

# OPTIMAL REGRET MINIMIZATION STRATEGIES IN MARKOV DECISION PROCESSES.

Odalric-Ambrym Maillard

JUNE 07, 2023

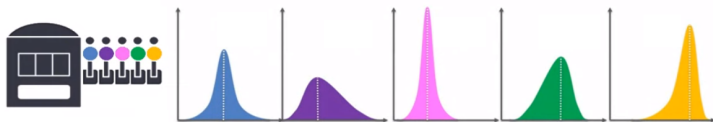
Inria Scool

- ▷ Pesquerel, F. and Maillard, O.A., 2022, November. *IMED-RL: Regret optimal learning of ergodic Markov decision processes*. In *NeurIPS 2022-Thirty-sixth Conference on Neural Information Processing Systems*.
- ▷ Honda, J. and Takemura, A., 2015. *Non-asymptotic analysis of a new bandit algorithm for semi-bounded rewards*. *J. Mach. Learn. Res.*, 16, pp.3721-3756.
- ▷ Agrawal, R., 1990, December. *Adaptive control of Markov chains under the weak accessibility*. In *29th IEEE Conference on Decision and Control* (pp. 1426-1431). IEEE.



# Warming-up

Model  $\mathbf{M} = (\mathcal{A}, \mathbf{r})$  where  $\mathcal{A} = \{1, \dots, 5\}$ ,  $\mathbf{r} : \mathcal{A} \rightarrow \mathcal{P}([0, 1])$ .



Consider a score, e.g. the mean,  $\mathbf{m} : \mathcal{A} \rightarrow [0, 1]$

$$\text{Let } \mathbf{m}_* = \max_{a \in \mathcal{A}} \mathbf{m}(a) \text{ and } \mathcal{O}(\mathbf{M}) = \text{Argmax}_{a \in \mathcal{A}} \mathbf{m}(a).$$

At each decision time  $t$ , choose  $A_t \in \mathcal{A}$ , receive  $X_t \sim \mathbf{r}(A_t)$ .

After a while

Form the trajectory  $\tau_t = (A_1, X_1, \dots, A_t, X_t)$ , and counts  $N_t(a) = \sum_{t'=1}^t \mathbb{I}\{A_{t'} = a\}$ .  
Try to guess e.g. the best distribution.



## Warming-up

**Log-likelihood** of model  $\mathbf{M} = (\mathcal{A}, \mathbf{r})$  versus  $\tilde{\mathbf{M}} = (\mathcal{A}, \tilde{\mathbf{r}})$  on  $\tau_T$ :

$$\log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}(\tau_T)} = \sum_{t=1}^T \log \frac{\mathbf{r}(A_t)(X_t)}{\tilde{\mathbf{r}}(A_t)(X_t)},$$
$$L_T(\mathbf{M}, \tilde{\mathbf{M}}) = \mathbb{E}_{\mathbf{M}} \left[ \log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}(\tau_T)} \right] = \sum_a \mathbb{E}_{\mathbf{M}} \left[ N_T(a) \right] \text{KL}(\mathbf{r}(a), \tilde{\mathbf{r}}(a)).$$



## Warming-up

**Log-likelihood** of model  $\mathbf{M} = (\mathcal{A}, \mathbf{r})$  versus  $\tilde{\mathbf{M}} = (\mathcal{A}, \tilde{\mathbf{r}})$  on  $\tau_T$ :

$$\log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}(\tau_T)} = \sum_{t=1}^T \log \frac{\mathbf{r}(A_t)(X_t)}{\tilde{\mathbf{r}}(A_t)(X_t)},$$
$$L_T(\mathbf{M}, \tilde{\mathbf{M}}) = \mathbb{E}_{\mathbf{M}} \left[ \log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}(\tau_T)} \right] = \sum_a \mathbb{E}_{\mathbf{M}} \left[ N_T(a) \right] \text{KL}(\mathbf{r}(a), \tilde{\mathbf{r}}(a)).$$

### Game

$\mathcal{B}(a, \mathbf{M})$ : For some  $a \notin \mathcal{O}(\mathbf{M})$ , find  $\tilde{\mathbf{M}} = \tilde{\mathbf{M}}_a$  such that

- ▶  $\forall a \in \mathcal{O}(\mathbf{M}), \mathbf{r}(a)$  and  $\tilde{\mathbf{r}}(a)$  undistinguishable
- ▶  $a \in \mathcal{O}(\tilde{\mathbf{M}})$

Say  $\mathcal{O}(\mathbf{M}) = \{5\}$ , then  $\tilde{\mathbf{r}}(5) = \mathbf{r}(5)$ ,  $\tilde{\mathbf{r}}(a)$  has mean  $\tilde{\mathbf{m}}(a) > \mathbf{m}(5)$ .

### Closest?



# Warming-up

**Log-likelihood** of model  $\mathbf{M} = (\mathcal{A}, \mathbf{r})$  versus  $\tilde{\mathbf{M}} = (\mathcal{A}, \tilde{\mathbf{r}})$  on  $\tau_T$ :

$$\log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}(\tau_T)} = \sum_{t=1}^T \log \frac{\mathbf{r}(A_t)(X_t)}{\tilde{\mathbf{r}}(A_t)(X_t)},$$
$$L_T(\mathbf{M}, \tilde{\mathbf{M}}) = \mathbb{E}_{\mathbf{M}} \left[ \log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}(\tau_T)} \right] = \sum_a \mathbb{E}_{\mathbf{M}} \left[ N_T(a) \right] \text{KL}(\mathbf{r}(a), \tilde{\mathbf{r}}(a)).$$

## Game

$\mathcal{B}(a, \mathbf{M})$ : For some  $a \notin \mathcal{O}(\mathbf{M})$ , find  $\tilde{\mathbf{M}} = \tilde{\mathbf{M}}_a$  such that

- ▶  $\forall a \in \mathcal{O}(\mathbf{M}), \mathbf{r}(a)$  and  $\tilde{\mathbf{r}}(a)$  undistinguishable
- ▶  $a \in \mathcal{O}(\tilde{\mathbf{M}})$

Say  $\mathcal{O}(\mathbf{M}) = \{5\}$ , then  $\tilde{\mathbf{r}}(5) = \mathbf{r}(5)$ ,  $\tilde{\mathbf{r}}(a)$  has mean  $\tilde{\mathbf{m}}(a) > \mathbf{m}(5)$ .

## Closest?

$$\inf_{\tilde{\mathbf{M}} \in \mathcal{B}(a, \mathbf{M})} L_T(\mathbf{M}, \tilde{\mathbf{M}}) = \mathbb{E}_{\mathbf{M}} \left[ N_T(a) \right] \inf_{\tilde{\mathbf{r}}} \{ \text{KL}(\mathbf{r}(a), \tilde{\mathbf{r}}) : \tilde{\mathbf{r}} \text{ has mean } > \mathbf{m}_* \}$$
$$= \mathbb{E}_{\mathbf{M}} \left[ N_T(a) \right] \underline{\mathbf{K}}(\mathbf{r}(a), \mathbf{m}_*)$$



How much does it cost to fool you and deviate from optimality?

**Fooling cost** The larger  $\inf_{\tilde{\mathbf{M}} \in \mathcal{B}(a, \mathbf{M})} L(\mathbf{M}, \tilde{\mathbf{M}})$  the most difficult to make arm  $a$  look optimal while it isn't.

**Lower bound** We consider algorithms that are uniformly “good” on all  $\mathbf{M} \in \mathfrak{M}$ . In particular on  $\mathbf{M}$ , and  $\tilde{\mathbf{M}}_a \in \mathcal{B}(a, \mathbf{M})$  for each  $a$ .

