# Online Configuration in Continuous Decision Space

**Davide Maran**[*]
Politecnico di Milano
davide.maran@polimi.it

**Pierriccardo Olivieri** [*]
Politecnico di Milano
pierriccardo.olivieri@polimi.it

**Francesco Emanuele Stradi**[*]
Politecnico di Milano
francescoemanuele.stradi@polimi.it

**Giuseppe Urso**
Politecnico di Milano
giuseppe.urso@mail.polimi.it

**Nicola Gatti**
Politecnico di Milano
nicola.gatti@polimi.it

**Marcello Restelli**
Politecnico di Milano
marcello.restelli@polimi.it

## Abstract

In this paper, we investigate the optimal *online configuration* of episodic Markov decision processes when the space of the possible configurations is continuous. Specifically, we study the interaction between a *learner* (referred to as the *configurator*) and an *agent* with a *fixed*, *unknown* policy, when the learner aims to minimize her losses by choosing transition functions in online fashion. The losses may be unrelated to the agent's rewards. This problem applies to many real-world scenarios where the learner seeks to manipulate the Markov decision process to her advantage. We study both *deterministic* and *stochastic* settings, where the losses are either fixed or sampled from an unknown probability distribution. We design two algorithms whose peculiarity is to rely on occupancy measures to explore with optimism the continuous space of transition functions, achieving constant regret in deterministic settings and $\tilde{\mathcal{O}}(\sqrt{T})$ regret in stochastic settings, respectively. Moreover, we prove that the regret bound is tight with respect to any constant factor in deterministic settings. Finally, we compare the empiric performance of our algorithms with a baseline in synthetic experiments.

## 1 Introduction

Reinforcement Learning (RL) investigates the sequential interaction between a learner and an environment, aiming at continually improving the learner's strategy (Sutton and Barto, 2018). In this context, the environment is customarily represented as a Markov Decision Process (MDP) with a fixed but unknown transition function. We study a general scenario where the interaction occurs in episodes, each with a predetermined length. Differently from the standard RL setting, we consider the learner not to be the agent playing the MDP, but the *configurator*. Precisely, at each episode, the learner picks the transition functions for the entire MDP (*i.e.*, a configuration) from a fixed continuous set. Next, she observes the loss suffered and the path traversed by the agent, which depend both on

---

[*]Equal Contribution.

the agent's fixed policy and the transition chosen for the specific episode. The aim of the configurator is to minimize her regret between her total loss and that provided by an optimal fixed configuration.

Our model represents various real-world situations where the learner aims to manipulate the stochastic nature of the MDP to her advantage. For example, consider the sale of hotel rooms, where the MDP states are characterized by the number of rooms booked in different categories each day, while the transitions depend on the hotel's pricing and user behavior. Customarily, the hotels use a fixed pricing strategy that is trained offline and implemented online. Given that users compare prices across hotels before booking rooms, a *competing* hotel (acting as the MDP configurator) can strategically adjust its pricing to influence user behavior and consequently alter the MDP transitions. Specifically, the competitor seeks to reduce the number of room reservations obtained by the agent to maximize her own. Although this example illustrates an adversarial setting, our model applies to general scenarios that do not require a relationship between the configurator's loss and the agent's reward.

## 1.1 Related Work

**Online learning in MDPs** Several works initially introduced for on online learning (Cesa-Bianchi and Lugosi, 2006; Hazan, 2019) have been subsequently extended to MDPs (Auer et al., 2008; Even-Dar et al., 2009; Neu et al., 2010). In particular, Azar et al. (2017) study the problem of optimal exploration in episodic MDPs with unknown transitions and stochastic losses when the feedback is bandit. Rosenberg and Mansour (2019b) study the online learning problem in episodic MDPs with adversarial losses and unknown transitions when the feedback is full information, presenting an online algorithm which provides a regret upper bound of $\tilde{\mathcal{O}}(\sqrt{T})$, where $T$ is the number of episodes. Rosenberg and Mansour (2019a) study the same setting when the feedback is bandit, providing a regret upper bound of $\tilde{\mathcal{O}}(T^{3/4})$, which was subsequently improved to $\tilde{\mathcal{O}}(\sqrt{T})$ by Jin et al. (2019).

**Configurable MDPs** In MDPs, the transition function is customarily assumed to be fixed, see, *e.g.*, Sutton and Barto (2018). However, various subsequent works represent environments with non-fixed transition probabilities, as provided in the works by Satia and Lave (1973), White and Eldeib (1994), and Bueno et al. (2017). Recently, the concept of Configurable Markov Decision Processes (Conf-MDPs) was formalized by Metelli et al. (2018). In particular, the authors propose an algorithm capable of optimizing, at the same time, the environment configuration, namely, the transition function and the policy of the learning agent. This line of research has been further expanded upon by Metelli et al. (2019) and Metelli et al. (2022). Moreover, Ramponi et al. (2021) extend the Conf-MDP setting to an online learning framework. This scenario involves a configurator who chooses online a transition function from a discrete set and aims to maximize her own reward, which is independent from the agent's one.

**Adversarial Attacks** Several works deal with adversarial attacks in MDPs, see, *e.g.*, Ilahi et al. (2021). In the *bounded state attacks* framework, the adversary can manipulate the current state of an MDP in order to force the learning agent to make suboptimal decisions, see, *e.g.*, Pattanaik et al. (2017), Korkmaz (2021), and Wu et al. (2022). Instead, in the *action attacks* setting, the adversary is capable of modifying the agent's actions, see, *e.g.*, Lee et al. (2019), Lee et al. (2021) and Tan et al. (2020). Finally, in the *model attacks* framework, the attack consists in a (bounded) perturbation of the transition function of the MDP performed by an adversary, see, *e.g.*, Rakhsha et al. (2020).

## 1.2 Original Contribution

We investigate the problem of *online configuration with continuous decision space* in MDPs, where the rewards may be both *deterministic* or *stochastic*. Precisely, we study the problem of an online configurator which chooses at any round a transition function from a continuous decision space and receives a loss which depends on both the configuration chosen and the fixed policy of the agent she is interacting with. First, we show that our setting can be seen as an instance of the well-known *Lipschitz bandit* framework, as well as a generalization of many *adversarial attacks* models. Then, we propose two algorithms, namely, O-DOSC (Online Deterministic Optimistic Configuration Search) for deterministic settings and O-SOSC (Online Stochastic Optimistic Configuration Search) for the stochastic ones. We prove that O-DOSC achieves constant regret, matching the lower bound that we provide for the deterministic setting. Then, we show that O-SOSC achieves a $\widetilde{\mathcal{O}}\left(\sqrt{T}\right)$ regret bound in stochastic settings. Finally, we empirically validate our results with synthetic simulations.

## 2 Problem Formulation

### 2.1 Online MDPs

We introduce *online episodic loop-free* MDPs $\mathcal{M} = (X, A, P, \mathcal{R})$ defined as follows.

- $T$ is the number of episodes, with $t \in [T]$ denoting a specific episode.

- $X$ and $A$ are the finite state and action spaces, respectively. By the loop-free property, $X$ is partitioned into $H$ layers $X_0, \ldots, X_H$ such that the first and the last layers are singletons, *i.e.*, $X_0 = \{x_0\}$ and $X_H = \{x_H\}$. We will refer to $H$ as the horizon. Moreover, we denote as $h(x)$ the layer of a specific state $x$.

- $P : X \times A \to \Delta(X)$ is the transition function, where, for ease of notation, we denote by $P(x'|x, a)$ the probability of going from state $x \in X$ to $x' \in X$ by taking action $a \in A$. By the loop-free property, it holds that $P(x'|x, a) > 0$ only if $x' \in X_{h+1}$ and $x \in X_h$ for some $h \in [0 .. H - 1]$.

- $\mathcal{R}$ is the reward function, which can be *deterministic*, that is, $\mathcal{R} : X \times A \to [0, 1]$, or *stochastic*, namely a distribution over $[0, 1]$ for every $(x, a)$. We refer to the reward of a specific state-action pair $x \in X, a \in A$ for a specific episode $t \in [T]$ as $r_t(x, a)$.

**Remark 1.** *Any episodic MDP with horizon $H$ that is* not *loop-free can be cast into a loop-free one by suitably duplicating the state space $H$ times,* i.e., *a state $x$ is mapped to a set of new states $(x, h)$, where $h \in [0 .. H]$.*

A *policy* $\pi : X \to \Delta(A)$ defines a probability distribution over actions at each state. For ease of notation, we denote by $\pi(\cdot|x)$ the probability distribution for a state $x \in X$, with $\pi(a|x)$ denoting the probability of action $a \in A$.

### 2.2 Continuous Configurable-MDPs

The framework we propose, called Continuous Configurable-MDPs, is characterized by:

- an *agent*, which knows the optimal policy $\pi^*$ of a fixed MDP $\mathcal{M}(X, A, \overline{P}, \mathcal{R})$. We assume, without loss of generality, that $\pi^*$ is deterministic, since it is well known that MDPs always admit an optimal deterministic policy;

- a *configurator*, which knows $X, A, \overline{P}, H, T$ and at every episode $t \in [T]$ can choose a configuration (*i.e.*, a transition function) $P_t$ from a bounded set $\mathcal{I}$, in order to minimize her loss $\mathcal{L}$. Similarly to the reward function, the loss function $\mathcal{L}$ can be *deterministic*, that is, $\mathcal{L} : X \times A \to [0, 1]$, or *stochastic*, namely a distribution over $(x, a)$, still bounded in $[0, 1]$. We refer to the loss of a specific state-action pair $x \in X, a \in A$ for a specific episode $t \in [T]$ as $\ell_t(x, a)$.

Customarily in the literature, it is assumed that the configurator's loss is directly tied to the agent's reward, namely $\mathcal{L} = \mathcal{R}$. Instead, in our setting, the two functions can be independent.

In Algorithm 1, we report the interaction between the agent and the configurator in the online MDP.

Precisely, at the beginning of each episode $t$, the loss function is either *deterministically* chosen (although this term may be slightly abused in this context) or *stochastically* chosen (refer to Line 2). Subsequently, the configurator chooses a transition function $P_t$ (as in Line 3), and the MDP is initialized in the state $x_0$ (as per Line 4). During the episode, the agent traverses all the layers based on her policy $\pi^*$ (as described in Line 6) and the transition $P_t$ (as per Line 7). Upon completion of the episode, the configurator observes the complete trajectory and losses (as stated in Line 9).

