably to high-dimensional problems.

7 Numerical Simulations

We present the numerical simulations, starting with an illustrative example and then moving to RL benchmarks. For the RL experiments, similarly to Papini et al. (2019), the evaluation is carried out in the parameter-based PO setting (Sehnke et al. 2008), where the policy parameters θ are sampled from a *hyperpolicy* ν_ξ and the optimization is performed in the space of *hyperparameters* Ξ (Appendix A). This setting is particularly convenient since the Rényi divergence between hyperpolicies can be computed exactly (at least for Gaussians). Details and an additional experiment on the Cartpole domain are reported in Appendix F.

Illustrative Problems The goal of this experiment is to show the advantages of the additional structure offered by the mediator feedback over the bandit feedback. We design a class of 5-policy PO problems, isomorphic to bandit problems, in which trajectories are collapsed to a single real action $\mathcal{T} = \mathbb{R}$ and $\mathcal{R}(\tau) = \max\{0, \min\{1, \tau/4\}\}$.

¹⁰ \mathfrak{g}_t^\dagger can be seen as obtained from \mathfrak{g}_t^* applying a quadrature with $\{\theta_1, \dots, \theta_{t-1}\}$ as nodes for the inner integral.

The policies are Gaussians ($\mathcal{N}(0, \sigma^2)$, $\mathcal{N}(1, \sigma^2)$, $\mathcal{N}(2, \sigma^2)$, $\mathcal{N}(2.95, \lambda^2)$, $\mathcal{N}(3, \sigma^2)$) defined in terms of the two values $\sigma, \lambda > 0$. The optimal policy is the fifth one and we have a near-optimal parameter, the fourth, with a different variance. Intuitively, we can tune the parameters σ and λ to vary the Rényi divergences. We compare RANDOMIST with $a = 8.1$ (as prescribed in Theorem 6.1) and $a = 1.1$, and $b = \sqrt{(\alpha \log t)/\eta_t(\theta)}$ for both cases, with OPTIMIST (Papini et al. 2019), FTL, UCB1 (Auer, Cesa-Bianchi, and Fischer 2002), PHE (Kveton et al. 2019a), and TS with Gaussian prior (Agrawal and Goyal 2013a). The cumulative regret is shown in Figure 2 for four combinations of σ and λ . In (a) and (d) we are in a perfect mediator feedback, but in (a) $\log v \simeq 2.25$ and (d) $\log v \simeq 900$. Instead, in (b) or (c), we have $v = \infty$. We notice that FTL displays a (near-)linear regret in (a) as expected since $v = \infty$ but also in (c) where v is finite but very large. RANDOMIST with theoretical value of $a = 8.1$ always displays a good behavior and better than OPTIMIST, except in (d) where the latter shows a remarkable constant regret. We also note that when the amount of information shared among parameters is small, UCB1 performs better than OPTIMIST as well as PHE over RANDOMIST. Furthermore, TS with Gaussian prior performs very well across the tasks, although it considers the bandit feedback. This can be explained since TS assumes the correct return distribution. It also suggests that RANDOMIST could be improved when coped with other perturbation distributions (e.g., Gaussian). Finally, we observe that RANDOMIST with $a = 1.1$, although violating the conditions of Theorem 6.1, keeps showing a sublinear regret even in (b) and (c) when $v = \infty$.

Linear Quadratic Gaussian Regulator The Linear Quadratic Gaussian Regulator (LQG, Curtain 1997) is a benchmark for continuous control. We consider the monodimensional case and a Gaussian hyperpolicy $\nu_\xi = \mathcal{N}(\xi, 0.15^2)$ where ξ is the learned parameter. From ν_ξ , we sample the gain θ of a deterministic linear policy: $a_h = \theta s_h$. This experiment aims at comparing RANDOMIST with UCB1 (Auer, Cesa-Bianchi, and Fischer 2002), GPUCB (Srinivas et al. 2010), and OPTIMIST (Papini et al. 2019) in a finite policy space by discretizing $[-1, 1]$ in $K = 100$ parameters. In Figure 3, we notice that OPTIMIST and RANDOMIST outperform UCB1. While RAN-

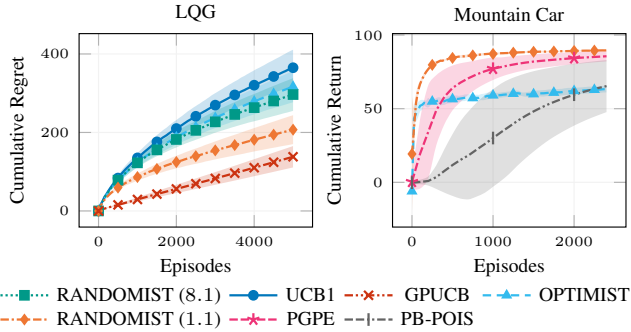


Figure 3: Cumulative regret in the LQG (30 runs, 95% c.i.) and cumulative return in the Mountain Car (5 runs, 95% c.i.).

DOMIST with $a = 8.1$ and OPTIMIST have similar performance, RANDOMIST improves significantly when setting a to 1.1. As in (Papini et al. 2019), the good performance of GPUCB is paired with a lack of theoretical guarantees due to the arbitrary choice of the GP kernel.

Mountain Car To test RANDOMIST in a continuous parameter space, we employ the approximation described above in the Mountain Car environment (Sutton and Barto 2018). We consider the setting of (Papini et al. 2019), employing PGPE (Sehnke et al. 2008) and PB-POIS (Metelli et al. 2018) as baselines. We use a Gaussian hyperpolicy $\nu_{\xi} = \mathcal{N}(\xi, \text{diag}(0.15, 3)^2)$ with learned mean ξ , from which we sample the parameters of a deterministic policy, linear in position and velocity. The exploration phase is performed by sampling from the approximate density g_t^{\dagger} , taking 10 steps of the Metropolis-Hastings algorithm (Owen 2013) with Gaussian proposal $q_m = \mathcal{N}(\theta_m, \text{diag}(0.15, 3)^2)$. Figure 3 shows that RANDOMIST outperforms both policy gradient baselines and OPTIMIST, in terms of learning speed and final performance.

8 Related Works

In this section, we revise the related literature, with attention to bandits with expert advice and to provably efficient PO. Additional comparisons are reported in Appendix B.

Mediator Feedback and Expert Advice A related formulation are the *Bandits with Expert Advice* (BEA, Bubeck and Cesa-Bianchi 2012, Section 4.2), introduced as an approach to adversarial contextual bandits. To draw a parallelism with PO, let \mathcal{T} be the set of arms and $\Theta = [K]$ the finite set of experts. At each step t , the agent receives *advice* $p_{\theta}^t \in \mathcal{P}(\mathcal{T})$ from each expert $\theta \in \Theta$, selects one expert θ_t , and pulls arm $\tau_t \sim p_{\theta_t}^t$. The goal is to minimize the *in-class* regret, competing with the best expert in hindsight. Differently from the trajectory distributions of PO, expert advice can change with time. A major concern of BEA, also relevant to PO, is the dependency of the regret on the number K of experts (resp. policies). A naïve application of Exp3 (Auer et al. 2002)

yields $\mathcal{O}(\sqrt{nK \log K})$ regret. Like our PO algorithms, this is impractical when the experts are exponentially many. Exp4 (Auer et al. 2002) achieves $\mathcal{O}(\sqrt{n|\mathcal{T}| \log K})$ regret, which scales well with K , but is vacuous in the case of infinite arms. McMahan and Streeter (2009) replace $|\mathcal{T}|$ with the *degree of agreement* of the experts, which has interesting similarities with our distributional-divergence approach. *Meta-bandit* approaches (Agarwal et al. 2017; Pacchiano et al. 2020) are so general that could be applied both to continuous-arm BEA and PO, but also exhibit a superlogarithmic dependence on K . Beygelzimer et al. (2011) obtain $\tilde{\mathcal{O}}(\sqrt{dn})$ regret competing with an infinite set of experts of VC-dimension d , mirrored in PO by OPTIMIST on compact spaces of dimension d (Papini et al. 2019, Theorem 3).

Provably Efficient PO Recently, a surge of approaches to deal with PO in a theoretically sound way, with both stochastic or adversarial environments, has emerged. These works consider either *full-information*, i.e., the agent observes the whole reward function $\{\mathcal{R}(s_h, a)\}_{a \in \mathcal{A}}$ regardless the played action (e.g., Rosenberg and Mansour 2019; Cai et al. 2019), or the *bandit feedback* (with a different meaning compared to the use we have made in this paper), in which only the reward of the chosen action is observed $\mathcal{R}(s_h, a_h)$ (e.g., Jin et al. 2019; Efroni et al. 2020). These methods are not directly comparable with the mediator feedback, although both settings exploit the structure of the PO problem. While with MF we explicitly model the policy space Π_{Θ} , these methods search in the space of all Markovian stationary policies. Furthermore, they are limited to tabular MDPs, while MF can deal natively with continuous state-action spaces.

9 Discussion and Conclusions

We have deepened the understanding of policy optimization as an online learning problem with additional feedback. We believe that mediator feedback has potential applications even beyond PO. Indeed, the problem of optimizing over probability distributions also encompasses GANs and variational inference (Chu, Blanchet, and Glynn 2019) and, more generally, MF emerges in any Bayesian network in which we control the conditional distributions on some vertices, via parameters θ , while the other are fixed and independent from θ . Furthermore, we have introduced a novel randomized algorithm, RANDOMIST, and we have shown its advantages both in terms of computational complexity and performance. The algorithm could be improved by adopting a different perturbation, e.g., Gaussian, as already hinted in (Kveton et al. 2019b). Further work is needed to match the theoretical regret lower bounds. Currently, a major discrepancy is the use of the KL-divergence in the lower bounds instead of the larger Rényi divergence required by algorithms based on IS. Moreover, the algorithm employs the ratio importance weight and, thus, it might suffer from the curse of horizon (Liu et al. 2018). Finally, the case of non-perfect mediator feedback could be related to graphical bandits (Alon et al. 2017), where finite Rényi divergences are the edges of a directed feedback graph, in order to capture the actual difficulty of this intermediate case.

Acknowledgments

This work has been partially supported by the Italian MIUR PRIN 2017 Project ALGADIMAR "Algorithms, Games, and Digital Markets".

References

- Abbasi-Yadkori, Y.; Pál, D.; and Szepesvári, C. 2011. Improved algorithms for linear stochastic bandits. In *NeurIPS*.
- Abeille, M.; and Lazaric, A. 2017a. Linear thompson sampling revisited. *Electronic Journal of Statistics* 11(2): 5165–5197.
- Abeille, M.; and Lazaric, A. 2017b. Thompson Sampling for Linear-Quadratic Control Problems. In *AISTATS*.
- Abeille, M.; and Lazaric, A. 2018. Improved regret bounds for thompson sampling in linear quadratic control problems. In *ICML*.
- Agarwal, A.; Luo, H.; Neyshabur, B.; and Schapire, R. E. 2017. Corraling a Band of Bandit Algorithms. In *COLT*.
- Agrawal, S.; and Goyal, N. 2012. Analysis of Thompson Sampling for the Multi-armed Bandit Problem. In *COLT*.
- Agrawal, S.; and Goyal, N. 2013a. Further Optimal Regret Bounds for Thompson Sampling. In *AISTATS*.
- Agrawal, S.; and Goyal, N. 2013b. Thompson sampling for contextual bandits with linear payoffs. In *ICML*.
- Alon, N.; Cesa-Bianchi, N.; Dekel, O.; and Koren, T. 2015. Online learning with feedback graphs: Beyond bandits. In *COLT*.
- Alon, N.; Cesa-Bianchi, N.; Gentile, C.; Mannor, S.; Mansour, Y.; and Shamir, O. 2017. Nonstochastic Multi-Armed Bandits with Graph-Structured Feedback. *SIAM J. Comput.* 46(6): 1785–1826.
- Auer, P.; Cesa-Bianchi, N.; and Fischer, P. 2002. Finite-time Analysis of the Multiarmed Bandit Problem. *Mach. Learn.* 47(2-3): 235–256.
- Auer, P.; Cesa-Bianchi, N.; Freund, Y.; and Schapire, R. E. 2002. The Nonstochastic Multiarmed Bandit Problem. *SIAM J. Comput.* 32(1): 48–77.
- Beskos, A.; and Stuart, A. 2009. Computational complexity of Metropolis-Hastings methods in high dimensions. In *Monte Carlo and Quasi-Monte Carlo Methods 2008*, 61–71.
- Beygelzimer, A.; Langford, J.; Li, L.; Reyzin, L.; and Schapire, R. E. 2011. Contextual Bandit Algorithms with Supervised Learning Guarantees. In *AISTATS*.
- Bubeck, S.; and Cesa-Bianchi, N. 2012. Regret Analysis of Stochastic and Nonstochastic Multi-armed Bandit Problems. *Foundations and Trends in Machine Learning* 5(1): 1–122.
- Bubeck, S.; Eldan, R.; and Lehec, J. 2018. Sampling from a log-concave distribution with projected langevin monte carlo. *Discrete & Computational Geometry* 59(4): 757–783.
- Bubeck, S.; Perchet, V.; and Rigollet, P. 2013. Bounded regret in stochastic multi-armed bandits. In *COLT*.
- Cai, Q.; Yang, Z.; Jin, C.; and Wang, Z. 2019. Provably Efficient Exploration in Policy Optimization. *arXiv preprint arXiv:1912.05830*.
- Caron, S.; Kveton, B.; Lelarge, M.; and Bhagat, S. 2012. Leveraging Side Observations in Stochastic Bandits. In *UAI*.
- Casella, G.; and George, E. I. 1992. Explaining the Gibbs sampler. *The American Statistician* 46(3): 167–174.
- Chapelle, O.; and Li, L. 2011. An Empirical Evaluation of Thompson Sampling. In *NeurIPS*.
- Chen, W.; Wang, Y.; Yuan, Y.; and Wang, Q. 2016. Combinatorial multi-armed bandit and its extension to probabilistically triggered arms. *JMLR* 17(1): 1746–1778.
- Cheng, X.; and Bartlett, P. 2018. Convergence of Langevin MCMC in KL-divergence. In *ALT*.
- Chowdhury, S. R.; and Gopalan, A. 2017. On kernelized multi-armed bandits. In *ICML*.
- Chowdhury, S. R.; and Gopalan, A. 2019. Online Learning in Kernelized Markov Decision Processes. In *AISTATS*.
- Chu, C.; Blanchet, J. H.; and Glynn, P. W. 2019. Probability Functional Descent: A Unifying Perspective on GANs, Variational Inference, and Reinforcement Learning. In Chaudhuri, K.; and Salakhutdinov, R., eds., *ICML*.
- Chung, F.; and Lu, L. 2006. Old and new concentration inequalities. *Complex Graphs and Networks* 107: 23–56.
- Cochran, W. G. 1977. *Sampling Techniques, 3rd Edition*. John Wiley. ISBN 0-471-16240-X.
- Combes, R.; Magureanu, S.; and Proutiere, A. 2017. Minimal exploration in structured stochastic bandits. In *NeurIPS*.
- Corless, R. M.; Gonnet, G. H.; Hare, D. E.; Jeffrey, D. J.; and Knuth, D. E. 1996. On the LambertW function. *Advances in Computational mathematics* 5(1): 329–359.
- Cortes, C.; Mansour, Y.; and Mohri, M. 2010. Learning Bounds for Importance Weighting. In *NeurIPS*.
- Curtain, R. F. 1997. Linear-quadratic control: An introduction. *Autom.* 33(5): 1004.
- Dani, V.; Hayes, T. P.; and Kakade, S. M. 2008. Stochastic Linear Optimization under Bandit Feedback. In *COLT*.
- Davis, W.; and Chatfield, J. 1970. Concerning product integrals and exponentials. *AMS*.
- Dean, S.; Mania, H.; Matni, N.; Recht, B.; and Tu, S. 2018. Regret bounds for robust adaptive control of the linear quadratic regulator. In *NeurIPS*.
- Deisenroth, M. P.; Neumann, G.; and Peters, J. 2013. A Survey on Policy Search for Robotics. *Foundations and Trends in Robotics* 2(1-2): 1–142.
- D’Eramo, C.; Nuara, A.; Pirota, M.; and Restelli, M. 2017. Estimating the maximum expected value in continuous reinforcement learning problems. In *AAAI*.
- Eckles, D.; and Kaptein, M. 2014. Thompson sampling with the online bootstrap. *arXiv preprint arXiv:1410.4009*.
- Efroni, Y.; Shani, L.; Rosenberg, A.; and Mannor, S. 2020. Optimistic Policy Optimization with Bandit Feedback. *arXiv preprint arXiv:2002.08243*.
- Gil, M.; Alajaji, F.; and Linder, T. 2013. Rényi divergence measures for commonly used univariate continuous distributions. *Information Sciences* 249: 124–131.
- Grant, J. A.; and Leslie, D. S. 2020. On Thompson Sampling for Smoother-than-Lipschitz Bandits. *arXiv preprint arXiv:2001.02323*.
- Hannan, J. 1957. Approximation to Bayes risk in repeated play. *Contributions to the Theory of Games* 3: 97–139.
- Ionides, E. L. 2008. Truncated importance sampling. *JCGS* 17(2): 295–311.
- Jaksch, T.; Ortner, R.; and Auer, P. 2010. Near-optimal Regret Bounds for Reinforcement Learning. *J. Mach. Learn. Res.* 11: 1563–1600.