Control achievable “performances”.



## Algorithm

Pick MLE  $\widehat{\mathbf{M}}_t$ , choose to sample  $A_t \in \mathcal{O}(\widehat{\mathbf{M}}_t)$

vs

Pick MLE  $\widehat{\mathbf{M}}_t$ , choose to “track”  $\inf_{\tilde{\mathbf{M}} \in \mathcal{B}(a, \widehat{\mathbf{M}}_t)} L(\widehat{\mathbf{M}}_t, \tilde{\mathbf{M}})$ .





## Algorithm

Pick MLE  $\widehat{\mathbf{M}}_t$ , choose to sample  $A_t \in \mathcal{O}(\widehat{\mathbf{M}}_t)$

vs

Pick MLE  $\widehat{\mathbf{M}}_t$ , choose to “track”  $\inf_{\tilde{\mathbf{M}} \in \mathcal{B}(a, \widehat{\mathbf{M}}_t)} L(\widehat{\mathbf{M}}_t, \tilde{\mathbf{M}})$ .

Ensure  $\log\left(\frac{N_t(a)}{t}\right) \simeq - \inf_{\tilde{\mathbf{M}} \in \mathcal{B}(a, \widehat{\mathbf{M}}_t)} L(\widehat{\mathbf{M}}_t, \tilde{\mathbf{M}})$  by choosing:

$$\begin{aligned} A_t &\in \underset{a \in \mathcal{A}}{\operatorname{argmin}} \inf_{\tilde{\mathbf{M}} \in \mathcal{B}(a, \widehat{\mathbf{M}}_t)} L(\widehat{\mathbf{M}}_t, \tilde{\mathbf{M}}) + \log\left(\frac{N_t(a)}{t}\right). \\ &= \underset{a \in \mathcal{A}}{\operatorname{argmin}} N_T(a) \underline{\mathbf{K}}(\hat{\mathbf{r}}_t(a), \hat{\mathbf{m}}_{*,t}) + \log\left(\frac{N_t(a)}{t}\right). \end{aligned}$$

▷ Provably optimal strategy for “regret minimization”



INTRODUCTION

ACHIEVABLE PERFORMANCES

LEARNING ALGORITHM

EXTENSIONS



$$\mathbf{M} = (\mathcal{S}, \mathcal{A}, \mathbf{p}_1, \mathbf{p}, \mathbf{r})$$

$\mathcal{S}$ : set of states,

$\mathcal{A} = (\mathcal{A}_s)_{s \in \mathcal{S}}$ : set of actions available in each state.

$\mathcal{X} = \{(s, a) : s \in \mathcal{S}, a \in \mathcal{A}_s\}$ : set of pairs.

$\mathbf{p}_1 \in \mathcal{P}(\mathcal{S})$ : initial state distribution,

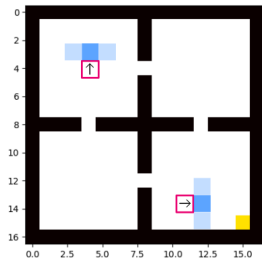
$\mathbf{p} : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{S})$ : transition distribution function.

$\mathbf{r} : \mathcal{X} \rightarrow \mathcal{P}(\mathbb{R})$ : reward distribution function, with mean reward function  $\mathbf{m} : \mathcal{X} \rightarrow \mathbb{R}$

In RL:  $\mathbf{p}, \mathbf{r}$  are **unknown**.



# Example: Grid-world



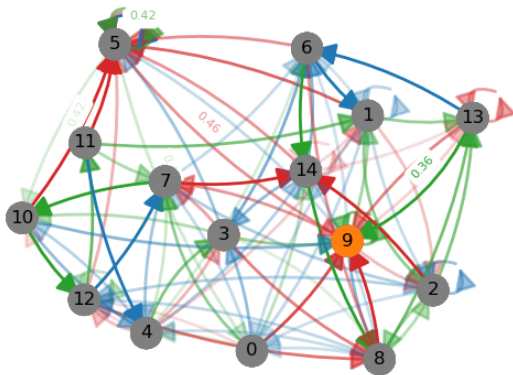
**Actions:**  $\{\leftarrow, \downarrow, \rightarrow, \uparrow\}$

**Rewards:** 1 when in goal state, 0 otherwise.

**Transitions:** "Frozen-lake". Reset to random initial state after reaching goal (yellow)



# Example: GARNET



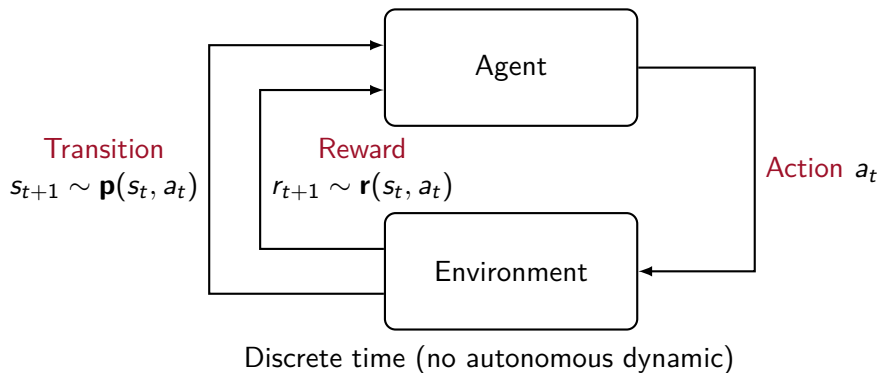
**Actions:** Colors

**Transitions:** Shaded arrows

**Rewards:** sparse.



# Reinforcement Learning



Goal: Maximize “cumulative return”.



Each stationary policy  $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$  acting in  $\mathbf{M}$  induces a Markov chain on  $\mathcal{S}$ .

$$\text{(Policy mean)} \quad \mathbf{m}_\pi(s) = \sum_{a \in \mathcal{A}_s} \mathbf{m}(s, a) \pi(a|s)$$

$$\text{(Policy transition)} \quad \mathbf{p}_\pi(s'|s) = \sum_{a \in \mathcal{A}_s} \mathbf{p}(s'|s, a) \pi(a|s)$$

For a sequence of policies  $\boldsymbol{\pi} = (\pi_t)_t$ , its value up to time  $T$  is

$$\mathbf{V}_{\boldsymbol{\pi}, T}(s_1) = \mathbb{E} \left[ \sum_{t=1}^T r_t \right] = \sum_{t=1}^T \left( \prod_{t'=1}^{t-1} \mathbf{p}_{\pi_{t'}} \mathbf{m}_{\pi_t} \right) (s_1).$$



# Optimality gap and Regret

Learning **with episodes** (of length  $H$ ):

- ▷ Learner wants to find a policy with **maximal value**  $\mathbf{V}_H^*(s) = \max_{\pi} \mathbf{V}_{\pi, H}(s)$ .
- ▷ At episode  $k$ , agent builds policy  $\pi_k$ . We measure the **gap to optimality**:

$$\Delta(\pi_k) = \mathbf{V}_H^*(s_1) - \mathbf{V}_{\pi_k, H}(s_1)$$

- ▷ As  $k \rightarrow \infty$ , we expect the gap to vanish, that is  $\Delta(\pi_k) \rightarrow 0$ .

vs





# Optimality gap and Regret

Learning **with episodes** (of length  $H$ ):

- ▶ Learner wants to find a policy with **maximal value**  $\mathbf{V}_H^*(s) = \max_{\pi} \mathbf{V}_{\pi, H}(s)$ .
- ▶ At episode  $k$ , agent builds policy  $\pi_k$ . We measure the **gap to optimality**:

$$\Delta(\pi_k) = \mathbf{V}_H^*(s_1) - \mathbf{V}_{\pi_k, H}(s_1)$$

- ▶ As  $k \rightarrow \infty$ , we expect the gap to vanish, that is  $\Delta(\pi_k) \rightarrow 0$ .

vs

Learning **within one episode** (this talk!)

