



# Efficient Regret Minimizing Strategies for Tabular Average-Reward MDPs

M. Sadegh Talebi  
Inria Lille - Nord Europe (SequeL)

Joint work with  
Odalric-Ambrym Maillard (Inria) and Hippolyte Bourel (Montanuniversität)

8ème Journée COSMOS, November 2019

# Undiscounted RL: MDP Model

We consider reinforcement learning (RL), where the environment is modeled as an undiscounted **Markov Decision Process (MDP)**.

**Undiscounted MDP**  $M = (\mathcal{S}, \mathcal{A}, p, \mu)$ :

- **State-space**  $\mathcal{S}$  with cardinality  $S$
- **Action-space**  $\mathcal{A}$  with cardinality  $A$
- **Transition kernel**  $p$ : Selecting  $a \in \mathcal{A}$  in  $s \in \mathcal{S}$  leads to a transition to  $s'$  with probability  $p(s'|s, a)$ .
- **Reward function**  $\mu$ : Selecting  $a \in \mathcal{A}$  in  $s \in \mathcal{S}$ , gives  $r(s, a)$  with mean  $\mu(s, a)$ .



$p$  and  $\mu$  are **unknown**, and the goal is to maximize  $\sum_{t=1}^T r_t$ .

# Undiscounted RL: MDP Model

We consider reinforcement learning (RL), where the environment is modeled as an undiscounted **Markov Decision Process (MDP)**.

**Undiscounted MDP**  $M = (\mathcal{S}, \mathcal{A}, p, \mu)$ :

- **State-space**  $\mathcal{S}$  with cardinality  $S$
- **Action-space**  $\mathcal{A}$  with cardinality  $A$
- **Transition kernel**  $p$ : Selecting  $a \in \mathcal{A}$  in  $s \in \mathcal{S}$  leads to a transition to  $s'$  with probability  $p(s'|s, a)$ .
- **Reward function**  $\mu$ : Selecting  $a \in \mathcal{A}$  in  $s \in \mathcal{S}$ , gives  $r(s, a)$  with mean  $\mu(s, a)$ .



$p$  and  $\mu$  are **unknown**, and the goal is to maximize  $\sum_{t=1}^T r_t$ .

# Undiscounted RL: Objective

**Goal:** To maximize the collected reward  $\sum_{t=1}^T r_t$ .

- A (Markov deterministic) **policy**  $\pi$  is a mapping from  $\mathcal{S}$  to  $\mathcal{A}$ .
- **Gain** (or long-term average reward) of a policy  $\pi$  is defined as

$$g^\pi(s_1) := \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T r_t(s_t, \pi(s_t)) \right]$$

- **Assumption:** We consider **communicating** MDPs in which every state is reachable from any other state by some appropriate policy. For communicating MDPs,  $g^\pi$  does not depend on  $s_1$ .

# Undiscounted RL: Objective

**Goal:** To maximize the collected reward  $\sum_{t=1}^T r_t$ .

- A (Markov deterministic) **policy**  $\pi$  is a mapping from  $\mathcal{S}$  to  $\mathcal{A}$ .
- **Gain** (or long-term average reward) of a policy  $\pi$  is defined as

$$g^\pi(s_1) := \liminf_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{t=1}^T r_t(s_t, \pi(s_t)) \right]$$

- **Assumption:** We consider **communicating** MDPs in which every state is reachable from any other state by some appropriate policy. For communicating MDPs,  $g^\pi$  does not depend on  $s_1$ .

# Undiscounted RL: Bellman's Equation

Any policy achieving  $g^* := \max_{\pi} g^{\pi}$  is called **gain-optimal**.

## Bellman's Optimality Equation (Poisson Equation)

$$g^* + b^*(s) = \max_{a \in \mathcal{A}} \left( \mu(s, a) + \sum_{s' \in \mathcal{S}} p(s'|s, a) b^*(s') \right), \quad \forall s$$

where  $g^*$  is called the **maximal gain** and  $b^*$  is called the **optimal bias function**.