**Algorithm 1** Agent-Configurator Interaction

1: **for** $t \in [T]$ **do**
2:     $\ell_t$ is chosen *deterministically* or *stochastically*
3:     configurator chooses $P_t \in \mathcal{I}$
4:     state is initialized to $x_0$
5:     **for** $h = 0, \dots, H-1$ **do**
6:         agent plays $a_h \sim \pi^*(\cdot|x_h)$
7:         environment evolves to $x_{h+1} \sim P_t(\cdot|x_h, a_h)$
8:     **end for**
9:     configurator observes $\{x_h, a_h\}_{h=0}^{H-1}$ and suffers $\{\ell_t(x_h, a_h)\}_{h=0}^{H-1}$
10: **end for**

## 2.3 Occupancy Measures

We introduce the notion of *occupancy measure*, see Rosenberg and Mansour (2019a). Given a transition function $P$ and a policy $\pi$, the occupancy measure $d^{P,\pi} \in [0,1]^{|X \times A \times X|}$ induced by $P$ and $\pi$ is such that, for every $x \in X_h$, $a \in A$, and $x' \in X_{h+1}$ with $h \in [0 .. H-1]$:

$$d^{P,\pi}(x, a, x') = \mathbb{P}[x_h = x, a_h = a, x_{h+1} = x' | P, \pi]. \tag{1}$$

Moreover, we also define:

$$d^{P,\pi}(x, a) = \sum_{x' \in X_{h+1}} d^{P,\pi}(x, a, x'), \tag{2}$$

$$d^{P,\pi}(x) = \sum_{a \in A} d^{P,\pi}(x, a). \tag{3}$$

Then, we can introduce the following lemma, which characterizes *valid* occupancy measures.

**Lemma 1** (Rosenberg and Mansour (2019b))**.** *For every $d \in [0,1]^{|X \times A \times X|}$, it holds that $d$ is a valid occupancy measure of an episodic loop-free MDP if and only if, for every $h \in [0 .. H-1]$, the following three conditions hold:*

$$\begin{cases} \sum\limits_{x \in X_h} \sum\limits_{a \in A} \sum\limits_{x' \in X_{h+1}} d(x, a, x') &= 1 \\ \sum\limits_{a \in A} \sum\limits_{x' \in X_{h+1}} d(x, a, x') &= \sum\limits_{x' \in X_{h-1}} \sum\limits_{a \in A} d(x', a, x) \qquad \forall x \in X_h \\ P^d &= P \end{cases}$$

*where $P$ is the transition function of the MDP and $P^d$ is the one induced by $d$ (see Equation* (4)*).*

Notice that any occupancy measure $d$ induces a transition function $P^d$ and a policy $\pi^d$ as:

$$P^d(x'|x, a) = \frac{d(x, a, x')}{d(x, a)}, \qquad \pi^d(a|x) = \frac{d(x, a)}{d(x)}. \tag{4}$$

## 2.4 Performance Metric

In order to have a proper performance metric for our algorithms, we introduce the notion of objective function of an MDP (in terms of loss).

**Definition 1** (Expected Loss)**.** *The expected loss suffered by the configurator at episode $t$ is defined as the expected value of the sum of the losses given the configuration chosen. Namely,*

$$J_t^\pi(P) := \mathbb{E}\left[\sum_{h=1}^{H} \ell_t(x_h, a_h)\Big|\pi, P\right].$$

*By definition of occupancy measure, this can be also written as*

$$J_t^\pi(P) = \sum_{x \in X, a \in A} \ell_t(x, a) d^{P,\pi}(x, a).$$

4

Thus, we define the cumulative regret as follows.

**Definition 2** (Cumulative Regret). *The cumulative regret is defined as*

$$R_T := \sum_{t=1}^{T} J_t^{\pi}(P_t) - J_t^{\pi}(P^*),$$

*where* $P^* := \arg\min_{P \in \mathcal{I}} J_t^{\pi}(P)$.

Following the formulation based on the occupancy measure, the cumulative regret can be written as $R_T := \sum_{t=1}^{T} \ell^{\top} d^{P_t, \pi^*} - \min_{P \in \mathcal{I}} \sum_{t=1}^{T} \ell^{\top} d^{P, \pi^*}$, or equivalently, $R_T := \sum_{t=1}^{T} \ell^{\top} d^{P_t, \pi^*} - \min_{d \in \Delta(\mathcal{I}, \pi^*)} \sum_{t=1}^{T} \ell^{\top} d$, where $d^{P, \pi}$ is the occupancy measure vector defined on the tuple $(x, a)$ given a transition function $P$ and a policy $\pi$, $\Delta(\mathcal{I}, \pi^*)$ is the space of occupancy measures built given the fixed policy $\pi^*$ and the transition function space $\mathcal{I}$, and $\ell$ is defined as:

- in the *deterministic* setting, $\ell$ is the loss vector composed by the loss values associated to each tuple $(x, a)$, namely $\mathcal{L}(x, a)$,

- in the *stochastic* setting, $\ell$, is the vector composed by the expected values of the loss distribution for every $(x, a)$, namely, $\mathbb{E}_{l \sim \mathcal{L}(x,a)}[l]$.

Given the definition of this setting, we aim that the regret is sublinear in $T$, namely $R_T = o(T)$.

The optimization problem described above is linear in the space of the occupancy measures, suggesting the potential adoption of online convex programming tools such as, *e.g.*, Bandit Linear Optimization (BLO) algorithms proposed by Abernethy et al. (2008). However, these methods cannot be adopted to our case. Indeed, without the knowledge of the agent's policy, the configurator cannot compute the exact occupancy measure corresponding to her transition and the agent's policy, thus precluding the design of online bandit linear optimization algorithms working on the occupancy measure space. In particular, the configurator can only choose a transition function $P_t$ and the objective function is highly nonlinear in the space of the transition functions.

## 3 Generality of the Setting and Interpretation

Our model captures various settings. In the following, we provide two different interpretations. The first focuses on MDPs with adversarial attacks, while the second focuses on Lipshitz bandits.

### 3.1 Interpreting Our Model as an MDP with Adversarial Attacks

We show that several forms of *adversarial attacks* in MDPs can be described by our model.

- *Bounded state attacks*. The adversary can modify the agent's state, substituting it with another state that is similar to the original one. This can be modeled by setting:

$$\mathcal{I} = \{P : \forall x \in X, a \in A, \exists x' \in B(x), \ P(\cdot|x, a) = \overline{P}(\cdot|x', a)\},$$

where $B(x) = \{x' : d(x, x') < \varepsilon\}$ for some distance function $d(\cdot)$ and $\varepsilon > 0$.

- *Action attacks*. Differently from the state attack scenarios, the adversary can perturb the action of the agent. This kind of attacks can be modeled by setting:

$$\mathcal{I} = \{P : \forall x \in X, a \in A, \exists a' \in B(a), \ P(\cdot|x, a) = \overline{P}(\cdot|x, a')\},$$

where the set $B(a)$ is defined as in the case of bounded state attacks.

- *Model attacks*. The adversary can change the transition probabilities and the amount of the change is upper bounded according to some metrics. In particular, we adopt the *total variation metric*, denoted with TV. Therefore, $\mathcal{I}$ can be defined as

$$\mathcal{I} = \{P : \text{TV}(P, \overline{P}) < \varepsilon\},$$

for some $\varepsilon > 0$, where $\text{TV}(P, P') := \sum_{x \in X, a \in A} \|P(\cdot|x, a) - P'(\cdot|x, a)\|_1$.

## 3.2 Interpreting Our Model as a Lipschitz Bandit

We can show that the optimization problem faced by the configurator can be seen as a *Lipschitz bandit*, namely, the objective function in the optimization problem is *Lipschitz* continuous.

**Theorem 2.** *Let $P, P'$ be two transition functions, and $\pi$ an arbitrary Markovian policy. Then,*

$$|J^\pi(P) - J^\pi(P')| \leq \frac{H^2}{2} TV(P, P'),$$

*where $TV(P, P') := \sum_{x \in X, a \in A} \|P(\cdot|x, a) - P'(\cdot|x, a)\|_1$.*

Theorem 2 suggests that algorithms for Lipschitz bandits can be used to solve our problem. In the specific case of the Zooming algorithm by Kleinberg et al. (2008)—one of the state-of-the-art algorithms for Lipschitz bandits—, we can derive the following upper regret bound.

**Corollary 3.** *The Zooming algorithm in our setting achieves a regret of $R_T \leq T^{\frac{1+\mathcal{D}(\mathcal{I})}{2+\mathcal{D}(\mathcal{I})}}$, where $\mathcal{D}(\mathcal{I})$ is the Zooming dimension of the space $\mathcal{I}$.*

When the decision space $\mathcal{I}$ depends on a family of $p$ continuous parameters, its Zooming dimension is exactly $p$, so that the regret becomes $T^{\frac{1+p}{2+p}}$. As we show in the following, this regret bound can be dramatically improved and therefore the Zooming algorithm is suboptimal for our problem.

## 4 Deterministic Settings

We focus on deterministic settings, and we present our algorithm and its theoretical guarantees. More precisely, we assume there is a fixed function $\ell : X \times A \to [0, 1]$, such that the configurator will always achieve the same loss whenever the agent chooses a particular action in a given state.

### 4.1 Algorithm

Algorithm 2 provides the pseudo-code of *Online Deterministic Optimistic Configuration Search* (O-DOSC), which tackles deterministic losses. As is customary in the online learning, the configurator needs to face an exploration-exploitation trade-off when searching for the optimal configuration. Specifically, the choice of $P_t$ needs to balance the exploration of unobserved states with the minimization of the configurator's losses.

As stated above, we assume that the optimal policy $\pi^*$ in the MDP is deterministic. Thus, our algorithm can safely keep track of the actions played and losses obtained. For this purpose, two sets are initialized: $\Pi$ contains all possible deterministic policies, while $\widehat{L}$ contains a loss value of $0$ for every tuple $(x, a)$ (Lines 1–2). Such an initialization for the set $\widehat{L}$ is chosen to guarantee optimism vs. uncertainty with respect to the actual loss function.