- Jin, C.; Allen-Zhu, Z.; Bubeck, S.; and Jordan, M. I. 2018. Is q-learning provably efficient? In *NeurIPS*.
- Jin, C.; Jin, T.; Luo, H.; Sra, S.; and Yu, T. 2019. Learning adversarial markov decision processes with bandit feedback and unknown transition. *arXiv preprint arXiv:1912.01192*.
- Kallus, N. 2018. Instrument-Armed Bandits. In *Algorithmic Learning Theory*, 529–546.
- Kaufmann, E.; Korda, N.; and Munos, R. 2012. Thompson Sampling: An Asymptotically Optimal Finite-Time Analysis. In *ALT*.
- Kveton, B.; Szepesvári, C.; Ghavamzadeh, M.; and Boutilier, C. 2019a. Perturbed-History Exploration in Stochastic Multi-Armed Bandits. In *IJCAI*.
- Kveton, B.; Szepesvári, C.; Vaswani, S.; Wen, Z.; Lattimore, T.; and Ghavamzadeh, M. 2019b. Garbage In, Reward Out: Bootstrapping Exploration in Multi-Armed Bandits. In *ICML*.
- Lai, T. L.; and Robbins, H. 1985. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics* 6(1): 4–22.
- Lattimore, T.; and Munos, R. 2014. Bounded Regret for Finite-Armed Structured Bandits. In *NeurIPS*.
- Lattimore, T.; and Szepesvári, C. 2018. Bandit algorithms.
- Liu, Q.; Li, L.; Tang, Z.; and Zhou, D. 2018. Breaking the Curse of Horizon: Infinite-Horizon Off-Policy Estimation. In *NeurIPS*, 5361–5371.
- Lu, X.; and Van Roy, B. 2017. Ensemble sampling. In *NeurIPS*.
- Magureanu, S.; Combes, R.; and Proutiere, A. 2014. Lipschitz bandits: Regret lower bounds and optimal algorithms. In *COLT*.
- Martino, L.; Elvira, V.; and Louzada, F. 2017. Effective sample size for importance sampling based on discrepancy measures. *Signal Processing* 131: 386–401.
- McMahan, H. B.; and Streeter, M. J. 2009. Tighter Bounds for Multi-Armed Bandits with Expert Advice. In *COLT*.
- Metelli, A. M.; Papini, M.; Faccio, F.; and Restelli, M. 2018. Policy Optimization via Importance Sampling. In *NeurIPS*.
- Metelli, A. M.; Papini, M.; Montali, N.; and Restelli, M. 2020. Importance Sampling Techniques for Policy Optimization. *JMLR* 21(141): 1–75.
- Osband, I.; Russo, D.; and Van Roy, B. 2013. (More) efficient reinforcement learning via posterior sampling. In *NeurIPS*.
- Osband, I.; Russo, D.; Wen, Z.; and Van Roy, B. 2017. Deep exploration via randomized value functions. *JMLR*.
- Owen, A. B. 2013. Monte Carlo theory, methods and examples. *Monte Carlo Theory, Methods and Examples*.
- Pacchiano, A.; Phan, M.; Abbasi-Yadkori, Y.; Rao, A.; Zimmert, J.; Lattimore, T.; and Szepesvári, C. 2020. Model Selection in Contextual Stochastic Bandit Problems. *CoRR* abs/2003.01704.
- Papini, M.; Metelli, A. M.; Lupo, L.; and Restelli, M. 2019. Optimistic Policy Optimization via Multiple Importance Sampling. In *ICML*.
- Peng, X. B.; Coumans, E.; Zhang, T.; Lee, T.-W.; Tan, J.; and Levine, S. 2020. Learning Agile Robotic Locomotion Skills by Imitating Animals. *arXiv preprint arXiv:2004.00784*.
- Peters, J.; and Schaal, S. 2008. Reinforcement learning of motor skills with policy gradients. *Neural Networks* 21(4): 682–697.
- Phan, M.; Yadkori, Y. A.; and Domke, J. 2019. Thompson Sampling and Approximate Inference. In *NeurIPS*.
- Puterman, M. L. 1994. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley.
- Rényi, A. 1961. On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*.
- Rosenberg, A.; and Mansour, Y. 2019. Online Convex Optimization in Adversarial Markov Decision Processes. In *ICML*.
- Rusmevichientong, P.; and Tsitsiklis, J. N. 2010. Linearly parameterized bandits. *Mathematics of Operations Research* 35(2): 395–411.
- Russo, D. J.; Van Roy, B.; Kazerouni, A.; Osband, I.; and Wen, Z. 2018. A Tutorial on Thompson Sampling. *Foundations and Trends® in Machine Learning* 11(1): 1–96.
- Sehnke, F.; Osendorfer, C.; Rückstieß, T.; Graves, A.; Peters, J.; and Schmidhuber, J. 2008. Policy Gradients with Parameter-Based Exploration for Control. In *ICANN*.
- Srinivas, N.; Krause, A.; Kakade, S. M.; and Seeger, M. W. 2010. Gaussian Process Optimization in the Bandit Setting: No Regret and Experimental Design. In *ICML*.
- Strehl, A. L.; and Littman, M. L. 2008. An analysis of model-based interval estimation for Markov decision processes. *JCSS* 74(8): 1309–1331.
- Sutton, R. S.; and Barto, A. G. 2018. *Reinforcement learning: An introduction*. MIT press.
- Tange, O. 2011. GNU Parallel - The Command-Line Power Tool. *login: The USENIX Magazine* 36(1): 42–47.
- Thompson, W. R. 1933. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika* 25(3/4): 285–294.
- Veach, E.; and Guibas, L. J. 1995. Optimally combining sampling techniques for Monte Carlo rendering. In Mair, S. G.; and Cook, R., eds., *SIGGRAPH*, 419–428. ACM.
- Vinyals, O.; Babuschkin, I.; Czarnecki, W.; Mathieu, M.; Dudzik, A.; Chung, J.; Choi, D.; Powell, R.; Ewalds, T.; Georgiev, P.; Oh, J.; Horgan, D.; Kroiss, M.; Danihelka, I.; Huang, A.; Sifre, L.; Cai, T.; Agapiou, J.; Jaderberg, M.; and Silver, D. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature* 575(7782): 350–354.

Index of the Appendix

In the following, we briefly recap the contents of the Appendix.

- Appendix A provides the formulation of the policy optimization problem with mediator feedback in the parameter-based setting.
- Appendix B completes the review of the relevant literature begun in Section 8, focusing on the approaches that share connections with the mediator feedback and RANDOMIST.
- Appendix C reports some central lemmas that are employed for the analysis of the TMIS estimator for expected return estimation.
- Appendix D provides the proofs that are omitted in the main paper.
- Appendix E reports some auxiliary lemmas that are employed in the analysis.
- Appendix F presents the experimental setting in more detail, additional experimental results and a discussion about implementation issues.

A Parameter-based PO and Mediator Feedback

In *parameter-based* policy optimization (PB-PO Sehnke et al. 2008) the policy parameters θ are sampled from a higher level distribution, called *hyperpolicy*, and the learning process occurs in the space of hyperpolicy parameters, named *hyperparameters*. More formally, we consider a space of parametric hyperpolicies $\mathcal{N}_\Xi = \{\nu_\xi \in \mathcal{P}(\Theta) : \xi \in \Xi\}$ where $\Xi \subseteq \mathbb{R}^d$ is the hyperparameter space. To each hyperpolicy ν_ξ we can associate an index of performance:

$$J(\xi) = \mathbb{E}_{\theta \sim \nu_\xi} [J(\theta)] = \mathbb{E}_{\theta \sim \nu_\xi} \left[\mathbb{E}_{\tau \sim p_\theta} [\mathcal{R}(\tau)] \right].$$

The goal consists in finding an optimal hyperparameter, i.e., any ξ^* maximizing $J(\xi)$. At each round $t \in [n]$, we evaluate a hyperparameter $\xi_t \in \Xi$ by sampling one (or more) policy parameters θ_t , running policy π_{θ_t} , collecting one (or more) trajectory τ_t and observing the corresponding return $\mathcal{R}(\tau_t)$. Then, based on the history of observations $\mathcal{H}_t = \{(\xi_i, (\theta_i, \tau_i), \mathcal{R}(\tau_i))\}_{i=1}^t$, we update the hyperparameter ξ_t to get ξ_{t+1} . Differently from the *action-based* paradigm (AB-PO), in PB-PO deterministic policies are typically employed, since the stochasticity of the hyperpolicy is a sufficient source of exploration (Sehnke et al. 2008).

From an *online learning* perspective, the goal of an agent consists in maximizing the sum of the expected payoffs over n rounds or, equivalently, minimize the cumulative regret $R(n)$ w.r.t. to an optimal decision:

$$\max_{\xi_1, \dots, \xi_n \in \Xi} \sum_{t=1}^n J(\xi_t) \quad \Leftrightarrow \quad \min_{\xi_1, \dots, \xi_n \in \Xi} R(n) = \sum_{t=1}^n \Delta(\xi_t), \quad (7)$$

where $\Delta(\xi) = J^* - J(\xi)$ is the optimality gap of $\xi \in \Xi$ and $J^* = \sup_{\xi \in \Xi} J(\xi)$.

B Related Works

In this appendix, we revise the additional relevant literature, with particular attention to structured bandits, approximate Thompson sampling, and RL approaches with regret guarantees.

Structured Bandits and Feedback Although the formulation is quite different, the mediator feedback can be thought of as a way to endow a bandit with a particular structure. Numerous works have studied different kinds of structure (e.g., linear (Abbasi-Yadkori, Pál, and Szepesvári 2011; Dani, Hayes, and Kakade 2008; Rusmevichientong and Tsitsiklis 2010), Lipschitz (Magureanu, Combes, and Proutiere 2014)). Of particular interest is (Lattimore and Munos 2014), in which *general structures* are considered and constant problem-dependent regret results of order $\mathcal{O}(\Delta^{-1} \log \Delta^{-1})$ are derived for specific cases. Concerning the regret lower bounds, in (Bubeck, Perchet, and Rigollet 2013; Lattimore and Munos 2014; Combes, Magureanu, and Proutiere 2017) several results are shown for different classes of structured bandits. Extensions of the bandit feedback in which, when an arm is pulled, the outcome of other arms is revealed (possibly with some probability) are typically based on a graph structure (Alon et al. 2015; Chen et al. 2016). More specifically, in (Caron et al. 2012) the notion of *side-observation* is introduced to consider free extra information (passive or active) that allows achieving constant problem-dependent regret. This is quite similar to our mediator feedback, although we do not receive further fresh samples but we employ a single sample to update the estimates of multiple arms. More generally, we can see the mediator feedback (but also the side observations) as something in between the bandit feedback and the full information (expert) feedback. Another related concept is that of instrument-armed bandits (Kallus 2018), where the reward of a decision (instrument) depends on an intermediate, observed variable (e.g., compliance to medical prescription). Different definitions of regret can be adopted depending on whether one is interested in the causal relationship between the instrument and the reward. We adopt here the *Intent-to-Treat Regret* formulation, since we are only interested in finding a good instrument (policy or hyperpolicy). In (Kallus 2018), this case is treated as a regular bandit problem, ignoring the mediator feedback.

Approximate Thompson Sampling Thompson Sampling (TS, Thompson 1933; Russo et al. 2018) is an effective methodology for randomized exploration in multi-armed bandits. The main bottlenecks of TS are the computation of and the sampling from the posterior distribution. Several works focused on the effect of sampling from *approximate posteriors* with guarantees on the degradation of the (Bayesian) regret (Lu and Van Roy 2017; Phan, Yadkori, and Domke 2019). Other works addressed the sampling issue by employing Monte Carlo Markov Chain (MCMC, Casella and George 1992) approaches with Laplace approximation (Chapelle and Li 2011), Langevin Monte Carlo (Bubeck, Eldan, and Lehec 2018; Cheng and Bartlett 2018), and bootstrapping (Eckles and Kaptein 2014). Apart from the contextual bandits with linear payoff (Agrawal and Goyal 2013b; Abeille and Lazaric 2017a), the case of infinite arm set has been addressed with TS only in a very limited number of works, deriving guarantees on the Bayesian regret under strong regularity conditions (Grant and Leslie 2020) (without proposing a sampling routine) or by employing adaptive discretizations in the context of GPs (Chowdhury and Gopalan 2017).

Reinforcement Learning We have already surveyed the approaches to PO with regret guarantees in Section 8. Here, we focus on other provably efficient RL approaches. The majority of RL methods with theoretical guarantees has been developed in the context of tabular RL and are based on optimistic exploration (e.g., Jaksch, Ortner, and Auer 2010; Jin et al. 2018; Strehl and Littman 2008). These approaches, typically, do not extend directly to continuous tasks and/or to a randomized form of exploration. Recently, a number of approaches have been proposed to apply *posterior sampling* for solving MDPs, mainly with guarantees on the Bayesian regret (e.g., Osband, Russo, and Van Roy 2013; Osband et al. 2017), but also on the frequentist regret for some specific classes of continuous problems (e.g., Abeille and Lazaric 2017b, 2018; Dean et al. 2018).

C Key Lemmas on Off-Distribution Payoff Estimation

In this appendix, we revise the key lemmas needed when using MIS and TMIS for off-policy expected return estimation. We start in Appendix C.1 with a result for bounding bias and variance for general truncation threshold M and then we focus, in Appendix C.2, on the specific threshold $M_t(\theta)$ used in the algorithms.

C.1 Lemmas for General Truncation

In this appendix, we consider the importance weights defined for general probability distributions $P, Q \in \mathcal{P}(\mathcal{T})$ with $P \ll Q$:

$$\omega_{P/Q}(\tau) = \frac{dP}{dQ}(\tau), \quad \check{\omega}_{P/Q}(\tau) = \min \{M, \omega_{P/Q}(\tau)\},$$

where $M < \infty$ is the truncation threshold. We start with an ancillary result, that extends Lemma 2 by Papini et al. (2019) in bounding the α -moments and the bias of the truncated weight.

Lemma C.1. *Let P and Q be probability measures on the measurable space $(\mathcal{T}, \mathcal{F})$ with $P \ll Q$. Then, for any $\alpha \in (0, \infty]$, the α -moment of the truncated importance weight $\check{\omega}_{P/Q}$ with threshold M can be bounded for any $\beta \in [0, 1]$ as:*

$$\mathbb{E}_{\tau \sim Q} [\check{\omega}_{P/Q}(\tau)^\alpha]^\frac{1}{\alpha} \leq M^{1-\beta} d_{\alpha\beta}(P\|Q)^{\beta-\frac{1}{\alpha}}. \quad (8)$$

Furthermore, the bias of the truncated weight $\check{\omega}_{P/Q}$ can be bounded for any $\beta \in [1, \infty]$ as:

$$\mathbb{E}_{\tau \sim Q} [\omega_{P/Q}(\tau) - \check{\omega}_{P/Q}(\tau)] \leq \left(\frac{d_\beta(P\|Q)}{M} \right)^{\beta-1}. \quad (9)$$

Proof. Let us start with the first result. Consider the following derivation with $\beta \in [0, 1]$:

$$\begin{aligned} \mathbb{E}_{\tau \sim Q} [\check{\omega}_{P/Q}(\tau)^\alpha] &= \mathbb{E}_{\tau \sim Q} [\min \{M, \omega_{P/Q}(\tau)\}^\alpha] \\ &= \mathbb{E}_{\tau \sim Q} \left[\min \{M, \omega_{P/Q}(\tau)\}^{\alpha(1-\beta)} \min \{M, \omega_{P/Q}(\tau)\}^{\alpha\beta} \right] \\ &\leq M^{\alpha(1-\beta)} \mathbb{E}_{\tau \sim Q} [\omega_{P/Q}(\tau)^{\alpha\beta}] \end{aligned} \quad (10)$$

$$= M^{\alpha(1-\beta)} d_{\alpha\beta}(P\|Q)^{\alpha\beta-1}, \quad (11)$$

where line (10) is obtained by bounding the minimum and line (11) comes from the definition of exponentiated Rényi divergence. The result is obtained by taking the $\frac{1}{\alpha}$ -power both sides.

For the second result, we consider the following derivation for $\beta \in [1, \infty]$:

$$\begin{aligned} \mathbb{E}_{\tau \sim Q} [\omega_{P/Q}(\tau) - \check{\omega}_{P/Q}(\tau)] &= \mathbb{E}_{\tau \sim Q} [(\omega_{P/Q}(\tau) - M) \mathbb{1} \{\omega_{P/Q}(\tau) > M\}] \\ &\leq \mathbb{E}_{\tau \sim Q} [\omega_{P/Q}(\tau) \mathbb{1} \{\omega_{P/Q}(\tau) > M\}] \end{aligned}$$

$$\leq \mathbb{E}_{\tau \sim Q} [\omega_{P/Q}(\tau)^\beta]^{\frac{1}{\beta}} \mathbb{E}_{\tau \sim Q} \left[\mathbb{1} \{ \omega_{P/Q}(\tau) > M \}^{\frac{\beta}{\beta-1}} \right]^{\frac{\beta-1}{\beta}} \quad (12)$$

$$= d_\beta(P\|Q)^{\frac{\beta-1}{\beta}} \mathbb{P}_{\tau \sim Q} (\omega_{P/Q}(\tau) > M)^{\frac{\beta-1}{\beta}}, \quad (13)$$

where line (12) is an application of Hölder's inequality and line (13) comes from the definition of exponentiated Rényi divergence. Considering the probability we have for any $\gamma \in (0, \infty)$:

$$\begin{aligned} \mathbb{P}_{\tau \sim Q} (\omega_{P/Q}(\tau) > M) &= \mathbb{P}_{\tau \sim Q} (\omega_{P/Q}(\tau)^\gamma > M^\gamma) \\ &\leq \frac{\mathbb{E}_{\tau \sim Q} [\omega_{P/Q}(\tau)^\gamma]}{M^\gamma} \end{aligned} \quad (14)$$

$$= \frac{d_\gamma(P\|Q)^{\gamma-1}}{M^\gamma}, \quad (15)$$

where line (14) follows from Markov's inequality and line (15) from the definition of exponentiated Rényi divergence. By taking $\gamma = \beta$, we get the result. \square

It is worth noting that, while for bounding the α -moment of the non-truncated weight we need the α -Rényi divergence to be finite, for the truncated weight we can employ any $\alpha\beta$ -Rényi divergence, where $\alpha\beta \leq \alpha$.

C.2 Lemmas for $M_t(\theta)$ Truncation

From now on, let \check{J}_t be the TMIS estimator for the expected return, in the case of finite policy space, at time t :

$$\begin{aligned} \check{J}_t(\theta) &= \frac{1}{t-1} \sum_{h=1}^K \sum_{l=1}^{T_t(\theta_h)} \min \left\{ M_t(\theta), \frac{p_\theta(\tau_{hl})}{\sum_{k=1}^K \frac{T_t(\theta_k)}{t-1} p_{\theta_h}(\tau_{hl})} \right\} \mathcal{R}(\tau_{hl}) \\ &= \sum_{i=1}^{t-1} \min \left\{ M_t(\theta), \frac{p_\theta(\tau_i)}{\sum_{j=1}^{t-1} p_{\theta_j}(\tau_i)} \right\} \mathcal{R}(\tau_i), \end{aligned} \quad (16)$$

where $K = |\Theta|$, $T_t(\theta)$ is the number of executions of policy π_θ previous to time $t-1$, τ_{hl} denotes the l -th trajectory from policy π_{θ_h} , and τ_i denotes the i -th trajectory overall. Note that, with slight abuse of notation, the underscript of trajectories is over the policy space in the first expression and over time in the second expression. Also, let us fix the truncation threshold:

$$M_t(\theta) = \sqrt{\frac{(t-1)d_2(p_\theta\|\Phi_t)}{\alpha \log t}} = d_2(p_\theta\|\Phi_t) \sqrt{\frac{\eta_t(\theta)}{\alpha \log t}}, \quad (17)$$

for some $\alpha > 1$, where $\eta_t(\theta) = \frac{t-1}{d_2(p_\theta\|\Phi_t)}$ is the effective number of trajectories (or effective sample size) and $\Phi_t = \sum_{h=1}^K \frac{T_t(\theta_h)}{t-1} p_{\theta_h} = \frac{1}{t-1} \sum_{i=1}^{t-1} p_{\theta_i}$ is the mixture of the distributions of trajectories sampled previous to t . We sometimes abbreviate $M_t(\theta)$ and $\eta_t(\theta)$ as M and s , respectively, when parameter and time are clear from context. When not specified, expected values are w.r.t. past history \mathcal{H}_{t-1} . We always assume $\mathcal{R}(\tau) \in [0, 1]$. With little abuse of language, we will sometimes say "policy θ " to mean "policy π_θ ".

