
Near Optimal Exploration-Exploitation in Non-Communicating Markov Decision Processes

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 While designing the state space of an MDP, it is common to include states that are
2 transient or not reachable by any policy (e.g., in mountain car, the product space
3 of speed and position contains configurations that are not physically reachable).
4 This leads to defining weakly-communicating or multi-chain MDPs. In this paper,
5 we introduce TUCRL, the first algorithm able to perform efficient exploration-
6 exploitation in any finite Markov Decision Process (MDP) without requiring any
7 form of prior knowledge. In particular, for any MDP with S^c communicating
8 states, A actions and $\Gamma^c \leq S^c$ possible communicating next states, we derive a
9 $\tilde{O}(D^c \sqrt{\Gamma^c S^c A T})$ regret bound, where D^c is the diameter (i.e., the longest shortest
10 path) of the communicating part of the MDP. This is in contrast with optimistic
11 algorithms (e.g., UCRL, Optimistic PSRL) that suffer linear regret in weakly-
12 communicating MDPs, as well as posterior sampling or regularised algorithms
13 (e.g., REGAL), which require prior knowledge on the bias span of the optimal policy
14 to bias the exploration to achieve sub-linear regret. We also prove that in weakly-
15 communicating MDPs, no algorithm can ever achieve a logarithmic growth of the
16 regret without first suffering a linear regret for a number of steps that is exponential
17 in the parameters of the MDP. Finally, we report numerical simulations supporting
18 our theoretical findings and showing how TUCRL overcomes the limitations of the
19 state-of-the-art.

20 1 Introduction

21 Reinforcement learning (RL) [1] studies the problem of learning in sequential decision-making
22 problems where the dynamics of the environment is unknown, but can be learnt by performing
23 actions and observing their outcome in an online fashion. A sample-efficient RL agent must trade
24 off the *exploration* needed to collect information about the environment, and the *exploitation* of
25 the experience gathered so far to gain as much reward as possible. In this paper, we focus on the
26 regret framework in *infinite-horizon average-reward* problems [2], where the exploration-exploitation
27 performance is evaluated by comparing the rewards accumulated by the learning agent and an optimal
28 policy. Jaksch et al. [2] showed that it is possible to efficiently solve the exploration-exploitation
29 dilemma using the *optimism in face of uncertainty* (OFU) principle. OFU methods build confidence
30 intervals on the dynamics and reward (i.e., construct a set of plausible MDPs), and execute the optimal
31 policy of the “best” MDP in the confidence region [e.g., 2, 3, 4, 5, 6]. An alternative approach is
32 posterior sampling (PS) [7], which maintains a posterior distribution over MDPs and, at each step,
33 samples an MDP and executes the corresponding optimal policy [e.g., 8, 9, 10, 11, 12].

34 **Weakly-communicating MDPs and misspecified states.** One of the main limitations of UCRL [2]
35 and optimistic PSRL [12] is that they require the MDP to be communicating so that its diameter
36 D (i.e., the longest shortest path) is finite. While assuming that all states are somehow reachable
37 may seem a reasonable assumption, it is rarely verified *in practice*. In fact, it requires a designer to
38 carefully define a state space \mathcal{S} that contains all reachable states (otherwise it may not be possible to

39 learn the optimal policy), but it excludes unreachable states (otherwise the resulting MDP would be
 40 non-communicating). This requires a considerable amount of prior knowledge about the environment.
 41 Consider a problem where we learn from images e.g., Breakout. The state space is the set of
 42 “plausible” configurations of the brick wall, ball and paddle positions. The situation in which the wall
 43 has an hole in the middle is a valid state (e.g., as an initial state) but it cannot be observed/reached
 44 starting from a dense wall. While it may be possible to design a suitable set of “reachable” states that
 45 define a communicating MDP, this is often a difficult and tedious task, sometimes even impossible.
 46 Whenever the state space is *misspecified* or the MDP is weakly communicating (i.e., $D = +\infty$),
 47 OFU-based algorithms (e.g., UCRL) optimistically attribute large reward and non-zero probability
 48 to reach states that have never been observed, and thus they tend to repeatedly attempt to *explore*
 49 unreachable states. This results in poor performance and linear regret. A first attempt to overcome
 50 this major limitation is REGAL.C [3] (Fruit et al. [6] recently proposed SCAL, an implementable
 51 efficient version of REGAL.C), which requires prior knowledge of an upper-bound H to the span (i.e.,
 52 range) of the optimal bias function h^* . The optimism of UCRL is then “constrained” to policies
 53 whose bias has span smaller than H . This implicitly “removes” non-reachable states, whose large
 54 optimistic reward would cause the span to become too large. Unfortunately, an accurate knowledge
 55 of the bias span may not be easier to obtain than designing a well-specified state space. Bartlett and
 56 Tewari [3] proposed an alternative algorithm – REGAL.D – that leverages the *doubling trick* to avoid
 57 any prior knowledge on the span. Nonetheless, we recently noticed a major flaw in the proof of [3,
 58 Theorem 3] that question the validity of the algorithm (see App. A for further details). PS-based
 59 algorithms also suffer from similar issues.¹ To the best of our knowledge, the only regret guarantees
 60 available in the literature for this setting are [14, 15, 16]. However, the counter-example of Osband
 61 and Roy [17] seems to invalidate the result of Abbasi-Yadkori and Szepesvári [14]. On the other
 62 hand, Ouyang et al. [15] and Theodorou et al. [16] present PS algorithms with expected *Bayesian*
 63 regret scaling linearly with H , where H is an upper-bound on the optimal bias spans of all the MDPs
 64 that can be drawn from the prior distribution ([15, Asm. 1] and [16, Sec. 5]). In [15, Remark 1], the
 65 authors claim that their algorithm does not require the knowledge of H to derive the regret bound.
 66 However, in App. B we show on a very simple example that for most continuous prior distributions
 67 (e.g., uninformative priors like Dirichlet), it is very likely that $H = +\infty$ implying that the regret
 68 bound may not hold (similarly for [16]). As a result, similarly to REGAL.D, the prior distribution
 69 should contain prior knowledge on the bias span to avoid poor performance.

70 In this paper, we present TUCRL, an algorithm designed to trade-off exploration and exploitation in
 71 weakly-communicating and multi-chain MDPs (e.g., MDPs with misspecified states) without any
 72 prior knowledge and under the only assumption that the agent starts from a state in a communicating
 73 subset of the MDP (Sec. 3). In communicating MDPs, TUCRL eventually (after a finite number
 74 of steps) performs as UCRL, thus achieving logarithmic regret. When the true MDP is weakly-
 75 communicating, we prove that TUCRL achieves a $\tilde{O}(\sqrt{T})$ regret that is polynomial in the MDP
 76 parameters. We also show that it is not possible to design an algorithm achieving logarithmic regret
 77 in weakly-communicating MDPs without having an exponential dependence on the MDP parameters
 78 (see Sec. 5). TUCRL is the first computationally tractable algorithm in the OFU literature that is able
 79 to adapt to the MDP nature without any prior knowledge. The theoretical findings are supported by
 80 experiments on several domains (see Sec. 4).

81 2 Preliminaries

82 We consider a finite *weakly-communicating* Markov decision process [18, Sec. 8.3] $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$
 83 with a set of states \mathcal{S} and a set of actions $\mathcal{A} = \bigcup_{s \in \mathcal{S}} \mathcal{A}_s$. Each state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}_s$
 84 is characterized by a reward distribution with mean $r(s, a)$ and support in $[0, r_{\max}]$ as well as a
 85 transition probability distribution $p(\cdot | s, a)$ over next states. In a weakly-communicating MDP, the set
 86 \mathcal{S} can be decomposed into two subsets: a *communicating* set (denoted by \mathcal{S}^c) in which for any pair
 87 of states s, s' there exists a policy that has a non-zero probability to reach s' starting from s , and a
 88 set of states that are *transient* under all policies (denoted by \mathcal{S}^T). We denote by $S = |\mathcal{S}|$, $S^c = |\mathcal{S}^c|$
 89 and $A = \max_{s \in \mathcal{S}} |\mathcal{A}_s|$ the number of states and actions, and by $\Gamma^c = \max_{s \in \mathcal{S}^c, a \in \mathcal{A}} \|p(\cdot | s, a)\|_0$ the
 90 maximum support of all transition probabilities $p(\cdot | s, a)$ with $s \in \mathcal{S}^c$. The sets \mathcal{S}^c and \mathcal{S}^T form a
 91 partition of \mathcal{S} i.e., $\mathcal{S}^c \cap \mathcal{S}^T = \emptyset$ and $\mathcal{S}^c \cup \mathcal{S}^T = \mathcal{S}$. A deterministic policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ maps states to

¹We notice that the problem of weakly-communicating MDPs and misspecified states does not hold in the
 more restrictive setting of finite horizon [e.g., 8] since exploration is directly tailored to the states that are
 reachable *within* the known horizon, or under the assumption of the existence of a recurrent state [e.g., 13].

92 actions and it has an associated *long-term average reward* (or *gain*) and a *bias function* defined as

$$g_M^\pi(s) := \lim_{T \rightarrow \infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T r(s_t, \pi(s_t)) \right]; \quad h_M^\pi(s) := C\text{-}\lim_{T \rightarrow \infty} \mathbb{E} \left[\sum_{t=1}^T (r(s_t, \pi(s_t)) - g_M^\pi(s_t)) \right],$$

93 where the bias $h_M^\pi(s)$ measures the expected total difference between the rewards accumulated by
 94 π starting from s and the stationary reward in *Cesaro-limit*² (denoted $C\text{-}\lim$). Accordingly, the
 95 difference of bias values $h_M^\pi(s) - h_M^\pi(s')$ quantifies the (dis-)advantage of starting in state s rather
 96 than s' . In the following, we drop the dependency on M whenever clear from the context and
 97 denote by $sp_S \{h^\pi\} := \max_{s \in \mathcal{S}} h^\pi(s) - \min_{s \in \mathcal{S}} h^\pi(s)$ the *span* of the bias function. In weakly
 98 communicating MDPs, any optimal policy $\pi^* \in \arg \max_{\pi} g^\pi(s)$ has *constant gain*, i.e., $g^{\pi^*}(s) = g^*$
 99 for all $s \in \mathcal{S}$. Finally, we denote by D , resp. D^C , the diameter of M , resp. the diameter of the
 100 communicating part of M (i.e., restricted to the set \mathcal{S}^C):

$$D := \max_{(s,s') \in \mathcal{S} \times \mathcal{S}, s \neq s'} \{\tau_M(s \rightarrow s')\}, \quad D^C := \max_{(s,s') \in \mathcal{S}^C \times \mathcal{S}^C, s \neq s'} \{\tau_M(s \rightarrow s')\}, \quad (1)$$

101 where $\tau_M(s \rightarrow s')$ is the expected time of the shortest path from s to s' in M .

102 **Learning problem.** Let M^* be the true (*unknown*) weakly-communicating MDP. We consider the
 103 learning problem where \mathcal{S} , \mathcal{A} and r_{\max} are *known*, while sets \mathcal{S}^C and \mathcal{S}^T , rewards r and transition
 104 probabilities p are *unknown* and need to be estimated on-line. We evaluate the performance of a
 105 learning algorithm \mathfrak{A} after T time steps by its cumulative *regret* $\Delta(\mathfrak{A}, T) = Tg^* - \sum_{t=1}^T r_t(s_t, a_t)$.
 106 Furthermore, we state the following assumption.

107 **Assumption 1.** *The initial state s_1 belongs to the communicating subset of states, i.e., $s_1 \in \mathcal{S}^C$.*

108 While this assumption somehow restricts the scenario we consider, it is fairly common in practice.
 109 For example, all the domains that are characterized by the presence of a resetting distribution (e.g.,
 110 episodic problems) satisfy this assumption (e.g., mountain car, cart pole, Atari games, taxi, etc.).

111 **Multi-chain MDPs.** While we consider weakly-communicating MDPs for ease of notation, all our
 112 results extend to the more general case of multi-chain MDPs.³ In this case, there may be multiple
 113 communicating and transient sets of states and the optimal gain g^* is different in each communicating
 114 subset. Under Asm. 1, we assume the learning agent starts from a state s_1 belonging to a specific
 115 communicating subset, which we denote by \mathcal{S}^C , while the set \mathcal{S}^T contains all other states, either
 116 transient or belonging to other (non-reachable from \mathcal{S}^C) communicating parts of the MDP. As a result,
 117 the regret is still defined as before, where the learning performance is compared to the optimal gain
 118 g^* related to the communicating part in $\mathcal{S}^C \ni s_1$.

119 3 Truncated Upper-Confidence for Reinforcement Learning

120 In this section we introduce TUCRL (Fig. 1), an optimistic online RL algorithm that efficiently
 121 balances exploration and exploitation to learn in non-communicating MDPs without prior knowledge.

122 Similar to UCRL, at the beginning of each episode k , TUCRL constructs confidence intervals for the
 123 reward and the dynamics of the MDP. Formally, for any $(s, a) \in \mathcal{S} \times \mathcal{A}$ we define

$$B_{p,k}(s, a) = \left\{ \tilde{p}(\cdot | s, a) \in \mathcal{C} : \forall s' \in \mathcal{S}, |\tilde{p}(s' | s, a) - \hat{p}(s' | s, a)| \leq \beta_{p,k}^{sas'} \right\}, \quad (2)$$

$$B_{r,k}(s, a) := [\hat{r}_k(s, a) - \beta_{r,k}^{sa}, \hat{r}_k(s, a) + \beta_{r,k}^{sa}] \cap [0, r_{\max}], \quad (3)$$

124 where $\mathcal{C} = \{p \in \mathbb{R}^{\mathcal{S}} | \forall s, p(s) \geq 0 \wedge \sum_s p(s) = 1\}$ is the $(S - 1)$ -probability simplex, while the
 125 size of the confidence intervals is constructed using the empirical Bernstein's inequality [19, 20] as

$$\beta_{r,k}^{sa} := \sqrt{\frac{14\hat{\sigma}_{r,k}^2(s, a)b_{k,\delta}}{N_k^+(s, a)}} + \frac{49}{3}r_{\max}b_{k,\delta}, \quad \beta_{p,k}^{sas'} := \sqrt{\frac{14\hat{\sigma}_{p,k}^2(s' | s, a)b_{k,\delta}}{N_k^+(s, a)}} + \frac{49}{3}b_{k,\delta}$$

126 where $N_k(s, a)$ is the number of visits in (s, a) before episode k , $N_k^+(s, a) := \max\{1, N_k(s, a)\}$,
 127 $N_k^\pm(s, a) := \max\{1, N_k(s, a) - 1\}$, $\hat{\sigma}_{r,k}^2(s, a)$ and $\hat{\sigma}_{p,k}^2(s' | s, a)$ are the empirical variances of $r(s, a)$

²For policies with an aperiodic chain, the standard limit exists.

³In the case of misspecified states, we implicitly define a multi-chain MDP, where each non-reachable state has a self-loop dynamics and it defines a ‘‘singleton’’ communicating subset.

| |
|--|
| <p>Input: Confidence $\delta \in]0, 1[$, r_{\max}, \mathcal{S}, \mathcal{A}</p> <p>Initialization: Set $N_0(s, a) := 0$ for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, $t := 1$ and observe s_1.</p> <p>For episodes $k = 1, 2, \dots$ do</p> <ol style="list-style-type: none"> 1. Set $t_k = t$ and episode counters $\nu_k(s, a) = 0$ 2. Compute estimates $\hat{p}_k(s' s, a)$, $\hat{r}_k(s, a)$ and a set $\mathcal{M}_k^\circ = \{\mathcal{M}_k, \text{ if } \mathcal{S}_k^T = \emptyset; \overline{\mathcal{M}}_k^+ \text{ otherwise}\}$. 3. Compute an $r_{\max}/\sqrt{t_k}$-approximation $\tilde{\pi}_k$ of $(\overline{\mathcal{M}}_k, \tilde{\pi}_k) = \arg \max_{M \in \mathcal{M}_k^\circ, \pi \in \Pi^{\text{SD}}} \{g_M^\pi\}$ 4. While $t_k == t$ or $(\sum_{a \in \mathcal{A}_{s_t}} N_k(s_t, a) > 0 \text{ and } \nu_k(s_t, \tilde{\pi}_k(s_t)) \leq \max\{1, N_k(s_t, \tilde{\pi}_k(s_t))\})$ do <ol style="list-style-type: none"> (a) Execute $a_t = \tilde{\pi}_k(s_t)$, obtain reward r_t, and observe s_{t+1} (b) Set $\nu_k(s_t, a_t) += 1$ and set $t += 1$ 5. Set $N_{k+1}(s, a) = N_k(s, a) + \nu_k(s, a)$ |
|--|

Figure 1: TUCRL algorithm.