- In the long run, maximal cumulative reward is achieved by following a gain-optimal policy.
- If MDP is known, one can find  $g^*$  and  $b^*$  by solving Bellman's optimality equation using numerical methods (e.g., **Value Iteration**).

# Undiscounted RL: Bellman's Equation

Any policy achieving  $g^* := \max_{\pi} g^{\pi}$  is called **gain-optimal**.

## Bellman's Optimality Equation (Poisson Equation)

$$g^* + b^*(s) = \max_{a \in \mathcal{A}} \left( \mu(s, a) + \sum_{s' \in \mathcal{S}} p(s'|s, a) b^*(s') \right), \quad \forall s$$

where  $g^*$  is called the **maximal gain** and  $b^*$  is called the **optimal bias function**.

- In the long run, maximal cumulative reward is achieved by following a gain-optimal policy.
- If MDP is known, one can find  $g^*$  and  $b^*$  by solving Bellman's optimality equation using numerical methods (e.g., **Value Iteration**).

# Undiscounted RL: Regret

**Goal:** To maximize the collected reward  $\sum_{t=1}^T r_t$ .

**Regret:** Defined as the difference between cumulative reward of the optimal policy  $\star$  and that gathered by the decision-maker (in expectation or w.h.p.):

$$\text{Regret}_T := \sum_{t=1}^T r_t^* - \sum_{t=1}^T r_t$$

Alternatively, the objective of the decision-maker is to **minimize the regret**.  
By Azuma-Hoeffding's inequality, with probability at least  $1 - \delta$ ,

$$\text{Regret}_T := Tg^* - \sum_{t=1}^T r_t + \mathcal{O}(\sqrt{T \log(2/\delta)})$$

So it makes sense to control the following notion of regret:

$$\mathfrak{R}_T := Tg^* - \sum_{t=1}^T r_t$$

# Undiscounted RL: Regret

**Goal:** To maximize the collected reward  $\sum_{t=1}^T r_t$ .

**Regret:** Defined as the difference between cumulative reward of the optimal policy  $\star$  and that gathered by the decision-maker (in expectation or w.h.p.):

$$\text{Regret}_T := \sum_{t=1}^T r_t^\star - \sum_{t=1}^T r_t$$

Alternatively, the objective of the decision-maker is to **minimize the regret**. By Azuma-Hoeffding's inequality, with probability at least  $1 - \delta$ ,

$$\text{Regret}_T := Tg^\star - \sum_{t=1}^T r_t + \mathcal{O}(\sqrt{T \log(2/\delta)})$$

So it makes sense to control the following notion of regret:

$$\mathfrak{R}_T := Tg^\star - \sum_{t=1}^T r_t$$

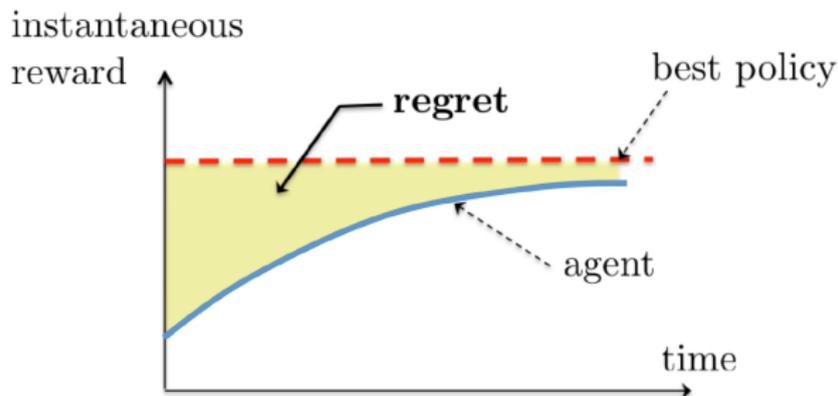
# Undiscounted RL: Regret

Alternatively, the objective of the decision-maker is to **minimize the regret**.

$$\mathfrak{R}_T := Tg^* - \sum_{t=1}^T r_t$$

The key difficulty to do so is to balance *exploration* vs. *exploitation*:

- Play the best action so far, ...
- ... or rather explore a different action?