In order to determine the transition function $P_t$ for each episode, an optimistic approach is adopted. In particular, we minimize the objective over the space of the occupancy measures, which is based on an estimate of the agent's policy (as reported in Line 4). This approach is optimistic with respect to both the policy and the loss function, which is set to be $0$ when non-visited. Additionally, it is possible to simplify the optimization over $\mathcal{I}$ and $\Pi$ by reducing to the optimization over the space $\Delta(\mathcal{I}, \Pi)$, where $d^{P, \pi} \in \Delta(\mathcal{I}, \Pi)$. For a detailed study of the computational complexity of the minimization update, please refer to Appendix A.

Then, once the agent's trajectory and losses suffered throughout the path have been observed (Line 5), the sets are updated as follows. For $\widehat{L}$, the $0$ values associated with the tuples $(x, a)$ visited during the episode are substituted with the observed losses (Line 6). Instead, for $\Pi$, the actions of the state traversed but not executed by the agent are discarded from the set (Line 7).

### 4.2 Upper and Lower Regret Bounds

In this section, we present the theoretical guarantees of our O-DOSC algorithm in deterministic settings. Initially, we state the regret bound achieved by our algorithm, and, subsequently, we show that the regret bound matches the lower bound for our specific setting.

---

**Algorithm 2** O-DOSC Algorithm

---

**Require:** $X, A, H, \mathcal{I}$
1: $\Pi \leftarrow$ set of all deterministic policies
2: $\widehat{L} \leftarrow \{0\}_{\forall (x,a) \in X \times A}$
3: **for** $t \in [T]$ **do**
4:    Choose

$$P_t = \underset{P \in \mathcal{I}, \pi \in \Pi}{\arg \min} \sum_{x,a} \left( \sum_{x' \in X_{h(x)+1}} d^{P,\pi}(x,a,x') \right) \widehat{\ell}(x,a) \quad \text{with} \quad \widehat{\ell}(x,a) \in \widehat{L}$$

5:    Observe $\{x_h, a_h, \ell(x_h, a_h)\}_{h=0}^{H-1}$
6:    $\widehat{L} \leftarrow \{\ell(x_h, a_h)\}_{h=0}^{H-1}$
7:    $\Pi \leftarrow \Pi \setminus \{x_h, a\}_{\forall a \neq a_h, \forall h \in [0\,..\,H-1]}$
8: **end for**

---

In deterministic settings, we show that Algorithm 2 achieves a constant regret bound.

**Theorem 4.** *In deterministic settings, Algorithm 2 guarantees a regret upper bound*

$$R_T \leq (H+1)|X|.$$

The previous result is rather intuitive. Indeed, since both the optimal policy and reward function are deterministic, once the configurator visited the entire MDP, the optimal configuration has been found.

The reader may wonder if the the regret bound shown in Theorem 4 is tight for the setting. In the following, we show that our result is the best *any* algorithm can achieve. Therefore, Algorithm 2 matches the lower bound of the deterministic setting. Indeed, we can show that,

**Theorem 5.** *In deterministic settings, any algorithm achieves a regret of order $\Omega(H|X|)$.*

## 5 Stochastic Settings

We focus on stochastic settings, and we present our algorithm and its theoretical guarantees. Precisely, we assume that there is a fixed probability distribution, denoted as $\mathcal{L}$, which drives the sampling of losses from the interval $[0,1]$ every time the agent chooses an action in a given state.

### 5.1 Algorithm

Algorithm 3 provides the pseudo-code of *Online Stochastic Optimistic Configuration Search* (O-SOSC) for stochastic losses. Similarly to what happens in deterministic settings, the configurator needs to address an exploration-exploitation trade-off when seeking for the optimal configuration. Again, the choice of $P_t$ is required to balance the exploration of non-visited states with the minimization of the configurator's losses. Furthermore, in this case we introduce an additional complexity, given by the way losses are chosen.

By the theory of MDPs, we can safely assume that the optimal policy $\pi^*$ for the MDP is deterministic. Algorithm 3 keeps track of the action played and the losses obtained by the configurator. For this purpose, we initialize two sets: $\Pi$ containing all possible deterministic policies, while $\widehat{L}$ contains a loss value of 0 for every tuple $(x,a)$ (Lines 1–2). We choose this initialization for the set $\widehat{L}$ to be optimistic with respect to the actual loss function.

To determine the transition function $P_t$ for each episode, we take an optimistic approach by minimizing the objective over the space of occupancy measures based on an estimate of the agent's policy (as reported in Line 4). It is worth noting that this update is optimistic with respect to both the policy and the loss function, which is set to 0 when non-visited, and is computed with UCB-like lower bound once traversed. Moreover, it is possible to simplify the optimization over $\mathcal{I}$ and $\Pi$ by reducing it to the optimization over the space $\Delta(\mathcal{I}, \Pi)$, where $d^{P,\pi} \in \Delta(\mathcal{I}, \Pi)$. For a detailed study of the computational complexity of the minimization update, please refer to Appendix A.

Once the agent's trajectory and losses suffered throughout the path have been observed (Line 5), the sets are updated as follows. For $\widehat{L}_t$, the values associated with the tuples $(x, a)$ visited during the episode are updated with a UCB-like term that depends on the number of visits of a specific state $N_t(x)$ (Line 6), which is subtracted to the empirical mean $\bar{\ell}(x, a)$ of the losses observed. For $\Pi$, the actions of the state traversed but not executed by the agent are discarded from the set (Line 7).

---

**Algorithm 3** O-SOSC Algorithm

---

**Require:** $X, A, H, \mathcal{I}, \delta, T$
1: $\Pi \leftarrow$ set of all deterministic policies
2: $\widehat{L}_0 \leftarrow \{0\}_{\forall (x,a) \in X \times A}$
3: **for** $t \in [T]$ **do**
4:     Choose

$$P_t = \underset{P \in \mathcal{I}, \pi \in \Pi}{\arg\min} \sum_{x,a} \left( \sum_{x' \in X_{h(x)+1}} d^{P,\pi}(x, a, x') \right) \widehat{\ell}(x, a) \quad \text{with } \widehat{\ell}(x, a) \in \widehat{L}_t$$

5:     Observe $\{x_h, a_h, \ell_t(x_h, a_h)\}_{h=0}^{H-1}$
6:     $\widehat{L}_{t+1} \leftarrow \left\{ \max\left( 0, \bar{\ell}(x_h, a_h) - \sqrt{\frac{-\log(\delta) + \log(N_t(x_h)(N_t(x_h)+1))}{2N_t(x_h)}} \right) \right\}_{h=0}^{H-1}$
7:     $\Pi \leftarrow \Pi \setminus \{x_h, a\}_{\forall a \neq a_h, \forall h \in [0 \,..\, H-1]}$
8: **end for**

---

## 5.2 Upper Regret Bound

In this section, we present the theoretical result for Algorithm 3. First of all, we can derive a simple lower bound of the regret in the stochastic case. Our model can be seen as a generalization of the multi-armed bandit setting. Specifically, given any multi-armed bandit problem, we can build an equivalent instance of our problem as follows. For every arm of the bandit problem, we have a transition function in $\mathcal{I}$ bringing deterministically from a common initial state to a different state. This implies that the number of transition functions in the MDP equals the number of arms of the bandit problem ($|\mathcal{I}| = |X|$). Therefore, the standard lower bound for multi-armed bandits with $|\mathcal{I}| = |X|$ number of arms leads to a regret of $R_T = \Omega(\sqrt{|X|T})$ which represents a lower bound for our problem. Now, we show that Algorithm 3 achieves a sublinear regret bound.

**Theorem 6.** *In the stochastic setting, for the choice $\delta = T^{-1/2}$, Algorithm 3 achieves a regret upper bounded as follows,*

$$R_T = \widetilde{\mathcal{O}}\left( |X|\sqrt{T} + H|X| \right).$$

We are interested in comparing our theoretical guarantees with the regret bounds of the algorithms available in the literature on online learning for adversarial Markov decision processes. It is well-established that in the online adversarial MDP setting, every algorithm achieves a regret bound of the order of $\Omega(H\sqrt{|X||A|T})$ (Jin et al., 2018). However, the current state-of-the-art result, achieved by Jin et al. (2019), provides a regret bound of $\widetilde{\mathcal{O}}(H|X|\sqrt{|A|T})$, leaving a gap of order $\mathcal{O}(\sqrt{|X|})$ open. In our setting, we observe a similar dependency: our regret bound depends linearly on the number of states, while the multi-armed bandits lower bound suggests that a dependency of order $\mathcal{O}(\sqrt{|X|})$ may be achievable.

## 6 Empirical Evaluation

In this section, we experimentally evaluate the performance of Algorithms 2 and 3 in terms of empiric regret. We describe the results obtained in the deterministic and stochastic settings separately. In each case, we conduct experiments with both discrete and continuous decision spaces $\mathcal{I}$.

As a baseline, we opt for UCB1 (Auer et al., 2002) since, in the case of discrete decision spaces, UCB1 is a standard baseline, while, in the case of continuous decision spaces, UCB1 can be preferred to Zooming (Kleinberg et al., 2008) for two reasons. The first reason is that the design of a suitable

Figure 1: Average cumulative regret with a 95% confidence interval over 10 experiments in deterministic settings with discrete (a, b) and continuous (c) decision spaces.



Figure 2: Average cumulative regret with a 95% confidence interval over 10 experiments in stochastic settings with discrete (a) and continuous (b) decision spaces.

covering oracle for Zooming raises several conceptual and computational issues due to the high number of dimensions whose solutions is open. The second reason is that, in our experimental settings, the optimal solution is one of the arms, and, in these cases, UCB1 is a more severe baseline than Zooming as it guarantees a much better regret bound.

In the following experiments, we consider a Markov decision process structured as follows. The MDP consists of four layers. As is standard in the loop-free model, the first and the last layers are singletons, while the second and third layers each comprise two states. Additionally, every state is associated with two actions. For reasons of space, the description of the experimental settings and additional details on the experimental results can be found in Appendix C.

**Deterministic Settings**   We report in Figure 1 the experimental results obtained with deterministic settings where the cumulative regret is averaged over 10 runs. In particular, Figure 1(a) shows the results with discrete settings, while Figure 1(c) shows the results with continuous settings. In both cases, O-DOSC dramatically outperforms UCB1. Figure 1(b) clearly shows that O-DOSC effectively computes the optimal transition function during the very initial rounds, and subsequently it ceases to explore. Indeed, once O-DOSC visited all the states, it can numerically compute the optimal transition. Instead, UBC1 keeps exploring for a long time.