128 and $p(s'|s, a)$ and $b_{k,\delta} = \ln(2SA t_k/\delta)$. The set of plausible MDPs associated with the confidence
129 intervals is then $\mathcal{M}_k = \{M = (\mathcal{S}, \mathcal{A}, \tilde{r}, \tilde{p}) : \tilde{r}(s, a) \in B_{r,k}(s, a), \tilde{p}(\cdot|s, a) \in B_{p,k}(s, a)\}$. UCRL
130 is optimistic w.r.t. the confidence intervals so that for all states s that have never been visited the
131 optimistic reward $\tilde{r}(s, a)$ is set to r_{\max} , while all transitions to s (i.e., $\tilde{p}(s|\cdot, \cdot)$) are set to the largest
132 value compatible with $B_{p,k}(\cdot, \cdot)$. Unfortunately, some of the states with $N_k(s, a) = 0$ may be actually
133 unreachable (i.e., $s \in \mathcal{S}^T$) and UCRL would uniformly explore the policy space with the hope that at
134 least one policy reaches those (optimistically desirable) states. TUCRL addresses this issue by first
135 constructing empirical estimates of \mathcal{S}^C and \mathcal{S}^T (i.e., the set of communicating and transient states in
136 M^*) using the states that have been visited so far, that is $\mathcal{S}_k^C := \{s \in \mathcal{S} \mid \sum_{a \in \mathcal{A}_s} N_k(s, a) > 0\} \cup$
137 $\{s_{t_k}\}$ and $\mathcal{S}_k^T := \mathcal{S} \setminus \mathcal{S}_k^C$, where t_k is the starting time of episode k .

138 In order to avoid optimistic exploration to unreachable states, we could simply execute UCRL on \mathcal{S}_k^C ,
139 which is guaranteed to contain only states in the communicating set (since $s_1 \in \mathcal{S}^C$ by Asm. 1, we
140 have that $\mathcal{S}_k^C \subseteq \mathcal{S}^C$). Nonetheless, this algorithm could *under-explore* state-action pairs that would
141 allow discovering other states in \mathcal{S}^C , thus getting stuck in a subset of the communicating states of the
142 MDP and suffering linear regret. While the states in \mathcal{S}_k^C are guaranteed to be in the communicating
143 subset, it is not possible to know whether states in \mathcal{S}_k^T are actually reachable from \mathcal{S}_k^C or not. Then
144 TUCRL first “guesses” a lower bound on the probability of transition from states $s \in \mathcal{S}_k^C$ to $s' \in \mathcal{S}_k^T$
145 and whenever the maximum probability in the confidence interval $\hat{p}_k(s'|s, a) + \beta_{p,k}^{sas'}$ is below the
146 lower bound, it assumes that such transition is actually not possible. The intuition behind this strategy
147 is that either the transition does not exist or it has a sufficiently “big” mass. However, these transitions
148 should be periodically reconsidered in order to avoid *under-exploration* issues. More formally, let
149 $(\rho_t)_{t \in \mathbb{N}}$ be a non-increasing sequence to be defined later, for all $s' \in \mathcal{S}_k^T$, $s \in \mathcal{S}_k^C$ and $a \in \mathcal{A}_s$, the
150 empirical mean $\hat{p}_k(s'|s, a)$ and variance $\hat{\sigma}_{p,k}^2(s'|s, a)$ are zero (i.e., this transition has never been
151 observed so far), so the largest probability (most optimistic) of transition from s to s' through a is
152 $\hat{p}_k^+(s'|s, a) = \frac{49}{3} \frac{b_{k,\delta}}{N_k^\pm(s, a)}$. TUCRL compares $\hat{p}_k^+(s'|s, a)$ to ρ_{t_k} and forces all transition probabilities
153 below the threshold to zero, while the confidence intervals of transitions to states that have already
154 been explored (i.e., in \mathcal{S}_k^C) are preserved unchanged. This corresponds to constructing the alternative
155 confidence interval

$$\overline{B}_{p,k}(s, a) = B_{p,k}(s, a) \cap \{\tilde{p}(\cdot|s, a) \in \mathcal{C} : \forall s' \in \mathcal{S}_k^T \text{ and } \tilde{p}_k^+(s'|s, a) < \rho_{t_k}, \tilde{p}(s'|s, a) = 0\}. \quad (4)$$

156 Given $\overline{B}_{p,k}$, TUCRL (implicitly) constructs the corresponding set of plausible MDPs $\overline{\mathcal{M}}_k$ and then
157 solves the optimistic optimization problem $(\overline{\mathcal{M}}_k, \tilde{\pi}_k) = \arg \max_{M \in \overline{\mathcal{M}}_k, \pi} \{g_M^\pi\}$.

158 In practice, we set $\rho_t = \frac{49b_{t,\delta}}{3} \sqrt{\frac{SA}{t}}$, so that the condition to remove transition reduces to $N_k^\pm(s, a) >$
159 $\sqrt{t_k/SA}$. This shows that only transitions from state-action pairs that have been poorly visited so far
160 are enabled, while if the state-action pair has already been tried often and yet no transition to $s' \in \mathcal{S}_k^T$
161 is observed, then it is assumed that s' is not reachable from s, a . We denote the set of state-action
162 pairs that are non-sufficiently explored by $\mathcal{K}_k = \{(s, a) \in \mathcal{S}_k^C \times \mathcal{A} : N_k^\pm(s, a) \leq \sqrt{t_k/SA}\}$.

163 For technical reasons (Sec. 3.1), we consider a *relaxation of the optimization problem* in which $\overline{\mathcal{M}}_k$ is
164 replaced by a relaxed extended MDP $\overline{\mathcal{M}}_k^+ \supseteq \overline{\mathcal{M}}_k$ defined by using ℓ_1 -norm concentration inequalities

165 for $p(\cdot|s, a)$. Let $B_{p,k}^+(s, a) = \{\tilde{p}(\cdot|s, a) \in \mathcal{C} : \|\tilde{p}(\cdot|s, a) - \hat{p}(\cdot|s, a)\|_1 \leq \sum_{s'} \beta_{p,k}^{sas'}\}$ (resp. $\overline{B}_{p,k}^+$) be
166 the relaxed confidence interval, then \mathcal{M}_k^+ (resp. $\overline{\mathcal{M}}_k^+$) is the corresponding (relaxed) set of plausible
167 MDPs. The optimistic optimization problem is solved by running extended value iteration (EVI) on
168 $\overline{\mathcal{M}}_k^+$ (up to accuracy $\epsilon_k = r_{\max}/\sqrt{t_k}$). Technically, we restrict EVI to work on the set of states $\mathcal{S}_k^{\text{EVI}}$
169 that are optimistically reachable from the communicating set \mathcal{S}_k^{C} . In practice, $\mathcal{S}_k^{\text{EVI}} = \mathcal{S}_k^{\text{C}}$ when $\mathcal{K}_k =$
170 \emptyset since all the transitions to \mathcal{S}_k^{T} are forbidden, otherwise $\mathcal{S}_k^{\text{EVI}} = \mathcal{S}$ (see Alg. 1 in App. C). Starting
171 from an initial vector $u_0 = 0$, EVI iteratively applies the optimal Bellman operator $\tilde{L}_{\overline{\mathcal{M}}_k^+}$ associated
172 to the (extended) MDP $\overline{\mathcal{M}}_k^+$ defined as: $\tilde{L}_{\overline{\mathcal{M}}_k^+} v(s) := \max_{a \in \mathcal{A}_s} \{ \max_{\tilde{r} \in B_{r,k}(s,a)} \tilde{r} + (\tilde{p}^{sa})^\top v \}$,
173 where $\tilde{p}^{sa} = \arg \max_{\tilde{p} \in \overline{B}_{p,k}^+(s,a)} \{ \tilde{p}^\top v_n \}$ can be solved using [2, Fig. 2], except for $(s, a) \notin \mathcal{K}_k$
174 for which we force $\tilde{p}^{sa}(s') := 0$ for any $s' \in \mathcal{S}_k^{\text{T}}$ (see Alg. 2 in App. C). If EVI is stopped when
175 $sp_{\mathcal{S}_k^{\text{EVI}}} \{u_{n+1} - u_n\} \leq \epsilon_k$ and the true MDP is sufficiently explored, then the greedy policy $\tilde{\pi}_k$ w.r.t.
176 u_n is ϵ_k -optimistic, i.e., $\tilde{g}_k \geq g_{M^*} - \epsilon_k$ (see Sec. 3.1 for details). The policy $\tilde{\pi}_k$ is then executed until
177 the number of visits to a state-action pair is doubled or a new state is “discovered” (i.e., $s_t \in \mathcal{S}_{k_t}^{\text{T}}$).
178 Finally, notice that when the true MDP M^* is communicating, there exists an episode \bar{k} s.t. for all
179 $k \geq \bar{k}$, $\mathcal{S}_k^{\text{T}} = \emptyset$ and TUCRL can be reduced to UCRL by considering \mathcal{M}_k in place of $\overline{\mathcal{M}}_k^+$.

180 3.1 Analysis of TUCRL

181 We prove that the regret of TUCRL is bounded as follows.

182 **Theorem 2.** *For any weakly communicating MDP M , with probability at least $1 - \delta$ it holds that for*
183 *any $T > 1$, the regret of TUCRL is bounded as*

$$\Delta(\text{TUCRL}, T) = O \left(r_{\max} D^c \sqrt{\Gamma^c S^c A T \ln \left(\frac{SAT}{\delta} \right)} + r_{\max} (D^c)^2 S^3 A \ln^2 \left(\frac{SAT}{\delta} \right) \right).$$

184 The first term in the regret shows the ability of TUCRL to adapt to the communicating part of the
185 true MDP M^* by scaling with the *communicating* diameter D^c and MDP parameters S^c and Γ^c . The
186 second term corresponds to the regret incurred in the early stage where the regret grows linearly.
187 When M^* is communicating, we match the square-root term of UCRL (first term), while the second
188 term is just slightly bigger than the one appearing in UCRL by a multiplicative factor S^2/Γ^c (ignoring
189 logarithmic terms, see Sec. 5).

190 We now provide a sketch of the main steps of the proof of Thm. 2 (the full proof is reported in
191 App. C). In order to preserve readability, in the following, all inequalities should be interpreted up to
192 minor approximations and in high probability.

193 Let $\Delta_k := \sum_{s,a} \nu_k(s, a)(g^* - r(s, a))$ be the regret incurred in episode k , where $\nu_k(s, a)$ is the
194 number of visits to s, a in episode k . We decompose the regret as

$$\Delta(\text{TUCRL}, T) \lesssim \sum_{k=1}^m \Delta_k \cdot \mathbb{1}\{M^* \in \mathcal{M}_k\} \lesssim \sum_{k=1}^m \Delta_k \cdot \mathbb{1}\{t_k < C(k)\} + \sum_{k=1}^m \Delta_k \cdot \mathbb{1}\{t_k \geq C(k)\}$$

195 where $C(k) = O((D^c)^2 S^3 A \ln^2(2SAT_k/\delta))$ defines the length of a full exploratory phase, where the
196 agent may suffer **linear regret**.

197 When $t_k \geq C(k)$, to further decompose Δ_k we have to prove that the solution found by EVI is
198 **optimistic**. The following lemma helps to identify the possible cases (see App. C.2).

199 **Lemma 3.** *Let episode k be such that $M^* \in \mathcal{M}_k$, $\mathcal{S}_k^{\text{T}} \neq \emptyset$ and $t_k \geq C(k)$. Then, either $\mathcal{S}_k^{\text{T}} = \mathcal{S}^{\text{T}}$*
200 *(case I) or $\mathcal{K}_k \neq \emptyset$, i.e., $\exists(s, a) \in \mathcal{S}_k^{\text{C}} \times \mathcal{A}$ for which transitions to \mathcal{S}_k^{T} are allowed (case II).*

201 We start noticing that when $\mathcal{S}_k^{\text{T}} = \emptyset$, the true MDP $M^* \in \mathcal{M}_k = \overline{\mathcal{M}}_k$ w.h.p. Similarly, if $\mathcal{S}_k^{\text{T}} = \mathcal{S}^{\text{T}}$
202 then $M^* \in \overline{\mathcal{M}}_k \subseteq \overline{\mathcal{M}}_k^+$ w.h.p., since TUCRL only truncates transitions that are indeed forbidden in
203 M^* itself. In both cases, we can use the same arguments in [2] to prove optimism. Finally, note that
204 in *case II* the gain of any state $s' \in \mathcal{S}_k^{\text{T}}$ is r_{\max} and, since there exists a path from \mathcal{S}_k^{C} to \mathcal{S}_k^{T} , the gain
205 of the solution returned by EVI is r_{\max} (i.e., optimistic).

206 Similarly to UCRL, by exploiting the optimism and the optimality equation, we further de-
207 compose the regret as $\Delta_k \cdot \mathbb{1}\{t_k \geq C(k)\} \lesssim \sum_{s,a} \nu_k(s, a)(\tilde{g}_k - \tilde{r}_k(s, a)) \mathbb{1}\{(s, a) \notin \mathcal{K}_k\} +$

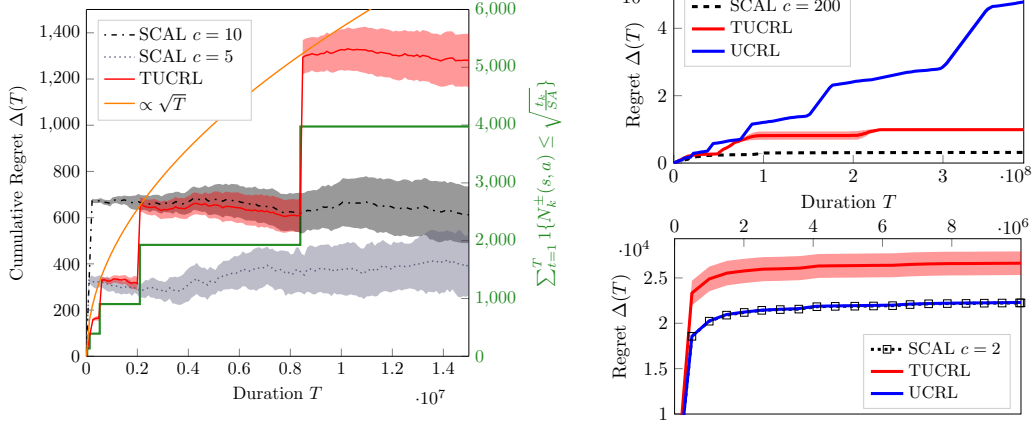


Figure 2: Cumulative regret in the weakly communicating three-states domain with $D = +\infty$ (left), in the taxi with misspecified states (right-top) and in the communicating taxi (right-bottom). Confidence intervals $\beta_{r,k}$ and $\beta_{p,k}$ are shrunk by a factor 0.05 and 0.01 for the three-states domain and taxi, respectively. Results are averaged over 20 runs and 95% confidence intervals are reported.

208 $r_{\max} \sum_{s,a} \nu_k(s,a) \mathbb{1}\{(s,a) \in \mathcal{K}_k\}$. We start focusing on the poorly visited state-action pairs, i.e.,
 209 $(s,a) \in \mathcal{K}_k$. In this case TUCRL may suffer the maximum per-step regret r_{\max} but the number of
 210 times this event happen is cumulatively “small” (see App. C.4.1 for the proof):

211 **Lemma 4.** For any $T \geq 1$ and any sequence of states and actions $\{s_1, a_1, \dots, s_T, a_T\}$ we have:

$$\sum_{k=1}^m \sum_{s,a} \nu_k(s,a) \mathbb{1}\left\{ \underbrace{N_k^\pm(s,a) \leq \sqrt{t_k/SA}}_{(s,a) \in \mathcal{K}_k} \right\} \leq \sum_{t=1}^T \mathbb{1}\left\{ N_{k_t}^\pm(s_t, a_t) \leq \sqrt{t/SA} \right\} \leq 2 \left(\sqrt{S^c AT} + S^c A \right)$$

212 When $(s,a) \notin \mathcal{K}_k$ (i.e., $N_k^\pm(s,a) > \sqrt{t_k/SA}$ holds), similar to UCRL, the “dominant” term in the
 213 upper-bound of $\sum_{s,a} \nu_k(s,a) (\bar{g}_k - \tilde{r}_k(s,a)) \cdot \mathbb{1}\{(s,a) \notin \mathcal{K}_k\}$ is:

$$\nu_k(\tilde{P}_k - I)\tilde{h}_k = \sum_{s \in \underline{\mathcal{S}}_k} \nu_k(s, \tilde{\pi}_k(s)) \left(\sum_{s' \in \underline{\mathcal{S}}_k} \tilde{p}_k(s'|s, \tilde{\pi}_k(s)) w_k(s') - w_k(s) \right),$$

214 where \tilde{h}_k is the vector returned by EVI and w_k is obtain by shifting \tilde{h}_k by a suitable constant.
 215 The critical step here is **bounding the optimistic vector** \tilde{h}_k (see Sec 4.3.2 in [2] and App. F.6
 216 in [6]). Using Holder inequality the sum can be decomposed into the ℓ_1 -norm of the transition
 217 probabilities (that can be bounded using a concentration argument) and $sp_{\underline{\mathcal{S}}_k}\{w_k\}$. In UCRL,
 218 $\underline{\mathcal{S}}_k = \mathcal{S}$ and $sp_{\mathcal{S}}\{w_k\} \leq D$ [2, Sec. 4.3]. Unfortunately, this result is uninformative when
 219 M^* is weakly communicating, since $D = +\infty$. However, using the new stopping condition, the
 220 truncation performed by $\overline{\mathcal{M}}_k^+$ and the fact that $(s,a) \notin \mathcal{K}_k$, in TUCRL we can prove that $\underline{\mathcal{S}}_k = \mathcal{S}_k^c$.
 221 Defining $w_k := \tilde{h}_k - \min_{s \in \mathcal{S}_k^c} \{\tilde{h}_k\} e$, we show that $sp_{\mathcal{S}_k^c}\{w_k\} = \max_{s \in \mathcal{S}_k^c} \{w_k\}$ is bounded by the
 222 communicating diameter D^c of M^* (see App. C.2). We first introduce the following useful lemma
 223 (see App. D).