- ▶ Learner wants to get **maximal reward** and maximize its **own** value.
- ▶ At decision step  $T$ , we measure the **regret** of the agent choosing  $\pi = (\pi_t)_t$ , accumulated **while learning**, compared to an optimal policy  $\star$  with

$$\mathfrak{R}_T(\pi, \mathbf{M}) = \mathbf{V}_{\star, T}(s_1) - \mathbf{V}_{\pi, T}(s_1)$$

- ▶ As  $T \rightarrow \infty$ , the regret cannot decrease, we still expect  $\mathfrak{R}_T(\pi, \mathbf{M})/T \rightarrow 0$ .



We are interested in policies that maximize rewards in the long run,

$$\lim_T \frac{1}{T} \mathbf{V}_{\pi, T}(s_1)$$

For stationary policy  $\pi$ , this is  $\left( \lim_T \frac{1}{T} \sum_{t=1}^T \mathbf{p}_{\pi}^{t-1} \mathbf{m}_{\pi} \right)(s_1)$ , hence we define

$$\text{(Average transition)} \quad \bar{\mathbf{p}}_{\pi} = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbf{p}_{\pi}^{t-1},$$

$$\text{(Average frequency)} \quad \mathbf{f}_{\pi, s_1}(s, a) = \bar{\mathbf{p}}_{\pi}(s|s_1)\pi(a|s)$$

The **average value** (aka gain) of a policy  $\pi$  starting in state  $s_1$ :

$$\mathbf{g}_{\pi}(s_1) = (\bar{\mathbf{p}}_{\pi} \mathbf{m}_{\pi})(s_1) = \langle \mathbf{f}_{\pi, s_1}, \mathbf{m}_{\pi} \rangle.$$



The **average value** (aka gain) of a policy  $\pi$  starting in state  $s_0$ :

$$\mathbf{g}_\pi(s_1) = (\bar{\mathbf{p}}_\pi \mathbf{m}_\pi)(s_1) = \langle \mathbf{f}_{\pi, s_1}, \mathbf{m}_\pi \rangle.$$

The **bias function** is given by  $\mathbf{b}_\pi(s) = \left( \sum_{t=1}^{\infty} (\mathbf{p}_\pi^{t-1} - \bar{\mathbf{p}}_\pi) \mathbf{m}_\pi \right)(s)$ .

**Fixed point property:**

$$(Poisson\ equation) \quad \mathbf{b}_\pi(s) = \mathbf{m}_\pi(s) - \mathbf{g}_\pi(s) + (\mathbf{p}_\pi \mathbf{b}_\pi)(s).$$

(Span  $\mathbb{S}(f) = \max_x f(x) - \min_x f(x)$  of function  $f : \mathcal{X} \rightarrow \mathbb{R}$ . The span operator induces a semi-norm on functions.)



An MDP  $M$  is **communicating** if every pair of state is communicating under some policy, that is

$$\forall s, s' \in \mathcal{S}, \exists \pi, \exists t < \infty : \mathbf{p}_{\pi}^t(s'|s) > 0.$$

An MDP  $M$  is **irreducible** if every policy  $\pi$  on  $M$  induces an irreducible Markov chain, that is

$$\forall \pi, \forall s, s' \in \mathcal{S}, \exists t < \infty : \mathbf{p}_{\pi}^t(s'|s) > 0.$$

An MDP  $M$  is **ergodic**<sup>1</sup> if every policy  $\pi$  on  $M$  induces an ergodic Markov chain, that is  $\forall \pi$ ,

(irreducibility)  $\forall s, s' \in \mathcal{S}, \exists t < \infty, \mathbf{p}_{\pi}^t(s'|s) > 0,$

(aperiodicity)  $\forall s \in \mathcal{S}, \text{GCD}(\{t : \mathbf{p}_{\pi}^t(s|s) > 0\}) = 1,$

(positive recurrence)  $\forall s \in \mathcal{S}, \mathbb{E}(\min\{t > 1 : s_t = s\} | s_1 = s) < \infty.$



INTRODUCTION

---

## Learning games

---

ACHIEVABLE PERFORMANCES

LEARNING ALGORITHM

EXTENSIONS



- ▷ Gain optimal policies:

$$\overline{\mathcal{O}}(\mathbf{M}, s_1) = \left\{ \pi : \mathbf{g}_\pi(s_1) \geq \mathbf{g}_{\pi'}(s_1) \forall \pi' \right\},$$

- ▷ In a communicating MDP with finite  $\mathcal{S}$  and  $\mathcal{A}$ ,  $\bigcap_{s_1 \in \mathcal{S}} \overline{\mathcal{O}}(\mathbf{M}, s_1) \neq \emptyset$  and contains a **unichain** policy.

- ▷ The **cumulative regret** up to time  $T \in \mathbb{N}$  of an algorithm executing policy  $\pi = (\pi_t)_t$  (that is, it plays policy  $\pi_t$  at time  $t \in \{1, \dots, T\}$ ) in MDP  $\mathbf{M}$  starting from initial state  $s_1$ , against optimal policies for the different optimality criteria are given by

$$\overline{\mathfrak{R}}_{T, s_1}(\pi; \mathbf{M}) = \mathbf{V}_{\star, T}(s_1) - \mathbf{V}_{\pi, T}(s_1) \text{ where } \star \in \overline{\mathcal{O}}(\mathbf{M}, s_1)$$



# Regret decomposition

The **cumulative regret** up to time  $T$  for the average gain and discounted criterion satisfy the following decomposition provided that an optimal policy  $\star$  is **unchain**,

$$\bar{\mathfrak{R}}_{T,s_1}(\pi; \mathbf{M}) = \sum_{x \in \mathcal{X}} \mathbb{E}_{\pi, s_1}[N_T(x)] \Delta(x) + \underbrace{\left( \left[ \prod_{t'=1}^T \mathbf{p}_{\pi_{t'}} - \mathbf{p}_{\star}^T \right] \mathbf{b}_{\star} \right)}_{\leq \mathbb{S}(\mathbf{b}_{\star})}(s_1),$$

where we introduced the **sub-optimality gap**

$$\Delta(s, a) = \mathbf{m}_{\star}(s) - \mathbf{m}(s, a) + ((\mathbf{p}_{\star} - \mathbf{p}_a) \mathbf{b}_{\star})(s).$$

Minimizing  $\Delta(s, a)$  is equivalent to maximizing  $\varphi_{\mathbf{M}}(s, a) = \mathbf{m}(s, a) + (\mathbf{p}_a \mathbf{b}_{\star})(s)$ .



# Regret decomposition

The **cumulative regret** up to time  $T$  for the average gain and discounted criterion satisfy the following decomposition provided that an optimal policy  $\star$  is **unchain**,

$$\bar{\mathfrak{R}}_{T,s_1}(\pi; \mathbf{M}) = \sum_{x \in \mathcal{X}} \mathbb{E}_{\pi, s_1}[N_T(x)] \Delta(x) + \underbrace{\left( \left[ \prod_{t'=1}^T \mathbf{p}_{\pi_{t'}} - \mathbf{p}_{\star}^T \right] \mathbf{b}_{\star} \right)}_{\leq \mathbb{S}(\mathbf{b}_{\star})}(s_1),$$

where we introduced the **sub-optimality gap**

$$\Delta(s, a) = \mathbf{m}_{\star}(s) - \mathbf{m}(s, a) + ((\mathbf{p}_{\star} - \mathbf{p}_a) \mathbf{b}_{\star})(s).$$

Minimizing  $\Delta(s, a)$  is equivalent to maximizing  $\varphi_{\mathbf{M}}(s, a) = \mathbf{m}(s, a) + (\mathbf{p}_a \mathbf{b}_{\star})(s)$ .