Lemma C.2. *The bias of \check{J}_t is bounded as follows:*

$$0 \leq J - \mathbb{E}[\check{J}_t(\theta)] \leq \sqrt{\frac{\alpha \log t}{\eta_t(\theta)}},$$

the variance:

$$\text{Var}[\check{J}_t(\theta)] \leq \frac{d_2(p_\theta\|\Phi_t)}{t-1},$$

and the estimator itself:

$$0 \leq \check{J}_t(\theta) \leq d_2(p_\theta\|\Phi_t) \sqrt{\frac{\eta_t(\theta)}{\alpha \log t}}.$$

Proof. The last property is evident from the chosen truncation threshold (17), and the first two can be easily deduced from (Papani et al. 2019), Lemma 2. \square

Lemma C.3. For all $\theta \in \Theta$, $t \geq 1$ and $\epsilon \geq 0$:

$$\mathbb{P}\left(\check{J}_t(\theta) - J(\theta) > \epsilon\right) \leq \exp\left[-\frac{\epsilon^2 \eta_t(\theta)}{2\left(1 + \frac{\epsilon}{3}\sqrt{\frac{\eta_t(\theta)}{\alpha \log t}}\right)}\right].$$

Moreover, if $\epsilon \geq \sqrt{\frac{\alpha \log t}{\eta_t(\theta)}}$:

$$\mathbb{P}\left(J(\theta) - \check{J}_t(\theta) > \epsilon\right) \leq \exp\left[-\frac{\eta_t(\theta)}{2}\left(\epsilon - \sqrt{\frac{\alpha \log t}{\eta_t(\theta)}}\right)^2\right].$$

Proof. For the first concentration inequality:

$$\mathbb{P}\left(\check{J}_t(\theta) - J(\theta) > \epsilon\right) \leq \mathbb{P}\left(\check{J}_t(\theta) - \mathbb{E}[\check{J}_t(\theta)] > \epsilon\right) \quad (18)$$

$$\leq \exp\left[\frac{-\epsilon^2(t-1)}{2\left(d_2(p_\theta \|\Phi_t) + \frac{\epsilon d_2(p_\theta \|\Phi_t)\sqrt{\frac{\eta_t(\theta)}{\alpha \log t}}}{3}\right)}\right] \quad (19)$$

$$= \exp\left[\frac{-\epsilon^2 \eta_t(\theta)}{2\left(1 + \frac{\epsilon}{3}\sqrt{\frac{\eta_t(\theta)}{\alpha \log t}}\right)}\right], \quad (20)$$

where we have used Lemma C.2 and (19) is from Theorem E.1. Similarly, for the second concentration inequality we still use Theorem E.1 together with Lemma C.2:

$$\begin{aligned} \mathbb{P}\left(J(\theta) - \check{J}_t(\theta) > \epsilon\right) &= \mathbb{P}\left(\mathbb{E}[\check{J}_t(\theta)] - \check{J}_t(\theta) > \epsilon + \mathbb{E}[\check{J}_t(\theta)] - J(\theta)\right) \\ &\leq \mathbb{P}\left(\mathbb{E}[\check{J}_t(\theta)] - \check{J}_t(\theta) > \epsilon - \sqrt{\frac{\alpha \log t}{\eta_t(\theta)}}\right) \\ &\leq \exp\left[-\frac{\left(\epsilon - \sqrt{\frac{\alpha \log t}{\eta_t(\theta)}}\right)^2 t}{2d_2(p_\theta \|\Phi_t)}\right] \\ &= \exp\left[-\frac{1}{2}\left(\epsilon - \sqrt{\frac{\alpha \log t}{\eta_t(\theta)}}\right)^2 \eta_t(\theta)\right]. \end{aligned}$$

□

Lemma C.4. The effective number of trajectories of a policy is always larger than the number of executions of that policy:

$$\eta_t(\theta) \geq T_t(\theta).$$

Moreover, if $v(\theta) = \sup_{x' \in \Theta} d_2(p_\theta \| p_{\theta'})$ is finite:

$$\eta_t(\theta) \geq \frac{t-1}{v(\theta)}.$$

Proof. The first inequality is trivial if $T_t(\theta) = 0$, so assume it is positive. From (Papini et al. 2019, Theorem 5) we know that $d_2(p_\theta \|\Phi_t)$ is bounded by the harmonic mean of pairwise divergences:

$$\begin{aligned} d_2(p_\theta \|\Phi_t) &\leq \frac{t-1}{\sum_{i=1}^{t-1} \frac{1}{d_2(p_\theta \| p_{\theta_i})}} \\ &= \frac{t-1}{T_t(\theta) + \sum_{i=1}^{t-1} \mathbb{1}\{\theta_i \neq \theta\} \frac{1}{d_2(p_\theta \| p_{\theta_i})}} \leq \frac{t-1}{T_t(\theta)}. \end{aligned}$$

Moreover, $d_2(p_\theta \|\Phi_t) \leq v(\theta)$ since the harmonic mean is never larger than the maximum. The claims follow by definition of $\eta_t(\theta)$. \square

Lemma C.5. *In PO, the optimality gap of policy θ is bounded as:*

$$\Delta(\theta) \leq \sqrt{2 \log v(\theta)}, \quad (21)$$

where $v(\theta) = \sup_{\theta' \in \Theta} d_2(p_\theta \| p_{\theta'})$.

Proof.

$$\Delta(\theta) = |J(\theta^*) - J(\theta)| \quad (22)$$

$$\begin{aligned} &= \left| \int_{\mathcal{T}} (p_{\theta^*}(z) - p_\theta(z)) \mathcal{R}(z) dz \right| \\ &\leq \int_{\mathcal{T}} |p_{\theta^*}(z) - p_\theta(z)| dz \\ &= 2D_{TV}(p_\theta, p_{\theta^*}) \\ &\leq \sqrt{2D_{KL}(p_\theta \| p_{\theta^*})} \quad (23) \end{aligned}$$

$$\leq \sqrt{2D_2(p_\theta \| p_{\theta^*})} \leq \sqrt{2 \log v(\theta)}, \quad (24)$$

where D_{TV} is the total variation distance, in (22) we use $\Delta(\theta) \geq 0$, (23) is from Pinsker's inequality, and (24) comes from the monotonicity of the Rényi divergence in the order. \square

D Proofs and Derivations

In this appendix, we report the proofs and the derivations we omitted in the main paper.

D.1 Proofs of Section 3

Theorem 3.1. *There exist an MDP and a parameter space $\Theta = \{\theta_1, \theta_2\}$ with $D_{KL}(p_{\theta_1} \| p_{\theta_2}) < \infty$, $D_{KL}(p_{\theta_2} \| p_{\theta_1}) < \infty$ and $J(\theta_1) - J(\theta_2) = \Delta$ such that, for sufficiently large n , all algorithms suffer regret $\mathbb{E} R(n) \geq \frac{1}{32\Delta}$.*

Proof. To prove the lower bound, we consider a pair of MDPs ν_1 and ν_2 with horizon 2 and $\mathcal{S} = \mathcal{A} = \mathbb{R}$. Thus, each trajectory is represented by the triple $\tau = (s, a, s')$. We take for the two MDPs the same reward function $\mathcal{R}(\tau) = s'$. The policy space is induced by $\Theta = \{\theta_1, \theta_2\}$. Let $\Delta \in [0, 1]$, for the first problem ν_1 we select the trajectory distributions as follows:¹¹

$$p_{\theta_1}^{\nu_1}(\tau) = \mu(s) \mathcal{N}(a|1, 1) \mathcal{N}(s'|a\Delta, 1), \quad p_{\theta_2}^{\nu_1}(\tau) = \mu(s) \mathcal{N}(a|0, 1) \mathcal{N}(s'|a\Delta, 1),$$

leading to the expected returns $J^{\nu_1}(\theta_1) = \Delta$ and $J^{\nu_1}(\theta_2) = 0$. Instead, for the second problem ν_2 we select:

$$p_{\theta_1}^{\nu_2}(\tau) = \mu(s) \mathcal{N}(a|1, 1) \mathcal{N}(s'| - a\Delta, 1), \quad p_{\theta_2}^{\nu_2}(\tau) = \mu(s) \mathcal{N}(a|0, 1) \mathcal{N}(s'| - a\Delta, 1),$$

leading to the expected returns $J^{\nu_2}(\theta_1) = -\Delta$ and $J^{\nu_2}(\theta_2) = 0$. For ν_1 the optimal decision is θ_1 , while for ν_2 the optimal policy is θ_2 and, for both, the gap is Δ . Furthermore, notice that:

$$D_{KL}(p_{\theta_1}^{\nu_1} \| p_{\theta_2}^{\nu_1}) = \int p_{\theta_1}^{\nu_1}(\tau) \log \frac{p_{\theta_1}^{\nu_1}(\tau)}{p_{\theta_2}^{\nu_1}(\tau)} d\tau = D_{KL}(\mathcal{N}(1, 1) \| \mathcal{N}(0, 1)) = \frac{1}{2}. \quad (25)$$

Similar derivations lead to $D_{KL}(p_{\theta_2}^{\nu_1} \| p_{\theta_1}^{\nu_1}) = D_{KL}(p_{\theta_1}^{\nu_2} \| p_{\theta_2}^{\nu_2}) = D_{KL}(p_{\theta_2}^{\nu_2} \| p_{\theta_1}^{\nu_2}) = \frac{1}{2}$.

Define a history of length t generated by the interaction of a policy with a problem ν as $\mathcal{H}_t^\nu = ((\theta_i, \tau_i, \mathcal{R}(\tau_i)))_{i=1}^t$. Given two problems we define $D_{KL}(\mathcal{H}_t^{\nu_1} \| \mathcal{H}_t^{\nu_2})$ as the KL-divergence between the distributions having generated the histories. Using standard derivations (Bubeck, Perchet, and Rigollet 2013) we have:

$$\begin{aligned} \max \left\{ \mathbb{E}_{\nu_1} R(n), \mathbb{E}_{\nu_2} R(n) \right\} &\geq \frac{1}{2} \left(\mathbb{E}_{\nu_1} R(n) + \mathbb{E}_{\nu_2} R(n) \right) \\ &= \frac{\Delta}{2} \sum_{t=1}^n \left(\mathbb{P}_{\nu_1}(\theta_t = \theta_2) + \mathbb{P}_{\nu_2}(\theta_t = \theta_1) \right) \\ &\geq \frac{\Delta}{4} \sum_{t=1}^n \exp[-D_{KL}(\mathcal{H}_t^{\nu_1} \| \mathcal{H}_t^{\nu_2})], \quad (26) \end{aligned}$$

¹¹The first factor is the initial-state distribution μ that is chosen equal for all the problems, the second factor is the policy π_θ , and the third factor is the transition model P .

where we denoted with \mathbb{E}_{ν_1} (resp. \mathbb{E}_{ν_2}) the expectation taken under the randomness of problem ν_1 (resp. ν_2) and we denoted with $\mathbb{P}_{\nu_1}(\boldsymbol{\theta}_t = \boldsymbol{\theta}_2)$ (resp. $\mathbb{P}_{\nu_2}(\boldsymbol{\theta}_t = \boldsymbol{\theta}_1)$) the probability of choosing decision $\boldsymbol{\theta}_2$ (resp. $\boldsymbol{\theta}_1$) at round t in the problem ν_1 (resp. ν_2). The last passage follows from Lemma 4 of (Bubeck, Perchet, and Rigollet 2013). Recalling that we have selected the same reward function for both problems, and again with standard derivations, we have:

$$D_{\text{KL}}(\mathcal{H}_t^{\nu_1} \|\mathcal{H}_t^{\nu_2}) = \mathbb{E}_{\nu_1}[T_t(\boldsymbol{\theta}_1)]D_{\text{KL}}(p_{\boldsymbol{\theta}_1}^{\nu_1} \| p_{\boldsymbol{\theta}_1}^{\nu_2}) + \mathbb{E}_{\nu_1}[T_t(\boldsymbol{\theta}_2)]D_{\text{KL}}(p_{\boldsymbol{\theta}_2}^{\nu_1} \| p_{\boldsymbol{\theta}_2}^{\nu_2}),$$

Let us now compute the divergences:

$$\begin{aligned} D_{\text{KL}}(p_{\boldsymbol{\theta}_1}^{\nu_1} \| p_{\boldsymbol{\theta}_1}^{\nu_2}) &= \int p_{\boldsymbol{\theta}_1}^{\nu_1}(\tau) \log \frac{p_{\boldsymbol{\theta}_1}^{\nu_1}(\tau)}{p_{\boldsymbol{\theta}_1}^{\nu_2}(\tau)} d\tau \\ &= \mathbb{E}_{a \sim \mathcal{N}(1,1)} [D_{\text{KL}}(\mathcal{N}(\Delta a, 1) \| \mathcal{N}(-\Delta a, 1))] \\ &= \mathbb{E}_{a \sim \mathcal{N}(1,1)} [2\Delta^2 a^2] = 4\Delta^2. \end{aligned}$$

In a similar way, we can derive $D_{\text{KL}}(p_{\boldsymbol{\theta}_2}^{\nu_1} \| p_{\boldsymbol{\theta}_2}^{\nu_2}) = \mathbb{E}_{a \sim \mathcal{N}(0,1)} [D_{\text{KL}}(\mathcal{N}(\Delta a, 1) \| \mathcal{N}(-\Delta a, 1))] = 2\Delta^2$. Thus, we have:

$$D_{\text{KL}}(\mathcal{H}_t^{\nu_1} \|\mathcal{H}_t^{\nu_2}) = 2\Delta^2 \left(2 \mathbb{E}_{\nu_1}[T_t(\boldsymbol{\theta}_1)] + \mathbb{E}_{\nu_1}[T_t(\boldsymbol{\theta}_2)] \right) \leq 4\Delta^2 t.$$

Plugging this result into Equation (26):

$$\max \left\{ \mathbb{E}_{\nu_1} R(n), \mathbb{E}_{\nu_2} R(n) \right\} \geq \frac{\Delta}{4} \sum_{t=1}^n \exp[-4\Delta^2 t] \geq \frac{1}{32\Delta},$$

where the last passage holds for sufficiently large n (Lattimore and Munos 2014). \square

Theorem 3.2. *There exist an MDP and a parameter space $\Theta = \{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2\}$ with $D_{\text{KL}}(p_{\boldsymbol{\theta}_1} \| p_{\boldsymbol{\theta}_2}) = \infty$ or $D_{\text{KL}}(p_{\boldsymbol{\theta}_2} \| p_{\boldsymbol{\theta}_1}) = \infty$, and $J(\boldsymbol{\theta}_1) - J(\boldsymbol{\theta}_2) = \Delta$ such that, for any $n \geq 1$, all algorithms suffer regret $\mathbb{E} R(n) \geq \frac{1}{8\Delta} \log(\Delta^2 n)$.*

Proof. The proofs follows the same steps of that of Theorem 3.1, but with a different construction of the trajectory distributions. We still consider a pair of MDPs ν_1 and ν_2 with horizon 2 defined over the policy space $\Theta = \{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2\}$ and having $\mathcal{S} = \mathcal{A} = \mathbb{R}$. We take for the two problems the same reward functions $\mathcal{R}(\tau) = a$. Let $\Delta \in [0, 1]$, for the first problem ν_1 we select the trajectory distributions as:

$$p_{\boldsymbol{\theta}_1}^{\nu_1}(\tau) = \mu(s)\mathcal{N}(a|1, 1)\mathcal{N}(s'|a\Delta, 1), \quad p_{\boldsymbol{\theta}_2}^{\nu_1}(\tau) = \mu(s)\delta_0(a)\mathcal{N}(s'|a\Delta, 1),$$

where δ_x is the Dirac measure centered in x , leading to the expected returns $J^{\nu_1}(\boldsymbol{\theta}_1) = \Delta$ and $J^{\nu_1}(\boldsymbol{\theta}_2) = 0$. Instead, for the second problem ν_2 we select:

$$p_{\boldsymbol{\theta}_1}^{\nu_2}(\tau) = \mu(s)\mathcal{N}(a|1, 1)\mathcal{N}(s'|-a\Delta, 1), \quad p_{\boldsymbol{\theta}_2}^{\nu_2}(\tau) = \mu(s)\delta_0(a)\mathcal{N}(s'|-a\Delta, 1),$$

leading to the expected returns $J^{\nu_2}(\boldsymbol{\theta}_1) = -\Delta$ and $J^{\nu_2}(\boldsymbol{\theta}_2) = 0$. For ν_1 the optimal policy is $\boldsymbol{\theta}_1$, while for ν_2 the optimal policy is $\boldsymbol{\theta}_2$ and, for both, the gap is Δ . Differently from the proof of Theorem 3.1, we considered Dirac deltas for the first factor of trajectory distribution of $\boldsymbol{\theta}_2$ instead of normal distributions. This leads to:

$$D_{\text{KL}}(p_{\boldsymbol{\theta}_1}^{\nu_1} \| p_{\boldsymbol{\theta}_2}^{\nu_1}) = \int p_{\boldsymbol{\theta}_1}^{\nu_1}(\tau) \log \frac{p_{\boldsymbol{\theta}_1}^{\nu_1}(\tau)}{p_{\boldsymbol{\theta}_2}^{\nu_1}(\tau)} d\tau = D_{\text{KL}}(\mathcal{N}(1, 1) \| \delta_0) = \infty. \quad (27)$$

Similar derivations lead to $D_{\text{KL}}(p_{\boldsymbol{\theta}_2}^{\nu_1} \| p_{\boldsymbol{\theta}_1}^{\nu_1}) = D_{\text{KL}}(p_{\boldsymbol{\theta}_1}^{\nu_2} \| p_{\boldsymbol{\theta}_2}^{\nu_2}) = D_{\text{KL}}(p_{\boldsymbol{\theta}_2}^{\nu_2} \| p_{\boldsymbol{\theta}_1}^{\nu_2}) = \infty$.

The analysis is now carried out w.r.t. the second problem ν_2 . First of all, we notice that:

$$\max \{R_{\nu_1}(n), R_{\nu_2}(n)\} \geq R_{\nu_2}(n) \geq \Delta \mathbb{E}_{\nu_2}[T_n(\boldsymbol{\theta}_1)].$$

Moreover, using standard derivations (Bubeck, Perchet, and Rigollet 2013) we have:

$$\begin{aligned} \max \left\{ \mathbb{E}_{\nu_1} R(n), \mathbb{E}_{\nu_2} R(n) \right\} &\geq \frac{1}{2} \left(\mathbb{E}_{\nu_1} R(n) + \mathbb{E}_{\nu_2} R(n) \right) \\ &= \frac{\Delta}{2} \sum_{t=1}^n \left(\mathbb{P}_{\nu_1}(\boldsymbol{\theta}_t = \boldsymbol{\theta}_2) + \mathbb{P}_{\nu_2}(\boldsymbol{\theta}_t = \boldsymbol{\theta}_1) \right) \\ &\geq \frac{\Delta}{4} \sum_{t=1}^n \exp[-D_{\text{KL}}(\mathcal{H}_t^{\nu_2} \|\mathcal{H}_t^{\nu_1})] \end{aligned} \quad (28)$$

$$\geq \frac{n\Delta}{4} \exp[-D_{\text{KL}}(\mathcal{H}_n^{\nu_2} \|\mathcal{H}_n^{\nu_1})], \quad (29)$$

where the only difference with the proof of Theorem 3.1 is that we switched the roles of ν_1 and ν_2 . Recalling that we have selected the same reward function for both problems, and again with standard derivations, we have:

$$D_{\text{KL}}(\mathcal{H}_t^{\nu_2} \|\mathcal{H}_t^{\nu_1}) = \mathbb{E}_{\nu_2}[T_t(\boldsymbol{\theta}_1)]D_{\text{KL}}(p_{\boldsymbol{\theta}_1}^{\nu_2} \| p_{\boldsymbol{\theta}_1}^{\nu_1}) + \mathbb{E}_{\nu_2}[T_t(\boldsymbol{\theta}_2)]D_{\text{KL}}(p_{\boldsymbol{\theta}_2}^{\nu_2} \| p_{\boldsymbol{\theta}_2}^{\nu_1}),$$

Let us now compute the divergences. For $\boldsymbol{\theta}_1$, $D_{\text{KL}}(p_{\boldsymbol{\theta}_1}^{\nu_2} \| p_{\boldsymbol{\theta}_1}^{\nu_1}) = 4\Delta^2$ as in Theorem 3.1. Instead, for the decision $\boldsymbol{\theta}_2$:

$$\begin{aligned} D_{\text{KL}}(p_{\boldsymbol{\theta}_2}^{\nu_2} \| p_{\boldsymbol{\theta}_2}^{\nu_1}) &= \int p_{\boldsymbol{\theta}_2}^{\nu_2}(\tau) \log \frac{p_{\boldsymbol{\theta}_2}^{\nu_2}(\tau)}{p_{\boldsymbol{\theta}_2}^{\nu_1}(\tau)} d\tau \\ &= \mathbb{E}_{a \sim \delta_0} [D_{\text{KL}}(\mathcal{N}(-a\Delta, 1) \|\mathcal{N}(a\Delta, 1))] \\ &= \mathbb{E}_{a \sim \delta_0} [2a^2\Delta^2] = 0. \end{aligned}$$