224 **Lemma 5.** For any pair $(s, \bar{s}) \in \mathcal{S} \times \mathcal{S}_k^c$, $\tau_{\mathcal{M}_k^+}(s \rightarrow \bar{s}) = \tau_{\overline{\mathcal{M}}_k^+}(s \rightarrow \bar{s})$. Moreover, EVI on $\overline{\mathcal{M}}_k^+$ is
 225 such that (starting from the vector $\mathbf{0}$): $\tilde{h}_k(s') - \tilde{h}_k(s) \leq r_{\max} \cdot \tau_{\overline{\mathcal{M}}_k^+}(s \rightarrow s')$, $\forall s, s' \in \mathcal{S}$.

226 Since $M^* \in \mathcal{M}_k^+$, for all $s, s' \in \mathcal{S}_k^c$, $\tau_{\mathcal{M}_k^+}(s \rightarrow s') \leq D^c$ (see App. D). Then, as a consequence of
 227 Lem. 5, for all $s, s' \in \mathcal{S}_k^c$, $w_k(s) - w_k(s') = \tilde{h}_k(s) - \tilde{h}_k(s') \leq r_{\max} D^c$. This is the reason why we
 228 had to consider the relaxed set of MDPs $\overline{\mathcal{M}}_k^+$ in the definition of TUCRL: Lem. 5 does not hold when
 229 we consider $\overline{\mathcal{M}}_k$ instead. The regret in Thm. 2 is proved by combining all the different regret terms.

230 4 Experiments

231 In this section, we present experiments to validate the theoretical findings of Sec. 3. We first consider
 232 the taxi problem [21] implemented in OpenAI Gym [22]. Even such a simple domain contains

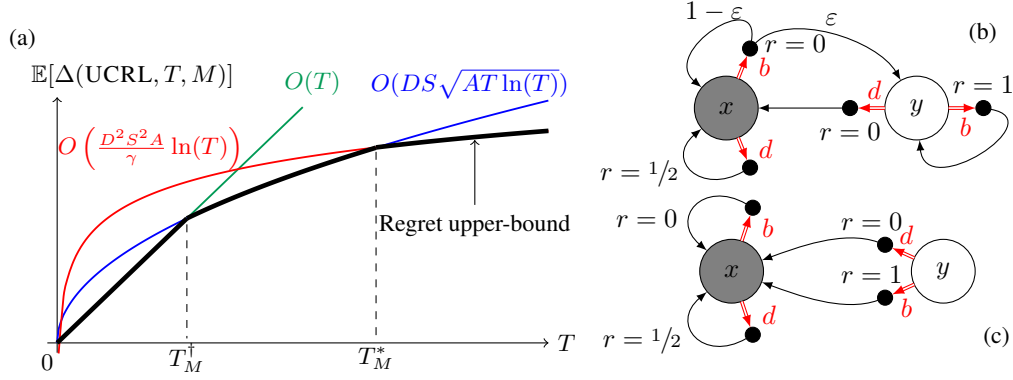


Figure 3: 3a Expected regret of UCRL (with known horizon T given as input) as a function of T . 3b 3c Toy example illustrating the difficulty of learning non-communicating MDPs. We represent a family of possible MDPs $\mathcal{M} = (M_\varepsilon)_{\varepsilon \in [0,1]}$ where the probability ε to go from x to y lies in $[0, 1]$.

233 misspecified states, since the state space is constructed as the outer product of the taxi position, the
 234 passenger position and the destination. This leads to states that cannot be reached from any possible
 235 starting configuration (all the starting states belong to \mathcal{S}^c). More precisely, out of 500 states in \mathcal{S} ,
 236 100 are non-reachable. On Fig. 2(right) we compare the regret of UCRL, SCAL and TUCRL when
 237 the misspecified states are present (top) and when they are removed (bottom). In the presence of
 238 misspecified states (top), the regret of UCRL clearly grows linearly with T while TUCRL is able to
 239 learn as expected. On the other hand, when the MDP is communicating (bottom) TUCRL performs
 240 similarly to UCRL. The small loss in performance is most likely due to the initial exploration phase
 241 during which the confidence intervals on the transition probabilities used by UCRL (see definition of
 242 $\overline{\mathcal{M}}_k$) are tighter than those used by TUCRL (see definition of $\overline{\mathcal{M}}_k^+$). TUCRL uses a “loose” bound on
 243 the ℓ_1 -norm while UCRL uses S different bounds, one for every possible next state. Finally, SCAL
 244 outperforms TUCRL by exploiting prior knowledge on the bias span.

245 We further study TUCRL regret in the simple three-state domain introduced in [6] (see App. F for
 246 details) with different reward distributions (uniform instead of Bernoulli). The environment is
 247 composed of only three states (s_0, s_1 and s_2) and one action per state, except in s_2 where two actions
 248 are available. As a result, the agent only has the choice between two possible policies. Fig. 2(left)
 249 shows the cumulative regret achieved by TUCRL and SCAL (with different upper-bounds on the
 250 bias span) when the diameter is infinite i.e., $\mathcal{S}^c = \{s_0, s_2\}$ and $\mathcal{S}^T = \{s_1\}$ (we omit UCRL, since
 251 it suffers linear regret). Both SCAL and TUCRL quickly achieve sub-linear regret as predicted by
 252 theory. However, SCAL and TUCRL seem to achieve different growth rates in regret: while SCAL
 253 appears to reach a logarithmic growth, the regret of TUCRL seems to grow as \sqrt{T} with periodic
 254 “jumps” that are increasingly distant (in time) from each other. This can be explained by the way the
 255 algorithm works: while most of the time TUCRL is optimistic on the restricted state space \mathcal{S}^c (i.e.,
 256 $\mathcal{S}_k^c = \mathcal{S}^c$), it periodically allows transitions to the set \mathcal{S}^T (i.e., $\mathcal{S}_k^c = \mathcal{S}$), which is indeed not reachable.
 257 Enabling these transitions triggers aggressive exploration during an entire episode. The policy played
 258 is then sub-optimal creating a “jump” in the regret. At the end of this exploratory episode, \mathcal{S}_k^c will be
 259 set again to \mathcal{S}^c and the regret will stop increasing until the condition $N_k^\pm \leq \sqrt{t_k}/SA$ occurs again
 260 (the time between two consecutive exploratory episodes grows quadratically). The cumulative regret
 261 incurred during exploratory episodes can be bounded by the term plotted in green on Fig. 2(left). In
 262 Lem. 4 we proved that this term is always bounded by $O(\sqrt{S^c AT})$. Therefore, it is not surprising
 263 to observe a \sqrt{T} increase of both the green and red curves. Unfortunately, the growth rate of the
 264 regret will keep increasing as \sqrt{T} and will never become logarithmic unlike SCAL (or UCRL when
 265 the MDP is communicating). This is because the condition $N_k^\pm \leq \sqrt{t_k}/SA$ will always be triggered
 266 $\Theta(\sqrt{T})$ times for any T . In Sec. 5 we show that this is not just a drawback specific to TUCRL, but it
 267 is rather an intrinsic limitation of learning in weakly-communicating MDPs.

268 5 Exploration-exploitation dilemma with infinite diameter

269 In this section we further investigate the empirical difference between SCAL and TUCRL and prove
 270 an impossibility result characterising the exploration-exploitation dilemma when the diameter is
 271 allowed to be infinite and no prior knowledge on the optimal bias span is available.

272 We first recall that the expected regret $\mathbb{E}[\Delta(\text{UCRL}, M, T)]$ of UCRL (with input parameter $\delta = 1/3T$)
 273 after $T \geq 1$ time steps and for any finite MDP M can be bounded in several ways:

$$\mathbb{E}[\Delta(\text{UCRL}, M, T)] \leq \begin{cases} r_{\max} T & \text{(by definition)} \\ C_1 \cdot r_{\max} D \sqrt{\Gamma S A T \ln(3T^2)} + \frac{1}{3} & \text{[2, Theorem 2]} \\ C_2 \cdot r_{\max} \frac{D^2 \Gamma S A}{\gamma} \ln(T) + C_3(M) & \text{[2, Theorem 4]} \end{cases} \quad (5)$$

274 where $\gamma = g_M^* - \max_{s, \pi} \{g_M^\pi(s) : g_M^\pi(s) < g_M^*\}$ is the gap in gain, $C_1 := 34$ and $C_2 := 34^2$
 275 are numerical constants independent of M , and $C_3(M) := O(\max_{\pi: \pi(s)=a} T_\pi)$ with T_π a measure of
 276 the ‘‘mixing time’’ of policy π . The three different bounds lead to three different *growth rates*
 277 for the function $T \mapsto \mathbb{E}[\Delta(\text{UCRL}, M, T)]$ as is illustrated on Fig. 3a: 1) for $T_M^\dagger \geq T \geq 0$, the
 278 expected regret is linear in T , 2) for $T_M^* \geq T \geq T_M^\dagger$ the expected regret grows as \sqrt{T} , 3) finally for
 279 $T \geq T_M^*$, the increase in regret is only logarithmic in T . These different ‘‘regimes’’ can be observed
 280 empirically (see [6, Fig. 5, 12]). Using (5), it is easy to show that the time it takes for UCRL to
 281 achieve sub-linear regret is at most $T_M^\dagger = \tilde{O}(D^2 \Gamma S A)$. We say that an algorithm is *efficient* when it
 282 achieves sublinear regret after a number of steps that is polynomial in the parameters of the MDP (i.e.,
 283 UCRL is then *efficient*). We now show with an example that *without prior knowledge*, any *efficient*
 284 learning algorithm must satisfy $T_M^* = +\infty$ when M has *infinite diameter* (i.e., it cannot achieve
 285 logarithmic regret).

286 **Example 1.** We consider a family of weakly-communicating MDPs $\mathcal{M} = (M_\varepsilon)_{\varepsilon \in [0,1]}$ represented
 287 on Fig. 3(right). Every MDP instance in \mathcal{M} is characterised by a specific value of $\varepsilon \in [0, 1]$ which
 288 corresponds to the probability to go from x to y . For $\varepsilon > 0$ (Fig. 3b), the optimal policy of M_ε is
 289 such that $\pi^*(x) = b$ and the optimal gain is $g_\varepsilon^* = 1$ while for $\varepsilon = 0$ (Fig. 3c) the optimal policy is
 290 such that $\pi^*(x) = d$ and the optimal gain is $g_0^* = 1/2$. We assume that the learning agent knows
 291 that the true MDP M^* belongs to \mathcal{M} but does not know the value ε^* associated to $M^* = M_{\varepsilon^*}$. We
 292 assume that all rewards are deterministic and that the agent starts in state x (coloured in grey).

293 **Lemma 6.** Let $C_1, C_2, \alpha, \beta > 0$ be positive real numbers and f a function defined for all $\varepsilon \in]0, 1]$
 294 by $f(\varepsilon) = C_1(1/\varepsilon)^\alpha$. There exists no learning algorithm \mathfrak{A}_T (with known horizon T) satisfying both
 295 1. for all $\varepsilon \in]0, 1]$, there exists $T_\varepsilon^\dagger \leq f(\varepsilon)$ such that $\mathbb{E}[\Delta(\mathfrak{A}_T, M_\varepsilon, x, T)] < 1/6 \cdot T$ for all $T \geq T_\varepsilon^\dagger$,
 296 2. and there exists $T_0^* < +\infty$ such that $\mathbb{E}[\Delta(\mathfrak{A}_T, M_0, x, T)] \leq C_2(\ln(T))^\beta$ for all $T \geq T_0^*$.

297 All the MDPs in \mathcal{M} share the same number of states $S = 2 \geq \Gamma$, number of actions $A = 2$, and gap
 298 in average reward $\gamma = 1/2$. As a result, any function of S, Γ, A and γ will be considered as constant.
 299 For $\varepsilon > 0$, the diameter coincides with the optimal bias span of the MDP and $D = sp_S \{h^*\} =$
 300 $1/\varepsilon < +\infty$, while for $\varepsilon = 0$, $D = +\infty$ but $sp_S \{h^*\} = 1/2$. As shown in Eq. 5 and Thm. 2, UCRL
 301 and TUCRL satisfy property 1. of Lem. 6 with $\alpha = 2$ and $C_1 = O(S^2 A)$ but do not satisfy 2. On the
 302 other hand, SCAL probably satisfies 2. with $\beta = 1$ and $C_2 = O(H^2 S A / \gamma)$ (it seems straightforward
 303 to adapt the proof of UCRL [2, Theorem 4] to SCAL) but since [6, Theorem 12] holds only when
 304 $H \geq sp_S \{h^*\}$, SCAL only satisfies 1. for $\varepsilon \geq 1/H$ and $\varepsilon = 0$ (not for $\varepsilon \in]0, 1/H[$). Lem. 6
 305 proves that no algorithm can actually achieve both 1. and 2. As a result, since TUCRL satisfies 1.,
 306 it cannot satisfy 2. This matches the empirical results presented in Sec. 4 where we observed that
 307 when the diameter is infinite, the growth rates of the regret of SCAL and TUCRL were respectively
 308 logarithmic and of order $\Theta(\sqrt{T})$. An algorithm that does not satisfy 1. could potentially satisfy 2.
 309 but, by definition of 1., it would suffer linear regret for a number of steps that is more than *polynomial*
 310 in the parameters of the MDP (more precisely, $e^{D^{1/\beta}}$). This is not a very desirable property and
 311 we claim that an *efficient* learning algorithm should always prefer *finite time guarantees* (1.) over
 312 *asymptotic guarantees* (2.) when they cannot be accommodated.

313 6 Conclusion

314 We introduced TUCRL, an algorithm that efficiently balances exploration and exploitation in weakly-
 315 communicating and multi-chain MDPs, when the starting state s_1 belongs to a communicating set
 316 (Asm. 1). We showed that TUCRL achieves a square-root regret bound and that, in the general case,
 317 it is not possible to design algorithm with logarithmic regret and polynomial dependence on the MDP
 318 parameters. Several questions remain open: **1)** relaxing Asm. 1 by considering a transient initial state
 319 (i.e., $s_1 \in \mathcal{S}^T$), **2)** refining the lower bound of Jaksch et al. [2] to finally understand whether it is
 320 possible to scale with $sp_S \{h^*\}$ (at least in communicating MDPs) instead of D without any prior
 321 knowledge (the flaw in REGAL.D may suggest it is indeed impossible).

References

- 322
- 323 [1] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1.
324 MIT press Cambridge, 1998.
- 325 [2] Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement
326 learning. *Journal of Machine Learning Research*, 11:1563–1600, 2010.
- 327 [3] Peter L. Bartlett and Ambuj Tewari. REGAL: A regularization based algorithm for reinforcement
328 learning in weakly communicating MDPs. In *UAI*, pages 35–42. AUAI Press, 2009.
- 329 [4] Ronan Fruit, Matteo Pirota, Alessandro Lazaric, and Emma Brunskill. Regret minimization in
330 mdps with options without prior knowledge. In *NIPS*, pages 3169–3179, 2017.
- 331 [5] Mohammad Sadegh Talebi and Odalric-Ambrym Maillard. Variance-aware regret bounds for
332 undiscounted reinforcement learning in mdps. In *ALT*, volume 83 of *Proceedings of Machine
333 Learning Research*, pages 770–805. PMLR, 2018.
- 334 [6] Ronan Fruit, Matteo Pirota, Alessandro Lazaric, and Ronald Ortner. Efficient bias-span-
335 constrained exploration-exploitation in reinforcement learning. *CoRR*, abs/1802.04020, 2018.
- 336 [7] William R. Thompson. On the likelihood that one unknown probability exceeds another in view
337 of the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- 338 [8] Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via
339 posterior sampling. In *NIPS*, pages 3003–3011, 2013.
- 340 [9] Yasin Abbasi-Yadkori and Csaba Szepesvári. Bayesian optimal control of smoothly parameter-
341 ized systems. In *UAI*, pages 1–11. AUAI Press, 2015.
- 342 [10] Ian Osband and Benjamin Van Roy. Why is posterior sampling better than optimism for
343 reinforcement learning? In *ICML*, volume 70 of *Proceedings of Machine Learning Research*,
344 pages 2701–2710. PMLR, 2017.
- 345 [11] Yi Ouyang, Mukul Gagrani, Ashutosh Nayyar, and Rahul Jain. Learning unknown markov
346 decision processes: A thompson sampling approach. In *NIPS*, pages 1333–1342, 2017.
- 347 [12] Shipra Agrawal and Randy Jia. Optimistic posterior sampling for reinforcement learning:
348 worst-case regret bounds. In *NIPS*, pages 1184–1194, 2017.
- 349 [13] Aditya Gopalan and Shie Mannor. Thompson sampling for learning parameterized markov
350 decision processes. In *COLT*, volume 40 of *JMLR Workshop and Conference Proceedings*,
351 pages 861–898. JMLR.org, 2015.
- 352 [14] Yasin Abbasi-Yadkori and Csaba Szepesvári. Bayesian optimal control of smoothly param-
353 eterized systems. In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial
354 Intelligence*, UAI’15, pages 2–11, Arlington, Virginia, United States, 2015. AUAI Press. ISBN
355 978-0-9966431-0-8.
- 356 [15] Yi Ouyang, Mukul Gagrani, Ashutosh Nayyar, and Rahul Jain. Learning unknown markov
357 decision processes: A thompson sampling approach. In *Advances in Neural Information
358 Processing Systems 30*, pages 1333–1342. Curran Associates, Inc., 2017.
- 359 [16] Georgios Theodorou, Zheng Wen, Yasin Abbasi-Yadkori, and Nikos Vlassis. Posterior
360 sampling for large scale reinforcement learning. *CoRR*, abs/1711.07979, 2017.
- 361 [17] Ian Osband and Benjamin Van Roy. Posterior sampling for reinforcement learning without
362 episodes. *CoRR*, abs/1608.02731, 2016.
- 363 [18] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*.
364 John Wiley & Sons, Inc., New York, NY, USA, 1994. ISBN 0471619779.
- 365 [19] Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. Tuning bandit algorithms in stochastic
366 environments. In *Algorithmic Learning Theory*, pages 150–165, Berlin, Heidelberg, 2007.
367 Springer Berlin Heidelberg.