▷ In Bandits  $|\mathcal{S}| = 1$ ,  $\Delta(a) = \mathbf{m}_{\star} - \mathbf{m}(a)$  and

$$\bar{\mathfrak{R}}_T(\pi, \mathbf{M}) = \sum_{a \in \mathcal{A}} \mathbb{E}_{\pi}[N_T(a)] \Delta(a).$$





INTRODUCTION

ACHIEVABLE PERFORMANCES

LEARNING ALGORITHM

EXTENSIONS



## Regular model $\mathfrak{M}$

For each  $x \in \mathcal{X}$ ,  $\mathbf{r}(x) \in \mathcal{D}_x \subset \mathcal{P}(\mathbb{R})$  where  $\mathcal{D}_x$  is known to the learner,  $\mathcal{D} = \otimes_x \mathcal{D}_x$ .

**Light-tail rewards** For all  $x \in \mathcal{X}$ , the moment generating function of reward  $\mathbf{r}(x)$  exists in a neighborhood of 0.

**Semi-bounded rewards** There exists a known quantity  $m_{\max}(x) \in \mathbb{R}$  such that  $\text{Supp}(\mathbf{r}(x)) \subset (-\infty, m_{\max}(x)]$  and  $\mathbf{m}(x) < m_{\max}(x)$ .

Good learner

## Definition (Uniformly Good strategies)

A agent is **uniformly-good** on  $\mathcal{D}$  if

$$\forall \mathbf{M} = (\mathbf{r}_a)_{a \in \mathcal{A}} \in \mathcal{D}, \forall a \notin \mathcal{O}(\mathbf{M}), \quad \mathbb{E}[N_T(a)] = o(T^\alpha) \quad \text{for all } \alpha \in (0, 1].$$



## Theorem (“Price for being uniformly-good”)

Any uniformly good strategy on  $\mathcal{D} = \text{Bern}^{\mathcal{A}}$  must satisfy (Lai & Robbins, 85)

$$\forall a \notin \mathcal{O}(\mathbf{M}) \quad \liminf_{T \rightarrow \infty} \frac{\mathbb{E}_{\mathbf{M}}[N_T(a)]}{\ln(T)} \geq \frac{1}{\text{kl}(\mathbf{m}_a, \mathbf{m}_*)}.$$

More generally (Burnetas & Katehakis, 96) for  $\mathbf{r} \in \otimes_{a \in \mathcal{A}} \mathcal{D}_a$

$$\forall a \notin \mathcal{O}(\mathbf{M}) \quad \liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\ln T} \geq \frac{1}{\underline{\mathbf{K}}_a(\mathbf{r}_a, \mathbf{m}_*)},$$

$$\underline{\mathbf{K}}_a(\mathbf{r}_a, \mathbf{m}_*) = \inf \{ \text{KL}(\mathbf{r}_a, \tilde{r}) : \tilde{r} \in \mathcal{D}_a, \mathbb{E}_{\tilde{r}}[X] > \mathbf{m}_* \}$$



For any **uniformly-good** strategy  $\pi$  on  $\mathfrak{M}$  and  $\mathbf{M} \in \mathfrak{M}$ ,

$$\liminf_{T \rightarrow \infty} \frac{\overline{\mathfrak{R}}_T(\pi, \mathbf{M})}{\ln T} \geq \sum_{a \in \mathcal{A}} \frac{\Delta_a}{\underline{\mathbf{K}}_a(\mathbf{r}_a, \mathbf{m}_*)}$$

where

$$\underline{\mathbf{K}}_a(\mathbf{r}_a, \mathbf{m}_*) = \inf \{ \text{KL}(\mathbf{r}_a, \tilde{r}) : \tilde{r} \in \mathcal{D}_a, \mathbb{E}_{\tilde{r}}[X] > \mathbf{m}_* \}$$



For any **uniformly-good** strategy  $\pi$  on  $\mathfrak{M}$  and  $\mathbf{M} \in \mathfrak{M}$ ,

$$\liminf_{T \rightarrow \infty} \frac{\overline{\mathfrak{R}}_T(\pi, \mathbf{M})}{\ln T} \geq \sum_{a \in \mathcal{A}} \frac{\Delta_a}{\underline{\mathbf{K}}_a(\mathbf{r}_a, \mathbf{m}_*)}$$

where

$$\underline{\mathbf{K}}_a(\mathbf{r}_a, \mathbf{m}_*) = \inf \{ \text{KL}(\mathbf{r}_a, \tilde{r}) : \tilde{r} \in \mathcal{D}_a, \mathbb{E}_{\tilde{r}}[X] > \mathbf{m}_* \}$$

▷ IMED strategy  $\pi_I$  is optimal in the sense that for each  $\mathbf{M} \in \mathfrak{M}$  (with regular  $\mathfrak{M}$ )

$$\limsup_{T \rightarrow \infty} \frac{\overline{\mathfrak{R}}_T(\pi_I, \mathbf{M})}{\ln T} \leq \sum_{a \in \mathcal{A}} \frac{\Delta_a}{\underline{\mathbf{K}}_a(\mathbf{r}_a, \mathbf{m}_*)}$$



▷ For any **test function**  $h$  bounded in  $[0, 1]$ ,

$$\mathbb{E}_{\mathbf{M}} \left[ \log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}(\tau_T)} \right] \geq \text{k1} (\mathbb{E}_{\mathbf{M}}[h(\tau_T)], \mathbb{E}_{\tilde{\mathbf{M}}}[h(\tau_T)]) ,$$

where  $\text{k1}(x, y) = x \log(x/y) + (1 - x) \log((1 - x)/(1 - y))$ .



# Fundamental inequality

▷ For any **test function**  $h$  bounded in  $[0, 1]$ ,

$$\mathbb{E}_{\mathbf{M}} \left[ \log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}(\tau_T)} \right] \geq \text{kl} (\mathbb{E}_{\mathbf{M}}[h(\tau_T)], \mathbb{E}_{\tilde{\mathbf{M}}}[h(\tau_T)]),$$

where  $\text{kl}(x, y) = x \log(x/y) + (1-x) \log((1-x)/(1-y))$ .