Thus, we have $D_{\text{KL}}(\mathcal{H}_t^{\nu_2} \|\mathcal{H}_t^{\nu_1}) = 4\Delta^2 \mathbb{E}_{\nu_2}[T_t(\boldsymbol{\theta}_2)]$. Plugging this result into Equation (29) and combining it with Equation 28:

$$\begin{aligned} \max_{\nu_1} \left\{ \mathbb{E}_{\nu_1} R(n), \mathbb{E}_{\nu_2} R(n) \right\} &\geq \max \left\{ \Delta \mathbb{E}_{\nu_2}[T_n(\boldsymbol{\theta}_1)], \frac{n\Delta}{4} \exp \left[-4\Delta^2 \mathbb{E}_{\nu_2}[T_n(\boldsymbol{\theta}_2)] \right] \right\} \\ &\geq \frac{\Delta}{2} \left(\mathbb{E}_{\nu_2}[T_n(\boldsymbol{\theta}_1)] + \frac{n}{4} \exp \left[-4\Delta^2 \mathbb{E}_{\nu_2}[T_n(\boldsymbol{\theta}_2)] \right] \right) \\ &\geq \frac{\Delta}{2} \min_{x \in [0, n]} \left\{ x + \frac{n}{4} \exp[-4\Delta^2 x] \right\} \\ &\geq \frac{1}{8\Delta} \log(\Delta^2 n), \end{aligned}$$

where the last line follows by solving the optimization problem over x , simply by zeroing the derivative. \square

D.2 Proofs of Section 5

Fact D.1. *There exist an MDP and a parameter space $\Theta = \{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2\}$ with $d_2(\boldsymbol{\theta}_1 \|\boldsymbol{\theta}_2) = \infty$ or $d_2(\boldsymbol{\theta}_2 \|\boldsymbol{\theta}_1) = \infty$ such that the expected regret of FTL is at least $\mathbb{E} R(n) \geq \frac{1}{16}(n-1)$.*

Proof. We consider a version of the FTL algorithm in which all policies are played once at the beginning. Let $\Theta = \{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2\}$ and $\mathcal{T} = \mathbb{R}$. We consider the following trajectory distributions:

$$p_{\boldsymbol{\theta}_1} = \text{Uni}([0, 1]), \quad p_{\boldsymbol{\theta}_2} = \delta_{1/4}, \quad (30)$$

where we denoted with Uni the uniform distribution. Finally, we select as reward function $\mathcal{R}(\tau) = \tau$. Clearly, the optimal policy is $\boldsymbol{\theta}_1$ having expected return $1/2$, while $\boldsymbol{\theta}_2$ has expected return $1/4$. Notice that $d_2(p_{\boldsymbol{\theta}_1} \| p_{\boldsymbol{\theta}_2}) = \infty$. Consider the bad event E in which, when pulled, $\boldsymbol{\theta}_1$ provides a reward $\mathcal{R}(\tau)$ that is smaller than $1/4$. This event has finite probability $\mathbb{P}(E) = \mathbb{P}(\tau < 1/4 | \tau \sim \text{Uni}([0, 1])) = 1/4$. After the initial play of the two policies, the estimates based on TMIS are just the on-policy ones. Indeed, on event E :

$$\begin{aligned} \check{J}_2(\boldsymbol{\theta}_1) &= \frac{1}{2} \left(\frac{p_{\boldsymbol{\theta}_1}(\tau)\tau}{\frac{1}{2}p_{\boldsymbol{\theta}_1}(\tau) + \frac{1}{2}p_{\boldsymbol{\theta}_2}(\tau)} + \frac{p_{\boldsymbol{\theta}_1}(1/4)1/4}{\frac{1}{2}p_{\boldsymbol{\theta}_1}(1/4) + \frac{1}{2}p_{\boldsymbol{\theta}_2}(1/4)} \right) = \tau \\ \check{J}_2(\boldsymbol{\theta}_2) &= \frac{1}{2} \left(\frac{p_{\boldsymbol{\theta}_2}(\tau)\tau}{\frac{1}{2}p_{\boldsymbol{\theta}_1}(\tau) + \frac{1}{2}p_{\boldsymbol{\theta}_2}(\tau)} + \frac{p_{\boldsymbol{\theta}_2}(1/4)1/4}{\frac{1}{2}p_{\boldsymbol{\theta}_1}(1/4) + \frac{1}{2}p_{\boldsymbol{\theta}_2}(1/4)} \right) = 1/4. \end{aligned}$$

Since $\check{J}_2(\boldsymbol{\theta}_1) < \check{J}_2(\boldsymbol{\theta}_2)$, FTL will play $\boldsymbol{\theta}_2$ at round 3. Moreover, since the samples from $\boldsymbol{\theta}_2$ do not change the estimate $\check{J}_2(\boldsymbol{\theta}_1)$, FTL will consistently play $\boldsymbol{\theta}_2$ suffering a regret of $\frac{1}{4}(n-1)$. Thus:

$$\mathbb{E} R(n) \geq \mathbb{E}[R(n)|E] \mathbb{P}(E) = \frac{1}{16}(n-1). \quad (31)$$

\square

We first prove Theorem 5.2 on OPTIMIST (Algorithm 1), then prove Theorem 5.1 on FTL as a variant. Before proceeding, a clarification on the initial executions performed by OPTIMIST is due.

Remark D.1. *We assume the expected reward estimators are initialized to an infinite value, i.e., $\check{J}_1(\boldsymbol{\theta}) \leftarrow +\infty$ for all $\boldsymbol{\theta} \in \Theta$, as is customary in OFU algorithms. Until there are infinite-valued estimates, one of the corresponding policies must necessarily*

be executed. We refer to this initial phase as Round-Robin regime. Note that the estimator of θ is updated (becomes finite) once a policy θ' at a finite Rényi divergence, i.e., $d_2(p_\theta \| p_{\theta'}) < \infty$, is executed (θ' can be θ itself). We distinguish two cases:

- (a) If all the pairwise Rényi divergences are finite (perfect mediator feedback), the initial policy θ_1 is executed, then all estimators are updated and the Round-Robin regime immediately ends.
- (b) If some Rényi divergences are infinite, let us call bad a policy such that $v(\theta) = \max_{\theta' \in \Theta} d_2(p_\theta \| p_{\theta'}) = \infty$. In this case, we first execute all the bad policies once in a Round-Robin fashion. After that, all expected return estimates must be finite. Notice that θ_1 need not be executed in this case unless it is itself a bad policy.

Hence, OPTIMIST only needs a *partial* initial Round-Robin, compared to the full Round-Robin of UCB1. Taking this into account, we now bound the expected regret:

Theorem 5.2. Let $\Theta = [K]$ and $v(\theta) = \max_{\theta' \in \Theta} d_2(p_\theta \| p_{\theta'})$ for all $\theta \in \Theta$ ($v(\theta)$ can be infinite). For any $\alpha > 1$, the expected regret of OPTIMIST with truncation $M_t(\theta) = \sqrt{\frac{td_2(p_\theta \| \Phi_t)}{\alpha \log t}}$ is bounded as:

(a) if $v := \max_{\theta \in \Theta} v(\theta) < \infty$:

$$\begin{aligned} \mathbb{E} R(n) &\leq \sum_{\theta \in \Theta: \Delta(\theta) > 0} \frac{48\alpha v(\theta)}{\Delta(\theta)} \log \frac{24\alpha v(\theta)}{\Delta(\theta)^2} \\ &\quad + \Delta(\theta_1) + \frac{2K}{\alpha - 1} \min \left\{ 1, \sqrt{2 \log v} \right\}; \end{aligned}$$

(b) in any case:

$$\mathbb{E} R(n) \leq \sum_{\theta \in \Theta: \Delta(\theta) > 0} \frac{24\alpha}{\Delta(\theta)} \log n + \frac{\alpha + 1}{\alpha - 1} K,$$

with an instance-independent expected regret of $\mathbb{E} R(n) \leq 4\sqrt{6\alpha K n \log n} + (\alpha + 1)K/(\alpha - 1)$.

Proof. We bound the expected number of executions $\mathbb{E} T_n(\theta)$ of policy $\theta \in \Theta$. The expected regret is then:

$$\mathbb{E} R(n) = \sum_{x \in \Theta: \Delta(x) > 0} \mathbb{E}[T_n(\theta)] \Delta(\theta). \quad (32)$$

Fix a policy $\theta \in \Theta$ and consider the following “good” events:

$$E_t = \left\{ \check{J}_t(\theta) \leq J(\theta) + (1 + \sqrt{2}) \sqrt{\frac{\alpha \log t}{\eta_t(\theta)}} \right\}, \quad F_t = \left\{ \check{J}_t(\theta^*) \geq J(\theta^*) - (1 + \sqrt{2}) \sqrt{\frac{\alpha \log t}{\eta_t(\theta^*)}} \right\}.$$

We will make sure that these events are well defined, i.e., $\eta_t(\theta) > 0$ always. By Lemma C.3:

$$\mathbb{P}(\overline{E}_t) = \mathbb{P} \left(\check{J}_t(\theta) - J(\theta) > (1 + \sqrt{2}) \sqrt{\frac{\alpha \log t}{\eta_t(\theta)}} \right) \quad (33)$$

$$\leq \exp \left[-\frac{3(8 + 5\sqrt{2})\alpha \log t}{28} \right] \leq t^{-\alpha}, \quad (34)$$

and also:

$$\mathbb{P}(\overline{F}_t) = \mathbb{P} \left(J(\theta^*) - \check{J}_t(\theta^*) > (1 + \sqrt{2}) \sqrt{\frac{\alpha \log t}{\eta_t(\theta^*)}} \right) \quad (35)$$

$$\leq \exp[-\alpha \log t] \leq t^{-\alpha}. \quad (36)$$

Under $E_t \cap F_t$:

$$J_t(\theta_t) + 2(1 + \sqrt{2}) \sqrt{\frac{\alpha \log t}{\eta_t(\theta_t)}} \geq \check{J}_t(\theta_t) + (1 + \sqrt{2}) \sqrt{\frac{\alpha \log t}{\eta_t(\theta_t)}} \quad (37)$$

$$\geq \check{J}_t(\theta^*) + (1 + \sqrt{2}) \sqrt{\frac{\alpha \log t}{\eta_t(\theta^*)}} \quad (38)$$

$$\geq J(\theta^*), \quad (39)$$

where (37) is from E_t , (38) is from the policy selection rule and (39) is from F_t . Rearranging:

$$\eta_t(\boldsymbol{\theta}_t) \leq \frac{4(1 + \sqrt{2})^2 \alpha \log t}{\Delta(\boldsymbol{\theta}_t)^2} \leq \frac{24\alpha \log t}{\Delta(\boldsymbol{\theta}_t)^2}. \quad (40)$$

Let $m = \max\{1, |\{\boldsymbol{\theta} \in \Theta \mid v(\boldsymbol{\theta}) = \infty\}|\}$. Hence:

$$\mathbb{E}[T_n(\boldsymbol{\theta})] = \mathbb{E} \left[\sum_{t=m+1}^n \mathbb{1}\{\boldsymbol{\theta}_t = \boldsymbol{\theta}, E_t \cap F_t\} + \sum_{t=m+1}^n \mathbb{1}\{\boldsymbol{\theta}_t = \boldsymbol{\theta}, \overline{E}_t \cup \overline{F}_t\} \right] + \mathbb{1}\{v(\boldsymbol{\theta}) = \infty \vee \boldsymbol{\theta} = \boldsymbol{\theta}_1\} \quad (41)$$

$$\leq \mathbb{E} \left[\sum_{t=m+1}^n \mathbb{1} \left\{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, \eta_t(\boldsymbol{\theta}) \leq \frac{24\alpha \log t}{\Delta(\boldsymbol{\theta})^2} \right\} \right] + \sum_{t=m+1}^n \mathbb{P}(\overline{E}_t \cup \overline{F}_t) + \mathbb{1}\{v(\boldsymbol{\theta}) = \infty \vee \boldsymbol{\theta} = \boldsymbol{\theta}_1\}, \quad (42)$$

where the third term is due to the initial round-robin (see Remark D.1). We first bound the second term:

$$\sum_{t=m+1}^n \mathbb{P}(\overline{E}_t \cup \overline{F}_t) \leq \sum_{t=m+1}^n [\mathbb{P}(\overline{E}_t) + \mathbb{P}(\overline{F}_t)] \quad (43)$$

$$\leq 2 \sum_{t=m+1}^n t^{-\alpha} \leq 2 \int_1^\infty t^{-\alpha} dt \leq \frac{2}{\alpha - 1}. \quad (44)$$

For the first term of (42), we bound it differently depending on whether $v(\boldsymbol{\theta})$ is finite or not.

If $v(\boldsymbol{\theta}) < \infty$, we use $\eta_t(\boldsymbol{\theta}) \geq \frac{t-1}{v(\boldsymbol{\theta})} > 0$ from Lemma C.4:

$$\sum_{t=m+1}^n \mathbb{1} \left\{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, \eta_t(\boldsymbol{\theta}) \leq \frac{24\alpha \log t}{\Delta(\boldsymbol{\theta})^2} \right\} \leq \sum_{t=m+1}^n \mathbb{1} \left\{ t \leq \frac{24\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log t + 1 \right\} \quad (45)$$

$$\leq \sum_{t=m+1}^n \mathbb{1} \left\{ t \leq \frac{48\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log \frac{24\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} + 1 \right\} \quad (46)$$

$$= \sum_{t=m}^{n-1} \mathbb{1} \left\{ t \leq \frac{48\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log \frac{24\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \right\} \quad (47)$$

$$\leq \sum_{t=1}^{n-1} \mathbb{1} \left\{ t \leq \frac{48\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log \frac{24\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \right\} \quad (48)$$

$$\leq \frac{48\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log \frac{24\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2}, \quad (49)$$

where (46) is from Lemma E.1.

Even if $v(\boldsymbol{\theta}) = \infty$, we can still use $\eta_t(\boldsymbol{\theta}) \geq T_t(\boldsymbol{\theta})$ from Lemma C.4 (in this case, $\eta_t(\boldsymbol{\theta}) > 0$ is guaranteed by the initial Round-Robin execution):

$$\sum_{t=m+1}^n \mathbb{1} \left\{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, \eta_t(\boldsymbol{\theta}) \leq \frac{24\alpha \log t}{\Delta(\boldsymbol{\theta})^2} \right\} \leq \sum_{t=1}^n \mathbb{1} \left\{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, T_t(\boldsymbol{\theta}) \leq \frac{24\alpha \log t}{\Delta(\boldsymbol{\theta})^2} \right\} \quad (50)$$

$$\leq \sum_{t=1}^n \mathbb{1} \left\{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, T_t(\boldsymbol{\theta}) \leq \frac{24\alpha \log n}{\Delta(\boldsymbol{\theta})^2} \right\} \quad (51)$$

$$\leq \frac{24\alpha}{\Delta(\boldsymbol{\theta})^2} \log n. \quad (52)$$

Statement (a) is obtained by using (49) for all policies. From (32):

$$\mathbb{E} R(n) \leq \sum_{x \in \Theta: \Delta(x) > 0} \frac{48\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})} \log \frac{24\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} + \Delta(\boldsymbol{\theta}_1) + \frac{2}{\alpha - 1} \sum_{x \in \Theta: \Delta(x) > 0} \Delta(\boldsymbol{\theta}) \quad (53)$$

$$\leq \sum_{x \in \Theta: \Delta(x) > 0} \frac{48\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})} \log \frac{24\alpha v(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} + \Delta(\boldsymbol{\theta}_1) + \frac{2}{\alpha - 1} \min \left\{ 1, \sqrt{2 \log v} \right\} K, \quad (54)$$

where the last inequality is by combining Lemma C.5 with the trivial $\Delta(\boldsymbol{\theta}) \leq 1$. Note that the $\mathbb{1}\{v(\boldsymbol{\theta}) = \infty \vee \boldsymbol{\theta} = \boldsymbol{\theta}_1\}$ terms from (42) amount to the unavoidable $\Delta(\boldsymbol{\theta}_1)$ in this case.

Similarly, (b) is obtained by using (52) for all policies. The $\mathbb{1}\{v(\boldsymbol{\theta}) = \infty \vee \boldsymbol{\theta} = \boldsymbol{\theta}_1\}$ terms from (42) amount to an additional K regret in the worst case.

The instance-independent regret is obtained from (b) by a standard reduction (see, e.g., Theorem 3 from (Kveton et al. 2019a)). \square

We now prove the regret bound for Follow The Leader (FTL), reported in Algorithm 3 for completeness.

Input: initial policy parameters $\boldsymbol{\theta}_1, \alpha > 1$
 Execute $\pi_{\boldsymbol{\theta}_1}$, observe $\tau_1 \sim p_{\boldsymbol{\theta}_1}$ and $\mathcal{R}(\tau_1)$
for $t = 2, \dots, n$ **do**
 Compute expected return estimate $\check{J}_t(\boldsymbol{\theta})$
 Select $\boldsymbol{\theta}_t \in \arg \max_{\boldsymbol{\theta} \in \Theta} \check{J}_t(\boldsymbol{\theta})$
 Execute $\pi_{\boldsymbol{\theta}_t}$, observe $\tau_t \sim p_{\boldsymbol{\theta}_t}$ and $\mathcal{R}(\tau_t)$
end for

Algorithm 3: Follow The Leader (FTL)

Theorem 5.1. *Let $\Theta = [K]$, $v(\boldsymbol{\theta}) = \max_{\boldsymbol{\theta}' \in \Theta} d_2(p_{\boldsymbol{\theta}} \| p_{\boldsymbol{\theta}'})$ for all $\boldsymbol{\theta} \in \Theta$ and $v^*(\boldsymbol{\theta}) = \max\{v(\boldsymbol{\theta}), v(\boldsymbol{\theta}^*)\}$, where $\pi_{\boldsymbol{\theta}^*}$ is an optimal policy. If $v := \max_{\boldsymbol{\theta} \in \Theta} v(\boldsymbol{\theta}) < \infty$, then, for any $\alpha > 1$, the expected regret of FTL using TMIS with truncation $M_t(\boldsymbol{\theta}) = \sqrt{\frac{td_2(p_{\boldsymbol{\theta}} \| \Phi_t)}{\alpha \log t}}$ is bounded as:*

$$\begin{aligned} \mathbb{E} R(n) &\leq \sum_{\boldsymbol{\theta} \in \Theta: \Delta(\boldsymbol{\theta}) > 0} \frac{48\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})} \log \frac{24\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \\ &\quad + \Delta(\boldsymbol{\theta}_1) + \frac{2K}{\alpha - 1} \min \left\{ 1, \sqrt{2 \log v} \right\}. \end{aligned} \quad (5)$$

Proof. The proof is similar to that of Theorem 5.2. We replace the argument in (39) with the following. Under $E_t \cap F_t$:

$$\Delta(\boldsymbol{\theta}_t) = J(\boldsymbol{\theta}^*) - J(\boldsymbol{\theta}_t) \quad (55)$$

$$= J(\boldsymbol{\theta}^*) - \check{J}(\boldsymbol{\theta}^*) + \check{J}(\boldsymbol{\theta}^*) - J(\boldsymbol{\theta}_t) \quad (56)$$

$$\leq J(\boldsymbol{\theta}^*) - \check{J}(\boldsymbol{\theta}^*) + \check{J}(\boldsymbol{\theta}_t) - J(\boldsymbol{\theta}_t) \quad (57)$$

$$\leq (1 + \sqrt{2}) \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta}^*)}} + (1 + \sqrt{2}) \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta}_t)}} \quad (58)$$

$$\leq (1 + \sqrt{2}) \sqrt{\frac{\alpha v(\boldsymbol{\theta}^*) \log t}{t - 1}} + (1 + \sqrt{2}) \sqrt{\frac{\alpha v(\boldsymbol{\theta}_t) \log t}{t - 1}}, \quad (59)$$

$$\leq 2(1 + \sqrt{2}) \sqrt{\frac{\alpha v^*(\boldsymbol{\theta}_t) \log t}{t - 1}}, \quad (60)$$

where (57) is by the policy selection rule, (58) is from $E_t \cap F_t$, and (59) is from Lemma C.4. Rearranging:

$$t \leq 4(1 + \sqrt{2})^2 \frac{\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log t + 1 \leq 24 \frac{\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log t + 1. \quad (61)$$

Hence:

$$\mathbb{E}[T_n(\boldsymbol{\theta})] = \mathbb{E} \left[\sum_{t=2}^n \mathbb{1}\{\boldsymbol{\theta}_t = \boldsymbol{\theta}, E_t \cap F_t\} + \sum_{t=2}^n \mathbb{1}\{\boldsymbol{\theta}_t = \boldsymbol{\theta}, \overline{E_t} \cup \overline{F_t}\} \right] + \mathbb{1}\{\boldsymbol{\theta} = \boldsymbol{\theta}_1\} \quad (62)$$

$$\leq \mathbb{E} \left[\sum_{t=2}^n \mathbb{1} \left\{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, t \leq 24 \frac{\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log t + 1 \right\} \right] + \sum_{t=2}^n \mathbb{P}(\overline{E_t} \cup \overline{F_t}) + \mathbb{1}\{\boldsymbol{\theta} = \boldsymbol{\theta}_1\}. \quad (63)$$

We bound the second term as in (44) from the proof of Theorem 5.2. For the first term:

$$\sum_{t=2}^n \mathbb{1} \left\{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, t \leq 24 \frac{\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log t + 1 \right\} = \sum_{t=2}^n \mathbb{1} \left\{ t \leq 24 \frac{\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log t + 1 \right\} \quad (64)$$

$$\leq \sum_{t=2}^n \mathbb{1} \left\{ t \leq 48 \frac{\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log 24 \frac{\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} + 1 \right\} \quad (65)$$

$$\leq \sum_{t=1}^{n-1} \mathbb{1} \left\{ t \leq 48 \frac{\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log 24 \frac{\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \right\} \quad (66)$$

$$\leq 48 \frac{\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \log 24 \frac{\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2}, \quad (67)$$

where (65) is from Lemma E.1. We then proceed as for the proof of statement (a) from Theorem 5.2. \square

D.3 Proofs of Section 6

In order to prove the results on the regret of RANDOMIST, we adopt an approach analogous to that of (Kveton et al. 2019b,a).