- 368 [20] Andreas Maurer and Massimiliano Pontil. Empirical bernstein bounds and sample-variance
369 penalization. In *COLT*, 2009.
- 370 [21] Thomas G. Dietterich. Hierarchical reinforcement learning with the MAXQ value function
371 decomposition. *J. Artif. Intell. Res.*, 13:227–303, 2000.
- 372 [22] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang,
373 and Wojciech Zaremba. Openai gym. *CoRR*, abs/1606.01540, 2016.
- 374 [23] Odalric-Ambrym Maillard, Phuong Nguyen, Ronald Ortner, and Daniil Ryabko. Optimal regret
375 bounds for selecting the state representation in reinforcement learning. In *Proceedings of the*
376 *30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine*
377 *Learning Research*, pages 543–551, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.

378 A Mistake in the regret bound of REGAL.D

379 A.1 Regularized optimistic RL (REGAL)

380 In weakly communicating MDPs, to avoid the over-optimism of UCRL, Bartlett and Tewari [3]
 381 proposed to penalise the optimism on g^* by the optimal bias span $sp_S \{h^*\}$. Formally, at each
 382 episode k , their algorithm --REGAL-- solves the following optimization problem:

$$\widetilde{M}_k = \arg \max_{M \in \mathcal{M}_k} \{g_M^* - C_k \cdot sp_S \{h_M^*\}\} \quad (6)$$

383 where $C_k \geq 0$ is a regularisation coefficient. Note that such optimization requires to first compute
 384 the optimal policy for a given MDP $M \in \mathcal{M}_k$ and then evaluate the regularized gain. Implicitly, this
 385 defines the optimistic policy $\widetilde{\pi}_k = \arg \max_{\pi \in \Pi^{\text{SD}}} \{g_{M_k}^{\pi}\}$. The term $sp_S \{h^*\}$ can be interpreted as a
 386 measure of the *complexity* of the environment: the bigger $sp_S \{h^*\}$, the more difficult it is to achieve
 387 the stationary reward g^* by following the optimal policy. In *supervised learning*, regularisation is
 388 often used to penalise the objective function by a measure of the complexity of the model so as to
 389 avoid *overfitting*. It is thus reasonable to expect that *over-optimism* in online RL can also be avoided
 390 through regularisation.

391 The regret bound of REGAL holds only when C_k is set to $\Theta(1/\sum_{s,a} \nu_k(s,a))$. This means that
 392 REGAL requires the knowledge of (future) visit counts $\nu_k(s,a)$ before episode k begins in order to
 393 tune the regularisation coefficient C_k . Unfortunately, an episode stops when the number of visits in a
 394 state-action pair $(s,a) \in \mathcal{S} \times \mathcal{A}$ has doubled and it is not possible to predict the future sequence of
 395 states of a given policy for two reasons: 1) the true MDP M^* is unknown and 2) what is observed
 396 is a random *sampled* trajectory (as opposed to *expected*). As a result, REGAL is not *implementable*.
 397 Bartlett and Tewari [3] proposed an alternative algorithm --REGAL.D-- that leverages on the *doubling*
 398 *trick* to guess the length of episode k (i.e., $\sum_{s,a} \nu_k(s,a)$) and proved a slightly worse regret bound
 399 than for REGAL. REGAL.D divides an episode k into sub-iterations where it applies the doubling trick
 400 techniques. At each sub-iteration j , REGAL.D guesses that the length of the episode will be at most
 401 2^j and it solves problem (6) with $C_{k,j} \propto 1/\sqrt{2^j}$. Then, it executes the optimistic policy $\widetilde{\pi}_{k,j}$ on the
 402 true MDP until the UCRL stopping condition is reached or 2^j steps are performed. In the first case
 403 the episode k ends since the guess was correct, while, in the second case, a new sub-iteration $j+1$ is
 404 started. This implies that for any k, j :

$$\sum_{s,a} \nu_{k,j}(s,a) \leq 2^j, \quad (7)$$

405 where $\nu_{k,j}(s,a)$ denotes the number of visits to (s,a) during episode k and sub-iteration j .

406 A.2 The doubling trick issue

407 The mistake in REGAL.D is located in the proof of the regret [3, Theorem 3] (see Sec. 6.3). Let $\widetilde{h}_{k,j}$
 408 denote the optimistic bias span at episode k and sub-iteration j induced by the doubling trick. At
 409 high level, the mistake comes from the attempt to upper-bound the term $x \cdot \sum_{s,a} \nu_{k,j}(s,a)$ by $x \cdot 2^j$
 410 using the fact the $\sum_{s,a} \nu_{k,j}(s,a) \leq 2^j$. Unfortunately, this is possible only under the assumption
 411 that $x \geq 0$ that does not hold in the case of REGAL.D.

412 Formally, while bounding $\sum_{k \in G} \Delta_k$, the authors have to deal with the term (derived by the combina-
 413 tion of [3, Eq. 15] and [3, Lem. 11] with [3, Eq. 14]):

$$U := \sum_{k \in G} \sum_j sp_S \{\widetilde{h}_{k,j}\} \left(c \sqrt{\sum_{s,a} \nu_{k,j}(s,a)} - C_{k,j} \sum_{s,a} \nu_{k,j}(s,a) \right)$$

414 where $c := 2S\sqrt{12\ln(2AT/\delta)} + \sqrt{2\ln(1/\delta)} \geq 0$ and recall that $\sum_{s,a} \nu_{k,j}(s,a)$ denotes the *actual*
 415 length of the episode k at sub-iteration j . In the REGAL.D proof the authors directly replaced
 416 the actual length of the episode with the guessed length $2^j := \ell_{k,j}$ showing that the first term
 417 can be upper-bounded by $c \cdot \sqrt{\sum_{s,a} \nu_{k,j}(s,a)} \leq c \cdot \sqrt{2^j}$ (using Eq. 7). Concerning the second
 418 term, they wrote $-C_{k,j} \sum_{s,a} \nu_{k,j}(s,a) \leq -C_{k,j} 2^j$ which is actually not true since $-C_{k,j} :=$

419 $-c/\sqrt{2^j} \leq 0$. Therefore, we can not guarantee that $U \leq 0$ which is the goal of Bartlett and
 420 Tewari [3]. To do this, we would need to *lower-bound* $\sum_{s,a} \nu_{k,j}(s,a)$. Unfortunately, the only
 421 lower bound with probability 1 available for that term is $\min_{s,a} \{N_k(s,a)\} + 2$. This is not big
 422 enough to cancel the term $c\sqrt{\sum_{s,a} \nu_{k,j}(s,a)}$ and $C_{k,j}$ needs to be increased. As a result, the term
 423 $sp_S \{h^*\} \sum_{k \in G} \sum_j C_{k,j} \sqrt{\sum_{s,a} \nu_{k,j}(s,a)}$ becomes too big and all the proof collapses.

424 Notice that a similar mistake is contained in the work by Maillard et al. [23] where they use a
 425 regularized approach to learn a state representation in online settings. Similarly to [3], the authors
 426 have to bound the term $\sum_{s,a} \nu_{k,j}(s,a)(g^* - \tilde{g}_{k,j})$. By exploiting the fact that $g^* - \tilde{g}_{k,j} \leq \alpha$ (we
 427 omit the full expression of α for sake of clarity) [23, Eq. 17 Sec. 5.2] the authors derived the bound
 428 $\sum_{s,a} \nu_{k,j}(s,a)(g^* - \tilde{g}_{k,j}) \leq 2^j \cdot \alpha$ [23, Eq. 18]. The difference $g^* - \tilde{g}_{k,j}$ might be negative in which
 429 case the result does not hold. Actually for the case in which there is no regularization $C_{k,j} = 0$,
 430 $g^* \leq \tilde{g}_{k,j}$ which is what is used in the regret proof of UCRL. Therefore, it is very likely that the sign
 431 of $g^* - \tilde{g}_{k,j}$ can sometimes be negative.

432 In conclusion, it seems unavoidable to use a *lower-bound* (and not an upper-bound) on $\sum_{s,a} \nu_{k,j}(s,a)$
 433 to derive a correct regret bound for REGAL.D. As already mentioned, given the current stopping
 434 condition of an episode, the only reasonable lower bound is $\min_{s,a} \{N_k(s,a)\} + 2$ and it does not
 435 seem sufficient to derive a sensible regret bound. Another research direction could be to change the
 436 stopping condition. However, one of the terms in the regret bound of REGAL (and of REGAL.D) scales
 437 as $m\sqrt{T} \log_2(T)$ where m is the number of episodes. The term m is highly sensitive to the stopping
 438 condition and there is very little margin if we want to avoid $m\sqrt{T} \log_2(T)$ to become the leading
 439 term in the regret bound. All the efforts we put in this direction were unsuccessful. We conjecture
 440 that regularising by the optimal bias span might not allow to learn MDPs with infinite diameter.

441 B Unbounded optimal bias span with continuous Bayesian priors/posteriors

442 Recently, Ouyang et al. [15] and Theodorou et al. [16] proposed posterior sampling algorithms and
 443 proved bounds on the expected Bayesian regret. The regret bounds that they derive scale linearly with
 444 H , where H is the highest optimal bias span of all the MDPs that can be drawn from the prior/posterior
 445 distribution. Formally, let $f(\theta)$ be the density function of the prior/posterior distribution over the
 446 family of MDPs (M_θ) parametrised by θ . Then:

$$H := \sup_{\theta: f(\theta) > 0} \{sp_S \{h_\theta^*\}\}.$$

447 In this section we present an example where H is *infinite* and argue that it is probably the case for
 448 most priors/posteriors used in practice.

449 **Example 2** (Unbounded optimal bias span with continuous prior/posterior). *Consider the example of*
 450 *Fig. 4. There is only one action in every state and so one optimal policy. The (unique) action that can*
 451 *be played in state s_0 loops on s_0 with probability $1 - \theta$ and goes to s_1 with probability θ . The reward*
 452 *associated to this action is 0. Symmetrically, the (unique) action that can be played in state s_1 loops*
 453 *on s_1 with probability $1 - \theta$ and goes to s_0 with probability θ . The reward associated to this action is*
 454 *1. This MDP is characterised by the parameter θ and we denote it by M_θ . For any $\theta \in [0, 1]$, we*
 455 *denote by g_θ^* (resp. h_θ^*) the optimal gain (resp. bias) of M_θ . Observe that when $\theta > 0$, M_θ is ergodic*
 456 *and therefore the optimal gain $g_\theta^* = 1/2$ is state-independent whereas when $\theta = 0$, M_θ is multichain*
 457 *and the optimal gain does depend on the initial state: $g_0^*(x) = 0 < 1 = g_0^*(y)$.*

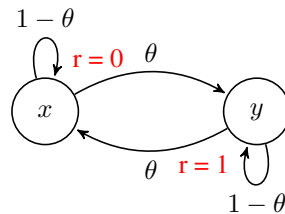


Figure 4: Toy example of a parametrised MDP M_θ with a single policy (one action per state).

458 Let's assume that the prior/posterior distribution we use on M_θ is characterised by a probability
 459 density function f satisfying $f(\theta) > 0$ for all $\theta > 0$ and $f(0) = 0$. Note that this assumption does
 460 not constrain the ‘‘smoothness’’ of f e.g., f can have continuous derivatives of all orders. Under this
 461 assumption, f is non-zero only for ergodic MDPs. It goes without saying that for all $\theta \in [0, 1]$ (0
 462 included), $sp_S \{h_\theta^*\} < +\infty$ by definition (the optimal bias span is always finite). More precisely we
 463 have:

$$g_\theta^* = \begin{cases} [1/2, 1/2]^T & \text{if } \theta > 0 \\ [0, 1]^T & \text{if } \theta = 0 \end{cases} \quad \text{and} \quad sp_S \{h_\theta^*\} = \begin{cases} \frac{1}{2\theta} & \text{if } \theta > 0 \\ 0 & \text{if } \theta = 0 \end{cases}$$

464 As a result, although $sp_S \{h_\theta^*\}$ is always *finite*, i.e., $\forall \theta \in [0, 1]$, $sp_S \{h_\theta^*\} < +\infty$, it is *unbounded*
 465 on the set of plausible MDPs $\theta \in]0, 1]$ satisfying $f(\theta) > 0$, i.e.,

$$H := \sup_{\theta \in]0, 1]} \{sp_S \{h_\theta^*\}\} = \lim_{\theta \rightarrow 0^+} \frac{1}{2\theta} = +\infty$$

466 Therefore, the regret bound $\tilde{O}(HS\sqrt{AT})$ proved by Ouyang et al. [15], Theocharous et al. [16] does
 467 not hold with prior/posterior f since $H = +\infty$. One might argue that the proofs in [15, 16] could
 468 be fixed by showing that H is bounded with probability 1. Unfortunately, for any $C \in [0, +\infty[$,
 469 the probability $\mathbb{P}(sp_S \{h_\theta^*\} \geq C) = \int_{\theta=0}^{\frac{1}{2C}} f(\theta)d\theta > 0$ of sampling an MDP with $sp_S \{h_\theta^*\} \geq C$ is
 470 strictly positive. We therefore conjecture that for this specific choice of priors/posteriors, the regret
 471 proof in [15, 16] cannot be fixed without major changes and new arguments. More generally, let's
 472 imagine that we have a prior/posterior distribution f satisfying:

- 473 • there exists θ_0 such that M_{θ_0} has non-constant gain i.e., $sp_S \{g_{\theta_0}^*\} > 0$,
- 474 • there exists an open neighbourhood of θ_0 denoted Θ_0 such that $\forall \theta \in \Theta_0$, M_θ has constant
 475 gain (e.g., M_θ is weakly-communicating) and $f(\theta) > 0$.

476 In this case we will face the same problem as in Ex. 2 i.e.,

$$\sup_{\theta: f(\theta) > 0} \{sp_S \{h_\theta^*\}\} = +\infty \quad \text{and} \quad \forall C \in [0, +\infty[, \mathbb{P}(sp_S \{h_\theta^*\} \geq C) > 0$$

477 When the set of plausible MDPs is *finite*, this problem cannot occur. But most priors/posteriors used
 478 in practice are *continuous* distributions. For instance, a Dirichlet distribution will most likely satisfy
 479 the above assumptions.

480 C Regret of TUCRL

481 We follow the proof structure of Jaksch et al. [2], Fruit et al. [6] and use similar notations. Nonetheless,
 482 several parts of the proof significantly differ from [2, 6]:

- 483 • in Sec. C.2 we prove that after a finite number of steps, TUCRL is *gain-optimistic* (which is
 484 not as straightforward as in the case of UCRL),
- 485 • in Sec. C.3 we show that the sums taken over the whole state space \mathcal{S} that appear in the main
 486 term of the regret decomposition of UCRL can be restricted to sums over \mathcal{S}_k^c thanks to the
 487 new stopping condition used for episodes and the use of the condition $N_k^\pm(s, a) > \sqrt{t_k/SA}$
 488 (see (15)),
- 489 • in Sec. C.4.1, we bound the number of time steps spent in ‘‘bad’’ state-action pairs (s, a)
 490 satisfying $N_k^\pm(s, a) \leq \sqrt{t_k/SA}$,
- 491 • in Sec. C.4.3, we bound the number of episodes with the new stopping condition.