**Stochastic Settings**   We report in Figure 2 the experimental results obtained with deterministic settings where the cumulative regret is averaged over 10 runs. In particular, Figure 2(a) shows the results with discrete settings, while Figure 2(b) shows the results with continuous settings. In both cases, O-SOSC outperforms UCB1. Differently from what happens in deterministic settings, O-SOSC does not find the optimal solution in the initial rounds, and additional exploration is required. However, the performance exhibited by O-SOSC in this setting is remarkably impressive.

# 7 Conclusions and Future Works

In this paper, we propose the problem of *online configuration* of Markov decision processes with *continuous decision spaces*. We study the problem both when the losses are deterministic and stochastic. We propose O-DOSC algorithm, which achieves constant regret in deterministic settings, and we show that this result is tight with respect to the lower bound. Then, we propose O-SOSC which achieves a sublinear regret bound when the losses are stochastic. Finally, we empirically validate our theoretical results with synthetic simulations.

In future work, we are interested in studying the problem when losses are *adversarial*, namely no statistical assumption are made. Furthermore, we aim to study the problem of *online configurations* against a learning agent, namely, when the policy of the agent is allowed to be dynamic.

# References

Jacob D. Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In *Annual Conference Computational Learning Theory*, 2008.

Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002.

Peter Auer, Thomas Jaksch, and Ronald Ortner. Near-optimal regret bounds for reinforcement learning. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems*, volume 21. Curran Associates, Inc., 2008. URL https://proceedings.neurips.cc/paper/2008/file/e4a6222cdb5b34375400904f03d8e6a5-Paper.pdf.

Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning, 2017. URL https://arxiv.org/abs/1703.05449.

Thiago P. Bueno, Denis D. Mauá, Leliane N. Barros, and Fabio G. Cozman. Modeling Markov decision processes with imprecise probabilities using probabilistic logic programming. In Alessandro Antonucci, Giorgio Corani, Inés Couso, and Sébastien Destercke, editors, *Proceedings of the Tenth International Symposium on Imprecise Probability: Theories and Applications*, volume 62 of *Proceedings of Machine Learning Research*, pages 49–60. PMLR, 10–14 Jul 2017. URL https://proceedings.mlr.press/v62/bueno17a.html.

Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.

Eyal Even-Dar, Sham M Kakade, and Yishay Mansour. Online markov decision processes. *Mathematics of Operations Research*, 34(3):726–736, 2009.

Elad Hazan. Introduction to online convex optimization. *CoRR*, abs/1909.05207, 2019. URL http://arxiv.org/abs/1909.05207.

Inaam Ilahi, Muhammad Usama, Junaid Qadir, Muhammad Umar Janjua, Ala Al-Fuqaha, Dinh Thai Hoang, and Dusit Niyato. Challenges and countermeasures for adversarial attacks on deep reinforcement learning, 2021.

Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably efficient? In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/d3b1fb02964aa64e257f9f26a31f72cf-Paper.pdf.

Chi Jin, Tiancheng Jin, Haipeng Luo, Suvrit Sra, and Tiancheng Yu. Learning adversarial mdps with bandit feedback and unknown transition, 2019. URL https://arxiv.org/abs/1912.01192.

Robert Kleinberg, Aleksandrs Slivkins, and Eli Upfal. Multi-armed bandits in metric spaces, 2008.

Ezgi Korkmaz. Investigating vulnerabilities of deep neural policies, 2021.

Xian Yeow Lee, Sambit Ghadai, Kai Liang Tan, Chinmay Hegde, and Soumik Sarkar. Spatiotemporally constrained action space attacks on deep reinforcement learning agents, 2019.

Xian Yeow Lee, Yasaman Esfandiari, Kai Liang Tan, and Soumik Sarkar. Query-based targeted action-space adversarial policies on deep reinforcement learning agents, 2021.

Alberto Maria Metelli, Mirco Mutti, and Marcello Restelli. Configurable markov decision processes, 2018.

Alberto Maria Metelli, Emanuele Ghelfi, and Marcello Restelli. Reinforcement learning in configurable continuous environments. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 4546–4555. PMLR, 09–15 Jun 2019. URL `https://proceedings.mlr.press/v97/metelli19a.html`.

Alberto Maria Metelli, Guglielmo Manneschi, and Marcello Restelli. Policy space identification in configurable environments. *Machine Learning*, 111(6):2093–2145, 2022. doi: 10.1007/s10994-021-06033-3. URL `https://doi.org/10.1007/s10994-021-06033-3`.

Gergely Neu, Andras Antos, András György, and Csaba Szepesvári. Online markov decision processes under bandit feedback. *Advances in Neural Information Processing Systems*, 23, 2010.

Anay Pattanaik, Zhenyi Tang, Shuijing Liu, Gautham Bommannan, and Girish Chowdhary. Robust deep reinforcement learning with adversarial attacks, 2017.

Amin Rakhsha, Goran Radanovic, Rati Devidze, Xiaojin Zhu, and Adish Singla. Policy teaching via environment poisoning: Training-time adversarial attacks against reinforcement learning, 2020.

Giorgia Ramponi, Alberto Maria Metelli, Alessandro Concetti, and Marcello Restelli. Learning in non-cooperative configurable markov decision processes. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 22808–22821. Curran Associates, Inc., 2021. URL `https://proceedings.neurips.cc/paper_files/paper/2021/file/c0f52c6624ae1359e105c8a5d8cd956a-Paper.pdf`.

Aviv Rosenberg and Yishay Mansour. Online stochastic shortest path with bandit feedback and unknown transition function. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019a. URL `https://proceedings.neurips.cc/paper/2019/file/a0872cc5b5ca4cc25076f3d868e1bdf8-Paper.pdf`.

Aviv Rosenberg and Yishay Mansour. Online convex optimization in adversarial Markov decision processes. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5478–5486. PMLR, 09–15 Jun 2019b. URL `https://proceedings.mlr.press/v97/rosenberg19a.html`.

Jay K. Satia and Roy E. Lave. Markovian decision processes with uncertain transition probabilities. *Operations Research*, 21(3):728–740, 1973. ISSN 0030364X, 15265463. URL `http://www.jstor.org/stable/169381`.

Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA, 2018. ISBN 0262039249.

Kai Liang Tan, Yasaman Esfandiari, Xian Yeow Lee, Aakanksha, and Soumik Sarkar. Robustifying reinforcement learning agents via action space adversarial training, 2020.

Chelsea C. White and Hany K. Eldeib. Markov decision processes with imprecise transition probabilities. *Operations Research*, 42(4):739–749, 1994. ISSN 0030364X, 15265463. URL `http://www.jstor.org/stable/171626`.

Fan Wu, Linyi Li, Zijian Huang, Yevgeniy Vorobeychik, Ding Zhao, and Bo Li. Crop: Certifying robust policies for reinforcement learning through functional smoothing, 2022.

# A Computational Complexity of the Minimization Problem

In this section we study the computational complexity of the minimization update performed by Algorithm 2 and Algorithm 3. Indeed, the optimization problem required to be solved strongly depends on how the decision space of the transition function is chosen beforehand. In the following, we show that the minimization update can be performed in polynomial time when the decision space $\mathcal{I}$ is:

- $\mathcal{I} = \big\{ P : |P(x'|x,a) - \overline{P}(x'|x,a)| \leq \epsilon(x,a,x'), \ \forall (x,a,x') \in X_h \times A \times X_{h+1} \big\}$
- $\mathcal{I} = \big\{ P : ||P(\cdot|x,a) - \overline{P}(\cdot|x,a)||_1 \leq \epsilon(x,a), \ \forall (x,a) \in X \times A \big\}$
- $\mathcal{I}$ is a discrete set.

Thus, we show that the *O-DOSC* and *O-SOSC* optimization problems applied to the previous decision spaces may be modeled as Linear Programs (or a combination of them), which implies that they can be solved in polynomial time.

Precisely, the optimization problem that has to be solved in Algorithm 2 is the following:

$$\underset{P \in \mathcal{I}, \pi \in \Pi}{\arg \min} \sum_{x,a} \left( \sum_{x' \in X_{h(x)+1}} d^{P,\pi}(x,a,x') \right) \widehat{\ell}(x,a)$$

The idea is to optimize on the occupancy space $\Delta(\mathcal{I}, \Pi)$, namely:

$$\underset{d \in \Delta(\mathcal{I}, \Pi)}{\arg \min} \sum_{x,a} \left( \sum_{x' \in X_{h(x)+1}} d(x,a,x') \right) \widehat{\ell}(x,a)$$

As previously stated, the optimization problems can be formulated as different LPs, depending on the choice of the set $\mathcal{I}$. Then the output of the LP, namely $d^*$, allows to compute the probability function $P$ (played by the algorithm) as:

$$P^{d^*}(x'|x,a) = \frac{d^*(x,a,x')}{\sum\limits_{y \in X_{h(x)+1}} d^*(x,a,y)}$$

In the rest of this section we will use the $\forall h$ term to identify $\forall h \in [0, \ldots, H-1]$. We start with the optimization problem for the decision space defined by the module of the difference between transition values for the triple $(x,a,x')$, namely:

- $\mathcal{I} = \big\{ P : |P(x'|x,a) - \overline{P}(x'|x,a)| \leq \epsilon(x,a,x'), \ \forall (x,a,x') \in X_h \times A \times X_{h+1} \big\}$

$$\arg \min \sum_{x,a} \left( \sum_{x' \in X_{h(x)+1}} d(x,a,x') \right) \widehat{\ell}(x,a) \tag{5}$$

s.t.

$$\sum_{x \in X_h, a \in A, x' \in X_{h+1}} d(x,a,x') = 1 \qquad \forall h$$

$$\tag{6}$$

$$\sum_{a \in A, x' \in X_{h+1}} d(x,a,x') = \sum_{x' \in X_{h-1}, a \in A} d(x',a,x) \qquad \forall h, \forall x \in X_h$$

$$\tag{7}$$

$$d(x,a,x') \leq [\overline{P}(x'|x,a) + \epsilon(x,a,x')] \cdot \sum_{y \in X_{h+1}} d(x,a,y) \qquad \forall h, \forall (x,a,x') \in X_h \times A \times X_{h+1}$$

$$\tag{8}$$

$$d(x, a, x') \geq [\overline{P}(x'|x, a) - \epsilon(x, a, x')] \cdot \sum_{y \in X_{h+1}} d(x, a, y) \quad \forall h, \forall (x, a, x') \in X_h \times A \times X_{h+1}$$

$$\tag{9}$$

$$d(x, a, x') \geq 0 \qquad \forall h, \forall (x, a, x') \in X_h \times A \times X_{h+1}$$

$$\tag{10}$$

$$\sum_{\bar{x} \in X_{h(x)-1}} \sum_{a' \in A} d(\bar{x}, a', x) = \sum_{x' \in X_{h(x)+1}} d(x, a, x') \qquad \forall x \in \overline{X}, \forall a \in \Pi(x)$$

$$\tag{11}$$

where $\overline{X}$ is the set of visited states, Constraints (6),(7),(10) define a valid occupancy measure, Constraints (8) and (9) define the space of the transition functions and finally Constraint (11) sets to 1 the probability that actions in $\Pi$ are played. It easy to check the previous optimization problem is a LP, which can be solved in polynomial time.