# Outline

- 1 UCRL2
- 2 UCRL3
- 3 KL-UCRL
- 4 Numerical Experiments
- 5 Technical Tools

# Outline

- 1 UCRL2
- 2 UCRL3
- 3 KL-UCRL
- 4 Numerical Experiments
- 5 Technical Tools

Two main approaches in RL:

- **Model-Based:** Consists in maintaining **an approximate MDP model** through estimating  $\mu$  and  $p$ , and deriving a **value function** from the *approximate* MDP.
  - Examples: UCB1, UCRL2.
- **Model-Free:** Directly learns a value function (without estimating  $\mu$  and  $p$ ).
  - Example: Variants of Q-learning.

In this talk we focus on model-based algorithms.

Two main approaches in RL:

- **Model-Based:** Consists in maintaining **an approximate MDP model** through estimating  $\mu$  and  $p$ , and deriving a **value function** from the *approximate* MDP.
  - Examples: UCB1, UCRL2.
- **Model-Free:** Directly learns a value function (without estimating  $\mu$  and  $p$ ).
  - Example: Variants of Q-learning.

In this talk we focus on model-based algorithms.

Under a given algorithm, we define:

- $N_t(s, a)$ : number of visits, up to time  $t$ , to  $(s, a)$ .
- $N_t(s, a, s')$ : number of visits, up to time  $t$ , to  $(s, a)$  followed by a visit to  $s'$ .
- Empirical estimates of transition probabilities and rewards:

$$\hat{\mu}_t(s, a) = \frac{\sum_{t'=0}^{t-1} r_{t'} \mathbb{I}\{s_{t'} = s, a_{t'} = a\}}{N_t(s, a)^+}$$

$$\hat{p}_t(s'|s, a) = \frac{N_t(s, a, s')}{N_t(s, a)^+}$$

with  $N_t(s, a)^+ := \max\{N_t(s, a), 1\}$ .

**UCRL2** (Jaksch et al., 2010): a **model-based** algorithm inspired by **UCB** for stochastic bandits:

- Maintains **confidence bounds** for  $\mu$  and  $p$ , and chooses an optimistic model that leads to the **highest gain**  $g$ .

Given  $\delta \in (0, 1)$ , UCRL2 defines a set  $\mathcal{M}_{t,\delta}$  of **plausible MDPs (models)** at time  $t$  as a collection of **candidate MDPs**  $M' = (\mathcal{S}, \mathcal{A}, \mu', p')$  satisfying:

For all  $s, a$ ,

$$\|\hat{p}_t(\cdot|s, a) - p'(\cdot|s, a)\|_1 \leq \sqrt{\frac{14S}{N_t(s, a)} \log\left(\frac{2At}{\delta}\right)}$$

$$|\hat{\mu}_t(s, a) - \mu'(s, a)| \leq \sqrt{\frac{7}{2N_t(s, a)} \log\left(\frac{2SA t}{\delta}\right)}$$

$\Rightarrow$  With **high probability**,  $M \in \mathcal{M}_{t,\delta}$ .

**UCRL2** (Jaksch et al., 2010): a **model-based** algorithm inspired by **UCB** for stochastic bandits:

- Maintains **confidence bounds** for  $\mu$  and  $p$ , and chooses an optimistic model that leads to the **highest gain**  $g$ .

Given  $\delta \in (0, 1)$ , UCRL2 defines a set  $\mathcal{M}_{t,\delta}$  of **plausible MDPs (models)** at time  $t$  as a collection of **candidate MDPs**  $M' = (\mathcal{S}, \mathcal{A}, \mu', p')$  satisfying:

For all  $s, a$ ,

$$\|\hat{p}_t(\cdot | s, a) - p'(\cdot | s, a)\|_1 \leq \sqrt{\frac{14S}{N_t(s, a)} \log\left(\frac{2At}{\delta}\right)}$$

$$|\hat{\mu}_t(s, a) - \mu'(s, a)| \leq \sqrt{\frac{7}{2N_t(s, a)} \log\left(\frac{2SA t}{\delta}\right)}$$

$\Rightarrow$  With **high probability**,  $M \in \mathcal{M}_{t,\delta}$ .