492 C.1 Splitting into episodes

493 The regret of TUCRL after T time steps is defined as: $\Delta(\text{TUCRL}, T) := Tg^* - \sum_{t=1}^T r_t(s_t, a_t)$.
 494 Defining $\Delta_k = \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \nu_k(s, a) (g^* - r(s, a))$ and using the same arguments as in [2, 6], it holds
 495 with probability $1 - \frac{\delta}{12T^{4/5}}$ that:

$$\Delta(\text{TUCRL}, T) \leq \sum_{k=1}^m \Delta_k + r_{\max} \sqrt{\frac{5}{2} T \ln \left(\frac{8T}{\delta} \right)} \quad (8)$$

Algorithm 1 TRUNCATED EXTENDED VALUE ITERATION (TEVI)

Input: value vector \tilde{h}_0 , extended MDP \mathcal{M} , set of states $\bar{\mathcal{S}}$, accuracy ϵ
Output: $\tilde{g}_k, \tilde{h}_k, \tilde{\pi}_k$
 $k := 0$
 $\tilde{h}_1(s) := \tilde{L}_{\mathcal{M}} \tilde{h}_0(s) := \max_{a \in \mathcal{A}_s} \left\{ \max_{\tilde{r} \in B_r(s,a)} \tilde{r} + \max_{\tilde{p} \in B_p(s,a)} \tilde{p}^\top \tilde{h}_0 \right\}, \forall s \in \bar{\mathcal{S}}$ (see Sec. 3)
while $\max_{s \in \bar{\mathcal{S}}} \left\{ \tilde{h}_{k+1}(s) - \tilde{h}_k(s) \right\} - \min_{s \in \bar{\mathcal{S}}} \left\{ \tilde{h}_{k+1}(s) - \tilde{h}_k(s) \right\} > \epsilon$ **do**
 $k := k + 1$
 $\tilde{h}_{k+1}(s) := \tilde{L}_{\mathcal{M}} \tilde{h}_k(s), \forall s \in \bar{\mathcal{S}}$
end while
 $\tilde{g}_k := \frac{1}{2} \left(\max_{s \in \bar{\mathcal{S}}} \left\{ \tilde{h}_{k+1}(s) - \tilde{h}_k(s) \right\} + \min_{s \in \bar{\mathcal{S}}} \left\{ \tilde{h}_{k+1}(s) - \tilde{h}_k(s) \right\} \right)$
 $\tilde{\pi}_k(s) \in \arg \max_{a \in \mathcal{A}_s} \left\{ \max_{\tilde{r} \in B_r(s,a)} \tilde{r} + \max_{\tilde{p} \in B_p(s,a)} \tilde{p}^\top \tilde{h}_k \right\}, \forall s \in \bar{\mathcal{S}}$

496 **C.2 Episodes with $M^* \in \mathcal{M}_k$**

497 We now assume that $M^* \in \mathcal{M}_k$. Let's denote by \tilde{g}_k, \tilde{h}_k and $\tilde{\pi}_k$ the outputs of TEVI($\mathbf{0}, \mathcal{M}_k^\circ, \mathcal{S}_k^{\text{EVI}}, \epsilon_k$)
498 (see Alg. 1) where $\epsilon_k := r_{\max}/\sqrt{t_k}$ and

$$\mathcal{S}_k^{\text{EVI}} = \begin{cases} \mathcal{S}_k^c & \text{if } \mathcal{K}_k = \emptyset \\ \mathcal{S} & \text{otherwise} \end{cases}, \quad \mathcal{M}_k^\circ = \begin{cases} \mathcal{M}_k = \overline{\mathcal{M}}_k & \text{if } \mathcal{S}_k^T = \emptyset \\ \overline{\mathcal{M}}_k^+ & \text{otherwise} \end{cases}. \quad (9)$$

499 In order to bound Δ_k we first show that $\tilde{g}_k \gtrsim g^*$ (up to $r_{\max}/\sqrt{t_k}$ -accuracy). If $\mathcal{S}_k^T = \emptyset$ then by
500 definition $\mathcal{M}_k^\circ = \overline{\mathcal{M}}_k = \mathcal{M}_k \ni M^*$ and so we can use the same argument as in [2, Sec. 4.3 & Thm.
501 7]. If $\mathcal{S}_k^T \neq \emptyset$, the true MDP M^* might not be ‘‘included’’ in the extended MDP $\overline{\mathcal{M}}_k^+$ considered by
502 EVI and we cannot use the same argument. To overcome this problem we first assume that t_k is big
503 enough which allows us to prove a useful lemma (Lem. 7):

$$t_k \geq \frac{2401}{9} (D^c)^2 SA \left(\mathcal{S}_k^T \ln \left(\frac{2SA t_k}{\delta} \right) \right)^2 := C(k) \quad (10)$$

504 where $\mathcal{S}_k^T := |\mathcal{S}_k^T|$ is the cardinal of \mathcal{S}_k^T .

505 **Lemma 7.** *Let episode k be such that $M^* \in \mathcal{M}_k, \mathcal{S}_k^T \neq \emptyset$ and (10) holds. Then,*

$$\left(\forall (s, a) \in \mathcal{S}_k^c \times \mathcal{A}, N_k^\pm(s, a) > \sqrt{\frac{t_k}{SA}} \right) \implies \mathcal{S}_k^T = \mathcal{S}^T$$

506 *Proof.* Assume that episode k is such that (10) holds and that for any state-action pair $(s, a) \in \mathcal{S}_k^c \times \mathcal{A}$

$$N_k^\pm(s, a) > \sqrt{\frac{t_k}{SA}} \geq \frac{49}{3} D^c \mathcal{S}_k^T \ln \left(\frac{2SA t_k}{\delta} \right)$$

507 Since $\mathcal{S}_k^T \neq \emptyset$ and $M^* \in \mathcal{M}_k$, for any $(s, a, s') \in \mathcal{S}_k^c \times \mathcal{A} \times \mathcal{S}_k^T$

$$\begin{aligned} \underbrace{p(s'|s, a)}_{\text{transition probability in } M^*} &\leq \underbrace{\hat{p}_k(s'|s, a)}_{=0} + \beta_k^{sas'} = \underbrace{\sqrt{\frac{14\hat{\sigma}_{p,k}^2(s'|s, a) \ln(2SA t_k/\delta)}{N_k^+(s, a)}}}_{=0} + \frac{49 \ln(2SA t_k/\delta)}{3N_k^\pm(s, a)} \\ &\leq \frac{49 \ln(2SA t_k/\delta)}{3N_k^\pm(s, a)} < \frac{1}{D^c \mathcal{S}_k^T} \end{aligned}$$

508 where we have exploited the fact that $\hat{p}(s'|s, a) = 0$ and $\hat{\sigma}_{p,k}^2(s'|s, a) = 0$ for any state $s' \in \mathcal{S}_k^T$
509 (remember that $N_k(s, a, s') = 0$).

510 We denote by $\tau_{M^*}(s \rightarrow s')$ the shortest path between any pair of states $(s, s') \in \mathcal{S} \times \mathcal{S}$ in the
511 true MDP M^* . Fix an arbitrary target state $\bar{s} \in \mathcal{S}_k^T$ and denote by $\tau(s) := \tau_{M^*}(s \rightarrow \bar{s})$ and

512 $\tau_{\min} := \min_{s \in \mathcal{S}_k^c} \{\tau(s)\}$. We have

$$\begin{aligned}
513 \quad & \tau(\bar{s}) = 0 \\
514 \quad & \forall s \in \mathcal{S}_k^c \quad \tau(s) = 1 + \min_{a \in \mathcal{A}_s} \left\{ \underbrace{\sum_{s' \in \mathcal{S}} p(s'|s, a) \tau(s')}_{\geq 0} \right\} \geq 1 + \min_{a \in \mathcal{A}_s} \left\{ \sum_{s' \in \mathcal{S}_k^c} p(s'|s, a) \underbrace{\tau(s')}_{\geq \tau_{\min}} \right\} \\
515 \quad & \geq 1 + \tau_{\min} \cdot \min_{a \in \mathcal{A}} \left\{ \sum_{s' \in \mathcal{S}_k^c} p(s'|s, a) \right\} = 1 + \tau_{\min} \cdot \min_{a \in \mathcal{A}} \left\{ 1 - \sum_{s' \in \mathcal{S}_k^T} p(s'|s, a) \right\} \\
516 \quad & > 1 + \tau_{\min} \left(1 - \sum_{s' \in \mathcal{S}_k^T} \frac{1}{D^c \mathcal{S}_k^T} \right) = 1 + \tau_{\min} \left(1 - \frac{1}{D^c} \right)
\end{aligned}$$

513 Applying the above inequality to $\tilde{s} \in \mathcal{S}_k^c$ achieving $\tau(\tilde{s}) = \tau_{\min}$ yields $\tau_{\min} > D^c$. This implies that
514 the shortest path in M^* between any state $s \in \mathcal{S}_k^c \subseteq \mathcal{S}^c$ and any state in $\bar{s} \in \mathcal{S}_k^T$ is strictly bigger
515 than D^c but by definition D^c is the longest shortest path between any pair of states in \mathcal{S}^c . Therefore,
516 $\bar{s} \in \mathcal{S}^T$. Since $\bar{s} \in \mathcal{S}_k^T$ was chosen arbitrarily, then $\mathcal{S}_k^T = \mathcal{S}^T$. \square

517 As a consequence of Lem. 7, under the assumptions that $M^* \in \mathcal{M}_k$, $\mathcal{S}_k^T \neq \emptyset$ and (10) holds, there
518 are only two possible cases:

- 519 1. Either $\mathcal{S}_k^T = \mathcal{S}^T$,
- 520 2. or $\exists (s, a) \in \mathcal{S}_k^c \times \mathcal{A} : N_k^\pm(s, a) \leq \sqrt{\frac{t_k}{SA}}$.

521 **Case 1:** $\mathcal{S}_k^T = \mathcal{S}^T$ implies that $M^* \in \overline{\mathcal{M}}_k^+$. This is because for any $(s, a, s') \in \mathcal{S}_k^c \times \mathcal{A} \times \mathcal{S}_k^T$ we have
522 $p(s'|s, a) = \tilde{p}_k(s'|s, a) = 0$ and for any $(s, a, s') \notin \mathcal{S}_k^c \times \mathcal{A} \times \mathcal{S}_k^T$ we have $|p(s'|s, a) - \tilde{p}_k(s'|s, a)| \leq$
523 $\beta_{p,k}^{sas'}$ and so $p(\cdot|s, a) \in \overline{B}_{p,k}^+(s, a)$. Since $M^* \in \overline{\mathcal{M}}_k^+$, we can use the same argument as Jaksch et al.
524 [2, Sec. 4.3 & Theorem 7] to prove $\tilde{g}_k \geq g^* - \frac{r_{\max}}{\sqrt{t_k}}$.

525 **Case 2:** For any $(s, a) \in \mathcal{S}_k^T \times \mathcal{A}$, $\overline{B}_{p,k}^+(s, a) = \mathcal{C}$ is the $(S-1)$ -simplex denoting the maximal
526 uncertainty about the transition probabilities, and $B_{r,k}(s, a) = [0, r_{\max}]$. We will now construct
527 an MDP $M' \in \overline{\mathcal{M}}_k^+$ with optimal gain r_{\max} . For all $(s, a) \in \mathcal{S}_k^T \times \mathcal{A}$, we set the transitions to
528 $p_{M'}(s'|s, a) = 1$ and rewards to $r_{M'}(s, a) = r_{\max}$. Let $(\bar{s}, \bar{a}) \in \mathcal{S}_k^c \times \mathcal{A}$ such that $N_k^\pm(\bar{s}, \bar{a}) \leq \sqrt{\frac{t_k}{SA}}$
529 (which exists by assumption). We set $p_{M'}(s'|\bar{s}, \bar{a}) > 0$ for all $s' \in \mathcal{S}_k^T$. This is possible because by
530 definition of $\overline{\mathcal{M}}_k^+$, the support of $p(\cdot|\bar{s}, \bar{a})$ is not restricted to \mathcal{S}_k^c . Finally, for all state-action pairs
531 $(s, a) \in \mathcal{S}_k^c \times \mathcal{A}$, we set $p_{M'}(\bar{s}|s, a) > 0$. This is possible because by definition of $\overline{\mathcal{M}}_k^+$, the support
532 of $p(\cdot|s, a)$ is only restricted to \mathcal{S}_k^c and $\bar{s} \in \mathcal{S}_k^c$. In M' , for all policies, all states in \mathcal{S}_k^T are absorbing
533 states (i.e., loop on themselves with probability 1) with maximal reward r_{\max} and all other states
534 $s \in \mathcal{S}_k^c$ are transient. The optimal gain of M' is thus r_{\max} and since $M' \in \overline{\mathcal{M}}_k^+$ we conclude that
535 $\tilde{g}_k \geq r_{\max} - \frac{r_{\max}}{\sqrt{t_k}} \geq g^* - \frac{r_{\max}}{\sqrt{t_k}}$.

536 In conclusion, TEVI is always returning an *optimistic* policy when the assumptions of Lem. 7 hold.
537 The regret Δ_k accumulated in episode k can thus be upper-bounded as:

$$\begin{aligned}
\Delta_k &= \sum_{s,a} \nu_k(s, a) (g^* - r(s, a)) = \sum_{s,a} \nu_k(s, a) \underbrace{(g^* - \tilde{r}_k(s, a))}_{\leq \tilde{g}_k - \frac{r_{\max}}{\sqrt{t_k}}} + \sum_{s,a} \nu_k(s, a) (\tilde{r}_k(s, a) - r(s, a)) \\
&\leq \underbrace{\sum_{s,a} \nu_k(s, a) (\tilde{g}_k - \tilde{r}_k(s, a))}_{:= \tilde{\Delta}_k} + \sum_{s,a} \nu_k(s, a) (\tilde{r}_k(s, a) - r(s, a)) + r_{\max} \sum_{s,a} \frac{\nu_k(s, a)}{\sqrt{t_k}}
\end{aligned}$$

538 To bound the difference between the optimistic reward \tilde{r}_k and the true reward r we introduce the
 539 estimated reward \hat{r}_k :

$$\forall s, a \in \mathcal{S} \times \mathcal{A}, \tilde{r}_k(s, a) - r(s, a) = \underbrace{\tilde{r}_k(s, a) - \hat{r}_k(s, a)}_{\leq \beta_{r,k}^{sa} \text{ by construction}} + \underbrace{\hat{r}_k(s, a) - r(s, a)}_{\leq \beta_{r,k}^{sa} \text{ since } M \in \mathcal{M}_k} \leq 2\beta_{r,k}^{sa}$$

540 and so in conclusion:

$$\Delta_k \leq \tilde{\Delta}_k + \underbrace{2 \sum_{s,a} \nu_k(s, a) \beta_{r,k}^{sa} + r_{\max} \sum_{s,a} \frac{\nu_k(s, a)}{\sqrt{t_k}}}_{:= U_k^1} \quad (11)$$

541 C.3 Bounding $\tilde{\Delta}_k$

542 The goal of this section is to bound the term $\tilde{\Delta}_k := \sum_{s,a} \nu_k(s, a) (\tilde{g}_k - \tilde{r}_k(s, a))$. We start by
 543 discarding the state-action pairs $(s, a) \in \mathcal{K}_k$ that have been poorly visited so far:

$$\begin{aligned} \tilde{\Delta}_k &= \sum_{s,a} \nu_k(s, a) (\tilde{g}_k - \tilde{r}_k(s, a)) \underbrace{\mathbb{1}\{(s, a) \notin \mathcal{K}_k\}}_{:= \mathbb{1}_k(s, a)} + \sum_{s,a} \nu_k(s, a) \underbrace{(\tilde{g}_k - \tilde{r}_k(s, a))}_{\leq r_{\max}} \mathbb{1}\{(s, a) \in \mathcal{K}_k\} \\ &\leq \underbrace{\sum_{s,a} \nu_k(s, a) (\tilde{g}_k - \tilde{r}_k(s, a)) \mathbb{1}_k(s, a)}_{:= \tilde{\Delta}'_k} + r_{\max} \sum_{s,a} \nu_k(s, a) \mathbb{1}\{(s, a) \in \mathcal{K}_k\} \end{aligned} \quad (12)$$

544 We will now bound the term $\tilde{\Delta}'_k = \sum_s \nu_k(s, \tilde{\pi}_k(s)) (\tilde{g}_k - \tilde{r}_k(s, \tilde{\pi}_k(s))) \mathbb{1}_k(s, \tilde{\pi}_k(s))$. We recall that
 545 the policy $\tilde{\pi}_k$ is obtained by executing TEVI(0, $\mathcal{M}_k^\circ, \mathcal{S}_k^{\text{EVI}}, \varepsilon_k$) 1 where $\varepsilon_k := r_{\max}/\sqrt{t_k}$ and $\mathcal{S}_k^{\text{EVI}}$
 546 and \mathcal{M}_k° are defined in (9). In all possible cases for both $\mathcal{S}_k^{\text{EVI}}$ and \mathcal{M}_k° , this amounts to applying
 547 value iteration to a *communicating* MDP with finite state space $\mathcal{S}_k^{\text{EVI}}$ and *compact* action space.
 548 By [18, Thm. 8.5.6], since the convergence criterion of value iteration is met we have:

$$\forall s \in \mathcal{S}_k^{\text{EVI}}, \left| \tilde{h}_k(s) + \tilde{g}_k - \tilde{r}_k(s, \tilde{\pi}_k(s)) - \sum_{s' \in \mathcal{S}} \tilde{p}_k(s'|s, \tilde{\pi}_k(s)) \tilde{h}_k(s') \right| \leq \frac{r_{\max}}{\sqrt{t_k}} \quad (13)$$

549 For all $s \notin \mathcal{S}_k^{\text{EVI}}$, $\nu_k(s, \tilde{\pi}_k(s)) = 0$ due to the stopping condition of episode k . Therefore we can
 550 plug (13) in $\tilde{\Delta}'_k$ and derive an upper bound restricted to the set $\mathcal{S}_k^{\text{C}} \subseteq \mathcal{S}_k^{\text{EVI}}$. Before to do that, we
 551 further decompose $\tilde{\Delta}'_k$ as:

$$\begin{aligned} \tilde{\Delta}'_k &\leq \sum_s \nu_k(s, \tilde{\pi}_k(s)) \left(\sum_{s' \in \mathcal{S}} \tilde{p}_k(s'|s, \tilde{\pi}_k(s)) \tilde{h}_k(s') - \tilde{h}_k(s) + \frac{r_{\max}}{\sqrt{t_k}} \right) \mathbb{1}_k(s, \tilde{\pi}_k(s)) \\ &= \nu'_k \left(\tilde{P}_k - I \right) \tilde{h}_k + r_{\max} \sum_{s,a} \frac{\nu_k(s, a)}{\sqrt{t_k}} \mathbb{1}_k(s, a) \end{aligned} \quad (14)$$