General Randomized Exploration with Shared History We start analyzing a more general algorithm that we call *General Randomized Exploration with Shared History* (GRE-SH, Algorithm 4). GRE-SH is the adaptation of the *General Randomized Exploration* (Algorithm 1 of (Kveton et al. 2019b)) to the mediator feedback setting.

Input: initial policy parameters $\boldsymbol{\theta}_1$
 Execute $\pi_{\boldsymbol{\theta}_1}$, observe $\tau_1 \sim p_{\boldsymbol{\theta}_1}$ and $\mathcal{R}(\tau_1)$
 Initialize $\mathcal{H}_1 = \{(\boldsymbol{\theta}_1, \tau_1, \mathcal{R}(\tau_1))\}$
for $t = 2, \dots, n$ **do**
 Draw $\theta_t(\boldsymbol{\theta}) \sim q_{\boldsymbol{\theta}}(\mathcal{H}_{t-1})$
 Select $\boldsymbol{\theta}_t \in \arg \max_{\boldsymbol{\theta} \in \Theta} \theta_t(\boldsymbol{\theta})$
 Execute $\pi_{\boldsymbol{\theta}_t}$, observe $\tau_t \sim p_{\boldsymbol{\theta}_t}$ and $y_t = \mathcal{R}(\tau_t)$
 Update $\mathcal{H}_t = \mathcal{H}_{t-1} \cup \{(\boldsymbol{\theta}_t, \tau_t, \mathcal{R}(\tau_t))\}$
end for

Algorithm 4: GRE-SH

For the sake of the analysis, let us define for any $t \in [n]$ and $\zeta \in \mathbb{R}$:

$$Q_t(\boldsymbol{\theta}, \zeta) = \mathbb{P}(\eta_t(\boldsymbol{\theta}) \geq \zeta | \eta_t(\boldsymbol{\theta}) \sim q_{\boldsymbol{\theta}}(\mathcal{H}_{t-1}), \mathcal{H}_{t-1}). \quad (68)$$

W.l.o.g. we will assume that the optimal policy $\boldsymbol{\theta}^*$ is unique. We can now provide the following result for Algorithm 4.

Theorem D.1. *For any tunable parameters $(\zeta(\boldsymbol{\theta}))_{\boldsymbol{\theta} \in \Theta \setminus \{\boldsymbol{\theta}^*\}} \in \mathbb{R}^{K-1}$ and $\alpha > 1$, the expected n -round regret of Algorithm 4 can be bounded from above as:*

$$\mathbb{E} R(n) = \sum_{\boldsymbol{\theta} \in \Theta \setminus \{\boldsymbol{\theta}^*\}} \Delta(\boldsymbol{\theta}) \mathbb{E}[T_n(\boldsymbol{\theta})] \leq \Delta(\boldsymbol{\theta}_1) + \sum_{\boldsymbol{\theta} \in \Theta \setminus \{\boldsymbol{\theta}^*\}} \Delta(\boldsymbol{\theta})(a(\boldsymbol{\theta}) + b(\boldsymbol{\theta}))$$

where:

$$a(\boldsymbol{\theta}) = \mathbb{E} \left[\sum_{t=2}^n \min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}^*, \zeta(\boldsymbol{\theta}))} - 1 \right) \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^* | \mathcal{H}_{t-1}), 1 \right\} \right]$$

$$b(\boldsymbol{\theta}) = \mathbb{E} \left[\sum_{t=2}^n \mathbb{1} \{ Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta})) > t^{-\alpha} \} \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta} | \mathcal{H}_{t-1}) \right] + \frac{1}{\alpha - 1}.$$

Proof. We extend the proof of Theorem 1 of (Kveton et al. 2019b). Our goal is to bound the expected number of execution for each suboptimal policy $\boldsymbol{\theta} \in \Theta \setminus \{\boldsymbol{\theta}^*\}$. Let us fix $\boldsymbol{\theta}$ and consider the event:

$$E_t(\boldsymbol{\theta}) = \{\eta_t(\boldsymbol{\theta}) \leq \zeta(\boldsymbol{\theta})\}. \quad (69)$$

We proceed to the decomposition:

$$\begin{aligned} \mathbb{E}[T_n(\boldsymbol{\theta})] &= \mathbb{E} \left[\sum_{t=1}^n \mathbb{1} \{ \boldsymbol{\theta}_t = \boldsymbol{\theta} \} \right] \\ &= \mathbb{E} \left[\sum_{t=2}^n \mathbb{1} \{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, E_t(\boldsymbol{\theta}) \} \right] + \mathbb{E} \left[\sum_{t=2}^n \mathbb{1} \{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, \overline{E_t(\boldsymbol{\theta})} \} \right] + \mathbb{1} \{ \boldsymbol{\theta} = \boldsymbol{\theta}_1 \}. \end{aligned}$$

To derive the expression of the term $b(\boldsymbol{\theta})$, let us consider the set of rounds $\mathcal{T} = \{t = 2, \dots, n : Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta})) > t^{-\alpha}\}$. We have:

$$\begin{aligned} \mathbb{E} \left[\sum_{t=2}^n \mathbb{1} \{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, \overline{E_t(\boldsymbol{\theta})} \} \right] &= \mathbb{E} \left[\sum_{t \in \mathcal{T}} \mathbb{1} \{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, \overline{E_t(\boldsymbol{\theta})} \} \right] + \mathbb{E} \left[\sum_{t \notin \mathcal{T}} \mathbb{1} \{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, \overline{E_t(\boldsymbol{\theta})} \} \right] \\ &\leq \mathbb{E} \left[\sum_{t \in \mathcal{T}} \mathbb{1} \{ \boldsymbol{\theta}_t = \boldsymbol{\theta} \} \right] + \mathbb{E} \left[\sum_{t \notin \mathcal{T}} \mathbb{1} \{ \overline{E_t(\boldsymbol{\theta})} \} \right] \\ &= \mathbb{E} \left[\sum_{t=2}^n \mathbb{1} \{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta})) > t^{-\alpha} \} \right] + \mathbb{E} \left[\sum_{t \notin \mathcal{T}} t^{-\alpha} \right] \\ &\leq \mathbb{E} \left[\sum_{t=2}^n \mathbb{1} \{ Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta})) > t^{-\alpha} \} \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta} | \mathcal{H}_{t-1}) \right] + \frac{1}{\alpha - 1}. \end{aligned}$$

where we note $\mathbb{E} \left[\mathbb{1} \{ \overline{E_t(\boldsymbol{\theta})} \} \right] = \mathbb{E} [\mathbb{E} [\mathbb{1} \{ E_t(\boldsymbol{\theta}) \} | \mathcal{H}_{t-1}]] = \mathbb{E} \left[\mathbb{P}(\overline{E_t(\boldsymbol{\theta})} | \mathcal{H}_{t-1}) \right] = \mathbb{E} [Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta}))]$ and that $\mathbb{E} [\mathbb{1} \{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta})) > t^{-\alpha} \}] = \mathbb{E} [\mathbb{1} \{ Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta})) > t^{-\alpha} \} \mathbb{E} [\mathbb{1} \{ \boldsymbol{\theta}_t = \boldsymbol{\theta} \} | \mathcal{H}_{t-1}]] = \mathbb{E} [\mathbb{1} \{ Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta})) > t^{-\alpha} \} \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta} | \mathcal{H}_{t-1})]$. Finally, we bounded the summation with the integral: $\sum_{t \notin \mathcal{T}} t^{-\alpha} \leq \sum_{t=2}^{\infty} t^{-\alpha} \leq \int_{x=1}^{\infty} x^{-\alpha} dx = \frac{1}{\alpha - 1}$ when $\alpha > 1$.

To derive the term $a(\boldsymbol{\theta})$, we need an auxiliary lemma, similar to Lemma 1 of (Agrawal and Goyal 2013a).

Lemma D.1. *For all $t \in [n]$ and for all $\boldsymbol{\theta} \in \Theta \setminus \{\boldsymbol{\theta}^*\}$ it holds that:*

$$\mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}, E_t(\boldsymbol{\theta}) | \mathcal{H}_{t-1}) \leq \left(\frac{1}{Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta}))} - 1 \right) \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^*, E_t(\boldsymbol{\theta}) | \mathcal{H}_{t-1}). \quad (70)$$

Proof. Let us consider the derivation:

$$\begin{aligned} \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta} | E_t(\boldsymbol{\theta}), \mathcal{H}_{t-1}) &= \mathbb{P}(\forall \boldsymbol{\theta}' \neq \boldsymbol{\theta} : \eta_t(\boldsymbol{\theta}) \geq \eta_t(\boldsymbol{\theta}') | E_t(\boldsymbol{\theta}), \mathcal{H}_{t-1}) \\ &\leq \mathbb{P}(\forall \boldsymbol{\theta}' \in \Theta : \eta_t(\boldsymbol{\theta}') \leq \zeta(\boldsymbol{\theta}) | E_t(\boldsymbol{\theta}), \mathcal{H}_{t-1}) \\ &= \mathbb{P}(\eta_t(\boldsymbol{\theta}^*) \leq \zeta(\boldsymbol{\theta}) | \mathcal{H}_{t-1}) \mathbb{P}(\forall \boldsymbol{\theta}' \neq \boldsymbol{\theta}^* : \eta_t(\boldsymbol{\theta}') \leq \zeta(\boldsymbol{\theta}) | E_t(\boldsymbol{\theta}), \mathcal{H}_{t-1}) \\ &= (1 - Q_t(\boldsymbol{\theta}^*, \zeta(\boldsymbol{\theta}))) \mathbb{P}(\forall \boldsymbol{\theta}' \neq \boldsymbol{\theta}^* : \eta_t(\boldsymbol{\theta}') \leq \zeta(\boldsymbol{\theta}) | E_t(\boldsymbol{\theta}), \mathcal{H}_{t-1}), \end{aligned}$$

where we exploited the fact that given \mathcal{H}_{t-1} the events $E_t(\boldsymbol{\theta})$ and $\eta_t(\boldsymbol{\theta}^*) \leq \zeta(\boldsymbol{\theta})$ are independent. Moreover, we have:

$$\begin{aligned} \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^* | E_t(\boldsymbol{\theta}), \mathcal{H}_{t-1}) &= \mathbb{P}(\forall \boldsymbol{\theta}' \neq \boldsymbol{\theta}^* : \eta_t(\boldsymbol{\theta}^*) \geq \eta_t(\boldsymbol{\theta}') | E_t(\boldsymbol{\theta}), \mathcal{H}_{t-1}) \\ &\geq \mathbb{P}(\forall \boldsymbol{\theta}' \neq \boldsymbol{\theta}^* : \eta_t(\boldsymbol{\theta}^*) > \zeta(\boldsymbol{\theta}) \geq \eta_t(\boldsymbol{\theta}') | E_t(\boldsymbol{\theta}), \mathcal{H}_{t-1}) \\ &= \mathbb{P}(\eta_t(\boldsymbol{\theta}^*) > \zeta(\boldsymbol{\theta}) | \mathcal{H}_{t-1}) \mathbb{P}(\forall \boldsymbol{\theta}' \neq \boldsymbol{\theta}^* : \eta_t(\boldsymbol{\theta}') \leq \zeta(\boldsymbol{\theta}) | E_t(\boldsymbol{\theta}), \mathcal{H}_{t-1}) \\ &= Q_t(\boldsymbol{\theta}^*, \zeta(\boldsymbol{\theta})) \mathbb{P}(\forall \boldsymbol{\theta}' \neq \boldsymbol{\theta}^* : \eta_t(\boldsymbol{\theta}') \leq \zeta(\boldsymbol{\theta}) | E_t(\boldsymbol{\theta}), \mathcal{H}_{t-1}). \end{aligned}$$

Putting together these two inequalities and using the rule of the conditional probability, we get the result. \square

Using Lemma D.1, we have:

$$\begin{aligned} \mathbb{E} \left[\sum_{t=2}^n \mathbb{1} \{ \boldsymbol{\theta}_t = \boldsymbol{\theta}, E_t(\boldsymbol{\theta}) \} \right] &= \mathbb{E} \left[\sum_{t=m+1}^n \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}, E_t(\boldsymbol{\theta}) | \mathcal{H}_{t-1}) \right] \\ &\leq \mathbb{E} \left[\sum_{t=2}^n \min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta}))} - 1 \right) \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^*, E_t(\boldsymbol{\theta}) | \mathcal{H}_{t-1}), 1 \right\} \right] \\ &\leq \mathbb{E} \left[\sum_{t=2}^n \min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta}))} - 1 \right) \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^* | \mathcal{H}_{t-1}), 1 \right\} \right], \end{aligned}$$

where we simply observed that $\mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^*, E_t(\boldsymbol{\theta}) | \mathcal{H}_{t-1}) \leq \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^* | \mathcal{H}_{t-1})$. \square

Proof of Theorem 6.1 Recalling that Algorithm 2 falls in the GRE-SH case, we can now proceed with the proof of Theorem 6.1.

Theorem 6.1. Let $\Theta = [K]$, $v(\boldsymbol{\theta}) = \max_{x' \in \Theta} d_2(p_{\boldsymbol{\theta}} \| p_{\boldsymbol{\theta}'})$ for all $\boldsymbol{\theta} \in \Theta$ ($v(\boldsymbol{\theta})$ can be infinite) and $v^*(\boldsymbol{\theta}) = \max\{v(\boldsymbol{\theta}), v(\boldsymbol{\theta}^*)\}$ where $\pi_{\boldsymbol{\theta}^*}$ is an optimal policy. For any $\alpha > 1$, the expected regret of RANDOMIST with truncation $M_t(\boldsymbol{\theta}) = \sqrt{\frac{td_2(p_{\boldsymbol{\theta}} \| \Phi_t)}{\alpha \log t}}$ is bounded as follows:

(a) if $v := \max_{\boldsymbol{\theta} \in \Theta} v(\boldsymbol{\theta}) < \infty$, $b \leq \sqrt{(\alpha \log t)/\eta_t(\boldsymbol{\theta})}$ and $a \geq 0$:

$$\begin{aligned} \mathbb{E} R(n) &\leq \frac{(188 + 32a)\alpha v^*(\boldsymbol{\theta})}{\sum_{\boldsymbol{\theta} \in \Theta: \Delta(\boldsymbol{\theta}) > 0} \Delta(\boldsymbol{\theta})} \log \frac{(94 + 16a)\alpha v^*(\boldsymbol{\theta})}{\Delta(\boldsymbol{\theta})^2} \\ &\quad + \Delta(\boldsymbol{\theta}_1) + \frac{\alpha + 3}{\alpha - 1} \min\left\{1, \sqrt{2 \log v}\right\} K; \end{aligned}$$

(b) no matter the value of v , if $a > 8$ and $J(\boldsymbol{\theta}) - \mathbb{E}[\check{J}_t(\boldsymbol{\theta})] \leq b \leq \sqrt{(\alpha \log t)/\eta_t(\boldsymbol{\theta})}$:

$$\mathbb{E} R(n) \leq \sum_{\boldsymbol{\theta} \in \Theta: \Delta(\boldsymbol{\theta}) > 0} \frac{(52 + 110a)c\alpha}{\Delta(\boldsymbol{\theta})} \log n + 2 \frac{\alpha + 1}{\alpha - 1} K,$$

where $c = 2 + \frac{e^2 \sqrt{a}}{\sqrt{2\pi}} \exp\left[\frac{16}{a-8}\right] \left(1 + \sqrt{\frac{\pi a}{a-8}}\right)$, with an instance-independent regret of $\mathbb{E} R(n) \leq 2\sqrt{(52 + 110a)c\alpha K n \log n} + 2 \frac{\alpha + 1}{\alpha - 1} K$.