**UCRL2** (Jaksch et al., 2010): a **model-based** algorithm inspired by **UCB** for stochastic bandits:

- Maintains **confidence bounds** for  $\mu$  and  $p$ , and chooses an optimistic model that leads to the **highest gain**  $g$ .

Given  $\delta \in (0, 1)$ , UCRL2 defines a set  $\mathcal{M}_{t,\delta}$  of **plausible MDPs (models)** at time  $t$  as a collection of **candidate MDPs**  $M' = (\mathcal{S}, \mathcal{A}, \mu', p')$  satisfying:

For all  $s, a$ ,

$$\|\hat{p}_t(\cdot|s, a) - p'(\cdot|s, a)\|_1 \leq \sqrt{\frac{14S}{N_t(s, a)} \log\left(\frac{2At}{\delta}\right)}$$

$$|\hat{\mu}_t(s, a) - \mu'(s, a)| \leq \sqrt{\frac{7}{2N_t(s, a)} \log\left(\frac{2SA t}{\delta}\right)}$$

$\Rightarrow$  With **high probability**,  $M \in \mathcal{M}_{t,\delta}$ .

---

## Algorithm 1 UCRL2

---

**Initialize:** For all  $(s, a)$ , set  $N_0(s, a) = 0$  and  $v_0(s, a) = 0$ . Set  $t_0 = 0$ ,  $t = 1$ ,  $k = 1$ , and observe the initial state  $s_1$ ;

**for** episodes  $k \geq 1$  **do**

    Set  $t_k = t$ ;

    Set  $N_{t_k}(s, a) = N_{t_{k-1}}(s, a) + v_k(s, a)$  for all  $(s, a)$ ;

    Compute  $\hat{\mu}_{t_k}(s, a)$  and  $\hat{p}_{t_k}(\cdot|s, a)$  for all  $(s, a)$ ;

    Compute  $\pi_{t_k}^+ = \text{EVI}\left(\hat{\mu}_{t_k}, \hat{p}_{t_k}, N_{t_k}, \frac{1}{\sqrt{t_k}}, \frac{\delta}{SA}\right)$ ;

**while**  $v_k(s_t, \pi_{t_k}^+(s_t)) < \max\{1, N_{t_k}(s_t, \pi_{t_k}^+(s_t))\}$  **do**

        Play  $a_t = \pi_{t_k}^+(s_t)$ , and observe  $s_{t+1}$  and  $r_t(s_t, a_t)$ ;

        Set  $v_k(s_t, a_t) = v_k(s_t, a_t) + 1$ ;

        Set  $t = t + 1$ ;

**end while**

**end for**

---

EVI stands for Extended Value Iteration

---

**Algorithm 2**  $\text{EVI}(\mu, p, N, \varepsilon, \delta)$

---

**Initialize:**  $u^{(0)} \equiv 0, u^{(-1)} \equiv -\infty, n = 0;$

**while**  $\max_s (u^{(n)}(s) - u^{(n-1)}(s)) - \min_s (u^{(n)}(s) - u^{(n-1)}(s)) > \varepsilon$  **do**

    For all  $(s, a)$ , set  $\mu'(s, a) = \mu(s, a) + \beta'_{N(s,a)}(\delta);$

    For all  $(s, a)$ , set  $p'(\cdot|s, a) \in \operatorname{argmax}_{q \in \mathcal{P}(s,a)} \sum_{x \in \mathcal{S}} q(x)u^{(n)}(x)$  where

$$\mathcal{P}(s, a) := \left\{ q \in \Delta^{\mathcal{S}} : \|q - p(\cdot|s, a)\|_1 \leq \beta_{N(s,a)}(\delta) \right\};$$

    For all  $s$ , update  $u^{(n+1)}(s) = \max_{a \in \mathcal{A}} \left( \mu'(s, a) + \sum_{x \in \mathcal{S}} p'(x|s, a)u^{(n)}(x) \right);$

    For all  $s$ , update  $\pi_{n+1}(s) \in \operatorname{argmax}_{a \in \mathcal{A}} \left( \mu'(s, a) + \sum_{x \in \mathcal{S}} p'(x|s, a)u^{(n)}(x) \right);$