552 where $\nu'_k = (\nu_k(s, \tilde{\pi}_k(s)) \mathbb{1}_k(s, \tilde{\pi}_k(s)))_{s \in \mathcal{S}}$ is the vector of visit counts for each state and the corre-
 553 sponding action chosen by $\tilde{\pi}_k$ multiplied by the indicator function $\mathbb{1}_k$, $\tilde{P}_k = (\tilde{P}_k(s'|s, \tilde{\pi}_k(s)))_{s, s' \in \mathcal{S}}$
 554 is transition matrix associated to $\tilde{\pi}_k$ in $\overline{\mathcal{M}}_k^+$ and I is the identity matrix. We now focus on the
 555 term $\nu'_k (\tilde{P}_k - I) \tilde{h}_k$. Since the rows of \tilde{P}_k sum to 1, $\forall \lambda \in \mathbb{R}$, $(\tilde{P}_k - I) \tilde{h}_k = (\tilde{P}_k - I) (\tilde{h}_k + \lambda e)$
 556 where $e = (1, \dots, 1)^\top$ is the vector of all ones. Let's take $\lambda := -\min_{s \in \mathcal{S}_k^{\text{C}}} \{\tilde{h}_k(s)\}$ and define
 557 $w_k := \tilde{h}_k + \lambda e$ so that for all $s \in \mathcal{S}_k^{\text{C}}$, $w_k(s) \geq 0$ and $\min_{s \in \mathcal{S}_k^{\text{C}}} \{w_k(s)\} = 0$. We have:

$$\nu'_k (\tilde{P}_k - I) \tilde{h}_k = \sum_{s \in \mathcal{S}} \nu_k(s, \tilde{\pi}_k(s)) \mathbb{1}_k(s, \tilde{\pi}_k(s)) \left(\sum_{s' \in \mathcal{S}} \tilde{p}_k(s'|s, \tilde{\pi}_k(s)) w_k(s') - w_k(s) \right)$$

558 We denote by $k_t := \sup\{k \geq 1 : t_k \leq t\}$ the current episode at time t . Whenever $s_t \in \mathcal{S}_{k_t}^{\text{T}}$, episode
 559 k_t stops before executing any action (see the stopping condition of TUCRL in Alg. 1) implying that

560 $\forall s \in \mathcal{S}_k^T, \nu_k(s, \tilde{\pi}_k(s)) = 0$. Therefore we have:

$$\nu'_k(\tilde{P}_k - I)\tilde{h}_k = \sum_{s \in \mathcal{S}_k^c} \nu_k(s, \tilde{\pi}_k(s)) \mathbb{1}_k(s, \tilde{\pi}_k(s)) \left(\sum_{s' \in \mathcal{S}} \tilde{p}_k(s'|s, \tilde{\pi}_k(s)) w_k(s') - w_k(s) \right)$$

561 For all states s such that $\mathbb{1}_k(s, \tilde{\pi}_k(s)) = 1$, i.e., satisfying $N_k^\pm(s, \tilde{\pi}_k(s)) > \sqrt{t_k/SA}$, we force TEVI
562 to set $\tilde{p}_k(s'|s, \tilde{\pi}_k(s)) = 0, \forall s' \in \mathcal{S}_k^T$, by construction of $\overline{\mathcal{M}}_k^+$ so that:

$$\nu'_k(\tilde{P}_k - I)\tilde{h}_k = \sum_{s \in \mathcal{S}_k^c} \nu_k(s, \tilde{\pi}_k(s)) \mathbb{1}_k(s, \tilde{\pi}_k(s)) \left(\sum_{s' \in \mathcal{S}_k^c} \tilde{p}_k(s'|s, \tilde{\pi}_k(s)) w_k(s') - w_k(s) \right) \quad (15)$$

563 We can now introduce p :

$$\sum_{s' \in \mathcal{S}_k^c} \tilde{p}_k(s'|s, \tilde{\pi}_k(s)) w_k(s') - w_k(s) = \sum_{s' \in \mathcal{S}_k^c} \tilde{p}_k(s'|s, \tilde{\pi}_k(s)) w_k(s') - p(s'|s, \tilde{\pi}_k(s)) w_k(s') \quad (16)$$

$$+ \left(\sum_{s' \in \mathcal{S}_k^c} p(s'|s, \tilde{\pi}_k(s)) w_k(s') - w_k(s) \right) \quad (17)$$

564 By definition $\mathcal{S}_k^c \subseteq \mathcal{S}^c$ and using $(1, \infty)$ -Hölder's inequality, the term (16) can be bounded as
565 (16) $\leq \|\tilde{p}_k(\cdot|s, \tilde{\pi}_k(s)) - p(\cdot|s, \tilde{\pi}_k(s))\|_{1, \mathcal{S}^c} \cdot \max_{s' \in \mathcal{S}_k^c} \{w_k(s')\}$ where for any vector $v \in \mathbb{R}^{\mathcal{S}}$,
566 $\|v\|_{1, \mathcal{S}^c} := \sum_{s \in \mathcal{S}^c} |v(s)|$. Define $\bar{s} \in \arg \max_{s \in \mathcal{S}_k^c} \{w_k(s)\}$ and $\tilde{s} \in \arg \min_{s \in \mathcal{S}_k^c} \{w_k(s)\}$. By
567 definition $\bar{s}, \tilde{s} \in \mathcal{S}_k^c$ and $w_k(\tilde{s}) = \min_{s \in \mathcal{S}_k^c} \{w_k(s)\} = 0$. By Lem. 9, we know that for all $s, s' \in \mathcal{S}_k^c$,
568 the difference $w_k(s') - w_k(s) = \tilde{h}_k(s') - \tilde{h}_k(s)$ is upper bounded by $r_{\max} \cdot \tau_{\overline{\mathcal{M}}_k^+}(s \rightarrow s')$. We also
569 know by Lem. 5 that for all $s, s' \in \mathcal{S}_k^c, \tau_{\mathcal{M}_k^+}(s \rightarrow s') = \tau_{\overline{\mathcal{M}}_k^+}(s \rightarrow s')$. Since $M^* \in \mathcal{M}_k^+$ (M^* is
570 the true MDP), we also have that for all $s, s' \in \mathcal{S}_k^c \subseteq \mathcal{S}^c, \tau_{\mathcal{M}_k^+}(s \rightarrow s') \leq \tau_{M^*}(s \rightarrow s') \leq D^c$.
571 In conclusion, $\forall s, s' \in \mathcal{S}_k^c, w_k(s') - w_k(s) \leq r_{\max} D^c$ and in particular $\max_{s' \in \mathcal{S}_k^c} \{w_k(s')\} =$
572 $w_k(\bar{s}) = w_k(\bar{s}) - w_k(\tilde{s}) \leq r_{\max} D^c$. Similarly to what we did to bound $|\tilde{r}_k - r|$ (11), we bound the
573 distance in ℓ_1 -norm between \tilde{p}_k and p by introducing \hat{p}_k :

$$\|\tilde{p}_k - p\|_{1, \mathcal{S}^c} \leq \|\tilde{p}_k - \hat{p}_k\|_{1, \mathcal{S}^c} + \|\hat{p}_k - p\|_{1, \mathcal{S}^c} \leq 2 \left(\sum_{s' \in \mathcal{S}^c} \beta_{p, k}^{s \tilde{\pi}_k(s) s'} \right) \quad (18)$$

574 We now bound the contribution of the term (17). Jaksch et al. [2] decompose this term into a
575 martingale difference sequence and a telescopic sum but due to the indicator function $\mathbb{1}_k$, in our case
576 the sum is not telescopic anymore and an additional term appears.

$$(17) = \sum_{s \in \mathcal{S}} \nu_k(s, \tilde{\pi}_k(s)) \mathbb{1}_k(s, \tilde{\pi}_k(s)) \left(\sum_{s' \in \mathcal{S}_k^c} \underbrace{p(s'|s, \tilde{\pi}_k(s)) w_k(s')}_{\geq 0, \forall s' \in \mathcal{S}_k^c} \mathbb{1}\{s \in \mathcal{S}_k^c\} - w_k(s) \mathbb{1}\{s \in \mathcal{S}_k^c\} \right)$$

$$\leq \sum_{s \in \mathcal{S}} \nu_k(s, \tilde{\pi}_k(s)) \mathbb{1}_k(s, \tilde{\pi}_k(s)) \left(\sum_{s' \in \mathcal{S}_k^c} p(s'|s, \tilde{\pi}_k(s)) w_k(s') - w_k(s) \mathbb{1}\{s \in \mathcal{S}_k^c\} \right)$$

$$= \sum_{t=t_k}^{t_{k+1}-1} \left(\sum_{s' \in \mathcal{S}_k^1} p(s'|s_t, \tilde{\pi}_k(s_t)) w_k(s') - w_k(s_t) \mathbb{1}\{s_t \in \mathcal{S}_k^c\} \right) \mathbb{1}_k(s_t, \tilde{\pi}_k(s_t))$$

$$= \sum_{t=t_k}^{t_{k+1}-1} \underbrace{\left(\sum_{s' \in \mathcal{S}_k^1} p(s'|s_t, \tilde{\pi}_k(s_t)) w_k(s') - w_k(s_{t+1}) \mathbb{1}\{s_{t+1} \in \mathcal{S}_k^c\} \right)}_{:= X_t} \mathbb{1}_k(s_t, \tilde{\pi}_k(s_t)) \quad (19)$$

$$+ \underbrace{\sum_{t=t_k}^{t_{k+1}-1} (w_k(s_{t+1}) \mathbb{1}\{s_{t+1} \in \mathcal{S}_k^c\} - w_k(s_t) \mathbb{1}\{s_t \in \mathcal{S}_k^c\}) \mathbb{1}_k(s_t, \tilde{\pi}_k(s_t))}_{\text{not telescopic due to } \mathbb{1}_k!} \quad (20)$$

577 Define the filtration $\mathcal{F}_t = \sigma(s_1, a_1, r_1, \dots, s_{t+1})$. Since k_t is \mathcal{F}_{t-1} -measurable:

$$\mathbb{E} \left[w_{k_t}(s_{t+1}) \mathbb{1}\{s_{t+1} \in \mathcal{S}_{k_t}^c\} \mathbb{1}_{k_t}(s_t, \tilde{\pi}_{k_t}(s_t)) | \mathcal{F}_{t-1} \right] = \underbrace{\sum_{s' \in \mathcal{S}_{k_t}^c} p(s' | s_t, \tilde{\pi}_{k_t}(s_t)) w_{k_t}(s') \mathbb{1}_{k_t}(s_t, \tilde{\pi}_{k_t}(s_t))}_{\mathcal{F}_{t-1}\text{-measurable}}$$

578 implying $\mathbb{E}[X_t | \mathcal{F}_{t-1}] = 0$ and so $(X_t, \mathcal{F}_t)_{t \geq 1}$ is a martingale difference sequence (MDS) with
 579 $|X_t| \leq r_{\max} D^c$. We will bound (19) in the next section (Sec. C.4) using Azuma's inequality. Using
 580 the fact that $\mathbb{1}_{k_t}(s_t, \tilde{\pi}_{k_t}(s_t)) = \mathbb{1}\{(s_t, \tilde{\pi}_{k_t}(s_t)) \notin \mathcal{K}_k\} = 1 - \mathbb{1}\{(s_t, \tilde{\pi}_{k_t}(s_t)) \in \mathcal{K}_k\}$ we can make a
 581 telescopic sum appear and rewrite (20) as:

$$\begin{aligned} (20) &= \underbrace{\sum_{t=t_k}^{t_{k+1}-1} w_k(s_{t+1}) \mathbb{1}\{s_{t+1} \in \mathcal{S}_k^c\} - w_k(s_t) \mathbb{1}\{s_t \in \mathcal{S}_k^c\}}_{=w_k(s_{t_{k+1}}) \mathbb{1}\{s_{t_{k+1}} \in \mathcal{S}_k^c\} - w_k(s_{t_k}) \mathbb{1}\{s_{t_k} \in \mathcal{S}_k^c\} \leq r_{\max} D^c} \quad (\text{telescopic sum}) \\ &+ \sum_{t=t_k}^{t_{k+1}-1} \underbrace{(w_k(s_t) \mathbb{1}\{s_t \in \mathcal{S}_k^c\} - w_k(s_{t+1}) \mathbb{1}\{s_{t+1} \in \mathcal{S}_k^c\})}_{\leq r_{\max} D^c} \mathbb{1}\{(s_t, \tilde{\pi}_k(s_t)) \in \mathcal{K}_k\} \\ &\leq r_{\max} D^c + r_{\max} D^c \sum_{s,a} \nu_k(s, a) \mathbb{1}\{(s_t, \tilde{\pi}_k(s_t)) \in \mathcal{K}_k\} \end{aligned} \quad (21)$$

582 By gathering (12), (14), (18), (19) and (21) we obtain the following bound for $\tilde{\Delta}_k$:

$$\begin{aligned} \tilde{\Delta}_k &\leq 2r_{\max} D^c \sum_{s,a} \sum_{s' \in \mathcal{S}^c} \underbrace{\mathbb{1}_k(s, a)}_{\leq 1} \underbrace{\nu_k(s, a) \beta_{p,k}^{sas'}}_{\geq 0} + \sum_{t=t_k}^{t_{k+1}-1} X_t + r_{\max} D^c \\ &+ r_{\max}(D^c + 1) \sum_{s,a} \nu_k(s, a) \mathbb{1}\{N_k^\pm(s, a) \leq \sqrt{t_k/S_A}\} + r_{\max} \sum_{s,a} \underbrace{\frac{\nu_k(s, a)}{\sqrt{t_k}}}_{\geq 0} \underbrace{\mathbb{1}_k(s, a)}_{\leq 1} \\ &\leq 2r_{\max} D^c \sum_{s,a} \sum_{s' \in \mathcal{S}^c} \nu_k(s, a) \beta_{p,k}^{sas'} + r_{\max}(D^c + 1) \sum_{s,a} \nu_k(s, a) \mathbb{1}\{(s, a) \in \mathcal{K}_k\} \\ &+ \sum_{t=t_k}^{t_{k+1}-1} X_t + r_{\max} D^c + r_{\max} \sum_{s,a} \frac{\nu_k(s, a)}{\sqrt{t_k}} := U_k^2 \end{aligned} \quad (22)$$

583 C.4 Summing over episodes with $M^* \in \mathcal{M}_k$ and $t_k \geq C(k)$

584 Denote by $\mathbb{1}(k) := \mathbb{1}\{t_k \geq C(k)\} \cdot \mathbb{1}\{M^* \in \mathcal{M}_k\}$ the indicator function taking value 1 only when
 585 both $M^* \in \mathcal{M}_k$ and $t_k \geq C(k)$. By gathering (11) and (22) we obtain:

$$\sum_{k=1}^m \Delta_k \cdot \mathbb{1}(k) \leq \sum_{k=1}^m \underbrace{U_k^1}_{\geq 0} \cdot \underbrace{\mathbb{1}(k)}_{\leq 1} + \sum_{k=1}^m \underbrace{U_k^2}_{\geq 0} \cdot \underbrace{\mathbb{1}(k)}_{\leq 1} \leq \sum_{k=1}^m U_k^1 + U_k^2 \quad (23)$$

586 and so

$$\begin{aligned} \sum_{k=1}^m \Delta_k \cdot \mathbb{1}(k) &\leq 2 \sum_{k=1}^m \sum_{s,a} \nu_k(s, a) \left(r_{\max} D^c \sum_{s' \in \mathcal{S}^c} \beta_{p,k}^{sas'} + \beta_{r,k}^{s,a} \right) + 2r_{\max} \sum_{k=1}^m \sum_{s,a} \frac{\nu_k(s, a)}{\sqrt{t_k}} \quad (24) \\ &+ r_{\max}(D^c + 1) \sum_{k=1}^m \sum_{s,a} \nu_k(s, a) \mathbb{1}\{(s, a) \in \mathcal{K}_k\} + \sum_{t=1}^T X_t \mathbb{1}(k_t) + r_{\max} m D^c \end{aligned}$$

587 We will now upper-bound the terms appearing in (24). The main novelty of (24) compared to UCRL
 588 is the term $\sum_{k=1}^m \sum_{s,a} \nu_k(s, a) \mathbb{1}\{(s, a) \in \mathcal{K}_k\}$ which is not present in the proof of Jaksch et al. [2].

589 We will show in the next section that this term is bounded by $O(\sqrt{S^c AT})$. All the other terms are
 590 similar to those found in UCRL.