We then focus on the case where the distance between transition functions for every tuple $(x, a)$ is computed by the $\|\cdot\|_1$-norm, namely:

- $\mathcal{I} = \left\{ P : \|P(\cdot|x, a) - \overline{P}(\cdot|x, a)\|_1 \leq \epsilon(x, a), \ \ \forall (x, a) \in X \times A \right\}$

$$\arg\min \sum_{x, a} \left( \sum_{x' \in X_{h(x)+1}} d(x, a, x') \right) \widehat{\ell}(x, a) \tag{12}$$

s.t.

$$\sum_{x \in X_h, a \in A, x' \in X_{h+1}} d(x, a, x') = 1 \qquad \forall h$$

$$\tag{13}$$

$$\sum_{a \in A, x' \in X_{h+1}} d(x, a, x') = \sum_{x' \in X_{h-1}, a \in A} d(x', a, x) \qquad \forall h, \forall x \in X_h$$

$$\tag{14}$$

$$d(x, a, x') - \overline{P}(x'|x, a) \cdot \sum_{y \in X_{h+1}} d(x, a, y) \leq \epsilon(x, a, x') \quad \forall h, \forall (x, a, x') \in X_h \times A \times X_{h+1}$$

$$\tag{15}$$

$$\overline{P}(x'|x, a) \cdot \sum_{y \in X_{h+1}} d(x, a, y) - d(x, a, x') \leq \epsilon(x, a, x') \quad \forall h, \forall (x, a, x') \in X_h \times A \times X_{h+1}$$

$$\tag{16}$$

$$d(x, a, x') \geq 0 \qquad \forall h, \forall (x, a, x') \in X_h \times A \times X_{h+1}$$

$$\tag{17}$$

$$\sum_{\bar{x} \in X_{h(x)-1}} \sum_{a' \in A} d(\bar{x}, a', x) = \sum_{x' \in X_{h(x)+1}} d(x, a, x') \qquad \forall x \in \overline{X}, \forall a \in \Pi(x)$$

$$\tag{18}$$

$$\sum_{x' \in X_{h+1}} \epsilon(x, a, x') \leq \epsilon(x, a) \cdot \sum_{x' \in X_{h+1}} d(x, a, x') \qquad \forall h, \forall (x, a) \in X_h \times A$$

$$\tag{19}$$

where Constraints (13),(14),(17) define a valid occupancy measure, Constraints (15), (16) and (19) define the space of the transition functions and finally Constraint (18) sets to 1 the probability that actions in $\Pi$ are played. It easy to check the previous optimization problem is a LP, which can be solved in polynomial time.

We conclude the section focusing on the case where transition functions are chosen from a discrete set. We show how the occupancy measure is computed for a fixed transition function $P_i$. Precisely, the occupancy measure can be obtained with a LP formulation, which implies that $|\mathcal{I}|$ LPs must be

solved to obtain the final result. Notice that, since a single LP can be solved in Polynomial time, performing it $|\mathcal{I}|$ times is still polynomial. Moreover, in the discrete case, the value of the occupancy measure given in output by the LP is not necessary; indeed, it is sufficient to obtain the optimal values of the objective function and then to minimize over those values.

- $\mathcal{I}$ is a discrete set

$$\arg\min \sum_{x,a} \left( \sum_{x' \in X_{h(x)+1}} d(x,a,x') \right) \widehat{\ell}(x,a) \tag{20}$$

s.t.

$$\sum_{x \in X_h, a \in A, x' \in X_{h+1}} d(x,a,x') = 1 \qquad\qquad \forall h \tag{21}$$

$$\sum_{a \in A, x' \in X_{h+1}} d(x,a,x') = \sum_{x' \in X_{h-1}, a \in A} d(x',a,x) \qquad\qquad \forall h, \forall x \in X_h \tag{22}$$

$$d(x,a,x') = P_i(x'|x,a) \sum_{y \in X_{h+1}} d(x,a,y) \qquad \forall h, \forall (x,a,x') \in X_h \times A \times X_{h+1} \tag{23}$$

$$\sum_{\bar{x} \in X_{h(x)-1}} \sum_{a' \in A} d(\bar{x},a',x) = \sum_{x' \in X_{h(x)+1}} d(x,a,x') \qquad\qquad \forall x \in \overline{X}, \forall a \in \Pi(x) \tag{24}$$

$$d(x,a,x') \geq 0 \qquad\qquad \forall h, \forall (x,a,x') \in X_h \times A \times X_{h+1} \tag{25}$$

where the meaning of the constraints is similar to the ones of the first LP.

In this section, we have shown that the computational complexity of the minimization problem in the *deterministic* setting, namely for Algorithm 2, is polynomial. Notice that, the same result holds in the *stochastic* setting as well. Specifically, it easy to check that by substituting the loss value with its lower-bound, which is not an optimization variable, the same results as in the *deterministic* setting can be obtained.

# B  Omitted Proofs

In the following, we provide the omitted proof of the theorems presented in the main paper, and the related lemmas. For the sake of clarity we name the following subsections as the main paper sections.

## B.1  Interpreting Our Model as a Lipschitz Bandit

**Theorem 2.** *Let $P, P'$ be two transition functions, and $\pi$ an arbitrary Markovian policy. Then,*

$$|J^\pi(P) - J^\pi(P')| \leq \frac{H^2}{2} TV(P,P'),$$

*where $TV(P,P') := \sum_{x \in X, a \in A} \|P(\cdot|x,a) - P'(\cdot|x,a)\|_1$.*

14

*Proof.* Let us denote as $d_h^{P,\pi}(\cdot) \in \Delta(X)$ the distribution of states of layer $h$ under configuration $P$. We have, for every $h > 1$,

$$\|d_h^{P,\pi}(\cdot) - d_h^{P',\pi}(\cdot)\|_1 = \sum_{x \in X} |d_h^{P,\pi}(x) - d_h^{P',\pi}(x)|$$

$$= \sum_{x \in X} \sum_{x_0 \in X, a \in A} |P(x|x_0,a)\pi(a|x_0)d_{h-1}^{P,\pi}(x_0) - P'(x|x_0,a)\pi(a|x_0)d_{h-1}^{P',\pi}(x_0)|$$

$$= \sum_{x \in X} \sum_{x_0 \in X, a \in A} \pi(a|x_0)|P(x|x_0,a)d_{h-1}^{P,\pi}(x_0) - P'(x|x_0,a)d_{h-1}^{P',\pi}(x_0)|$$

$$\leq \sum_{x \in X} \sum_{x_0 \in X, a \in A} \pi(a|x_0)|P(x|x_0,a)d_{h-1}^{P,\pi}(x_0) - P(x|x_0,a)d_{h-1}^{P',\pi}(x_0)|$$

$$+ \pi(a|x_0)|P(x|x_0,a)d_{h-1}^{P',\pi}(x_0) - P'(x|x_0,a)d_{h-1}^{P',\pi}(x_0)|$$

$$= \sum_{x \in X} \sum_{x_0 \in X, a \in A} \pi(a|x_0)P(x|x_0,a)|d_{h-1}^{P,\pi}(x_0) - d_{h-1}^{P',\pi}(x_0)|$$

$$+ \pi(a|x_0)d_{h-1}^{P',\pi}(x_0)|P(x|x_0,a) - P'(x|x_0,a)|.$$

Here, we can swap the order of the two sums, having, for the first,

$$\sum_{x_0 \in X, a \in A} \pi(a|x_0)|d_{h-1}^{P,\pi}(x_0) - d_{h-1}^{P',\pi}(x_0)| \sum_{x \in X} P(x|x_0,a)$$

$$= \sum_{x_0 \in X, a \in A} \pi(a|x_0)|d_{h-1}^{P,\pi}(x_0) - d_{h-1}^{P',\pi}(x_0)|$$

$$\leq \|d_{h-1}^{P,\pi}(\cdot) - d_{h-1}^{P',\pi}(\cdot)\|_1.$$

and, for the second,

$$\sum_{x_0 \in X, a \in A} \pi(a|x_0)d_{h-1}^{P',\pi}(x_0) \sum_{x \in X} |P(x|x_0,a) - P'(x|x_0,a)|$$

$$= \sum_{x_0 \in X, a \in A} \pi(a|x_0)d_{h-1}^{P',\pi}(x_0)\|P(\cdot|x_0,a) - P'(\cdot|x_0,a)\|_1$$

$$\leq \sum_{x_0 \in X, a \in A} \|P(\cdot|x_0,a) - P'(\cdot|x_0,a)\|_1 = \text{TV}(P,P').$$

It follows that $\|d_h^{P,\pi}(\cdot) - d_h^{P',\pi}(\cdot)\|_1 \leq \|d_{h-1}^{P,\pi}(\cdot) - d_{h-1}^{P',\pi}(\cdot)\|_1 + \text{TV}(P,P')$. Thus, applying the induction, we get

$$\|d_h^{P,\pi}(\cdot) - d_h^{P',\pi}(\cdot)\|_1 \leq h\text{TV}(P,P').$$

Now we focus on the quantity,

$$J^\pi(P) - J^\pi(P') := \sum_{h=1}^{H} \ell(x_h, a_h)\pi(a_h|x_h)d_h^{P,\pi}(x_h).$$

Since the loss is bounded by 1, we get

$$|J^\pi(P) - J^\pi(P')| = \left|\sum_{h=1}^{H} \ell(x_h, a_h)\pi(a_h|x_h)(d_h^{P,\pi}(x_h) - d_h^{P',\pi}(x_h))\right| \leq \sum_{h=1}^{H} \|d_h^{P,\pi}(\cdot) - d_h^{P',\pi}(\cdot)\|_1,$$

which, applying the previous relation, is bounded by

$$\frac{H^2}{2}\text{TV}(P, P').$$

$\square$

## B.2 Deterministic setting

Before being able to prove our main result, let us focus on a simple proposition that will help in the next.