Proof. For the sake of the proof, we denote $\eta_t(\boldsymbol{\theta}) = \check{J}_t(\boldsymbol{\theta}) + U_t(\boldsymbol{\theta})$. We apply Theorem D.1 with the choice $\zeta(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \frac{a}{2} + \frac{\Delta(x)}{2}$

Upper Bound on $b(\boldsymbol{\theta})$ Let us start with rewriting $Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta}))$ as:

$$Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta})) = \mathbb{P}\left(\eta_t(\boldsymbol{\theta}) > J(\boldsymbol{\theta}) + \frac{a}{2} + \frac{\Delta(x)}{2} \mid \mathcal{H}_{t-1}\right).$$

We ignore the dependence on $\zeta(\boldsymbol{\theta})$ whenever clear from the context, thus $Q_t(\boldsymbol{\theta}) = Q_t(\boldsymbol{\theta}, \zeta(\boldsymbol{\theta}))$. We start with bounding the term $b(\boldsymbol{\theta})$. Let us consider the event:

$$E_t(\boldsymbol{\theta}) = \left\{ \check{J}_t(\boldsymbol{\theta}) - J(\boldsymbol{\theta}) \leq \frac{\Delta(\boldsymbol{\theta})}{4} \right\}.$$

We can bound the probability that event $E_t(\boldsymbol{\theta})$ does not occur, by means of the inequalities of Lemma C.3:

$$\mathbb{P}\left(\overline{E_t(\boldsymbol{\theta})}\right) = \mathbb{P}\left(\check{J}_t(\boldsymbol{\theta}) - J(\boldsymbol{\theta}) > \frac{\Delta(\boldsymbol{\theta})}{4}\right) \leq \exp\left[-\frac{\Delta(\boldsymbol{\theta})^2 \eta_t(\boldsymbol{\theta})}{32 \left(1 + \frac{\Delta(\boldsymbol{\theta})}{12} \sqrt{\frac{\eta_t(\boldsymbol{\theta})}{\alpha \log t}}\right)}\right] \leq t^{-\alpha},$$

provided that $\eta_t(\boldsymbol{\theta}) \geq \frac{32(\sqrt{19}+10)\alpha}{9\Delta(\boldsymbol{\theta})^2} \log t \simeq \frac{52\alpha}{\Delta(\boldsymbol{\theta})^2} \log t$. Under event $E_t(\boldsymbol{\theta})$, we can bound the probability $Q_t(\boldsymbol{\theta})$ by means of Hoeffding's inequality and recalling that in any case $b \leq \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta})}}$:

$$\begin{aligned} Q_t(\boldsymbol{\theta}) &= \mathbb{P}\left(\check{J}_t(\boldsymbol{\theta}) + U_t(\boldsymbol{\theta}) > J(\boldsymbol{\theta}) + \frac{a}{2} + \frac{\Delta(\boldsymbol{\theta})}{2}\right) \\ &\leq \mathbb{P}\left(U_t(\boldsymbol{\theta}) - \frac{a}{2} > \frac{\Delta(\boldsymbol{\theta})}{4}\right) \\ &= \mathbb{P}\left(\frac{1}{a\eta_t(\boldsymbol{\theta})} \sum_{l=1}^{a\eta_t(\boldsymbol{\theta})} z_l - \frac{1}{2} > \frac{\Delta(\boldsymbol{\theta})}{4a} - \frac{b}{a}\right) \\ &\leq \mathbb{P}\left(\frac{1}{a\eta_t(\boldsymbol{\theta})} \sum_{l=1}^{a\eta_t(\boldsymbol{\theta})} z_l - \frac{1}{2} > \frac{\Delta(\boldsymbol{\theta})}{4a} - \frac{1}{a} \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta})}}\right) \\ &\leq \exp\left[-\frac{2}{a} \left(\frac{\Delta(\boldsymbol{\theta})}{4} - \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta})}}\right)^2 \eta_t(\boldsymbol{\theta})\right] \leq t^{-\alpha}, \end{aligned}$$

where we have to enforce the following two conditions:

$$\begin{aligned} \frac{\Delta(\boldsymbol{\theta})}{4} - \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta})}} &> 0 \\ \frac{2}{a} \left(\frac{\Delta(\boldsymbol{\theta})}{4} - \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta})}} \right)^2 \eta_t(\boldsymbol{\theta}) &\geq \alpha \log t. \end{aligned}$$

The second condition leads to:

$$\begin{aligned} \frac{2}{a} \left(\frac{\Delta(\boldsymbol{\theta})}{4} - \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta})}} \right)^2 s &\geq \alpha \log t \implies \frac{\Delta(\boldsymbol{\theta})}{4} - \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta})}} > \sqrt{\frac{a\alpha \log t}{2\eta_t(\boldsymbol{\theta})}} \\ \implies \eta_t(\boldsymbol{\theta}) &> \frac{16\alpha}{\Delta(\boldsymbol{\theta})^2} \left(1 + \sqrt{\frac{a}{2}} \right)^2 \log t \simeq \frac{16(2+a)\alpha}{\Delta(\boldsymbol{\theta})^2} \log t. \end{aligned}$$

Combining the two conditions, we require $\eta_t(\boldsymbol{\theta}) \geq \frac{(52+16a)\alpha}{\Delta(\boldsymbol{\theta})^2} \log t$. If $v(\boldsymbol{\theta}) < \infty$, we use $\eta_t(\boldsymbol{\theta}) \geq \frac{t-1}{v(\boldsymbol{\theta})}$ and we apply Lemma E.1 to get the following condition on the number of rounds that we denote with t^\dagger :

$$t \geq v(\boldsymbol{\theta}) \frac{2(52+16a)\alpha}{\Delta(\boldsymbol{\theta})^2} \log \left[v(\boldsymbol{\theta}) \frac{(52+16a)\alpha}{\Delta(\boldsymbol{\theta})^2} \right] + 1 := t^\dagger. \quad (71)$$

Now, we bound the term $b(\boldsymbol{\theta})$ when $v(\boldsymbol{\theta}) < \infty$:

$$b(\boldsymbol{\theta}) = \mathbb{E} \left[\sum_{t=2}^n \mathbb{1} \{Q_t(\boldsymbol{\theta}) > t^{-\alpha}\} \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta} | \mathcal{H}_{t-1}) \right] + \frac{1}{\alpha-1} \quad (72)$$

$$\leq \mathbb{E} \left[\sum_{t=2}^n \mathbb{1} \{Q_t(\boldsymbol{\theta}) > t^{-\alpha}\} \right] + \frac{1}{\alpha-1} \quad (73)$$

$$= \sum_{t=2}^n \mathbb{P}(Q_t(\boldsymbol{\theta}) > t^{-\alpha}) + \frac{1}{\alpha-1}. \quad (74)$$

If $t \leq t^\dagger$, we bound trivially $\mathbb{P}(Q_t(\boldsymbol{\theta}) > t^{-\alpha}) \leq 1$. Otherwise, we have:

$$\begin{aligned} \mathbb{P}(Q_t(\boldsymbol{\theta}) > t^{-\alpha}) &= \mathbb{P}(Q_t(\boldsymbol{\theta}) > t^{-\alpha} | E_t(\boldsymbol{\theta})) \mathbb{P}(E_t(\boldsymbol{\theta})) + \mathbb{P}(Q_t(\boldsymbol{\theta}) > t^{-\alpha} | \overline{E_t(\boldsymbol{\theta})}) \mathbb{P}(\overline{E_t(\boldsymbol{\theta})}) \\ &\leq \mathbb{P}(Q_t(\boldsymbol{\theta}) > t^{-\alpha} | E_t(\boldsymbol{\theta})) + \mathbb{P}(\overline{E_t(\boldsymbol{\theta})}) \leq 0 + t^{-\alpha}, \end{aligned}$$

where we exploited that under $E_t(\boldsymbol{\theta})$, we have $Q_t(\boldsymbol{\theta}) \leq t^{-\alpha}$. Thus, we have:

$$b(\boldsymbol{\theta}) \leq \sum_{t=2}^{\lfloor t^\dagger \rfloor} 1 + \sum_{t=\lfloor t^\dagger \rfloor+1}^n t^{-\alpha} + \frac{1}{\alpha-1} \leq t^\dagger - 1 + \frac{2}{\alpha-1},$$

where we bounded the summation with the integral, recalling that it must be $\lfloor t^\dagger \rfloor + 1 \geq 2$.

Instead, when $v(\boldsymbol{\theta}) = \infty$, we use $\eta_t(\boldsymbol{\theta}) \geq T_t(\boldsymbol{\theta})$, leading to the condition on the number of executions of policy $\boldsymbol{\theta}$, that we denote with s^\ddagger :

$$T_t(\boldsymbol{\theta}) \geq \frac{(52+16a)\alpha}{\Delta(\boldsymbol{\theta})^2} \log n := s^\ddagger. \quad (75)$$

To get the second bound on $b(\boldsymbol{\theta})$, i.e., the bound when $v(\boldsymbol{\theta}) = \infty$, we need some further manipulations. We denote with $t_s(\boldsymbol{\theta})$ the random round in which policy $\boldsymbol{\theta}$ is executed for the s -th time, with $t_0(\boldsymbol{\theta}) = 0$.

$$\begin{aligned} b(\boldsymbol{\theta}) &= \mathbb{E} \left[\sum_{t=2}^n \mathbb{1} \{Q_t(\boldsymbol{\theta}) > t^{-\alpha}\} \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta} | \mathcal{H}_{t-1}) \right] + \frac{1}{\alpha-1} \\ &= \mathbb{E} \left[\sum_{s=0}^{n-1} \sum_{t=t_s(\boldsymbol{\theta})+1}^{t_{s+1}(\boldsymbol{\theta})} \mathbb{1} \{Q_t(\boldsymbol{\theta}) > t^{-\alpha}\} \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta} | \mathcal{H}_{t-1}) \right] + \frac{1}{\alpha-1} \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} \left[\sum_{s=0}^{\lfloor s^\dagger \rfloor} \sum_{t=t_s(\boldsymbol{\theta})+1}^{t_{s+1}(\boldsymbol{\theta})} \mathbb{1} \{Q_t(\boldsymbol{\theta}) > t^{-\alpha}\} \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta} | \mathcal{H}_{t-1}) \right] \\
&\quad + \mathbb{E} \left[\sum_{s=\lfloor s^\dagger \rfloor+1}^{n-1} \sum_{t=t_s(\boldsymbol{\theta})+1}^{t_{s+1}(\boldsymbol{\theta})} \mathbb{1} \{Q_t(\boldsymbol{\theta}) > t^{-\alpha}\} \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta} | \mathcal{H}_{t-1}) \right] + \frac{1}{\alpha-1} \\
&\leq \mathbb{E} \left[\sum_{s=0}^{\lfloor s^\dagger \rfloor} \sum_{t=t_s(\boldsymbol{\theta})+1}^{t_{s+1}(\boldsymbol{\theta})} \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta} | \mathcal{H}_{t-1}) \right] + \mathbb{E} \left[\sum_{t=t_{\lfloor s^\dagger \rfloor+1}(\boldsymbol{\theta})+1}^n \mathbb{1} \{Q_t(\boldsymbol{\theta}) > t^{-\alpha}\} \right] + \frac{1}{\alpha-1} \\
&= \sum_{s=0}^{\lfloor s^\dagger \rfloor} 1 + \sum_{t=t_{\lfloor s^\dagger \rfloor+1}(\boldsymbol{\theta})+1}^n \mathbb{P}(Q_t(\boldsymbol{\theta}) > t^{-\alpha}) + \frac{1}{\alpha-1}.
\end{aligned}$$

Now, for $t \geq t_{\lfloor s^\dagger \rfloor+1}(\boldsymbol{\theta}) + 1$ we know that policy $\boldsymbol{\theta}$ was executed at least s^\dagger times, i.e., $T_t(\boldsymbol{\theta}) \geq s^\dagger$. Thus, we have, similarly as before:

$$\mathbb{P}(Q_t(\boldsymbol{\theta}) > t^{-\alpha}) \leq \mathbb{P}(Q_t(\boldsymbol{\theta}) > t^{-\alpha} | E_t(\boldsymbol{\theta})) + \mathbb{P}(\overline{E_t(\boldsymbol{\theta})}) \leq 0 + t^{-\alpha},$$

where we exploited that under event $E_t(\boldsymbol{\theta})$ we have that $Q_t(\boldsymbol{\theta}) \leq t^{-\alpha}$ and the upper bound on the probability that event $E_t(\boldsymbol{\theta})$ does not occur. Thus, we have:

$$b(\boldsymbol{\theta}) \leq \sum_{s=0}^{\lfloor s^\dagger \rfloor} 1 + \sum_{t=t_{\lfloor s^\dagger \rfloor+1}(\boldsymbol{\theta})+1}^n t^{-\alpha} + \frac{1}{\alpha-1} \leq s^\dagger + \frac{2}{\alpha-1},$$

where we bounded the summation with the integral, recalling that it must be $t_{\lfloor s^\dagger \rfloor+1}(\boldsymbol{\theta}) + 1 \geq 2$.

Upper Bound on $a(\boldsymbol{\theta})$ Fix a suboptimal policy $\boldsymbol{\theta} \neq \boldsymbol{\theta}^*$. We can rewrite the expression of $Q_t(\boldsymbol{\theta}^*, \zeta(\boldsymbol{\theta}))$ for the choice $\zeta(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \frac{a}{2} + \frac{\Delta(x)}{2} = J(\boldsymbol{\theta}^*) - \frac{a}{2} + \frac{\Delta(x)}{2}$:

$$Q_t(\boldsymbol{\theta}^*, \zeta(\boldsymbol{\theta})) = \mathbb{P} \left(\eta_t(\boldsymbol{\theta}^*) > J(\boldsymbol{\theta}^*) + \frac{a}{2} - \frac{\Delta(x)}{2} | \mathcal{H}_{t-1} \right) \quad (76)$$

Again, we discard the dependence on $\zeta(\boldsymbol{\theta})$, i.e., $Q_t(\boldsymbol{\theta}^*) = Q_t(\boldsymbol{\theta}^*, \zeta(\boldsymbol{\theta}))$. Large part of the derivation exploits tools similar to those employed for $b(\boldsymbol{\theta})$. We redefine event $E_t(\boldsymbol{\theta})$ as:

$$E_t(\boldsymbol{\theta}) = \left\{ J(\boldsymbol{\theta}^*) - \check{J}_t(\boldsymbol{\theta}^*) \leq \frac{\Delta(\boldsymbol{\theta})}{4} \right\}.$$

We now bound the probability that event $E_t(\boldsymbol{\theta})$ does not occur:

$$\mathbb{P}(\overline{E_t(\boldsymbol{\theta})}) = \mathbb{P} \left(J(\boldsymbol{\theta}^*) - \check{J}_t(\boldsymbol{\theta}^*) > \frac{\Delta(\boldsymbol{\theta})}{4} \right) \leq \exp \left[-\frac{1}{2} \left(\frac{\Delta(\boldsymbol{\theta})}{4} - \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta}^*)}} \right)^2 \eta_t(\boldsymbol{\theta}^*) \right] \leq t^{-\alpha}.$$

For this, we have to enforce two conditions. The first one for the bias and the second one for fulfilling the inequality.

$$\begin{aligned}
\frac{\Delta(\boldsymbol{\theta})}{4} - \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta}^*)}} &> 0 \\
\frac{1}{2} \left(\frac{\Delta(\boldsymbol{\theta})}{4} - \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta}^*)}} \right)^2 \eta_t(\boldsymbol{\theta}^*) &> \alpha \log t.
\end{aligned}$$

Let us consider the following derivation in which we retain only the positive solution:

$$\begin{aligned}
\frac{1}{2} \left(\frac{\Delta(\boldsymbol{\theta})}{4} - \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta}^*)}} \right)^2 \eta_t(\boldsymbol{\theta}^*) > \alpha \log t &\implies \frac{\Delta(\boldsymbol{\theta})}{4} - \sqrt{\frac{\alpha \log t}{\eta_t(\boldsymbol{\theta}^*)}} > \sqrt{\frac{2\alpha \log t}{\eta_t(\boldsymbol{\theta}^*)}} \\
\implies \eta_t(\boldsymbol{\theta}^*) > \frac{16(1 + \sqrt{2})^2 \alpha}{\Delta(\boldsymbol{\theta})^2} \log t.
\end{aligned}$$

Notice that this condition implies the first one on the bias. Consequently, we enforce $\eta_t(\boldsymbol{\theta}^*) > \frac{16(1+\sqrt{2})^2\alpha}{\Delta(\boldsymbol{\theta}^*)^2} \log t \simeq \frac{94\alpha}{\Delta(\boldsymbol{\theta}^*)^2} \log t$. Similarly for $Q_t(\boldsymbol{\theta}^*)$ under event $E_t(\boldsymbol{\theta})$:

$$\begin{aligned}
Q_t(\boldsymbol{\theta}^*) &= \mathbb{P} \left(\check{J}_t(\boldsymbol{\theta}^*) + U_t(\boldsymbol{\theta}^*) > J(\boldsymbol{\theta}) + \frac{a}{2} - \frac{\Delta(x)}{2} \right) \\
&\geq \mathbb{P} \left(U_t(\boldsymbol{\theta}^*) > \frac{a}{2} - \frac{\Delta(x)}{4} \right) \\
&= 1 - \mathbb{P} \left(\frac{a}{2} - U_t(\boldsymbol{\theta}^*) > \frac{\Delta(x)}{4} \right) \\
&= 1 - \mathbb{P} \left(\frac{1}{2} - \frac{1}{a\eta_t(\boldsymbol{\theta}^*)} \sum_{l=1}^{a\eta_t(\boldsymbol{\theta}^*)} z_l - \frac{1}{a} \sqrt{\frac{\alpha \log t}{s}} > \frac{\Delta(x)}{4a} \right) \\
&\geq 1 - \mathbb{P} \left(\frac{1}{2} - \frac{1}{a\eta_t(\boldsymbol{\theta}^*)} \sum_{l=1}^{a\eta_t(\boldsymbol{\theta}^*)} z_l > \frac{\Delta(x)}{4a} \right) \\
&\geq 1 - t^{-\alpha},
\end{aligned}$$

provided that $\eta_t(\boldsymbol{\theta}^*) \geq \frac{8a\alpha}{\Delta(\boldsymbol{\theta}^*)^2} \log t$ (using Hoeffding's inequality). Moreover, we have:

$$\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \leq \frac{t^\alpha}{t^\alpha - 1} - 1 = \frac{1}{t^\alpha - 1} \leq \frac{1}{(t-1)^\alpha},$$

for $t \geq 2$. Putting together these conditions we require $\eta_t(\boldsymbol{\theta}^*) \geq \frac{94a\alpha}{\Delta(\boldsymbol{\theta}^*)^2} \log t$. Those conditions, lead to the very similar requirements on the rounds and on the number of executions:

$$t \geq v(\boldsymbol{\theta}^*) \frac{188a\alpha}{\Delta(\boldsymbol{\theta}^*)^2} \log \left[v(\boldsymbol{\theta}^*) \frac{94a\alpha}{\Delta(\boldsymbol{\theta}^*)^2} \right] + 1 := t^\dagger, \quad (77)$$

$$T_t(\boldsymbol{\theta}^*) \geq \frac{94a\alpha}{\Delta(\boldsymbol{\theta}^*)^2} \log n := s^\dagger. \quad (78)$$

We now proceed at bounding $a(\boldsymbol{\theta})$. If $v(\boldsymbol{\theta}^*) < \infty$, we have:

$$a(\boldsymbol{\theta}) = \mathbb{E} \left[\sum_{t=2}^n \min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \right) \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^* | \mathcal{H}_{t-1}), 1 \right\} \right] \leq \mathbb{E} \left[\sum_{t=2}^n \min \left\{ \frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1, 1 \right\} \right].$$

If the round index is smaller than $t \leq t^\dagger$, we bound the min with 1. Otherwise, we proceed to the following decomposition, based on whether event $E_t(\boldsymbol{\theta})$ occurs:

$$\begin{aligned}
\mathbb{E} \left[\min \left\{ \frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1, 1 \right\} \right] &= \mathbb{E} \left[\min \left\{ \frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1, 1 \right\} | E_t(\boldsymbol{\theta}) \right] \mathbb{P}(E_t(\boldsymbol{\theta})) \\
&\quad + \mathbb{E} \left[\min \left\{ \frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1, 1 \right\} | \overline{E_t(\boldsymbol{\theta})} \right] \mathbb{P}(\overline{E_t(\boldsymbol{\theta})}) \\
&\leq \mathbb{E} \left[\min \left\{ \frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1, 1 \right\} | E_t(\boldsymbol{\theta}) \right] + \mathbb{P}(\overline{E_t(\boldsymbol{\theta})}) \\
&\leq \frac{1}{(t-1)^\alpha} + t^{-\alpha}.
\end{aligned}$$

Putting all together, we have:

$$a(\boldsymbol{\theta}) \leq \sum_{t=2}^{\lfloor t^\dagger \rfloor} 1 + \sum_{\lfloor t^\dagger \rfloor + 1}^n \left(\frac{1}{(t-1)^\alpha} + t^{-\alpha} \right) \leq t^\dagger - 1 + \frac{\alpha + 1}{\alpha - 1},$$

where we bounded the summations with the integrals, recalling that $\lfloor t^\dagger \rfloor + 1 \geq 2$.