591 **C.4.1 Poorly visited state-action pairs**

592 We first notice that by definition $t_{k_t} \leq t$ where $k_t := \sup\{k \geq 1 : t_k \leq t\}$ is the current episode at
 593 time t . As a result,

$$\mathbb{1}\{(s, a) \in \mathcal{K}_{k_t}\} := \mathbb{1}\left\{N_{k_t}^\pm(s_t, a_t) \leq \sqrt{t_{k_t}/SA}\right\} \leq \mathbb{1}\left\{N_{k_t}^\pm(s_t, a_t) \leq \sqrt{t/SA}\right\}$$

594 Instead of directly bounding $\sum_{k=1}^m \sum_{s,a} \nu_k(s, a) \mathbb{1}\{(s, a) \in \mathcal{K}_k\}$ we will bound the number of visits
 595 Z_T in state-action pairs that have been visited less than $\sqrt{t/SA}$ times

$$Z_T := \sum_{t=1}^T \mathbb{1}\left\{N_{k_t}^\pm(s_t, a_t) \leq \sqrt{t/SA}\right\} \geq \sum_{k=1}^m \sum_{s,a} \nu_k(s, a) \mathbb{1}\{(s, a) \in \mathcal{K}_k\}$$

596 Note that the quantity $N_k(s, a)$ is updated only after the end of episode k and the stopping condition
 597 of episodes used by TUCRL implies that (see Alg. 1):

$$\forall k \geq 1, \forall (s, a) \in \mathcal{S} \times \mathcal{A}, \nu_k(s, a) \leq N_k^+(s, a) \quad (25)$$

598 Moreover, for all $(s, a) \notin \mathcal{S}^c \times \mathcal{A}$, $\nu_k(s, a) = 0$ implying that only the states $s \in \mathcal{S}^c$ should be
 599 considered in the above sums. Using (25), we prove the following lemma:

600 **Lemma 8.** For any $T \geq 1$ and any sequence of states and actions $\{s_1, a_1, \dots, s_T, a_T\}$ we have:

$$Z_T \leq 2\sqrt{S^c AT} + 2S^c A.$$

601 *Proof.* For any episode k starting at time t_k , and for any state-action pair (s, a) we recall that $N_k(s, a)$
 602 denotes the number of visits in (s, a) prior to episode k (k not included) and by $\nu_k(s, a)$ the number
 603 of visits in (s, a) during episode k :

$$N_k(s, a) := \sum_{t=1}^{t_k-1} \mathbb{1}\{(s_t, a_t) = (s, a)\} \quad \text{and} \quad \nu_k(s, a) := \sum_{t=t_k}^{t_{k+1}-1} \mathbb{1}\{(s_t, a_t) = (s, a)\}$$

604 and so $N_k(s, a) = \sum_{i=1}^{k-1} \nu_i(s, a)$. By convention, we denote by $N_{k_T+1}(s, a) := \sum_{t=1}^T \mathbb{1}\{(s_t, a_t) =$
 605 $(s, a)\}$ the total number of visits in (s, a) after T time steps (T included). We first decompose Z_T as:

$$Z_T := \sum_{s,a} \sum_{t=1}^T \mathbb{1}\left\{\max\{1, N_{k_t}(s, a) - 1\} \leq \sqrt{t/SA}\right\} \cdot \mathbb{1}\{(s_t, a_t) = (s, a)\} = \sum_{s \in \mathcal{S}^c} \sum_a Z_T(s, a)$$

$$\text{where } Z_T(s, a) := \sum_{t=1}^T \mathbb{1}\left\{\max\{1, N_{k_t}(s, a) - 1\} \leq \sqrt{t/SA}\right\} \cdot \mathbb{1}\{(s_t, a_t) = (s, a)\}$$

606 Using the fact that for all $t \geq 1$, $t_{k_t} \leq t \leq t_{k_t+1} - 1$ we have:

$$\begin{aligned} \forall T \geq \tau \geq 1, \quad Z_\tau(s, a) &= \sum_{t=1}^{\tau} \underbrace{\mathbb{1}\left\{\max\{1, N_{k_t}(s, a) - 1\} \leq \sqrt{t/SA}\right\}}_{\leq 1} \cdot \underbrace{\mathbb{1}\{(s_t, a_t) = (s, a)\}}_{\geq 0} \\ &\leq \sum_{t=1}^{\tau} \mathbb{1}\{(s_t, a_t) = (s, a)\} \leq \sum_{t=1}^{t_{k_\tau+1}-1} \mathbb{1}\{(s_t, a_t) = (s, a)\} \\ &= N_{k_\tau+1}(s, a) \end{aligned} \quad (26)$$

607 Let's define $t_{s,a}$ as the last time that $Z_t(s, a)$ was incremented by 1:

$$\begin{aligned} t_{s,a} &:= \max\left\{T \geq t \geq 1 : \max\{1, N_{k_t}(s, a) - 1\} \leq \sqrt{t/SA} \text{ and } (s_t, a_t) = (s, a)\right\} \\ &= \min\left\{T \geq t \geq 1 : Z_t(s, a) = Z_T(s, a)\right\} \end{aligned}$$

608 We denote by $m_{s,a} := k_{t_{s,a}}$ the corresponding episode. By definition,

$$Z_T(s, a) = Z_{t_{s,a}}(s, a) \quad (27)$$

609 and

$$\max\{1, N_{m_{s,a}}(s, a) - 1\} \leq \sqrt{t_{s,a}/SA} \quad (28)$$

610 Using (26) with $\tau = t_{s,a}$ we obtain:

$$Z_{t_{s,a}} \leq N_{m_{s,a}+1}(s, a) \quad (29)$$

611 Moreover, by definition of $N_k(s, a)$ and (25):

$$\begin{aligned} N_{m_{s,a}+1}(s, a) &= N_{m_{s,a}}(s, a) + \underbrace{\nu_{m_{s,a}}(s, a)}_{\leq N_{m_{s,a}}^+(s, a)} \leq 2 \underbrace{\max\{1, N_{m_{s,a}}(s, a)\}}_{\leq \max\{1, N_{m_{s,a}}(s, a) - 1\} + 1} \\ &\implies N_{m_{s,a}+1}(s, a) \leq 2 \cdot \max\{1, N_{m_{s,a}}(s, a) - 1\} + 2 \end{aligned} \quad (30)$$

612 Gathering (27), (28), (29), and (30) we obtain:

$$\begin{aligned} Z_T(s, a) &= Z_{t_{s,a}}(s, a) \leq \max\{1, N_{m_{s,a}+1}(s, a) - 1\} + 1 \leq 2 \cdot \max\{1, N_{m_{s,a}}(s, a) - 1\} + 2 \\ &\leq 2\sqrt{t_{s,a}/SA} + 2 \\ &\leq 2\sqrt{T/SA} + 2 \\ &\implies Z_T = \sum_{s \in S^c} \sum_a Z_T(s, a) \leq 2\sqrt{S^c AT} + 2S^c A \end{aligned}$$

613 where for the last inequality we used the fact that $S^c \leq S$ (by definition) implying $S^c/\sqrt{S} =$
614 $\sqrt{S^c/S} \cdot \sqrt{S} \leq \sqrt{S^c}$. This concludes the proof. \square

615 As a consequence of Lem. 8:

$$\sum_{k=1}^m \sum_{s,a} \nu_k(s, a) \mathbb{1}\{N_k^\pm(s, a) \leq \sqrt{t_k/S^c A}\} \leq Z_T \leq 2\sqrt{S^c AT} + 2S^c A \quad (31)$$

616 C.4.2 Confidence bounds $\beta_{r,k}^{sa}$ and $\beta_{p,k}^{sas'}$

617 Since (25) holds, Lemma 19 of Jaksch et al. [2] can still be applied. Moreover, exploiting again the
618 fact that for all $(s, a) \notin S^c \times \mathcal{A}$, $\nu_k(s, a) = 0$ we obtain

$$\sum_{k=1}^m \sum_{s,a} \frac{\nu_k(s, a)}{\sqrt{t_k}} \leq \sum_{k=1}^m \sum_{s,a} \frac{\nu_k(s, a)}{\sqrt{N_k^+(s, a)}} \leq (\sqrt{2} + 1) \sqrt{S^c AT} \quad (32)$$

619 and as shown in [6, Appendix F.7] (with the difference that S is restricted to S^c) we have:

$$\sum_{k=1}^m \sum_{s,a} \frac{\nu_k(s, a)}{N_k^\pm(s, a)} \leq 6S^c A + 2S^c A \ln(T) \quad (33)$$

620 The terms $\sum_{k=1}^m \sum_{s,a} \nu_k(s, a) \beta_{r,k}^{sa}$ and $\sum_{k=1}^m \sum_{s,a,s' \in S^c} \nu_k(s, a) \beta_{p,k}^{sas'}$ can then be bounded exactly
621 as in [6, App. F.7] with S replaced by S^c (except in the logarithm).

622 C.4.3 Number of episodes

623 The stopping condition of episodes used by TUCRL (see Alg. 1) combines the original stopping
624 condition of UCRL with the condition $s_t \in \mathcal{S}_{k_t}^T$. Using only inequality (25), Jaksch et al. [2, Figure
625 1] proved that for any any sequence $\{s_1, a_1, \dots, s_T, a_T\}$, the number of episodes is bounded by
626 $1 + 2SA + SA \log_2 \left(\frac{T}{SA}\right)$. Since (25) also holds in our case, the total number of episodes m after T
627 time steps can be bounded by the same quantity (with S replaced by S^c since states in \mathcal{S}^T will never
628 be visited) plus the number of times the event $s_t \in \mathcal{S}_{k_t}^T$ occurs. Since whenever $s_t \in \mathcal{S}_{k_t}^T$ state s_t is
629 removed from $\mathcal{S}_{k_{t+1}}^T$ and s_t necessarily belongs to S^c (by definition), this event can happen at most
630 S^c times. By Proposition 18 in [2] we thus have:

$$m \leq 1 + 2S^c A + S^c A \log_2 \left(\frac{T}{S^c A}\right) + S^c \quad (34)$$

631 **C.4.4 Martingale Difference Sequence $X_t \cdot \mathbb{1}(k_t)$**

632 In Sec. C.3 we already proved that $(X_t, \mathcal{F}_t)_{t \geq 1}$ is an MDS i.e., for all $t \geq 1$, $\mathbb{E}[X_t | \mathcal{F}_{t-1}] = 0$.
 633 Since k_t is \mathcal{F}_{t-1} -measurable, we also have $\mathbb{E}[X_t \mathbb{1}(k_t) | \mathcal{F}_{t-1}] = \mathbb{1}(k_t) \cdot \mathbb{E}[X_t | \mathcal{F}_{t-1}] = 0$ with
 634 $|X_t \mathbb{1}(k_t)| \leq r_{\max} D^c$. Therefore, $(X_t \mathbb{1}(k_t), \mathcal{F}_t)_{t \geq 1}$ is also an MDS. By Azuma's inequality (see for
 635 example [2, Lemma 10]):

$$\sum_{t=1}^T X_t \mathbb{1}(k_t) \leq r_{\max} D^c \sqrt{\frac{5}{2} T \ln \left(\frac{8T}{\delta} \right)} \quad \text{w.p.} \geq 1 - \frac{\delta}{12T^{5/4}} \quad (35)$$

636 **C.5 Completing the regret bound**

637 By gathering (24), (31), (32), (33), (35) and (34) we conclude that with probability at least $1 - \frac{\delta}{12T^{5/4}}$:

$$\begin{aligned} \sum_{k=1}^m \Delta_k \cdot \mathbb{1}(k) &\leq 2 \left(\sqrt{28} + \sqrt{14} \right) r_{\max} \sqrt{S^c A T \ln \left(\frac{2SAT}{\delta} \right)} \left(D^c \sqrt{\Gamma^c - 1} + 1 \right) \\ &\quad + \frac{196}{3} r_{\max} S^c A \ln \left(\frac{2SAT}{\delta} \right) (3 + \ln(T)) (D^c S^c + 1) \\ &\quad + 2r_{\max} (D^c + 1) (\sqrt{S^c A T} + S^c A) \\ &\quad + r_{\max} D^c \sqrt{\frac{5}{2} T \ln \left(\frac{8T}{\delta} \right)} + 2 \left(\sqrt{2} + 1 \right) r_{\max} \sqrt{S^c A T} \\ &\quad + r_{\max} D^c \left(1 + 2S^c A + S^c A \log_2 \left(\frac{T}{SA} \right) + S^c \right) \\ &\leq C \cdot \left(r_{\max} D^c \sqrt{\Gamma^c S^c A T \ln \left(\frac{SAT}{\delta} \right)} + r_{\max} D^c (S^c)^2 A \ln^2 \left(\frac{SAT}{\delta} \right) \right) \end{aligned} \quad (36)$$

638 where C is a numerical constant independent of the MDP instance.

639 From (8), with probability at least $1 - \frac{\delta}{12T^{5/4}}$:

$$\begin{aligned} \Delta(\text{TUCRL}, T) &\leq \sum_{k=1}^m \Delta_k + r_{\max} \sqrt{\frac{5}{2} T \ln \left(\frac{8T}{\delta} \right)} \\ &= \underbrace{\sum_{k=1}^m \Delta_k \mathbb{1}(k)}_{\text{see (36)}} + \sum_{k=1}^m \Delta_k \cdot (1 - \mathbb{1}(k)) + r_{\max} \sqrt{\frac{5}{2} T \ln \left(\frac{8T}{\delta} \right)} \end{aligned}$$

640 where $1 - \mathbb{1}(k)$ is the complement of $\mathbb{1}(k)$ i.e., takes value 1 only when either $t_k < C(k)$ (see (10)
 641 for the definition of $C(k)$) or $M^* \notin \mathcal{M}_k$. As is proved in Appendix F.2 of [6], since both (25) and
 642 Theorem 1 of Fruit et al. [6] hold, we have that with probability at least $1 - \frac{\delta}{20T^{5/4}} \geq 1 - \frac{\delta}{12T^{5/4}}$:

$$\sum_{k=1}^m \Delta_k \mathbb{1}\{M^* \notin \mathcal{M}_k\} \leq r_{\max} \sqrt{T} \quad (37)$$

643 As a consequence of (25) $t_{k+1} \leq 2t_k$. Thus, by definition of the condition $t_k < C(k)$ we have

$$\sum_{k=1}^m \Delta_k \cdot \underbrace{\mathbb{1}\{t_k < C(k)\}}_{\geq 0} \leq 2r_{\max} C(k) \leq \frac{4802}{9} r_{\max} (D^c)^2 S^3 A \ln^2 \left(\frac{2SAT}{\delta} \right) \quad (38)$$

644 Finally, by Boole's inequality: $1 - \mathbb{1}(k) \leq \mathbb{1}\{M^* \notin \mathcal{M}_k\} + \mathbb{1}\{t_k < C(k)\}$ and so

$$\sum_{k=1}^m \Delta_k \cdot (1 - \mathbb{1}(k)) \leq \underbrace{\sum_{k=1}^m \Delta_k \cdot \mathbb{1}\{M^* \notin \mathcal{M}_k\}}_{\text{see (37)}} + \underbrace{\sum_{k=1}^m \Delta_k \cdot \mathbb{1}\{t_k < C(k)\}}_{\text{see (38)}}$$

Algorithm 2 OPTIMISTIC TRANSITION PROBABILITIES (OTP) [2]

Input: Probability estimate $\hat{p} \in \mathbb{R}^n$, confidence interval $\beta \in \mathbb{R}$, value vector $v \in \mathbb{R}^n$, subset of states $\mathcal{I} \subseteq \{s_1, \dots, s_m\}$, $m \leq n$, such that $\sum_{s \in \mathcal{I}} \hat{p}(s) = 1$
Output: Optimistic probabilities $\tilde{p} \in \mathbb{R}^n$

Let $\mathcal{I} = \{s_1, s_2, \dots, s_m\}$ such that $v(s_1) \geq v(s_2) \geq \dots \geq v(s_m)$

$$\tilde{p}_1(s_1) = \min \left\{ 1, \hat{p}(s_1) + \frac{\beta}{2} \right\}$$

$$\tilde{p}_1(s_j) = \hat{p}(s_j), \quad \forall 1 < j \leq m$$

$$j = m$$

$$i = 1$$

while $\sum_{s \in \mathcal{I}} \tilde{p}_i(s) > 1$ **do**

$$i = i + 1$$

$$\tilde{p}_i(s) = \tilde{p}_{i-1}(s), \quad \forall s \neq s_j$$

$$\tilde{p}_i(s_j) = \max \left\{ 0, 1 - \sum_{s \in \mathcal{I} \setminus \{s_j\}} \tilde{p}_{i-1}(s) \right\}$$

$$j = j - 1$$

end while

$$\tilde{p}_i(s) := 0, \quad \forall s \in \mathcal{S} \setminus \mathcal{I}$$

$$\tilde{p} := \tilde{p}_i$$

645 In conclusion, there exists a numerical constant C independent of the MDP instance such that for any
646 MDP and any $T > 1$, with probability at least $1 - \frac{\delta}{12T^{5/4}} - \frac{\delta}{12T^{5/4}} - \frac{\delta}{12T^{5/4}} = 1 - \frac{\delta}{4T^{5/4}}$ we have:

$$\Delta(\text{TUCRL}, T) \leq C \cdot \left(r_{\max} D^c \sqrt{\Gamma S^c A T \ln \left(\frac{SAT}{\delta} \right)} + r_{\max} (D^c)^2 S^3 A \ln^2 \left(\frac{SAT}{\delta} \right) \right) \quad (39)$$

647 Since $\sum_{T=2}^{+\infty} \frac{\delta}{4T^{5/4}} = \delta$, by taking a union bound we have that the regret bound (39) holds with
648 probability at least $1 - \delta$ for all $T > 1$.