**Proposition 7.** *Let $\pi_1, \pi_2$ be two policies. Then,*

$$TV(d^{P,\pi_1}, d^{P,\pi_2}) \leq H d^{P,\pi_2}(\{\pi_1(x) \neq \pi_2(x)\})$$

*Proof.* Let us suppose $\{\pi_1(x) \neq \pi_2(x)\}$ corresponds to a single state $x_\star$ belonging to layer $h_\star$, in the opposite case we can simply use linearity and sum their visiting distributions. Then,

$$\text{TV}(d^{P,\pi_1}, d^{P,\pi_2}) \leq \sum_{h=1}^{H} \text{TV}(d_h^{P,\pi_1}, d_h^{P,\pi_2})$$

$$= \sum_{h=h_\star}^{H} \text{TV}(d_h^{P,\pi_1}, d_h^{P,\pi_2}).$$

This is true since, for $h < h_\star$ the effect of $x_\star$ is null. In the opposite case, we have

$$\text{TV}(d_h^{P,\pi_1}, d_h^{P,\pi_2}) = \sup_{S \subset X} \mathbb{P}_{\pi_1}(x_h \in S) - \mathbb{P}_{\pi_2}(x_h \in S)$$

$$= \sup_{S \subset X} \mathbb{P}_{\pi_1}(x_h \in S) - \mathbb{P}_{\pi_2}(x_h \in S | x_{h_\star} = x_\star)\mathbb{P}_{\pi_2}(x_{h_\star} = x_\star)$$

$$- \mathbb{P}_{\pi_2}(x_h \in S | x_{h_\star} \neq x_\star)\mathbb{P}_{\pi_2}(x_{h_\star} \neq x_\star),$$

where the last step holds due to the law of total probabilities. Moreover, under the event $\{s_{h_\star} \neq s_\star\}$, the two process are the same, so that

$$\mathbb{P}_{\pi_1}(x_h \in S | x_{h_\star} \neq s_\star)\mathbb{P}_{\pi_2}(x_{h_\star} \neq x_\star) = \mathbb{P}_{\pi_2}(x_h \in S | x_{h_\star} \neq x_\star)\mathbb{P}_{\pi_2}(x_{h_\star} \neq x_\star).$$

This leads to, for all $h \geq h_\star$,

$$\text{TV}(d_h^{P,\pi_1}, d_h^{P,\pi_2}) = \sup_{S \subset X} \mathbb{P}_{\pi_1}(x_h \in S) - \mathbb{P}_{\pi_2}(x_h \in S)$$

$$= \sup_{S \subset X} d^{P,\pi_2}(x_\star)(\mathbb{P}_{\pi_1}(x_h \in S | x_{h_\star} = x_\star) - \mathbb{P}_{\pi_2}(x_h \in S | x_{h_\star} = x_\star))$$

$$\leq d^{P,\pi_2}(x_\star).$$

Summing over $h$ concludes the proof. $\square$

**Theorem 4.** *In deterministic settings, Algorithm 2 guarantees a regret upper bound*

$$R_T \leq (H+1)|X|.$$

*Proof.* Since the policy is fixed and deterministic, the loss in a given state is always the same, and the dependence on the action can be omitted. For this reason we write

$$\ell(x) := \ell(x, \pi(x)).$$

Using algorithm 2, at any timestep we play the configuation $P_t \in \mathcal{I}$ minimizing the following quantity

$$LB_t(P_t) := \min_{P \in \mathcal{I}, \pi \in \Pi} \sum_x d^{P,\pi}(x)\widehat{\ell}(x),$$

16

where $\widehat{\ell}$ is the loss estimated by the algorithm which, due to the determinism of the loss, is always a lower bound for the true loss. This means being optimistic on the actions of the policy in unknown states, assuming they have loss of $0$ (the best possible).

In this way, we have, $LB_t(P) \leq J^\pi(P)$ at any time step $t$ and for any $P \in \mathcal{I}$. From now on, denote as $\overline{X}_t$ the set of unknown state at time $t$. We can underline some crucial facts about the algorithm:

1. If we have visited all the states we play the optimal configuration $P_t = P_\star$

2. Let us call $\varepsilon_t := J^\pi(P_t) - LB_t(P_t)$. We can note that, at any time step $t$, we must have

$$\sum_{x \in \overline{X}_t} d^{P_t, \pi}(x) \geq \frac{\varepsilon_t}{H+1}.$$

Indeed,

$$
\begin{aligned}
J^\pi(P_t) - LB_t(P_t) &= \sum_{x \in X} d^{P_t, \pi}(x)\ell(x) - \sum_{x \in X \setminus \overline{X}_t} d^{P_t, \widehat{\pi}_t}(x)\ell(x) \\
&= \sum_{x \in X} d^{P_t, \pi}(x)\ell(x) - \sum_{x \in X \setminus \overline{X}_t} d^{P_t, \pi}(x)\ell(x) \\
&\quad + \sum_{x \in X \setminus \overline{X}_t} d^{P_t, \pi}(x)\ell(x) - \sum_{x \in X \setminus \overline{X}_t} d^{P_t, \widehat{\pi}}(x)\ell(x) \\
&= \sum_{x \in \overline{X}_t} d^{P_t, \pi}(x)\ell(x) + \sum_{x \in X \setminus \overline{X}_t} (d^{P_t, \pi}(x) - d^{P_t, \widehat{\pi}_t}(x))\ell(x).
\end{aligned}
$$

Since the loss is in $[0,1]$, the first term is bounded by the sum of the visiting distribution of the unknown states

$$\sum_{\overline{X}_t} d^{P_t, \pi}(x)\ell(x) \leq \sum_{\overline{X}_t} d^{P_t, \pi}(x),$$

while the second one is bounded by $\mathrm{TV}(d^{P_t, \pi_1}, d^{P_t, \pi_2})$, again since the loss is in $[0,1]$. Therefore, we can use proposition 7 to have

$$\mathrm{TV}(d^{P_t, \pi_1}, d^{P_t, \pi_2}) \leq H \sum_{\overline{X}_t} d^{P_t, \pi}(x).$$

Therefore, substituting in the previous formula for the lower bound we get

$$J^\pi(P_t) - LB_t(P_t) \leq \sum_{\overline{X}_t} d^{P_t, \pi}(x) + H \sum_{\overline{X}_t} d^{P_t, \pi}(x) = (H+1) \sum_{\overline{X}_t} d^{P_t, \pi}(x),$$

from which

$$\sum_{\overline{X}_t} d^{P_t, \pi}(x) \geq \frac{J^\pi(P_t) - LB_t(P_t)}{H+1} = \frac{\varepsilon_t}{H+1}.$$

17

3. Our regret is bounded by $\mathbb{E}\left[\sum_{t=1}^{T} \varepsilon_t\right]$. Indeed,

$$R_T = \mathbb{E}\left[\sum_{t=1}^{T} J^\pi(P_t) - J^\pi(P_\star)\right]$$

$$\leq \mathbb{E}\left[\sum_{t=1}^{T} J^\pi(P_t) - LB_t(P_t)\right]$$

$$= \mathbb{E}\left[\sum_{t=1}^{T} \varepsilon_t\right].$$

Now, let us define the following sequence of random variables $N_t$ for every $t \in 1, \ldots T$.

$$N_t := \text{number of new states discovered at step } t.$$

With this definition, we can also define the number of states visited up to any time $t$, which corresponds to the size of $\overline{X}_t$ at that time $t$,

$$V_t := \sum_{\tau=1}^{t} N_\tau = |\overline{X}_t|.$$

Also, we will define

$$T_X := \inf\{t \in 1, \ldots T : V_t = |X|\}.$$

With this definitions, we have indeed

$$\sum_{t=1}^{T_X} N_t = |S| \qquad \text{a.s.} \tag{26}$$

Now, recall that, by points $2, 3$ we have

$$R_T \leq \mathbb{E}\left[\sum_{t=1}^{T} \varepsilon_t\right]$$

$$\leq (H+1)\mathbb{E}\left[\sum_{t=1}^{T} \sum_{x \in \overline{X}_t} d^{P_t, \pi}(x)\right]$$

$$= (H+1)\mathbb{E}\left[\sum_{t=1}^{T_X} \sum_{x \in \overline{X}_t} d^{P_t, \pi}(x)\right].$$

Moreover, since the MDP is assumed without loss of generality to be loop free, the quantity $\sum_{x \in \overline{X}_t} d^{P_t, \pi}(x)$ corresponding to the expected time spent in the set $\overline{X}_t$ at time $t$, also corresponds to the expected value of the number of states in $\overline{X}_t$ visited, as no state can be visited multiple times in the same episode. This quantity was called $N_t$ in the previous steps. Therefore,

$$R_T \leq (H+1)\mathbb{E}\left[\sum_{t=1}^{T_X} \sum_{x \in \overline{X}_t} d^{P_t, \pi}(x)\right]$$

$$= (H+1)\mathbb{E}\left[\sum_{t=1}^{T_X} \mathbb{E}[N_t]\right].$$

To conclude, we have only to derive a bound on this quantity based on equation (26). Indeed, we have

$$R_T \leq (H+1)\mathbb{E}\left[\sum_{t=1}^{T_X}\mathbb{E}[N_t]\right]$$

$$= (H+1)\mathbb{E}\left[\sum_{t=1}^{T_X}N_t\right]$$

$$\overset{eq.26}{=} (H+1)|X|.$$

which concludes the proof.