    Set  $n = n + 1;$

**end while**

**Output:**  $\pi_{n+1}$

---

## Definition (Diameter (Jaksch et al., 2010))

Let  $T_\pi(s'|s)$  denote the first hitting time of state  $s'$  when following stationary policy  $\pi$  from initial state  $s$ . The diameter  $D$  of an MDP  $M$  is defined as

$$D := \max_{s \neq s'} \min_{\pi} \mathbb{E}[T_\pi(s'|s)].$$

For any communicating MDP, under UCRL2, with probability at least  $1 - \delta$ ,

$$\mathfrak{R}_T \leq 34DS \sqrt{AT \log(T/\delta)}$$

Minimax lower bound (Jaksch et al., 2010):  $\Omega(\sqrt{DSAT})$

## Definition (Diameter (Jaksch et al., 2010))

Let  $T_\pi(s'|s)$  denote the first hitting time of state  $s'$  when following stationary policy  $\pi$  from initial state  $s$ . The diameter  $D$  of an MDP  $M$  is defined as

$$D := \max_{s \neq s'} \min_{\pi} \mathbb{E}[T_\pi(s'|s)].$$

For any communicating MDP, under UCRL2, with probability at least  $1 - \delta$ ,

$$\mathfrak{R}_T \leq 34DS\sqrt{AT \log(T/\delta)}$$

Minimax lower bound (Jaksch et al., 2010):  $\Omega(\sqrt{DSAT})$

## Definition (Diameter (Jaksch et al., 2010))

Let  $T_\pi(s'|s)$  denote the first hitting time of state  $s'$  when following stationary policy  $\pi$  from initial state  $s$ . The diameter  $D$  of an MDP  $M$  is defined as

$$D := \max_{s \neq s'} \min_{\pi} \mathbb{E}[T_\pi(s'|s)].$$

For any communicating MDP, under UCRL2, with probability at least  $1 - \delta$ ,

$$\mathfrak{R}_T \leq 34DS \sqrt{AT \log(T/\delta)}$$

Minimax lower bound (Jaksch et al., 2010):  $\Omega(\sqrt{DSAT})$

Despite its strong regret guarantee, UCRL2 does not perform well in practice (even in small environments) – In particular, it suffers from a long burn-in phase.

Drawbacks of UCRL2:

- (i) Loose and polytopic set of models
- (ii) Conservative optimistic policy
- (iii) Inefficient stopping criterion for internal episodes

We discuss two variants of UCRL2 aiming to remove (i) and (ii).

Despite its strong regret guarantee, UCRL2 does not perform well in practice (even in small environments) – In particular, it suffers from a long burn-in phase.

Drawbacks of UCRL2:

- (i) Loose and polytopic set of models
- (ii) Conservative optimistic policy
- (iii) Inefficient stopping criterion for internal episodes

We discuss two variants of UCRL2 aiming to remove (i) and (ii).

# Outline

1 UCRL2

2 UCRL3

3 KL-UCRL

4 Numerical Experiments

5 Technical Tools

UCRL3 is a variant of UCRL2, with the following key differences:

- Uses **tight element-wise** confidence intervals for  $p$ 
  - Defined for individual transition probabilities  $p(s'|s, a)$ , in contrast to UCRL2 that does for  $p(\cdot|s, a)$ .
  - Intersection of **time-uniform** Bernstein and Bernoulli concentration for each  $p(s'|s, a)$
- Computes a **less conservative** optimistic policy.

To simplify the presentation, we assume that  $\mu$  is known.

UCRL3 is a variant of UCRL2, with the following key differences:

- Uses **tight element-wise** confidence intervals for  $p$ 
  - Defined for individual transition probabilities  $p(s'|s, a)$ , in contrast to UCRL2 that does for  $p(\cdot|s, a)$ .
  - Intersection of **time-uniform** Bernstein and Bernoulli concentration for each  $p(s'|s, a)$
- Computes a **less conservative** optimistic policy.

To simplify the presentation, we assume that  $\mu$  is known.