▷ For  $h(\tau_T) = 1 - N_T(a)/T$ , using **uniformly-good** assumption, we can show

$$\mathbb{E}_{\mathbf{M}} \left[ \log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}_a(\tau_T)} \right] \geq (1 - \alpha) \log(T) + o(\log(T)) \text{ for all } \alpha,$$

from which we get for each  $a \notin \mathcal{O}(\mathbf{M})$

$$\liminf_T \frac{\mathbb{E}_{\mathbf{M}} \left[ \log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}_a(\tau_T)} \right]}{\log(T)} \geq 1 \text{ i.e. } \liminf_T \frac{\mathbb{E}_{\mathbf{M}} [N_T(a)]}{\log(T)} \geq \frac{1}{\underline{\mathbf{K}}(\mathbf{r}(a), \mathbf{m}_*)}.$$



INTRODUCTION

ACHIEVABLE PERFORMANCES

---

**Back to MDPs**

---

LEARNING ALGORITHM

EXTENSIONS





For a given trajectory  $\tau_T = (s_1, a_1, r_1, s_2, \dots, s_{T+1})$  and two MDPs  $\mathbf{M} = (\mathcal{S}, \mathcal{A}, \mathbf{p}_1, \mathbf{p}, \mathbf{r})$ ,  $\tilde{\mathbf{M}} = (\mathcal{S}, \mathcal{A}, \tilde{\mathbf{p}}_1, \tilde{\mathbf{p}}, \tilde{\mathbf{r}})$ :

$$\log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}(\tau_T)} = \log \frac{p_1(s_1)}{\tilde{p}_1(s_1)} + \sum_{t=1}^T \log \frac{p(s_t, a_t)(s_{t+1})}{\tilde{p}(s_t, a_t)(s_{t+1})} + \log \frac{r(s_t, a_t)(r_t)}{\tilde{r}(s_t, a_t)(r_t)}$$

$$\mathbb{E}_{\mathbf{M}} \left[ \log \frac{\mathbf{M}(\tau_T)}{\tilde{\mathbf{M}}(\tau_T)} \right] = \text{KL}(\mathbf{p}_1, \tilde{\mathbf{p}}_1) + \sum_x \mathbb{E}_{\mathbf{M}} \left[ N_T(x) \right] \left[ \text{KL}(\mathbf{p}(x), \tilde{\mathbf{p}}(x)) + \text{KL}(\mathbf{r}(x), \tilde{\mathbf{r}}(x)) \right]$$

▷ For a policy  $\pi$  generating  $\tau_T$  from starting state  $s_1$  (i.e.  $\mathbf{p}_1 = \tilde{\mathbf{p}}_1 = \delta_{s_1}$ ), we let

$$L_T(\mathbf{M}, \tilde{\mathbf{M}}) = \sum_x \mathbb{E}_{\mathbf{M}, \pi, s_1} \left[ N_T(x) \right] \text{KL}_x(\mathbf{M}, \tilde{\mathbf{M}}).$$

where  $\text{KL}_x(\mathbf{M}, \tilde{\mathbf{M}}) = \text{KL}(\mathbf{p}(x), \tilde{\mathbf{p}}(x)) + \text{KL}(\mathbf{r}(x), \tilde{\mathbf{r}}(x))$ .



## ▷ Regret decomposition

$$\overline{\mathfrak{R}}_{T,s_1}(\pi; \mathbf{M}) = \sum_{x \in \mathcal{X}} \mathbb{E}_{\pi} [N_T(x)] \Delta(x) + \text{constant}$$

with  $\Delta(s, a) = \varphi_{\mathbf{M}}(s, \star(s)) - \varphi_{\mathbf{M}}(s, a)$ .

## ▷ Optimal policies

$$\overline{\mathcal{O}}(\mathbf{M}, s_1) = \left\{ \pi : \mathbf{g}_{\pi}(s_1) \geq \mathbf{g}_{\pi'}(s_1) \forall \pi' \right\},$$

▷ For each  $\pi$  deviating from  $\overline{\mathcal{O}}(\mathbf{M}, s_1)$ , find **confusing**  $\tilde{\mathbf{M}}_{\pi}$  that minimizes

$$L_T(\mathbf{M}, \tilde{\mathbf{M}}) = \sum_{x \in \mathcal{X}} \mathbb{E}_{\mathbf{M}} [N_T(x)] \text{KL}_x(\mathbf{M}, \tilde{\mathbf{M}}).$$



**Uniformly good strategies** An algorithm is said uniformly  $\eta$ -good on  $\mathfrak{M}$ , where  $\eta : \mathbb{N} \rightarrow \mathbb{R}^+$  that satisfies  $\lim_{t \rightarrow \infty} \eta(t) = 0$  specifies a rate function, if for all  $T$ , for all  $\mathbf{M} \in \mathfrak{M}$ , for all  $x \in \mathcal{X}$  such that  $\Delta_{\mathbf{M}}(x) > 0$ , then

$$\mathbb{E}_{\mathbf{M}, \pi, s_1} [N_T(x) / T] \leq \eta(T).$$

(e.g.  $\eta(T) = O(T^{\alpha-1})$ )



**Uniformly good strategies** An algorithm is said uniformly  $\eta$ -good on  $\mathfrak{M}$ , where  $\eta : \mathbb{N} \rightarrow \mathbb{R}^+$  that satisfies  $\lim_{t \rightarrow \infty} \eta(t) = 0$  specifies a rate function, if for all  $T$ , for all  $\mathbf{M} \in \mathfrak{M}$ , for all  $x \in \mathcal{X}$  such that  $\Delta_{\mathbf{M}}(x) > 0$ , then

$$\mathbb{E}_{\mathbf{M}, \pi, s_1} [N_T(x) / T] \leq \eta(T).$$

(e.g.  $\eta(T) = O(T^{\alpha-1})$ )

**Beneficial MDPs** for policy  $\pi$ , are  $\tilde{\mathbf{M}} \in \overline{\mathcal{B}}_{s_1}(\pi, \mathbf{M}; \mathfrak{M})$  such that

- ▶  $\forall \star \in \mathcal{O}(\mathbf{M})$ ,  $(\mathbf{r}_\star, \mathbf{p}_\star)$  and  $(\tilde{\mathbf{r}}_\star, \tilde{\mathbf{p}}_\star)$  are undistinguishable.
- ▶  $\forall \star \in \mathcal{O}(\mathbf{M})$ ,  $\tilde{\mathbf{g}}_\pi(s_1) > \tilde{\mathbf{g}}_\star(s_1)$

Note: So  $\tilde{\mathbf{g}}_\star = \mathbf{g}_\star$  and  $\varphi_{\tilde{\mathbf{M}}}(s, \star(s)) = \varphi_{\mathbf{M}}(s, \star(s)) = \varphi_{\mathbf{M}}^\star(s)$



# Confusing instance for Ergodic MDPs

$$\inf_{\tilde{\mathbf{M}} \in \bar{\mathcal{B}}_{s_1}(\pi, \mathbf{M})} \sum_{x \in \mathcal{X}} \mathbb{E}_{\mathbf{M}} [N_T(x)] \text{KL}_x(\mathbf{M}, \tilde{\mathbf{M}})$$

- ▷ **Perturbation of optimal policy** For  $\star \in \mathcal{O}(\mathbf{M})$ , consider  $\pi = \star_{s,a}$  that plays action  $a$  in state  $s$  and otherwise  $\star$ .
- ▷ In ergodic MDPs, all states are recurrent under **each policy**.

In particular  $\mathcal{X}_{\star} \Delta \mathcal{X}_{\star_{s,a}} = \{(s, a)\}$ , where  $\mathcal{X}_{\pi} = \{x : \mathbf{f}_{\pi, s_1}(x) > 0\}$ .



$$\inf_{\tilde{M} \in \bar{\mathcal{B}}_{s_1}(\pi, \mathbf{M})} \sum_{x \in \mathcal{X}} \mathbb{E}_{\mathbf{M}} [N_T(x)] \text{KL}_x(\mathbf{M}, \tilde{\mathbf{M}})$$

- ▷ **Perturbation of optimal policy** For  $\star \in \mathcal{O}(\mathbf{M})$ , consider  $\pi = \star_{s,a}$  that plays action  $a$  in state  $s$  and otherwise  $\star$ .
- ▷ In ergodic MDPs, all states are recurrent under **each policy**.

In particular  $\mathcal{X}_{\star} \Delta \mathcal{X}_{\star_{s,a}} = \{(s, a)\}$ , where  $\mathcal{X}_{\pi} = \{x : \mathbf{f}_{\pi, s_1}(x) > 0\}$ .