Remark D.2. *It is worth noting that in this derivation of the bound on the term $a(\boldsymbol{\theta})$ we did not exploit the properties of the perturbation distribution. This is justified by the fact that we are considering the case $v(\boldsymbol{\theta}) < \infty$ and each sample for each policy is informative for all the policies. Indeed, FTL enjoys constant regret in this setting.*

For the case $v(\boldsymbol{\theta}^*) = \infty$, we need additional manipulations on the term $a(\boldsymbol{\theta})$:

$$\begin{aligned}
a(\boldsymbol{\theta}) &= \mathbb{E} \left[\sum_{t=2}^n \min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \right) \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^* | \mathcal{H}_{t-1}), 1 \right\} \right] \\
&= \mathbb{E} \left[\sum_{s=0}^{n-1} \sum_{t=t_s(\boldsymbol{\theta}^*)+1}^{t_{s+1}(\boldsymbol{\theta}^*)} \min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \right) \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^* | \mathcal{H}_{t-1}), 1 \right\} \right] \\
&= \mathbb{E} \left[\sum_{s=0}^{\lfloor s^\dagger \rfloor} \sum_{t=t_s(\boldsymbol{\theta}^*)+1}^{t_{s+1}(\boldsymbol{\theta}^*)} \min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \right) \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^* | \mathcal{H}_{t-1}), 1 \right\} \right] \\
&\quad + \mathbb{E} \left[\sum_{s=\lfloor s^\dagger \rfloor+1}^{n-1} \sum_{t=t_s(\boldsymbol{\theta}^*)+1}^{t_{s+1}(\boldsymbol{\theta}^*)} \min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \right) \mathbb{P}(\boldsymbol{\theta}_t = \boldsymbol{\theta}^* | \mathcal{H}_{t-1}), 1 \right\} \right] \\
&\leq \mathbb{E} \left[\mathbb{E} \left[\sum_{s=0}^{\lfloor s^\dagger \rfloor} \sum_{t=t_s(\boldsymbol{\theta}^*)+1}^{t_{s+1}(\boldsymbol{\theta}^*)} \left(\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \right) \mathbb{1}_{\{\boldsymbol{\theta}_t = \boldsymbol{\theta}^*\}} | \mathcal{H}_{t-1} \right] \right] \\
&\quad + \mathbb{E} \left[\sum_{t=t_{\lfloor s^\dagger \rfloor}(\boldsymbol{\theta}^*)+1}^n \min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \right), 1 \right\} \right] \\
&= \mathbb{E} \left[\sum_{s=0}^{\lfloor s^\dagger \rfloor} \left(\frac{1}{Q_{t_{s+1}(\boldsymbol{\theta}^*)}(\boldsymbol{\theta}^*)} - 1 \right) \right] + \mathbb{E} \left[\sum_{t=t_{\lfloor s^\dagger \rfloor}(\boldsymbol{\theta}^*)+1}^n \min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \right), 1 \right\} \right].
\end{aligned}$$

Now, for $s \leq s^\dagger$ we apply the upper bound that will be proved later in Lemma D.2:

$$\mathbb{E} \left[\frac{1}{Q_t(\boldsymbol{\theta}^*)} \right] = \mathbb{E} \left[\frac{1}{\mathbb{P}(\eta_t(\boldsymbol{\theta}^*) \geq J(\boldsymbol{\theta}^*) + \frac{a}{2} - \frac{\Delta(\boldsymbol{\theta})}{2} | \mathcal{H}_{t-1})} \right] \leq c, \quad \forall t \geq 1. \quad (79)$$

Instead, for $t \geq t_{\lfloor s^\dagger \rfloor}(\boldsymbol{\theta}^*) + 1$, we know that the optimal policy $\boldsymbol{\theta}^*$ was executed at least s^\dagger times, i.e., $T_t(\boldsymbol{\theta}) \geq s^\dagger$. Therefore:

$$\begin{aligned}
\mathbb{E} \left[\min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \right), 1 \right\} \right] &= \mathbb{E} \left[\min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \right), 1 \right\} | E_t(\boldsymbol{\theta}^*) \right] \mathbb{P}(E_t(\boldsymbol{\theta}^*)) \\
&\quad + \mathbb{E} \left[\min \left\{ \left(\frac{1}{Q_t(\boldsymbol{\theta}^*)} - 1 \right), 1 \right\} | \overline{E}_t(\boldsymbol{\theta}^*) \right] \mathbb{P}(\overline{E}_t(\boldsymbol{\theta}^*)) \\
&\leq \frac{1}{(t-1)^\alpha} + t^\alpha.
\end{aligned}$$

Putting all together, we have:

$$\begin{aligned}
a(\boldsymbol{\theta}) &\leq 1 + \sum_{s=0}^{\lfloor s^\dagger \rfloor} c + \sum_{t=t_{\lfloor s^\dagger \rfloor}(\boldsymbol{\theta}^*)+1}^n \left(\frac{1}{(t-1)^\alpha} + t^{-\alpha} \right) \\
&\leq cs^\dagger + \frac{\alpha+1}{\alpha-1}.
\end{aligned}$$

Putting together $a(\boldsymbol{\theta})$ and $b(\boldsymbol{\theta})$ Consider the case in which $v(\boldsymbol{\theta}) < \infty$:

$$\begin{aligned}
a(\boldsymbol{\theta}) + b(\boldsymbol{\theta}) &\leq v(\boldsymbol{\theta}^*) \frac{188a\alpha}{\Delta(\boldsymbol{\theta})^2} \log \left[v(\boldsymbol{\theta}^*) \frac{94a\alpha}{\Delta(\boldsymbol{\theta})^2} \right] + \frac{\alpha+1}{\alpha-1} \\
&\quad + v(\boldsymbol{\theta}) \frac{2(52+16a)\alpha}{\Delta(\boldsymbol{\theta})^2} \log \left[v(\boldsymbol{\theta}) \frac{(52+16a)\alpha}{\Delta(\boldsymbol{\theta})^2} \right] + \frac{2}{\alpha-1} \\
&\leq v^*(\boldsymbol{\theta}) \frac{(188+32a)\alpha}{\Delta(\boldsymbol{\theta})^2} \log \left[v^*(\boldsymbol{\theta}) \frac{(94+16a)\alpha}{\Delta(\boldsymbol{\theta})^2} \right] + \frac{\alpha+3}{\alpha-1},
\end{aligned}$$

where we exploited the definition of $v^*(\boldsymbol{\theta}) = \max\{v(\boldsymbol{\theta}), v(\boldsymbol{\theta}^*)\}$. The result is obtained by recalling that each $\Delta(\boldsymbol{\theta}) \leq \min\{1, \sqrt{2 \log v}\}$ (Lemma C.5). Instead, for the case $v(\boldsymbol{\theta}) = \infty$ we have:

$$\begin{aligned} a(\boldsymbol{\theta}) + b(\boldsymbol{\theta}) &\leq c \frac{94a\alpha}{\Delta(\boldsymbol{\theta})^2} \log n + \frac{2}{\alpha - 1} + \frac{(52 + 16a)\alpha}{\Delta(\boldsymbol{\theta})^2} \log n + \frac{\alpha + 1}{\alpha - 1} \\ &\leq \frac{(110a + 52)c\alpha}{\Delta(\boldsymbol{\theta})^2} \log n + \frac{\alpha + 3}{\alpha - 1}, \end{aligned}$$

where we simply exploited that $c > 1$. The result follows by trivially bounding each $\Delta(\boldsymbol{\theta}) \leq 1$. \square

Derivation of the c term We now explicitly derive the value of c bounding the expectation $\mathbb{E}\left[\frac{1}{Q_t(\boldsymbol{\theta}^*)}\right]$ for all $t \in [n]$. For the sake of the analysis, we will replace a with $2a$. Let us define the following symbols:

$$\bar{X} = s \mathbb{E}[\check{J}_t(\boldsymbol{\theta}^*)] \in [0, s], \quad X = s \check{J}_t(\boldsymbol{\theta}^*) \in [0, sM_t], \quad \bar{Y} = as, \quad Y = \sum_{l=1}^{2as} z_l,$$

where $z_l \sim \text{Ber}(1/2)$ and i.i.d.. First of all, we need to take into account the beneficial effect of our correction, that leads to the following derivation recalling that $J(\boldsymbol{\theta}^*) - \mathbb{E}[\check{J}_t(\boldsymbol{\theta}^*)] \leq b$:

$$\begin{aligned} \mathbb{P}\left(\eta_t(\boldsymbol{\theta}^*) \geq J(\boldsymbol{\theta}^*) + a - \frac{\Delta(\boldsymbol{\theta})}{2} \mid \mathcal{H}_{t-1}\right) &= \mathbb{P}\left(\check{J}_t(\boldsymbol{\theta}^*) + U_t(\boldsymbol{\theta}^*) \geq J(\boldsymbol{\theta}^*) + a - \frac{\Delta(\boldsymbol{\theta})}{2} \mid \mathcal{H}_{t-1}\right) \\ &= \mathbb{P}\left(\check{J}_t(\boldsymbol{\theta}^*) + \frac{1}{s} \sum_{l=1}^{2as} z_l \geq \underbrace{J(\boldsymbol{\theta}^*) - b + a - \frac{\Delta(\boldsymbol{\theta})}{2}}_{\leq \mathbb{E}[\check{J}_t(\boldsymbol{\theta}^*)]} \mid \mathcal{H}_{t-1}\right) \\ &\geq \mathbb{P}\left(\check{J}_t(\boldsymbol{\theta}^*) + \frac{1}{s} \sum_{l=1}^{2as} z_l \geq \mathbb{E}[\check{J}_t(\boldsymbol{\theta}^*)] + a - \frac{\Delta(\boldsymbol{\theta})}{2} \mid \mathcal{H}_{t-1}\right) \\ &= \mathbb{P}\left(X + Y \geq \bar{X} + \bar{Y} - \frac{\Delta(\boldsymbol{\theta})s}{2} \mid \mathcal{H}_{t-1}\right) \\ &\geq \mathbb{P}(X + Y \geq \bar{X} + \bar{Y} \mid \mathcal{H}_{t-1}). \end{aligned}$$

The following result bounds the probability $\mathbb{P}(X + Y \geq \bar{X} + \bar{Y} \mid \mathcal{H}_{t-1})$.

Lemma D.2. *For any $a > 4$, it holds that:*

$$\mathbb{E}\left[\frac{1}{\mathbb{P}(X + Y \geq \bar{X} + \bar{Y} \mid \mathcal{H}_{t-1})}\right] \leq 2 + \frac{2e^2\sqrt{a}}{\sqrt{\pi}} \left[\frac{8}{a-4}\right] \left(1 + \sqrt{\frac{\pi a}{2(a-4)}}\right).$$

Proof. The proof puts together some of the results presented in Appendix A of (Kveton et al. 2019a). Let $f(X)$ be defined as follows:

$$f(X) = \left[\sum_{y=\lceil \bar{X}-X+as \rceil}^{2as} g_{\text{Bin}(2as, 1/2)}(y) \right]^{-1},$$

where $g_{\text{Bin}(n,p)}$ is the p.d.f. of a Binomial distribution of parameters n and p . Given the definition of f , the following identity holds: $\mathbb{E}\left[\frac{1}{\mathbb{P}(X+Y \geq \bar{X} + \bar{Y} \mid \mathcal{H}_{t-1})}\right] = \mathbb{E}[f(X)]$. Let us define the partitioning of the interval $[0, sM_t]$, with i_0 the smallest integer s.t. $(i_0 + 1)\sqrt{s} \geq \bar{X}$:

$$\mathcal{P}_i = \begin{cases} (s, sM_t] & \text{if } i = -1, \\ (\max\{\bar{X} - \sqrt{s}, s\}, s] & \text{if } i = 0, \\ (\max\{\bar{X} - (i+1)\sqrt{s}, \bar{X} - i\sqrt{s}\}, \max\{\bar{X} - (i+1)\sqrt{s}, \bar{X} - i\sqrt{s}\}] & \text{if } i \in \{1, \dots, i_0\}, \end{cases} \quad (80)$$

We can now decompose the expectation $\mathbb{E}[f(X)]$ over the partitioning:

$$\mathbb{E}[f(X)] = \sum_{i=-1}^{i_0} \mathbb{E}[\mathbb{1}\{X \in \mathcal{P}_i\} f(X)]$$

$$\leq \mathbb{P}(X \in \mathcal{P}_{-1})f(s) + \sum_{i=0}^{i_0-1} \mathbb{P}(X \in \mathcal{P}_i)f(\bar{X} - (i+1)\sqrt{s}) + \mathbb{P}(X \in \mathcal{P}_{i_0})f(0),$$

where we simply observed that f is a decreasing function of X . We now proceed at bounding the probabilities for $i \geq 1$:

$$\begin{aligned} \mathbb{P}(X \in \mathcal{P}_i) &\leq \mathbb{P}(X \leq \bar{X} - i\sqrt{s}) \\ &\leq \mathbb{P}\left(\mathbb{E}[\check{J}_t(\boldsymbol{\theta}^*)] - \check{J}_t(\boldsymbol{\theta}^*) \geq \frac{i}{\sqrt{s}}\right) \\ &\leq \exp\left[-\frac{i^2}{2s}\right] = \exp\left[-\frac{i^2}{2}\right], \end{aligned}$$

where we applied the concentration inequalities in Lemma C.3. For $i \in \{-1, 0\}$ we bound trivially $\mathbb{P}(X \in \mathcal{P}_i) \leq 1$. Notice that we are satisfying the constraint on the bias, thanks to the bias correction.

We now analyze the terms $f(\bar{X} - (i+1)\sqrt{s})$. First of all, recall that $f(s) \leq 2$ since $[\bar{X} - s + as] \leq as$. For the other terms we apply Lemma 2 of (Kveton et al. 2019a) for $\delta \in [0, as]$:

$$\sum_{y=\lceil as+\delta \rceil}^{2as} g_{\text{Bin}(2as, 1/2)}(y) \geq \frac{\sqrt{\pi}}{e^2\sqrt{a}} \exp\left[-\frac{2(\delta + \sqrt{s})^2}{as}\right],$$

Specifically, in our case, we have for $\delta = (i+1)\sqrt{s}$:

$$\sum_{y=\lceil as+(i+1)\sqrt{s} \rceil}^{2as} g_{\text{Bin}(2as, 1/2)}(y) \geq \frac{\sqrt{\pi}}{e^2\sqrt{a}} \exp\left[-\frac{2(i+2)^2}{a}\right].$$

Instead, for $\delta = \bar{X}$:

$$\sum_{y=\lceil as+\bar{X} \rceil}^{2as} g_{\text{Bin}(2as, 1/2)}(y) \geq \frac{\sqrt{\pi}}{e^2\sqrt{a}} \exp\left[-\frac{2(\bar{X} + \sqrt{s})^2}{as}\right] \geq \frac{\sqrt{\pi}}{e^2\sqrt{a}} \exp\left[-\frac{2(i_0+2)^2}{a}\right].$$

As a consequence, the expectation of f can be rewritten as:

$$\mathbb{E}[f(X)] \leq 2 + \frac{e^2\sqrt{a}}{\sqrt{\pi}} \sum_{i=0}^{i_0} \exp\left[-\frac{ai^2 - 4(i+2)^2}{2a}\right].$$

To get to a result, we complete the square:

$$ai^2 - 4(i+2)^2 = (a-4) \left(i - \frac{8}{a-4}\right)^2 - \frac{16a}{a-4}.$$

It follows that, by bounding the summation with the integral under the assumption that $a > 4$:

$$\begin{aligned} \mathbb{E}[f(X)] &\leq 2 + \frac{e^2\sqrt{a}}{\sqrt{\pi}} \sum_{i=0}^{i_0} \exp\left[-\frac{a-4}{2a} \left(i - \frac{8}{a-4}\right)^2 + \frac{8}{a-4}\right] \\ &\leq 2 + \frac{2e^2\sqrt{a}}{\sqrt{\pi}} \exp\left[\frac{8}{a-4}\right] \sum_{i=0}^{\infty} \exp\left[-\frac{a-4}{2a} i^2\right] \\ &\leq 2 + \frac{2e^2\sqrt{a}}{\sqrt{\pi}} \left[\frac{8}{a-4}\right] \left(1 + \int_{x=0}^{\infty} \exp\left[-\frac{a-4}{2a} x^2\right] dx\right) \\ &\leq 2 + \frac{2e^2\sqrt{a}}{\sqrt{\pi}} \left[\frac{8}{a-4}\right] \left(1 + \sqrt{\frac{\pi a}{2(a-4)}}\right). \end{aligned}$$

□

E Auxiliary Lemmas

In this appendix, we provide some auxiliary lemmas that are employed to prove the main results.

Theorem E.1. *Let X_1, \dots, X_N are independent random variables satisfying $0 \leq X_i \leq M$ and $\mathbb{E}[X_i^2] \leq v^2$. Let $\bar{X} =$*

$\frac{1}{n} \sum_{i=1}^N X_i$. Then for any $\epsilon \geq 0$, we have:

$$\begin{aligned} \mathbb{P}(\bar{X} - \mathbb{E}[\bar{X}] \geq \epsilon) &\leq \exp\left[\frac{-\epsilon^2 N}{2\left(v^2 + \frac{M\epsilon}{3}\right)}\right], \\ \mathbb{P}(\mathbb{E}[\bar{X}] - \bar{X} \geq \epsilon) &\leq \exp\left[\frac{-\epsilon^2 N}{2v^2}\right]. \end{aligned}$$

Proof. These are just Theorems 2.7 and 2.8 of (Chung and Lu 2006) suitably rephrased. \square

Lemma E.1. Let $f(t) = a \log t - t + 1$ for $t \geq 1$ and $a \geq 1$. Then, the following statements hold:

1. f is concave;
2. the maximum of f is attained by $t_{\max} = a$;
3. the derivative of f in $t = 1$ is $a - 1$;
4. f admits two zeros: one in $t_1 = 1$ and the other in:

$$t_2 = -aW_{-1}\left(-\frac{1}{a}e^{-\frac{1}{a}}\right) \leq 2a \log a + 1, \quad (81)$$

where W_{-1} is the secondary component of the Lambert function (Corless et al. 1996).

Proof. The first three points are trivial. We just prove 4. Consider the equation and the subsequent derivation:

$$\begin{aligned} a \log t &= t - 1 \\ t^a &= e^{t-1} \\ t &= e^{t/a} e^{-1/a} \\ -\frac{t}{a} &= -\frac{1}{a} e^{\frac{t}{a}} e^{-\frac{1}{a}} \\ -\frac{t}{a} e^{-\frac{t}{a}} &= -\frac{1}{a} e^{-\frac{1}{a}}. \end{aligned}$$

Now we can apply the Lambert function yielding two solutions since the right-hand-side is in $[-1/e, 0]$, one for W_0 (the principal component) and one for W_{-1} (the secondary component):

$$t_{1,2} = -aW_{\star}\left(-\frac{1}{a}e^{-\frac{1}{a}}\right), \quad \star \in \{0, -1\}. \quad (82)$$

Moreover, since $-1/a > -1$, the first solution simplifies into $t_1 = -aW_0\left(-\frac{1}{a}e^{-\frac{1}{a}}\right) = 1$, which was already clear from the definition of function $f(t)$. The other solution t_2 cannot be further simplified. We proceed to bound its value. Notice that, in order to prove that $t_2 \leq 2a \log a + 1$, it is sufficient to show that $f(2a \log a + 1) \leq 0$, since f is concave:

$$f(2a \log a + 1) = a \log(2a \log a + 1) - 2a \log a \leq 0 \quad \text{if} \quad 2a \log a + 1 \leq a^2.$$

So it is enough to prove that function $g(a) := a^2 - 2a \log a - 1 \geq 0$ for all $a \geq 1$. But $g(1) = 0$ and g is monotonically increasing in a . \square

F Experimental Details

In this appendix, we present the practical aspects of RANDOMIST omitted in the main paper. In particular, we discuss the computation of Rényi divergences (Appendix F.1), we provide a detailed description of the adaptation of MCMC used in compact parameter spaces (Appendix F.2), we discuss the caching strategy that allows saving a factor t in the computational complexity (Appendix F.3) and provide additional experimental results (Appendix F.4).