At time  $t$ , UCRL3 considers the set  $\mathcal{M}_{t,\delta}$  of plausible MDPs

$$\mathcal{M}_{t,\delta} = \left\{ M' = (\mathcal{S}, \mathcal{A}, p', \mu) : p'(\cdot | s, a) \in \mathcal{C}_{t,\delta}(s, a), \quad \forall s, a, s' \right\}$$

where for each  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\mathcal{C}_{t,\delta}(s, a) := \left\{ q \in \Delta_{\mathcal{S}} : \forall s', q(s') \in \underbrace{\mathcal{C}_{t,\delta}^1(s', s, a)}_{\text{Bernstein}} \cap \underbrace{\mathcal{C}_{t,\delta}^2(s', s, a)}_{\text{sub-Gaussian}} \right\}$$

- $\mathcal{C}_{t,\delta}^1(s', s, a)$  is defined using Bernstein concentration inequality, modified using **a peeling technique**.
- $\mathcal{C}_{t,\delta}^2(s', s, a)$  is obtained by applying **the method of mixture** (a.k.a. **the Laplace method**) for sub-Gaussian random variables.

## UCRL3: Set of Models

At time  $t$ , UCRL3 considers the set  $\mathcal{M}_{t,\delta}$  of plausible MDPs

$$\mathcal{M}_{t,\delta} = \left\{ M' = (\mathcal{S}, \mathcal{A}, p', \mu) : p'(\cdot | s, a) \in \mathcal{C}_{t,\delta}(s, a), \quad \forall s, a, s' \right\}$$

where for each  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\mathcal{C}_{t,\delta}(s, a) := \left\{ q \in \Delta_{\mathcal{S}} : \forall s', q(s') \in \underbrace{\mathcal{C}_{t,\delta}^1(s', s, a)}_{\text{Bernstein}} \cap \underbrace{\mathcal{C}_{t,\delta}^2(s', s, a)}_{\text{sub-Gaussian}} \right\}$$

$$\mathcal{C}_{t,\delta}^1(s', s, a) = \left\{ \lambda : |\hat{p}_t(s' | s, a) - \lambda| \leq \sqrt{\frac{2\lambda(1-\lambda)\ell_{N_t(s,a)}\left(\frac{\delta}{2SA}\right)}{N_t(s,a)}} + \frac{\ell_{N_t(s,a)}\left(\frac{\delta}{2SA}\right)}{3N_t(s,a)} \right\}$$

where  $\ell_n(\delta) = \eta \log\left(\frac{\log(n)\log(\eta n)}{\log(\eta^2)\delta}\right)$  with  $\eta = 1.12$  (an arbitrary choice).

## UCRL3: Set of Models

At time  $t$ , UCRL3 considers the set  $\mathcal{M}_{t,\delta}$  of plausible MDPs

$$\mathcal{M}_{t,\delta} = \left\{ M' = (\mathcal{S}, \mathcal{A}, p', \mu) : p'(\cdot | s, a) \in \mathcal{C}_{t,\delta}(s, a), \quad \forall s, a, s' \right\}$$

where for each  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\mathcal{C}_{t,\delta}(s, a) := \left\{ q \in \Delta_{\mathcal{S}} : \forall s', q(s') \in \underbrace{\mathcal{C}_{t,\delta}^1(s', s, a)}_{\text{Bernstein}} \cap \underbrace{\mathcal{C}_{t,\delta}^2(s', s, a)}_{\text{sub-Gaussian}} \right\}$$

$$\mathcal{C}_{t,\delta}^1(s', s, a) = \left\{ \lambda : |\hat{p}_t(s' | s, a) - \lambda| \leq \sqrt{\frac{2\lambda(1-\lambda)\ell_{N_t(s,a)}\left(\frac{\delta}{2SA}\right)}{N_t(s,a)}} + \frac{\ell_{N_t(s,a)}\left(\frac{\delta}{2SA}\right)}{3N_t(s,a)} \right\}$$

where  $\ell_n(\delta) = \eta \log\left(\frac{\log(n)\log(\eta n)}{\log(\eta^2)\delta}\right)$  with  $\eta = 1.12$  (an arbitrary choice).