$\square$

**Proof of the lower bound (deterministic setting)** . To prove the lower bound, we propose to use a family of MDPs which is represented in the following figure:



**Theorem 8.** *In deterministic settings, any algorithm achieves a regret of order $\Omega(H|X|)$.*

*Proof.* We use the family of MDPs defined in the previous figure. Formally, the state space is defined in this way

19

1. The first and last layers are trivial.

2. The second layer is made by $N$ states.

3. The layers $h = 3, ...H - 1$ are made by 2 states.

Instead, the action set corresponds to $\{1, 2\}$. The loss is defined as

$$\ell_h(x, a) = \begin{cases} 1 & x = x_{h,1} \ \ h \in \{3, \ldots H - 1\} \\ 0 & \text{otherwise} \end{cases}.$$

This means that the loss is only distributed on the first column (the states of the form $x_{h,1}$ and is constant $+1$). The set $\mathcal{I}$ of possible transition is defined by the set of transitions $P$ satisfiying the following conditions

- $P_h(\cdot|x, a)$ is always deterministic

- $h = 2$ : for any $x \in X_2, a \in A$

$$P_2(x_{3,1}|x, a) = \begin{cases} 1 & a = 1 \\ 0 & a = 2 \end{cases} \qquad P_2(x_{3,2}|x, a) = \begin{cases} 0 & a = 1 \\ 1 & a = 2 \end{cases}.$$

  In other words, at layer 2, the next state is only decided by the action of the agent.

- $h = 3, \ldots H - 2$ for any $x \in X_2, a \in A$, we have

$$P_h(x_{h+1,1}|x, a) = \begin{cases} 1 & x = x_{h,1} \\ 0 & x = x_{h,2} \end{cases} \qquad P_h(x_{h+1,2}|x, a) = \begin{cases} 0 & x = x_{h,1} \\ 1 & x = x_{h,2} \end{cases}.$$

  Roughly speaking, this tells us that after the second layer, the process proceeds on the same vertical line regardless of the action of the agent.

From these two condition, we can see that the only transition that is not fixed is the one from state $x_{1,1}$ to the second layer, which can be arbitrary, until it is deterministic. This means that we, as configurator can choose arbitrarily the second state of the agent. It is then able to choose the state $x_{3,1}$ or $x_{3,2}$, and proceed on all the states $x_{h,1}$ in the first case or $x_{h,2}$ in the former.

By definition $|\mathcal{I}| = N$, since it corresponds to the possible choice if the state in $h = 2$. At this point, we want to show that for any algorithm there is a problem instance (which in this case is given by the agent policy $\pi$, the only element unknown to the configurator) where it cannot achieve expected regret less than $N(H - 3)$.

First, note that by Yao's principle it suffices to show that there exist a distribution over the problem instances such that any *deterministic* algorithm suffers at least $N/2(H - 3)$ regret when the instance that the algorithm runs on is chosen randomly from the distribution. As distribution of instances we simply choose the uniform distribution over the set $\Pi$ of the policies $\pi$ such that

$$\pi_2(2|y) = \begin{cases} 1 & y = x_{2,n} \\ 0 & \text{otherwise} \end{cases} \qquad n \in [1, \ldots N].$$

Of course, this set has exactly cardinality $N$. Indeed, any deterministic algorithm can be viewed as a sequence of permutation of the indexes $1, \ldots N$, which are repeated until loss $0$ is found. In each round where this is not found, the loss is instead $H - 3$. Therefore, by the expected regret of such algorithms can be computed exactly as

$$\mathbb{E}[R_T] \geq (H - 3)\frac{N}{2},$$

since, whichever the permutation choosen, the expected order of a random element is $N/2$.

Now, we can rewrite with the substitution $|X| = N + 2(H - 2)$, which gives

$$\mathbb{E}[R_T] \geq (H-3)\frac{|X| - 2(H-2)}{2} = \frac{H|X| - 2H^2 - 3|X| + 10H - 12}{2}.$$

$\square$

## B.3 Stochastic setting

In this section we focus on the more challenging version where the reward is stochastic. Before the main theorem, we have to prove a minor result.

**Lemma 2.** *Let us consider a sequence of i.i.d. random variables $Y_t$ for $t = 1, \ldots T$ of mean $\mu$ and bounded in $[0, 1]$. For every $\delta > 0$, we have*

$$\mathbb{P}\left(\exists t : \ \overline{|Y_t} - \mu| > \sqrt{\frac{-\log(\delta) + \log(t(t+1))}{2t}}\right) \leq \delta.$$

*where*

$$\bar{Y}_t := \frac{1}{t}\sum_{i=1}^{t} Y_i.$$

*Proof.* By Hoeffding's bound, we have, for every $t$,

$$\mathbb{P}\left(\bar{Y}_t - \mu > \sqrt{\frac{-\log(\delta) + \log(t(t+1))}{2t}}\right) \leq e^{-2t\frac{-\log(\delta)+\log(t(t+1))}{2t}}$$

$$= e^{\log(\delta) - \log(t(t+1))}$$

$$= \frac{\delta}{t(t+1)}.$$

Now, we can just use the union bound:

$$\mathbb{P}\left(\exists t : \ \bar{Y}_t - \mu > \sqrt{\frac{-\log(\delta) + \log(t(t+1))}{2t}}\right) \leq \sum_{t=1}^{T} \mathbb{P}\left(\exists t : \ \bar{Y}_t - \mu > \sqrt{\frac{-\log(\delta) + \log(t(t+1))}{2t}}\right)$$

$$\leq \sum_{t=1}^{T} \frac{\delta}{t(t+1)}$$

$$= \delta.$$

At the same way, we can prove

$$\mathbb{P}\left(\exists t : \ \mu - Y_t > \sqrt{\frac{-\log(\delta) + \log(t(t+1))}{2t}}\right) \leq \delta,$$

The two results being equivalent to the thesis.

$\square$

**Theorem 6.** *In the stochastic setting, for the choice $\delta = T^{-1/2}$, Algorithm 3 achieves a regret upper bounded as follows,*

$$R_T = \widetilde{\mathcal{O}}\left(|X|\sqrt{T} + H|X|\right).$$

*Proof.* Since the policy is fixed and deterministic, the reward in a given state is always the same, and the dependence on the action can be omitted. For this reason we write

$$\ell(x) := \ell(x, \pi(x)).$$

Our algorithm plays (Line 4), at any time $t$, the configuration $P \in \mathcal{I}$ minimizing the following lower bound

$$LB_t(P) = \underset{P \in \mathcal{I}, \pi \in \Pi}{\arg\min} \sum_{x,a} \left( \sum_{x' \in X_{h(x)+1}} d^{P,\pi}(x,a,x') \right) \widehat{\ell}(x,a)$$

$$\widehat{\ell}(x,a) = \max \left( 0, \overline{\ell}(x_h, a_h) - \sqrt{\frac{-\log(\delta) + \log(N_t(x_h)(N_t(x_h)+1))}{2N_t(x_h)}} \right).$$

We will call $P_t, \pi_t$ the couple configuration, policy attaining the minimum.

Define, for every $t = 1, \dots T$,

$$\varepsilon_t := J^\pi(P_t) - LB_t(P_t).$$

**(Part 1)** *Failure probability.*

Let us note

$$E := \left\{ \exists x \in X, t \in [T] : |\overline{\ell}_t(x) - \ell(x)| > \sqrt{\frac{-\log(\delta) + \log(N_t(x)(N_t(x)+1))}{2N_t(x)}} \right\},$$

where $N_t(x)$ denotes the number of visits of state $x$ at time $t$. By lemma 2, we have $\mathbb{P}(E) \leq 2|X|\delta$.

**(Part 2)** *Decomposition of the regret.* Let us suppose at time $t$ we have pulled a suboptimal configuration $P_t$. Assume that we are under the event $E^c$: we have that all lower bounds are respected, so that at any time step $t$, $LB_t(P^*) \leq J^\pi(P^*)$. This fact allows the following inequality

$$R_T = \mathbb{E}\left[ \sum_{t=1}^T J^\pi(P_t) - J^\pi(P_\star) \right]$$

$$\leq \mathbb{E}\left[ \sum_{t=1}^T J^\pi(P_t) - LB_t(P_t) \right]$$

$$= \mathbb{E}\left[ \sum_{t=1}^T \varepsilon_t \right].$$

In this way we have proved that, under the event $E^c$, our regret is bounded by $\sum_{t=1}^T \varepsilon_t$.

**(Part 3)** *From regret to visiting state distribution.* By definition, we have at any time $t$

$$\varepsilon_t = J^\pi(P_t) - LB_t(P_t)$$
$$= \sum_{x \in X} d^{P_t,\pi}(x)\ell(x) - \sum_{x \in X} d^{P_t,\pi}(x)\widehat{\ell}(x)$$
$$= \sum_{x \in X} d^{P_t,\pi}(x)\ell(x) - \sum_{x \in X} d^{P_t,\pi}(x)\widehat{\ell}(x)$$
$$+ \sum_{x \in X} d^{P_t,\pi}(x)\widehat{\ell}(x) - \sum_{x \in X} d^{P_t,\pi}(x)\widehat{\ell}(x)$$
$$= \sum_{x \in X} d^{P_t,\pi}(x)(\ell(x) - \widehat{\ell}(x)) + \sum_{x \in X} (d^{P_t,\pi}(x) - d^{P_t,\pi}(x))\widehat{\ell}(x).$$

22

Note that, under the event $E^c$, we have, for any $x \in X$,

$$\ell(x) - \widehat{\ell}_t(x) = \ell(x) - \bar{\ell}_t(x) + \sqrt{\frac{-\log(\delta) + \log(N_t(x)(N_t(s) + 1))}{2t}}$$

$$\leq \ell(x) - \bar{\ell}_t(x) + 2\sqrt{\frac{-\log(\delta) + \log(N_t(x)(N_t(x) + 1))}{2N_t(x)}}$$

$$= \underbrace{2\sqrt{\frac{-\log(\delta) + \log(N_t(x)(N_t(x) + 1))}{2N_t(x)}}}_{L(N_t(x),\delta)}.$$

This ensures that

$$\sum_{x \in X} d^{P_t,\pi}(x)(\ell(x) - \widehat{\ell}(x)) \leq \sum_{x \in X} d^{P_t,\pi}(x)L(N_t(x), \delta). \qquad (27)$$