- ▷ Also, since  $\star$  and  $\pi$  only differ in  $(s, a)$ ,  $\tilde{\mathbf{g}}_{\star_{s,a}}(s_1) > \tilde{\mathbf{g}}_{\star}(s_1)$  iff  $\varphi_{\tilde{\mathbf{M}}}(s, a) > \varphi_{\tilde{\mathbf{M}}}(s, \star(s)) = \varphi_{\mathbf{M}}^*(s)$ , that is  $\tilde{\mathbf{m}}(s, a) + (\tilde{\mathbf{p}}_a \mathbf{b}_{\star})(s) > \varphi_{\mathbf{M}}^*(s)$ .



$$\inf_{\tilde{M} \in \tilde{\mathcal{B}}_{s_1}(\pi, \mathbf{M})} \sum_{x \in \mathcal{X}} \mathbb{E}_{\mathbf{M}} [N_T(x)] \text{KL}_x(\mathbf{M}, \tilde{\mathbf{M}})$$

- ▷ **Perturbation of optimal policy** For  $\star \in \mathcal{O}(\mathbf{M})$ , consider  $\pi = \star_{s,a}$  that plays action  $a$  in state  $s$  and otherwise  $\star$ .
- ▷ In ergodic MDPs, all states are recurrent under **each policy**.

In particular  $\mathcal{X}_{\star} \Delta \mathcal{X}_{\star_{s,a}} = \{(s, a)\}$ , where  $\mathcal{X}_{\pi} = \{x : \mathbf{f}_{\pi, s_1}(x) > 0\}$ .

- ▷ Also, since  $\star$  and  $\pi$  only differ in  $(s, a)$ ,  $\tilde{\mathbf{g}}_{\star_{s,a}}(s_1) > \tilde{\mathbf{g}}_{\star}(s_1)$  iff  $\varphi_{\tilde{\mathbf{M}}}(s, a) > \varphi_{\tilde{\mathbf{M}}}(s, \star(s)) = \varphi_{\mathbf{M}}^{\star}(s)$ , that is  $\tilde{\mathbf{m}}(s, a) + (\tilde{\mathbf{p}}_a \mathbf{b}_{\star})(s) > \varphi_{\mathbf{M}}^{\star}(s)$ .

## Ergodic simplification

For each  $(s, a)$ , the confusing cost becomes

$$\mathbb{E}_{\mathbf{M}} [N_T(s, a)] \underline{\mathbf{K}}_{s,a}(\mathbf{M}, \varphi_{\mathbf{M}}^{\star}(s))$$

$$\begin{aligned} \text{where } \underline{\mathbf{K}}_x(\mathbf{M}, \varphi) &= \inf \{ \text{KL}(\mathbf{r}(x), \tilde{\mathbf{r}}(x)) + \text{KL}(\mathbf{p}(x), \tilde{\mathbf{p}}(x)) : \varphi_{\tilde{\mathbf{M}}}(x) > \varphi \} \\ &= \inf \left\{ \text{KL}(\mathbf{r}(x) \otimes \mathbf{p}(x), \tilde{\mathbf{r}} \otimes \tilde{\mathbf{p}}) : \tilde{\mathbf{m}} + \tilde{\mathbf{q}} \mathbf{b}_{\star} > \varphi \right\} \end{aligned}$$



For all  $\eta$ -uniformly-good strategy with  $\eta(T) = O(T^{\alpha-1})$  for all  $\alpha \in (0, 1)$ ,

$$\begin{aligned} \liminf_T \frac{\overline{\mathfrak{R}}_{T, s_1}(\pi, \mathbf{M})}{\ln(T)} &\geq \inf_{\kappa \in \mathbb{R}_+^{|\mathcal{X}|}} \left\{ \sum_{x \in \mathcal{X}} \kappa_x \Delta(x) : \forall (s, a), \kappa_{s,a} \underline{\mathbf{K}}_{s,a}(\mathbf{M}, \varphi_{\mathbf{M}}^*(s)) \geq 1 \right\}, \\ &\geq \sum_{(s,a) \in \mathcal{X}} \frac{\Delta(s, a)}{\underline{\mathbf{K}}_{s,a}(\mathbf{M}, \varphi_{\mathbf{M}}^*(s))} \end{aligned}$$





INTRODUCTION

ACHIEVABLE PERFORMANCES

LEARNING ALGORITHM

EXTENSIONS



**Light-tail rewards** For all  $x \in \mathcal{X}$ , the moment generating function of reward  $\mathbf{r}(x)$  exists in a neighborhood of 0.

**Semi-bounded rewards** For all  $x \in \mathcal{X}$ ,  $\mathbf{r}(x)$  belongs to a subset  $\mathcal{F}_x \subset \mathcal{P}(\mathbb{R})$  known to the learner.

There exists a known quantity  $m_{\max}(x) \in \mathbb{R}$  such that  $\text{Supp}(\mathbf{r}(x)) \subset (-\infty, m_{\max}(x))$  and its mean satisfies  $\mathbf{m}(x) < m_{\max}(x)$ .

**Ergodicity** The MDP is ergodic,  $\forall s, s', \forall \pi, \exists t \in \mathbb{N} : \mathbf{p}_{\pi}^t(s'|s) > 0$ .



IMED-RL is a **model-based** algorithm that keeps empirical estimates of the transitions  $\mathbf{p}$  and rewards  $\mathbf{r}$ .

While **policy iteration** constructs a sequence of policies that are increasingly better, IMED-RL constructs a sequence of **sub-MDPs** of the original MDP that are increasingly better with high probability.



**Empirical MDP**  $\hat{r}(s, a)$  and  $\hat{p}_t(s, a)$  the empirical reward and transition distributions after  $t$  time steps, *i.e.* using  $N_{s,a}(t)$  samples from the distributions.  $\widehat{\mathbf{M}}_t$  is the empirical MDP built from those estimations.

**Skeleton** The skeleton at time  $t$ , is defined by

$$\mathcal{A}_s(t) = \{a \in \mathcal{A}_s : N_{s,a}(t) \geq \log^2(\max_{a' \in \mathcal{A}_s} N_{s,a'}(t))\}.$$

It is **homogeneous** in the sense that it does not depend explicitly on  $t$  nor  $N_s(t)$ , only  $N_{s,a}(t)$  and  $\max_{a'} N_{s,a'}(t)$ .

**Restricted MDP**  $\widehat{\mathbf{M}}_t|_t = \widehat{\mathbf{M}}_t(\mathcal{A}(t))$  is the empirical MDP whose action space is restricted to the skeleton.



A value iteration scheme (or Policy iteration scheme) is applied to  $\widehat{\mathbf{M}}_t$ , to find a policy  $\widehat{\star}_t$ .

▷ This enables to build **estimate**  $\varphi_{\widehat{\mathbf{M}}_t}$  of:  $\varphi_{\mathbf{M}}(s, a) = \mathbf{m}(s, a) + (\mathbf{p}_a \mathbf{b}_{\star})(s)$ .

Using  $\widehat{\star}_t, \widehat{\mathbf{m}}_t, \widehat{\mathbf{p}}_t, \widehat{\mathbf{b}}_{\widehat{\star}_t}$ .



A value iteration scheme (or Policy iteration scheme) is applied to  $\widehat{\mathbf{M}}|_t$ , to find a policy  $\widehat{\star}_t$ .

▷ This enables to build **estimate**  $\varphi_{\widehat{\mathbf{M}}|_t}$  of:  $\varphi_{\mathbf{M}}(s, a) = \mathbf{m}(s, a) + (\mathbf{p}_a \mathbf{b}_{\star})(s)$ .