Infrastructure The experiments have been run on a machine with two CPUs Intel(R) Xeon(R) CPU E7-8880 v4 @ 2.20GHz (22 cores, 44 thread, 55 MB cache) and 128 GB RAM.

F.1 Computing Rényi Divergences

To generate the perturbation $U_t(\theta)$ of RANDOMIST (Algorithm 2), we need to compute the effective number of trajectories $\eta_t(\theta)$, which in turn requires the Rényi divergence $d_2(p_\theta \| \Phi_t)$ between the candidate and the mixture of previously executed

policies (see Section 3). This latter quantity can be challenging to compute since the mixture Φ_t is typically difficult to characterize. From (Papini et al. 2019) (Theorem 5), this quantity is upper bounded by the harmonic mean of pairwise divergences:

$$d_2(p_\theta \|\Phi_t) \leq \frac{t-1}{\sum_{i=1}^{t-1} \frac{1}{d_2(p_\theta \| p_{\theta_i})}}. \quad (83)$$

In practice, we replace all occurrences of $d_2(p_\theta \|\Phi_t)$ in RANDOMIST with the harmonic mean. It is easy to show that Theorem 6.1 still holds for this modified version, since the proof is already based on the above upper bound. The same argument holds for the exploration bonus of OPTIMIST, as already observed in (Papini et al. 2019).

To compute $d_2(p_\theta \|\Phi_t)$, we just have to compute the pairwise Rényi divergences $d_2(p_\theta \| p_{\theta'})$ for each θ' previously executed. This is straightforward in the PB-PO framework, where the divergence is between hyperpolicies:

$$\int_{\Theta} \int_{\mathcal{T}} \nu_{\xi'}(\theta) p_\theta(\tau) \left(\frac{\nu_\xi(\theta) p_\theta(\tau)}{\nu_{\xi'}(\theta) p_\theta(\tau)} \right)^2 d\tau d\theta = \int_{\Theta} \nu_{\xi'}(\theta) \left(\frac{\nu_\xi(\theta)}{\nu_{\xi'}(\theta)} \right)^2 d\theta = d_2(\nu_\xi \|\nu_{\xi'}), \quad (84)$$

and the hyperpolicies are perfectly known, often Gaussian.¹²

In the AB-PO framework, the outcome distributions are trajectory distributions, which are unknown (although we can easily compute probability ratios). Possible estimators for $d_2(p_\theta \| p_{\theta'})$ are discussed in (Metelli et al. 2018). It is an open problem whether this approximation affects the regret. In this work, we only experiment with the easier PB-PO framework.

F.2 RANDOMIST in compact parameter spaces

In Section 6, we presented a version of RANDOMIST for infinite policy spaces. We now first provide more details on the algorithm, then we present a more practical version with reduced computational complexity. For both the versions, implemented in the parameter-based setting (PB-PO), we improve the hyperpolicy parameters by taking $M = 10$ steps of the Metropolis-Hastings algorithm with Gaussian proposal. The pseudocode of our method is presented in Algorithm 5. The two versions differ for how g^\dagger is computed.

```

Input: initial policy parameters  $\theta_1$ , kernel covariance  $\Sigma$ 
Execute  $\pi_{\theta_1}$ , observe  $\tau_1 \sim p_{\theta_1}$  and  $\mathcal{R}(\tau_1)$ 
for  $t = 2, \dots, n$  do
  Initialize  $\theta_m = \theta_{t-1}$ 
  for  $m = 1, \dots, M-1$  do
    Select proposed point  $\theta_p \sim \mathcal{N}(\theta_m, \Sigma)$ 
    Sample  $\epsilon \sim \mathcal{U}(0, 1)$ 
    if  $\epsilon < \frac{g_t^\dagger(\theta_p | \theta_m)}{g_t^\dagger(\theta_m | \theta_p)}$  then
      Set  $\theta_{m+1} = \theta_p$ 
    else
      Set  $\theta_{m+1} = \theta_m$ 
    end if
  end for
  Set  $\theta_t = \theta_M$ 
  Execute  $\pi_{\theta_t}$ , observe  $\tau_t \sim p_{\theta_t}$  and  $\mathcal{R}(\tau_t)$ 
end for

```

Algorithm 5: MCMC-RANDOMIST

Full-Density MCMC-RANDOMIST As discussed in Section 6, it is possible to approximately compute the probability for a policy θ of being the one with maximum perturbed estimated expected return. Since the perturbation is drawn from a binomial distribution that sums $as(\theta)$ Bernoulli samples, we compute the outer integral as a summation over the $as(\theta)$ possible values, considering the generic probability density function $g_{\text{Bin}(as(\theta), \frac{1}{2})}$ of a binomial distribution. To evaluate the probability of being the best for each one of these possible perturbed estimated expected return for a fixed θ , we would like to multiply the probabilities for it of being larger than the estimated expected return of any other policy: this amounts to multiplying the cumulative density functions of the distribution of the estimated expected return of other policies θ' , evaluated in the particular candidate expected return value for θ . Since, in a compact parameter space, the number of θ' s is infinite, a *product integral*, i.e., the multiplication equivalent of the standard integral, must be considered. We compute this integral by a kind of numerical

¹²Closed-form expressions for the Rényi divergence are available for several common distributions, including Gaussians (Gil, Alajaji, and Linder 2013).

quadrature, only considering policies taken in previous timesteps $t' < t$, leading to the following expression for the density:

$$\mathbf{g}_t^\dagger(\boldsymbol{\theta}) \propto \int_{\mathbb{R}} g_{\boldsymbol{\theta}}(y) \prod_{t'=1}^{t-1} G_{\boldsymbol{\theta}_{t'}}(y) dy \quad (85)$$

$$= \sum_{i=0}^{\#} g_{\text{Bin}(\#, \frac{1}{2})}(i) \prod_{t'=1}^{t-1} G_{\text{Bin}(\#, \frac{1}{2})} \left(\left[\frac{s(\boldsymbol{\theta}_{t'})}{s(\boldsymbol{\theta})} i + s(\boldsymbol{\theta}_{t'}) \left(\hat{J}_t(\boldsymbol{\theta}) - \hat{J}_t(\boldsymbol{\theta}_{t'}) \right) \right] \right), \quad (86)$$

where $\# = \lceil as(\boldsymbol{\theta}) \rceil$. Note that, for practical and efficiency purposes, we shift the argument of the cumulative density function $G_{(\#, \frac{1}{2})}$, rather than modifying its parameters for the different $\boldsymbol{\theta}_{t'}$. Thus, given that it can only take a discrete number of values, we employ memoization, by caching the first 10^5 values for each combination of $\#$ and c.d.f. argument after their first computation.

1-Step-Density MCMC-RANDOMIST The MCMC-based RANDOMIST needs, for each call to the density function \mathbf{g}^\dagger , to have an updated estimate of the expected returns of each policy taken in the previous timesteps. This is particularly expensive, being based on the balance heuristic estimator, leading to a $\mathcal{O}(dt^3)$ per-iteration complexity, obviously computationally demanding. Therefore, we propose an alternative heuristic for the computation of the density, based on the idea of comparing the estimated expected return of a proposed point just with the one of the last node in the chain constructed by the MCMC algorithm:

$$\mathbf{g}_t^\dagger(\boldsymbol{\theta}|\boldsymbol{\theta}') = \sum_{i=0}^{\#} g_{\text{Bin}(\#, \frac{1}{2})}(i) G_{\text{Bin}(\#, \frac{1}{2})} \left(\left[\frac{s(\boldsymbol{\theta}')}{s(\boldsymbol{\theta})} i + s(\boldsymbol{\theta}') \left(\hat{J}_t(\boldsymbol{\theta}) - \hat{J}_t(\boldsymbol{\theta}') \right) \right] \right), \quad (87)$$

where $\# = \lceil as(\boldsymbol{\theta}) \rceil$. In other words, we employ $\frac{\mathbb{P}(\hat{J}_t(\boldsymbol{\theta}') \geq \hat{J}_t(\boldsymbol{\theta}))}{\mathbb{P}(\hat{J}_t(\boldsymbol{\theta}) \geq \hat{J}_t(\boldsymbol{\theta}'))}$ as a ratio for the Metropolis-Hastings algorithm. Thus, the estimation of the expected return for all the policies played in the previous rounds can be avoided, gaining a factor of t in the asymptotic time complexity. Despite its potentially greedy aptitude, this version enjoys very similar performance to the algorithm that uses the full density.

We leave as a future work the study of the convergence properties of these approximations of the density.

F.3 Reducing Complexity via Caching

The main computational overhead for the class of algorithms presented in this paper is the update of the estimated expected return of a policy. Indeed, for both OPTIMIST and RANDOMIST, both the importance weights (through the balance heuristic) and the Rényi divergence between the policy and the mixture of the other policies must be computed, leading to a time complexity of the order of $\mathcal{O}(t^2)$. In fact, for each policy, potentially, every other policy should be queried in order to obtain the correct MIS denominator for a given sample and the correct Rényi divergence.

Nonetheless, it is possible to reduce the computational complexity by noting at time t that, for all the samples collected up to time $t - 1$, all the probabilities relative to those samples have already been computed. Therefore, by storing a $(t - 1) \times (t - 1)$ matrix that holds the probability of each previous sample under each previous policy, one can tradeoff memory for time, and avoid the repeated computations of these values. A similar reasoning applies to the Rényi divergences. As discussed in Appendix F.1, we follow the same approach as in (Papini et al. 2019), and bound the Rényi divergence $d_2(p_{\boldsymbol{\theta}} \parallel \Phi_t)$ with the harmonic mean of the pairwise divergences between $\boldsymbol{\theta}$ and the other policies. At time step t , instead of recomputing from scratch the divergences for an old policies in order to update its TMIS estimator, we use a cached sum of the previous divergences, that we plug into the computation of the harmonic mean. Overall, the time complexity of these computations is reduced to $\mathcal{O}(t)$.

F.4 Tasks, additional experiments, hyperparameters

We now give additional details on the tasks used to evaluate RANDOMIST, also presenting additional results. For all the tasks, we parallelized the execution of different runs of our experiments by using the GNU Parallel tool (Tange 2011). We employ $\alpha = 2$ for the importance weight truncation and computation of the exploration bonus.

LQG For the LQG task, we obtain the cumulative regret by computing, for all the $K = 100$ hyperpolicies generated by the discretization, the expected return in closed form, then comparing it, at each round, with the one obtained using the selected hyperpolicy in a deterministic version of the environment.

Mountain Car To test the MCMC-RANDOMIST algorithm, we employ the standard Mountain Car (Sutton and Barto 2018) task. We compute the *cumulative return* at time t as the the average of the return obtained across all iterations for $t' \leq t$, following the same approach of (Papini et al. 2019). We normalize the returns from the interval $[-5, 95]$ to $[0, 1]$ whenever it is appropriate for our algorithm, for instance inside the argument of the cumulative density function of the binomial distribution.

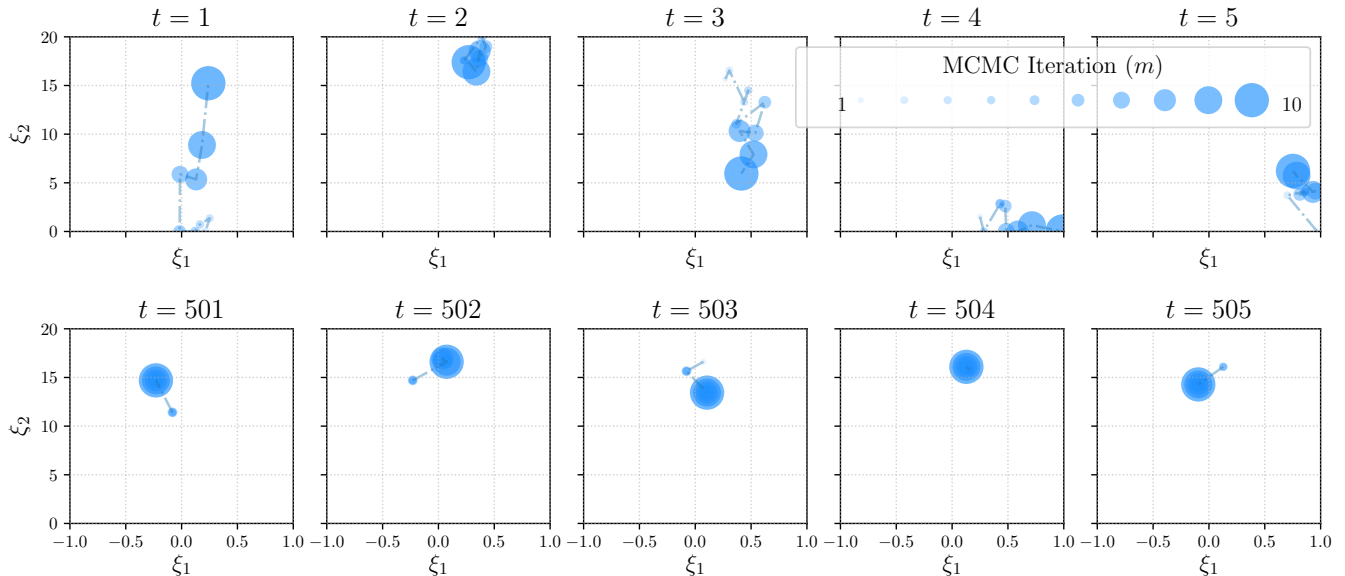


Figure 4: Trajectories traced by the $M = 10$ Metropolis-Hastings steps in the two-dimensional space Ξ of the hyperparameters $\xi = [\xi_1, \xi_2]$ of the hyperpolicy ν_ξ during different iterations. Top row: first 5 iterations. Bottom row: iterations from 501 to 505.

We restrict our search for the hyperparameters inside the box $[-1, 1] \times [0, 20]$ and heuristically set the covariance matrix of the Gaussian proposal used to construct the MCMC chain (we use the full-density approximation) to be equal to the one of the hyperpolicy, i.e., $\Sigma = \text{diag}(0.15, 3)^2$. We employ linear policies in the state. In addition to the experiment discussed in Section 7, Figure 4 reports an investigation of how the behavior of MCMC-RANDOMIST changes during its training. In the first iterations of the algorithm, when very few samples have been collected, no hyperpolicy is able to obtain a high probability of leading to the maximum expected return. Therefore, few proposed points are rejected and the Markov Chains constructed by the Metropolis-Hastings algorithm cross a considerable portion of the Ξ space. By contrast, in later iterations of the algorithm, when it is converging towards its maximum performance, acceptance of new proposed hyperparameters becomes rare and the steps that compose the trajectories are very small. This behavior can be naturally explained as a manifestation, on the way MCMC trajectories are traced, of an initial exploration phase of the space Ξ , followed by an exploitation phase.

Continuous Cartpole An interesting feature of the continuous version of RANDOMIST is its scalability compared to existing approaches (i.e., the discretization-based version of OPTIMIST, called OPTIMIST2 in (Papini et al. 2019)). To assess it in practice, we employ a continuous Cartpole environment, with $\dim(\mathcal{S}) = 4$ and $\dim(\mathcal{A}) = 1$ and an horizon $H = 200$. We employ the most efficient version of RANDOMIST, by only considering proposed and current points during the computation of the density ratio in the MCMC step and the hyperparameter $a = 0.1$. We again employ the parameter-based setting, by leveraging a Gaussian hyperpolicy of learned mean and constant covariance $\Sigma = \text{diag}(1, 1, 1, 1)$, together with policies linear in the state features. We run, at each iteration, $M = 10$ steps of MCMC (we use the one-step approximation), with a Gaussian proposal centered in the current point and covariance equal to the one of the hyperpolicy, and search for the optimal hyperparameters in the box $[-2, 2] \times [0, 4] \times [0, 10] \times [0, 12]$. Figure 5 (left) shows the cumulative return, computed as in the mountain car experiment, obtained by RANDOMIST, OPTIMIST2¹³ and PGPE (Sehnke et al. 2008) run with different step sizes. RANDOMIST is able to scale to this task and to obtain satisfying performance, comparable to the one of a policy gradient approach. It is interesting to observe that OPTIMIST2 completely fails to learn the task, or, at least, displays a very slow improving behavior. We tested OPTIMIST2 with different discretizations, all fulfilling Theorem 4 of (Papini et al. 2019) with $\kappa \in \{2, 3, 4\}$. Increasing the value of κ leads to coarser discretization and, consequently, reduces the computational complexity at the cost of a degradation of the regret guarantee (that remains sublinear anyway). For clarity, since all versions performed almost the same, in the left plot we report the case of $\kappa = 4$ only. The plot on the right shows the number of points in the grid employed by OPTIMIST2 as a function of the number of iterations. We observe that the number of points soon becomes intractable. This example shows how a discretization-based approach, like OPTIMIST2, despite its theoretical guarantees, does not scale in practice as the number of dimensions increases.

¹³For running OPTIMIST2, we employ the publicly-available official implementation.

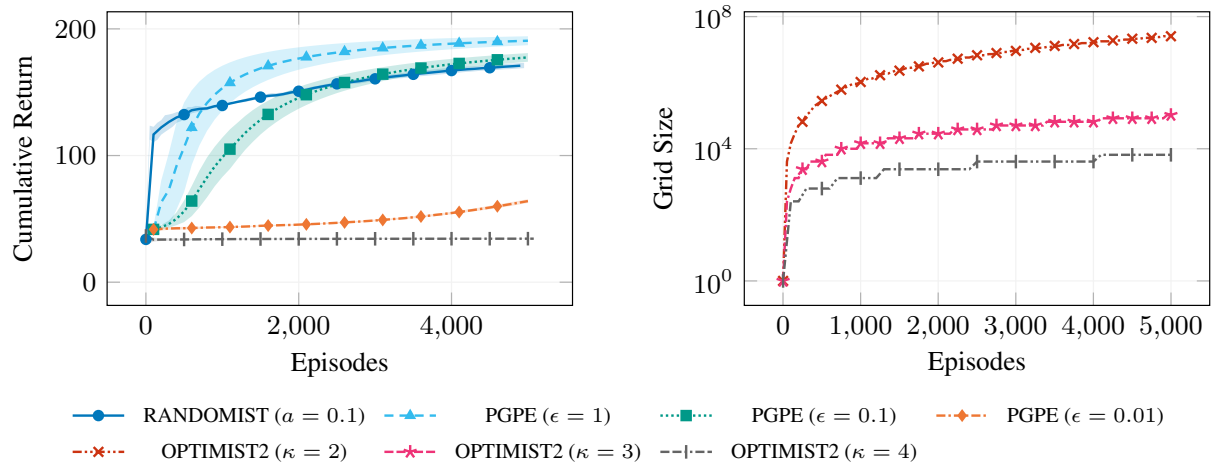


Figure 5: Cumulative Return in the continuous cartpole task (5 runs, 95% c.i.) and number of points on the discretization grid of OPTIMIST2.