About the second term, we can say that it is bounded by $\mathrm{TV}(d^{P,\pi_1}, d^{P,\pi_2})$, since the reward is in $[0, 1]$. Therefore, we can use proposition 7 to have

$$\sum_{x \in X} (d^{P_t,\pi}(x) - d^{P_t,\pi_t}(x))\widehat{\ell}_t(x) \leq \mathrm{TV}(d^{P,\pi}, d^{P,\pi_t}) \leq H \sum_{\overline{X}_t} d^{P_t,\pi}(x),$$

where, as in the previous proofs, $\overline{X}_t$ indicates the set of unknown states at time $t$. If we define the function

$$G(N_t(x)) = \begin{cases} H & N_t(x) = 0 \\ 0 & N_t(x) \geq 1 \end{cases}$$

The previous can be rewritten as

$$\sum_{x \in X} (d^{P_t,\pi}(x) - d^{P_t,\pi_t}(x))\widehat{\ell}_t(x) \leq \sum_X d^{P_t,\pi}(x)G(N_t(x)),$$

which, together with equation (27), gives

$$\varepsilon_t \leq \sum_X d^{P_t,\pi}(x)(L(N_t(x), \delta) + G(N_t(x))).$$

**(Part 4)** *Rewriting the regret.* From the previous results, we have

$$R_T \leq \mathbb{E}\left[\sum_{t=1}^T \varepsilon_t\right]$$

$$\leq \mathbb{E}\left[\sum_{t=1}^T \sum_{x \in X} d^{P_t,\pi}(x)(L(N_t(x), \delta) + G(N_t(x)))\right]$$

$$= \mathbb{E}\left[\sum_{x \in X} \sum_{t=1}^T d^{P_t,\pi}(x)(L(N_t(x), \delta) + G(N_t(x)))\right].$$

Which, noting as $1_{P_t,t}(x)$ the indicator function of state $x$ being visited at step $t$ by configuration $P_t$, can also be written as

$$R_T \leq \mathbb{E}\left[\sum_X \sum_{t=1}^{T} d_{P_t}^{\pi}(x)(L(N_t(x), \delta) + G(N_t(x)))\right]$$

$$= \mathbb{E}\left[\sum_X \sum_{t=1}^{T} \mathbf{1}_{P_t,t}(x)(L(N_t(x), \delta) + G(N_t(x)))\right]. \qquad (28)$$

the last step being valid due to the fact that $\mathbb{E}[\mathbf{1}_{P_t,t}(x)|\mathcal{F}_{t-1}] = d^{P_t,\pi}(x)$, which is true thanks to the loop-free assumption, and the fact that the other two random quantities $N_t(x), P_t$ are $\mathcal{F}_{t-1}-$measurable. Therefore, we need to bound the two sums

$$\sum_{t=1}^{T} \mathbf{1}_{P_t,t}(x) L(N_t(x), \delta) + \sum_{t=1}^{T} \mathbf{1}_{P_t,t}(x) G(N_t(x)).$$

**(Part 5)** *Bounding the two sums.* Due to the fact that $N_t(x) = \sum_{\tau=1}^{t} \mathbf{1}_{P_t,t}(x)$, we have

$$\sum_{t=1}^{T} \mathbf{1}_{P_t,t}(x) L(N_t(x), \delta) \leq \sum_{n=1}^{T} L(n, \delta),$$

$$\sum_{t=1}^{T} \mathbf{1}_{P_t,t}(x) G(N_t(x)) \leq \sum_{n=1}^{T} G(n).$$

The second sum is trivial: by definition of $G$ we get exactly $H$. About the first one we can say that

$$\sum_{n=1}^{T} L(n, \delta) = \sum_{n=1}^{T} 2\sqrt{\frac{-\log(\delta) + \log(n(n+1))}{2n}}$$

$$\leq \sum_{n=1}^{T} 2\sqrt{\frac{-\log(\delta)}{2n}} + 2\sqrt{\frac{\log(n(n+1))}{2n}},$$

by convexity. The first part is

$$\sum_{n=1}^{T} 2\sqrt{\frac{-\log(\delta)}{2n}} = \sqrt{-2\log(\delta)} \sum_{n=1}^{T} \frac{1}{\sqrt{n}} \leq \sqrt{-2\log(\delta)}(1 + 2\sqrt{T}).$$

While the second is

$$\sqrt{2} \sum_{n=1}^{T} \sqrt{\frac{\log(n(n+1))}{n}} \leq \sqrt{2\log(T(T+1))} \sum_{n=1}^{T} \sqrt{\frac{1}{n}}$$

$$\leq 2\sqrt{\log(T+1)} \sum_{n=1}^{T} \sqrt{\frac{1}{n}}$$

$$\leq 2\sqrt{\log(T+1)}(1 + 2\sqrt{T}).$$

Putting all the parts together we have that the sum of all the terms is bounded by

$$H + (2\sqrt{\log(T+1)} + \sqrt{-2\log(\delta)})(1 + 2\sqrt{T}).$$

**(Part 6)** *Final considerations.*

As pointed out, the expected regret is bounded by the expected value of the quantity

$$S_T := \sum_{t=1}^{T} \mathbf{1}_{P_t,t}(x)L(N_t(x),\delta) + \sum_{t=1}^{T} \mathbf{1}_{P_t,t}(x)G(N_t(x)),$$

that was bounded in the previous step. The expected regret is then bounded as follows, for every $\delta > 0$:

1. Under $E$, which has probability $2\delta|X|$, the regret is bounded by $T$.

2. Under $E^c$, by the previous point

$$S_T \leq \sum_{x \in X} H + (2\sqrt{\log(T+1)} + \sqrt{-2\log(\delta)})(1 + 2\sqrt{T})$$
$$\leq |X|H + |X|(2\sqrt{\log(T+1)} + \sqrt{-2\log(\delta)})(1 + 2\sqrt{T}).$$

Therefore, choosing $\delta = T^{-1/2}$, we get

$$R_T = T\mathbb{P}(E) + S_T\mathbb{P}(E^c)$$
$$\leq 2|X|\sqrt{T} + |X|H + |X|(2\sqrt{\log(T+1)} + \sqrt{\log(T)})(1 + 2\sqrt{T}).$$

The final expected regret is then bounded by

$$R_T \leq 2|X|\sqrt{T} + |X|H + |X|(2\sqrt{\log(T+1)} + \sqrt{\log(T)})(1 + 2\sqrt{T}).$$

$\square$

## C  Experiments

For the sake of clarity, we report in the followings additional details on the five instances presented in Figures 1,2. Each instance was tested on the MDP presented in Figure 3. We report the original configuration of the MDP (**config 0**) and all the configurations used for the discrete case in Table 1:

| State | Action | State | Config 0 | Config 1 | Config 2 | Config 3 |
|-------|--------|-------|----------|----------|----------|----------|
| S0 | A1 | S1 | 0.1 | 0.9 | 0.5 | 0.1 |
| S0 | A1 | S2 | 0.9 | 0.1 | 0.5 | 0.9 |
| S0 | A0 | S1 | 0.1 | 0.9 | 0.5 | 1.0 |
| S0 | A0 | S2 | 0.9 | 0.1 | 0.5 | 0.0 |
| S1 | A3 | S3 | 0.1 | 0.9 | 0.5 | 1.0 |
| S1 | A3 | S4 | 0.9 | 0.1 | 0.5 | 0.0 |
| S1 | A2 | S3 | 1.0 | 1.0 | 1.0 | 1.0 |
| S2 | A5 | S3 | 0.1 | 0.9 | 0.5 | 0.1 |
| S2 | A5 | S4 | 0.9 | 0.1 | 0.5 | 0.9 |
| S2 | A4 | S4 | 1.0 | 1.0 | 1.0 | 1.0 |
| S3 | A6 | E | 1.0 | 1.0 | 1.0 | 1.0 |
| S4 | A7 | E | 1.0 | 1.0 | 1.0 | 1.0 |

Table 1: Tabular representation of the transition function for each configuration.

- *Instance of Figure 1a*:
    - number of rounds $T = 1000$
    - number of experiments $Exp = 10$
    - arms $n = 4$
    - transition functions described in Table 1
    - loss vector $\ell = [0.58, 0.42, 0.5, 0.4]$,

- *Instance of Figure 1b*:
    - number of rounds $T = 15$
    - number of experiments $Exp = 10$
    - arms $n = 4$
    - transition functions described in Table 1
    - loss vector $\boldsymbol{\ell} = [0.58, 0.42, 0.5, 0.4]$,

- *Instance of Figure 1c*:
    - number of rounds $T = 1000$
    - number of experiments $Exp = 10$
    - arms $n = 4$
    - transition functions $\mathcal{I} = \left\{ P : ||P(\cdot|x,a) - \overline{P}(\cdot|x,a)||_1 \leq \epsilon(x,a), \ \forall(x,a) \in X \times A \right\}$
    - $\epsilon = 5$
    - loss vector $\boldsymbol{\ell} = [0.58, 0.42, 0.5, 0.4]$,

- *Instance of Figure 2a*:
    - number of rounds $T = 100000$
    - number of experiments $Exp = 10$
    - arms $n = 4$
    - transition functions described in 1
    - mean loss vector $\boldsymbol{\ell} = [0.58, 0.42, 0.5, 0.4]$,
    - unitary variance for each arm

- *Instance of Figure 2b*:
    - number of rounds $T = 100000$
    - number of experiments $Exp = 10$
    - arms $n = 4$
    - transition functions $\mathcal{I} = \left\{ P : ||P(\cdot|x,a) - \overline{P}(\cdot|x,a)||_1 \leq \epsilon(x,a), \ \forall(x,a) \in X \times A \right\}$
    - $\epsilon = 5$
    - mean loss vector $\boldsymbol{\ell} = [0.58, 0.42, 0.5, 0.4]$,
    - unitary variance for each arm

**Training Details** In the main paper we have presented five experiments, each corresponding to a different setting. Each experiment is performed with a fixed random seed. The computational time for one experiment depends on the setting. We run the experiments of each setting in parallel with a total computational time of approximately 12 hours.

**Compute** We run the numerical simulations on a server with the following specifications:

- CPU: `128x Intel(R) Xeon(R) Platinum 8358 CPU @ 2.60GHz`
- RAM: `512,0 GB`
- Operating system: `Ubuntu 20.04.5 LTS`
- System type: `64 bit`

**Reproducibility** We have performed every experiment with a fixed seed. The seed influences the loss generation by the environment and the transitions to the next states.

Figure 3: Graphical representation of the MDP used for our experiments