649 D Shortest Path Analysis

650 We are interesting in comparing the shortest path of any pair $(s, \bar{s}) \in \mathcal{S} \times \mathcal{S}_k^c$ in \mathcal{M}_k^+ and $\overline{\mathcal{M}}_k^+$.
651 Formally, given a target state \bar{s} , the stochastic shortest path $\tau_M(s) := \tau_M(s \rightarrow \bar{s})$ of an (extended)
652 MDP M is the (negation) solution of the following Bellman equation

$$\begin{aligned} \tau_M(s) &= -1 + \max_{a \in \mathcal{A}_s, p \in B_p(s, a)} \{p^T \tau_M\}, \quad \forall s \neq \bar{s} \\ \tau_M(\bar{s}) &= 0 \end{aligned} \quad (40)$$

653 D.1 Proof of Lem. 5

654 In order to analyse the properties of the stochastic shortest path we need to investigate the maximiza-
655 tion over the confidence interval $B_p(s, a)$ either in \mathcal{M}_k^+ or $\overline{\mathcal{M}}_k^+$. This problem can be solved using
656 Alg. 2. For any state-action pair (s, a) , we define $\tilde{p}_{\mathcal{M}_k^+}(\cdot | s, a) = \text{OTP}(\hat{p}(\cdot | s, a), B_{p,k}^+(s, a), \tau, \mathcal{S})$ and
657 $\tilde{p}_{\overline{\mathcal{M}}_k^+}(\cdot | s, a) = \text{OTP}(\hat{p}(\cdot | s, a), \overline{B}_{p,k}^+(s, a), \tau, \mathcal{S})$. It is easy to notice that the optimistic probability
658 vectors built by Alg. 2 satisfy (either in \mathcal{M}_k^+ or in $\overline{\mathcal{M}}_k^+$)

$$\forall i \in \{1, \dots, n\}, \quad \tilde{p}_i(s_1) \geq \hat{p}(s_1)$$

$$\begin{aligned} \forall i \in \{2, \dots, n\}, \forall l \in \{n - i + 2, n\}, \quad \tilde{p}_i(s_l) &= \max \left\{ 0, 1 - \sum_{s' \neq s_l} \tilde{p}_{i-1}(s') \right\} \\ &= \max \left\{ 0, \hat{p}(s_l) - \left(\sum_{s'} \tilde{p}_{i-1}(s') - 1 \right) \right\} \\ &\leq \hat{p}(s_l) \end{aligned}$$

659 where s_1, \dots, s_n are such that $\tau(s_1) \geq \dots \geq \tau(s_n)$. The algorithm may stop before n iterations but
 660 this means that the states not processed are kept at \hat{p} .

661 We start considering the case in which $(s, a) \in \mathcal{K}_k$. Recall that $\forall s' \in \mathcal{S}_k^T$, $\hat{p}(s'|s, a) = 0$ by
 662 definition since s' is not reachable from \mathcal{S}_k^C (i.e., $N_k(s, a, s') = 0$) and that \mathcal{M}_k^+ and $\overline{\mathcal{M}}_k^+$ consider
 663 the same empirical average for the transition probabilities (i.e., \hat{p}). The shortest path to \bar{s} is such that
 664 $\max_s \{\tau(s)\} = \tau(\bar{s}) = 0$ and $\tau(s) \leq -1$ for any state $s \in \mathcal{S} \setminus \{\bar{s}\}$ (either in \mathcal{M}_k^+ or $\overline{\mathcal{M}}_k^+$). As a
 665 consequence, $s_1 = \bar{s}$ and for any $s' \in \mathcal{S}_k^T$,

$$\tilde{p}_{\mathcal{M}_k^+}(s') \leq \hat{p}(s') = 0, \text{ and } \tilde{p}_{\overline{\mathcal{M}}_k^+}(s') \leq \hat{p}(s') = 0$$

666 which ensures that $\forall (s, a) \in \mathcal{K}_k$ the constraints in $\overline{\mathcal{M}}_k^+$ hold. This results is independent from the
 667 vector v provided to OTP. Then, for any vector $v \in V = \{v \in \mathbb{R}^{\mathcal{S}} | v(\bar{s}) = 0 \wedge v(s) \leq -1, \forall s \in$
 668 $\mathcal{S} \setminus \{\bar{s}\}\}$, we have that $\mathcal{I}^1 = \mathcal{I}^2$, since $\beta^1 = \beta^2$ and $\tilde{p}_{\mathcal{M}_k^+}(s) = \tilde{p}_{\overline{\mathcal{M}}_k^+}(s) = 0$ for any $s \in \mathcal{S}_k^2$ then:
 669 $\tilde{p}_{\mathcal{M}_k^+}(s') = \tilde{p}_{\overline{\mathcal{M}}_k^+}(s')$, $\forall s' \in \mathcal{S}$. Finally, $\forall (s, a) \in (\mathcal{S} \times \mathcal{A}) \setminus \mathcal{K}_k$ it is trivial to notice that: $\forall s' \in \mathcal{S}$,
 670 $\forall v \in V$, $\tilde{p}_{\mathcal{M}_k^+}(s') = \tilde{p}_{\overline{\mathcal{M}}_k^+}(s')$ since $B_{p,k}^+(s, a) = \overline{B}_{p,k}^+(s, a)$.

671 The proof follows by noticing that $\tau_{\mathcal{M}_k^+} \in V$ and $\tau_{\overline{\mathcal{M}}_k^+} \in V$.

672 D.2 Bounding the bias span

673 **Lemma 9.** Consider an (extended) MDP M and define L_M as the associated optimal (extended)
 674 Bellman operator. Given $h_0 = \mathbf{0}$, and $h_i = (L_M)^i h_0$ we have that

$$\forall s, s' \in \mathcal{S}, \quad h_i(s') - h_i(s) \leq r_{\max} \tau_M(s \rightarrow s')$$

675 where $\tau_M(s \rightarrow s')$ is the minimum expected shortest path from s to s' in M .

676 *Proof.* The proof follows from the application of the argument in [2, Sec. 4.3.1]. \square

677 E Proof of Lem. 6

678 We prove the statement by contradiction: we assume that there exists a learning algorithm denoted
 679 \mathfrak{A}_T satisfying

- 680 1. for all $\varepsilon \in]0, 1]$, there exists $T_\varepsilon^\dagger \leq f(\varepsilon)$ such that $\mathbb{E}[\Delta(\mathfrak{A}_T, M_\varepsilon, x, T)] < 1/6 \cdot T$ for all $T \geq T_\varepsilon^\dagger$,
- 681 2. there exists $T_0^* < +\infty$ such that $\mathbb{E}[\Delta(\mathfrak{A}_T, M_0, x, T)] \leq C_2(\ln(T))^\beta$ for all $T \geq T_0^*$.

682 Any randomised strategy for choosing an action at time t is equivalent to an (a priori) random
 683 choice from the set of all deterministic strategies. Thus, it is sufficient to show a contradiction
 684 when the action played by \mathfrak{A}_T at any time t is a deterministic function of the past trajectory $h_t :=$
 685 $\{s_1, a_1, r_1, \dots, s_t\}$. In the rest of the proof we assume that \mathfrak{A}_T maps any sequence of observations
 686 $h_t = \{s_1, a_1, r_1, \dots, s_t\}$ to a (single) action a_t .

687 By trivial induction it is easy to see that as long as state y has not been visited, the history h_t is
 688 independent of ε (\mathfrak{A}_T can not distinguish between different values of ε and plays exactly the same
 689 action when the past history is the same).

690 Let's define $N_T^0(x, b) := \sum_{t=1}^T \mathbb{1}\{(s_t, a_t) = (x, b)\}$ the number of visits in (x, b) with $a_t = \mathfrak{A}_T(h_t)$
 691 and $\varepsilon = 0$. Note that $N_T^0(x, b)$ is not random since when $\varepsilon = 0$ both action b and action d loop on x
 692 with probability 1. For any $\varepsilon \in [0, 1]$ and any horizon T define the event:

$$F(T, \varepsilon) := \bigcap_{1 \leq t \leq T} \{s_t \neq y\}$$

693 where the sequence of states s_t is obtained by executing \mathfrak{A}_T on MDP M_ε . We will denote by $\overline{F(T, \varepsilon)}$
 694 the complement of $F(T, \varepsilon)$.

695 For any horizon T , and independently of ε , there is only one possible trajectory $h_T =$
 696 $\{s_1, a_1, r_1, \dots, s_T\}$ that never goes to y and which corresponds to the trajectory observed when

697 $\varepsilon = 0$. When $\varepsilon = 0$, the probability of this trajectory is 1 and so $\mathbb{P}(F(T, 0)) = 1$ (recall that
 698 everything is deterministic in this case) while in general we have:

$$\forall T \geq 1, \forall \varepsilon \in [0, 1], \mathbb{P}(F(T, \varepsilon)) = (1 - \varepsilon)^{N_T^0(x, b)} \quad (41)$$

699 We now prove by contradiction that

$$\lim_{T \rightarrow +\infty} N_T^0(x, b) = +\infty \quad (42)$$

700 Let's assume that $C := \max\{10, \max_{T \geq 1}\{N_T^0(x, b)\}\} < +\infty$. Taking $\varepsilon = 1/C$ and applying the
 701 law of total expectation we obtain:

$$\begin{aligned} \forall T \geq 1, \mathbb{E}[\Delta(\mathfrak{A}_T, M_{1/C}, x, T)] &= \underbrace{\mathbb{E}[\Delta(\mathfrak{A}_T, M_{1/C}, x, T) | F(T, 1/C)]}_{=T/2+1/2 \cdot N_T^0(x, b) \geq T/2} \cdot \underbrace{\mathbb{P}(F(T, 1/C))}_{=(1-1/C)^{N_T^0(x, b)}} \\ &\quad + \underbrace{\mathbb{E}[\Delta(\mathfrak{A}_T, M_{1/C}, x, T) | \overline{F(T, 1/C)}]}_{\geq 0} \cdot \mathbb{P}(\overline{F(T, 1/C)}) \\ &\geq \frac{T}{2} \cdot \left(1 - \frac{1}{C}\right)^{N_T^0(x, b)} \geq \frac{T}{2} \cdot \underbrace{\left(1 - \frac{1}{C}\right)^C}_{\geq 1/3 \text{ by Lem. 10}} \geq \frac{T}{6} \end{aligned}$$

702 where we used the fact that

- 703 • $N_T^0(x, b) \leq C$ and $(1 - 1/C) \in [0, 1]$ by definition, implying $(1 - \frac{1}{C})^{N_T^0(x, b)} \leq (1 - \frac{1}{C})^C$,
- 704 • since $C \geq 10$ we have $(1 - \frac{1}{C})^C \geq 1/3$ by Lem. 10 applied to $x = 1/C$,
- 705 • and finally under event $F(T, 1/C)$, the regret incurred is exactly $T/2 + 1/2 \cdot N_T^0(x, b) \geq$
 706 $T/2$.

707 This contradicts our assumption that there exists $T_{1/C}^\dagger < +\infty$ such that for all $T \geq T_{1/C}^\dagger$,
 708 $\mathbb{E}[\Delta(\mathfrak{A}_T, M_{1/C}, x, T)] < T/6$ and so (42) holds.

709 Since $\lim_{T \rightarrow +\infty} N_T^0(x, b) = +\infty$, it is possible to construct a strictly increasing sequence $(T_n)_{n \in \mathbb{N}}$
 710 such that:

$$\forall n \in \mathbb{N}, N_{T_{n+1}}^0(x, b) > N_{T_n}^0(x, b), T_0 = T_0^*, T_1 \geq C_2, T_1 \geq C_2(\ln(T_1))^\beta \text{ and } N_{T_1}^0(x, b) \geq 10$$

711 We also define the (strictly decreasing) sequence: $\varepsilon_n := 1/N_{T_n}^0(x, b)$, $\forall n \geq 1$. By the law of total
 712 expectation:

$$\begin{aligned} \mathbb{E}[\Delta(\mathfrak{A}_{T_n}, M_{\varepsilon_n}, x, T_n)] &= \underbrace{\mathbb{E}[\Delta(\mathfrak{A}_{T_n}, M_{\varepsilon_n}, x, T_n) | F(T_n, \varepsilon_n)]}_{\geq T_n/2} \cdot \underbrace{\mathbb{P}(F(T_n, \varepsilon_n))}_{=(1-\varepsilon_n)^{N_{T_n}^0(x, b)}} \\ &\quad + \underbrace{\mathbb{E}[\Delta(\mathfrak{A}_{T_n}, M_{\varepsilon_n}, x, T_n) | \overline{F(T_n, \varepsilon_n)}]}_{\geq 0} \cdot \mathbb{P}(\overline{F(T_n, \varepsilon_n)}) \\ &\geq \frac{T_n}{2} \cdot (1 - \varepsilon_n)^{N_{T_n}^0(x, b)} = \frac{T_n}{2} \cdot \underbrace{(1 - \varepsilon_n)^{1/\varepsilon_n}}_{\geq 1/3 \text{ by Lem. 10}} \geq \frac{T_n}{6} \end{aligned} \quad (43)$$

713 where we applied Lem. 10 to $x = \varepsilon_n \leq 1/10$ since $N_{T_n}^0(x, b) \geq 10$ for all $n \geq 1$. Moreover, since
 714 by construction for all $n \geq 1$, $T_n > T_0 = T_0^*$ we have by assumption that

$$\begin{aligned} \forall n \geq 1, \mathbb{E}[\Delta(\mathfrak{A}_{T_n}, M_0, x, T_n)] &= \frac{1}{2} N_{T_n}^0(x, b) = \frac{1}{2\varepsilon_n} \leq C_2(\ln(T_n))^\beta \\ \implies T_n &\geq \exp\left(\frac{1}{(2C_2 \cdot \varepsilon_n)^{1/\beta}}\right) \end{aligned}$$

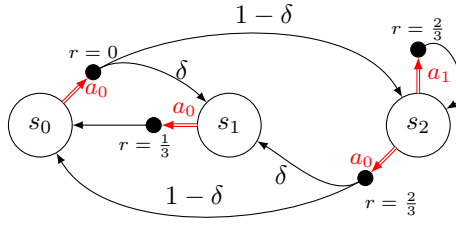


Figure 5: Three-state domain introduced in [6]

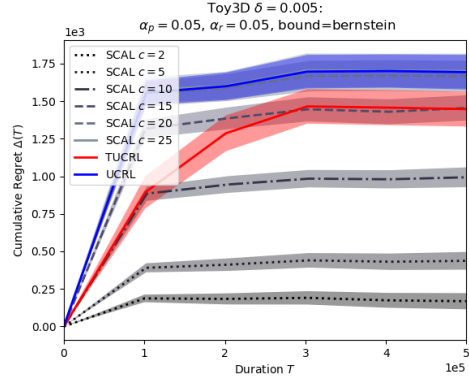


Figure 6: Communicating three-state domain ($\delta = 0.005$)

715 Since $\lim_{n \rightarrow +\infty} 1/\varepsilon_n = +\infty$ and $\lim_{x \rightarrow +\infty} \exp(x^{1/\beta})/x^\alpha = +\infty$ there exists $N \in \mathbb{N}$ such that
 716 for all $n \geq N$, $T_n \geq f(\varepsilon_n)$. By assumption, for all $n \geq N$,

$$\mathbb{E}[\Delta(\mathfrak{A}_{T_n}, M_{\varepsilon_n}, x, T_n)] < \frac{T_n}{6}$$

717 which contradicts (43) therefore concluding the proof.

718 **Lemma 10.** For all $x \in]0, 1/10]$, we have $(1-x)^{1/x} \geq 1/3$.

719 *Proof.* It is easy to verify that the derivative of $x \mapsto (1-x)^{1/x}$ is:

$$\forall x \in]0, 1/10], \frac{d}{dx} \left((1-x)^{1/x} \right) = - \underbrace{\frac{(1-x)^{1/x-1}}{x^2}}_{\geq 0} \cdot ((1-x) \ln(1-x) + x)$$

720 It is well known that for all $x \in]0, 1[$, $x < -\ln(1-x) < \frac{x}{1-x}$ implying that $(1-x) \ln(1-x) + x$
 721 is positive. Therefore, $\frac{d}{dx} \left((1-x)^{1/x} \right)$ is negative on $]0, 1/10]$ implying that $x \mapsto (1-x)^{1/x}$ is
 722 decreasing. As a result: $\forall x \in]0, 1/10]$, $(1-x)^{1/x} \geq 0.9^{10} > 1/3$. \square

723 F Experiments - Three-State Domain

724 This domain was introduced in [6] in order to show the inability of UCRL to learn in weakly
 725 communicating MDPs. The graphical representation of the domain is reported in Fig. 5. We keep the
 726 same means for the rewards (reported on Fig. 5) but we change the distributions: uniform distributions
 727 with range $1/5$ instead of Bernouillis. In the main paper we showed how the algorithms behave when
 728 $\delta = 0$. Here we consider the case the MDP is communicating by defining $\delta = 0.005$. Fig. 6 shows
 729 that, as expected, TUCRL behaves similarly to UCRL. In this example it is able to outperform UCRL
 730 since the preliminary phase in which transitions to non-observed states are forbidden leads to a more
 731 conservative exploration that, due to the structure of the problem (s_1 is difficult to reach but it is also
 732 non-optimal), results in a smaller regret.