Using  $\widehat{\star}_t, \widehat{\mathbf{m}}_t, \widehat{\mathbf{p}}_t, \widehat{\mathbf{b}}_{\widehat{\star}_t}$ .

▷ By **ergodicity**, the pairs in the skeleton have enough visits ( $\simeq \log^2(t)$ ), so well-estimated.

▷ In **ergodic** MDPs, for each  $\pi \notin \mathcal{O}(\mathbf{M})$ , there exists a policy improving over  $\pi$  at Hamming distance 1 of  $\pi$ .



The the IMED-RL index of  $(s, a)$  at time  $t$ ,  $\mathbf{H}_{s,a}(t)$ , is defined as

$$\mathbf{H}_{s,a}(t) = N_t(s, a) \underline{\mathbf{K}}_{s,a}(\widehat{\mathbf{M}}_{|t}, \varphi_{\widehat{\mathbf{M}}_{|t}}^*(s)) + \log N_t(s, a),$$

$$\underline{\mathbf{K}}_{s,a}(\mathbf{M}, \varphi) = \inf \left\{ \text{KL}(\mathbf{r}(x) \otimes \mathbf{p}(x), \tilde{\mathbf{r}} \otimes \tilde{\mathbf{p}}) : \tilde{m} + \tilde{\mathbf{p}}\mathbf{b}_* > \varphi \right\}$$

where  $\varphi_{\mathbf{M}}^*(s) = \max_{a \in \mathcal{A}_s} \varphi_{\mathbf{M}}(s, a)$  with  $\varphi_{\mathbf{M}}(x) = \mathbf{m}(s, a) + (\mathbf{p}_a \mathbf{b}_*)(s)$ .



---

## Algorithm 1 IMED-RL

---

**Initialisation** State  $s_1$

**for**  $t \geq 1$  **do**

    Compute the skeleton  $\mathcal{A}_s(t)$  for each  $s$  and set  $\widehat{\mathbf{M}}|_t$

    Run value iteration on  $\widehat{\mathbf{M}}|_t$  to get  $\varphi_{\widehat{\mathbf{M}}|_t}^*$ .

    Sample  $A_t \in \arg \min_{a \in \mathcal{A}_s} N_t(s, a) \mathbf{K}_{s,a}(\widehat{\mathbf{M}}|_t, \varphi_{\widehat{\mathbf{M}}|_t}^*(s)) + \log N_t(s, a)$

---

▷ Note:  $\varphi_{\widehat{\mathbf{M}}|_t}^*$  should be good estimate of  $\varphi_{\mathbf{M}}^*$ .





## Theorem (Asymptotic Optimality)

*IMED-RL* is **asymptotically optimal**, that is,

$$\lim_{T \rightarrow +\infty} \frac{\bar{\mathbf{R}}_{T, s_1}(\text{IMED-RL}; \mathbf{M})}{\log T} \leq \sum_{s, a \in \mathcal{X}} \frac{\Delta_{s, a}(\mathbf{M})}{\underline{\mathbf{K}}_{s, a}(\mathbf{M})}.$$



- ▶ Ergodic Environment
- ▶ Communicating Environment

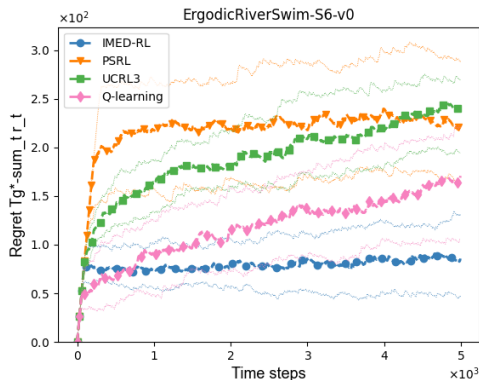
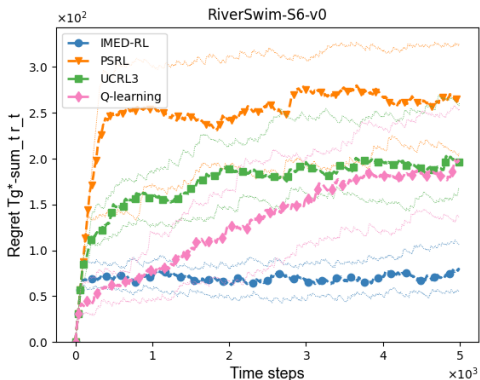
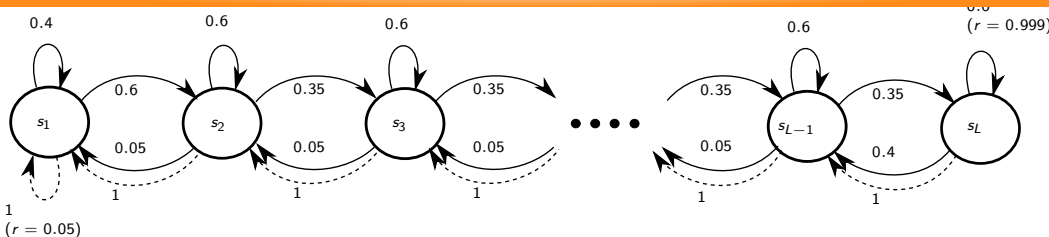
From a communicating only MDP, by mixing its transition  $\mathbf{p}$  with that obtained from playing a uniform policy, formally

$$\mathbf{p}_\varepsilon(\cdot|s, a) = (1 - \varepsilon)\mathbf{p}(\cdot|s, a) + \varepsilon \sum_{a' \in \mathcal{A}_s} \mathbf{p}(\cdot|s, a')/|\mathcal{A}_s|,$$

for an arbitrarily small  $\varepsilon > 0$  one obtain an ergodic MDP.  
Experimentally, we take  $\varepsilon = 10^{-3}$ .

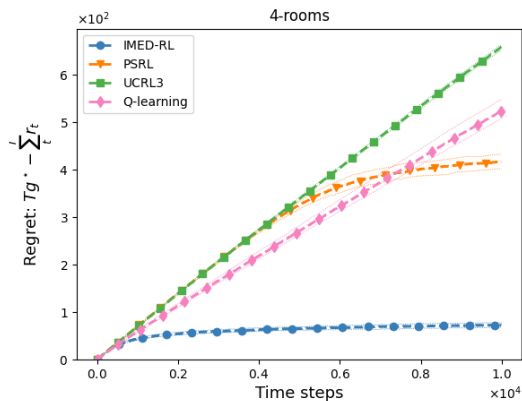
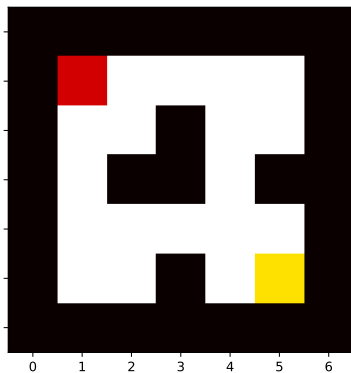


# RiverSwim



Average regret and quantiles (0.1 ; 0.9) on a communicating (left) and an ergodic (right) 6-states RiverSwim