# UCRL3: Set of Models

At time  $t$ , UCRL3 considers the set  $\mathcal{M}_{t,\delta}$  of plausible MDPs

$$\mathcal{M}_{t,\delta} = \left\{ M' = (\mathcal{S}, \mathcal{A}, p', \mu) : p'(\cdot | s, a) \in \mathcal{C}_{t,\delta}(s, a), \quad \forall s, a, s' \right\}$$

where for each  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\mathcal{C}_{t,\delta}(s, a) := \left\{ q \in \Delta_{\mathcal{S}} : \forall s', q(s') \in \underbrace{\mathcal{C}_{t,\delta}^1(s', s, a)}_{\text{Bernstein}} \cap \underbrace{\mathcal{C}_{t,\delta}^2(s', s, a)}_{\text{sub-Gaussian}} \right\}$$

$$\mathcal{C}_{t,\delta}^2(s', s, a) = \left\{ \lambda : -\sqrt{\underline{g}(\lambda)} \beta_{N_t(s,a)} \left( \frac{\delta}{2SA} \right) \leq \hat{p}_t(s' | s, a) - \lambda \leq \sqrt{g(\lambda)} \beta_{N_t(s,a)} \left( \frac{\delta}{2SA} \right) \right\}$$

where  $\beta_n(\delta) := \sqrt{\frac{(1+\frac{1}{n}) \log(\sqrt{n+1}/\delta)}{2n}}$ , and where

$$\underline{g}(\lambda) = \begin{cases} g(\lambda) & \text{if } \lambda < 1/2 \\ \lambda(1-\lambda) & \text{else} \end{cases}, \text{ and } g(\lambda) = \frac{1/2 - \lambda}{\log(1/\lambda - 1)}.$$

# UCRL3: Set of Models

At time  $t$ , UCRL3 considers the set  $\mathcal{M}_{t,\delta}$  of plausible MDPs

$$\mathcal{M}_{t,\delta} = \left\{ M' = (\mathcal{S}, \mathcal{A}, p', \mu) : p'(\cdot | s, a) \in \mathcal{C}_{t,\delta}(s, a), \quad \forall s, a, s' \right\}$$

where for each  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\mathcal{C}_{t,\delta}(s, a) := \left\{ q \in \Delta_{\mathcal{S}} : \forall s', q(s') \in \underbrace{\mathcal{C}_{t,\delta}^1(s', s, a)}_{\text{Bernstein}} \cap \underbrace{\mathcal{C}_{t,\delta}^2(s', s, a)}_{\text{sub-Gaussian}} \right\}$$

$$\mathcal{C}_{t,\delta}^2(s', s, a) = \left\{ \lambda : -\sqrt{\underline{g}(\lambda)} \beta_{N_t(s,a)} \left( \frac{\delta}{2SA} \right) \leq \hat{p}_t(s' | s, a) - \lambda \leq \sqrt{g(\lambda)} \beta_{N_t(s,a)} \left( \frac{\delta}{2SA} \right) \right\}$$

where  $\beta_n(\delta) := \sqrt{\frac{(1+\frac{1}{n}) \log(\sqrt{n+1}/\delta)}{2n}}$ , and where

$$\underline{g}(\lambda) = \begin{cases} g(\lambda) & \text{if } \lambda < 1/2 \\ \lambda(1-\lambda) & \text{else} \end{cases}, \text{ and } g(\lambda) = \frac{1/2 - \lambda}{\log(1/\lambda - 1)}.$$

## UCRL3: Set of Models

At time  $t$ , UCRL3 considers the set  $\mathcal{M}_{t,\delta}$  of plausible MDPs

$$\mathcal{M}_{t,\delta} = \left\{ M' = (\mathcal{S}, \mathcal{A}, p', \mu) : p'(\cdot | s, a) \in \mathcal{C}_{t,\delta}(s, a), \quad \forall s, a, s' \right\}$$

where for each  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\mathcal{C}_{t,\delta}(s, a) := \left\{ q \in \Delta_{\mathcal{S}} : \forall s', q(s') \in \underbrace{C_{t,\delta}^1(s', s, a)}_{\text{Bernstein}} \cap \underbrace{C_{t,\delta}^2(s', s, a)}_{\text{sub-Gaussian}} \right\}$$