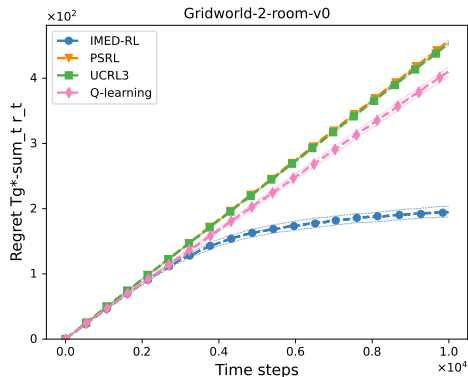
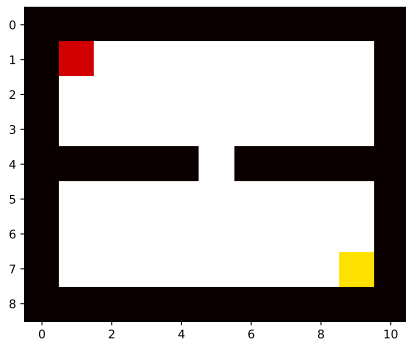




Average *regret* in the 4-rooms environment



# 2-rooms



Average **regret** and quantiles (0.1 and 0.9) curves in the 2-rooms environment



INTRODUCTION

ACHIEVABLE PERFORMANCES

LEARNING ALGORITHM

EXTENSIONS



INTRODUCTION

ACHIEVABLE PERFORMANCES

LEARNING ALGORITHM

EXTENSIONS

---

**Beyond discrete MDPs**

Beyond ergodic MDPs



---

## Algorithm 2 IMED-RL

---

**Initialisation** State  $s_1$

**for**  $t \geq 1$  **do**

    Compute the skeleton  $\mathcal{A}_s(t)$  for each  $s$  and set  $\widehat{\mathbf{M}}|_t$

    Run value iteration on  $\widehat{\mathbf{M}}|_t$  to get  $\varphi_{\widehat{\mathbf{M}}|_t}^*$ .

    Sample  $A_t \in \arg \min_{a \in \mathcal{A}_s} N_t(s, a) \underline{\mathbf{K}}_{s,a}(\widehat{\mathbf{M}}|_t, \varphi_{\widehat{\mathbf{M}}|_t}^*(s)) + \log N_t(s, a)$

---

- ▷ Optimization problem  $\underline{\mathbf{K}}_{s,a}$  computed at each step: find an **iterative approach** (reusing previous computation)?
- ▷ Compute  $N_t(s, a)$ ,  $\underline{\mathbf{K}}_{s,a}$  for continuous  $s, a$ : Introduce encoder for  $s, a$ , density estimation, etc.





INTRODUCTION

ACHIEVABLE PERFORMANCES

LEARNING ALGORITHM

EXTENSIONS

Beyond discrete MDPs

---

**Beyond ergodic MDPs**

---



## Recurrent classes

- ▷ Perturbation of  $\star$  at  $(s, a)$  may completely change the recurrent class:  $\mathcal{X}_\star$  vs  $\mathcal{X}_{\star, s, a}$
- ▷  $(s, a)$  may not be recurrent: find confusing instance by changing  $\mathbf{M}$  at other pairs.

$$L_T(\mathbf{M}, \tilde{\mathbf{M}}) = \sum_{x \in \mathcal{X}_{\star, s, a} \setminus \mathcal{X}_\star} \mathbb{E}_{\mathbf{M}} \left[ N_T(x) \right] \text{KL}_x(\mathbf{M}, \tilde{\mathbf{M}}).$$

- ▷ NO clear simplification in general.

## Estimation

- ▷ Some states may not be visited frequently, even by an optimal policy
- ▷ Estimates of  $\mathbf{r}(x)$ ,  $\mathbf{p}(x)$  may not converge for such  $x \notin \mathcal{X}_\star!$
- ▷ Study case with known support of  $\mathbf{p}$ ?



▷ **Finite time lower bound** For any **uniformly  $\eta$ -good** strategy on  $\mathfrak{M}$ , and any MDP  $\mathbf{M} \in \mathfrak{M}$ , it holds for all  $T$  such that  $\eta_T \leq \frac{1}{2|\mathcal{X}|}$ ,

$$\bar{\mathfrak{R}}_{T,s_1}(\pi; \mathbf{M}) \geq \bar{\mathbb{C}}_{T,s_1}(\mathbf{M}, \mathfrak{M}),$$

where we introduced

$$\bar{\mathbb{C}}_{T,s_1}(\mathbf{M}, \mathfrak{M}) = \inf_{n \in T \cdot \mathcal{P}_{|\mathcal{X}|}} \left\{ \sum_{x \in \mathcal{X}} n_x \Delta(x) - \mathbb{S}(b_*) : \forall \pi, \inf_{\tilde{\mathbf{M}} \in \bar{\mathcal{B}}_{s_1}(\pi, \mathbf{M}; \mathfrak{M})} \sum_{x \in \mathcal{X}} \frac{n_x \text{KL}_x(\mathbf{M}, \tilde{\mathbf{M}})}{\mathbf{k}1(1 - \eta_T |\mathcal{X}|, \eta_T |\mathcal{X}|)} \right\}$$



▷ **Finite time lower bound** For any **uniformly  $\eta$ -good** strategy on  $\mathfrak{M}$ , and any MDP  $\mathbf{M} \in \mathfrak{M}$ , it holds for all  $T$  such that  $\eta_T \leq \frac{1}{2^{|\mathcal{X}|}}$ ,

$$\bar{\mathfrak{R}}_{T,s_1}(\pi; \mathbf{M}) \geq \bar{\mathfrak{C}}_{T,s_1}(\mathbf{M}, \mathfrak{M}),$$

where we introduced

$$\bar{\mathfrak{C}}_{T,s_1}(\mathbf{M}, \mathfrak{M}) = \inf_{n \in T \cdot \mathcal{P}_{|\mathcal{X}|}} \left\{ \sum_{x \in \mathcal{X}} n_x \Delta(x) - \mathbb{S}(b_*) : \forall \pi, \inf_{\tilde{\mathbf{M}} \in \bar{\mathcal{B}}_{s_1}(\pi, \mathbf{M}; \mathfrak{M})} \sum_{x \in \mathcal{X}} \frac{n_x \text{KL}_x(\mathbf{M}, \tilde{\mathbf{M}})}{\mathbf{k}1(1 - \eta_T |\mathcal{X}|, \eta_T |\mathcal{X}|)} \right\}$$

▷ **Asymptotic lower bound** Furthermore, for  $\eta_T = O(T^{\alpha-1})$  for all  $\alpha \in (0, 1)$ , asymptotically, we get

$$\liminf_T \frac{\bar{\mathfrak{R}}_{T,s_1}(\pi, \mathbf{M})}{\ln(T)} \geq \inf_{\kappa \in \mathbb{R}_+^{|\mathcal{X}|}} \left\{ \sum_{x \in \mathcal{X}} \kappa_x \Delta(x) : \forall \pi, \inf_{\tilde{\mathbf{M}} \in \bar{\mathcal{B}}_{s_1}(\pi, \mathbf{M}; \mathfrak{M})} \sum_{x \in \mathcal{X}} \kappa_x \text{KL}_x(\mathbf{M}, \tilde{\mathbf{M}}) \geq 1 \right\}$$



- ▷ Computing **confusing MDP** hence cost implies searching over **all places** where MDP could be modified ( $|\mathcal{X}_{*s,a} \setminus \mathcal{X}_*|$ , possibly  $O(SA)$  pairs). Conjecture: cost is exponential in  $|\mathcal{X}_{*s,a} \setminus \mathcal{X}_*|$ .
- ▷ In non-unichain MDPs, there may be no improving policies over  $\pi$  at Hamming distance 1 of  $\pi$ .



*"The more applied you go, the stronger theory you need"*

# MERCI



odalricambrym.maillard@inria.fr  
odalricambrymmaillard.wordpress.com