Lemma (Time-uniform confidence bounds)

For any MDP with transition function  $p$ , for all  $\delta \in (0, 1)$ , it holds

$$\mathbb{P}\left(\exists t \in \mathbb{N}, \exists (s, a) \in \mathcal{S} \times \mathcal{A}, p(\cdot | s, a) \notin \mathcal{C}_{t,\delta}(s, a)\right) \leq \delta$$

$$\implies \mathbb{P}(\exists t \in \mathbb{N}, M \notin \mathcal{M}_{t,\delta}) \leq \delta.$$

## UCRL3: Set of Models

At time  $t$ , UCRL3 considers the set  $\mathcal{M}_{t,\delta}$  of plausible MDPs

$$\mathcal{M}_{t,\delta} = \left\{ M' = (\mathcal{S}, \mathcal{A}, p', \mu) : p'(\cdot | s, a) \in \mathcal{C}_{t,\delta}(s, a), \quad \forall s, a, s' \right\}$$

where for each  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\mathcal{C}_{t,\delta}(s, a) := \left\{ q \in \Delta_{\mathcal{S}} : \forall s', q(s') \in \underbrace{C_{t,\delta}^1(s', s, a)}_{\text{Bernstein}} \cap \underbrace{C_{t,\delta}^2(s', s, a)}_{\text{sub-Gaussian}} \right\}$$

### Lemma (Time-uniform confidence bounds)

For any MDP with transition function  $p$ , for all  $\delta \in (0, 1)$ , it holds

$$\mathbb{P}\left(\exists t \in \mathbb{N}, \exists (s, a) \in \mathcal{S} \times \mathcal{A}, p(\cdot | s, a) \notin \mathcal{C}_{t,\delta}(s, a)\right) \leq \delta$$

$$\implies \mathbb{P}(\exists t \in \mathbb{N}, M \notin \mathcal{M}_{t,\delta}) \leq \delta.$$

## UCRL3: Revisiting EVI

- To compute an optimistic policy, UCRL2 uses EVI as a subroutine, which involves computing

$$p_n^+ : s, a \mapsto \operatorname{argmax}\{p'u_n, p' \in \mathcal{C}_{t,\delta}(s, a)\}$$

at iteration  $n$  of EVI.

- EVI outputs a conservative policy, in particular when transition function  $p$  has a sparse support.
- UCRL3 remedies this issue by combining EVI with an **adaptive support selection**.

# Adaptive Support Selection

Given  $\tilde{\mathcal{S}} \subset \mathcal{S}$ , a pair  $(s, a)$ , and a function  $f : \mathcal{S} \rightarrow \mathbb{R}$ , define:

$$\bar{f}_{s,a}(\tilde{\mathcal{S}}) = \max \left\{ \sum_{s' \in \tilde{\mathcal{S}}} f(s')q(s') : q \text{ s.t. } \forall s' \in \tilde{\mathcal{S}}, q(s') \in \mathcal{C}_{t,\delta}(s', s, a) \text{ and } \sum_{s' \in \tilde{\mathcal{S}}} q(s') \leq 1 \right\}$$

$$\bar{p}_{s,a} = \operatorname{argmax} \left\{ \sum_{s' \in \tilde{\mathcal{S}}} f(s')q(s') : q \text{ s.t. } \forall s' \in \tilde{\mathcal{S}}, q(s') \in \mathcal{C}_{t,\delta}(s', s, a) \text{ and } \sum_{s' \in \tilde{\mathcal{S}}} q(s') \leq 1 \right\}$$

---

## Algorithm 3 Adaptive Support Selection (for $(s, a)$ )

---

**Input:** Target function  $f$ , parameter  $\kappa \in (0, 1)$ .

Let  $\tilde{\mathcal{S}} = \operatorname{supp}(\hat{p}_t(\cdot|s, a)) \cup \operatorname{argmax}_{s \in \mathcal{S}} f(s)$

**while**  $\bar{f}_{s,a}(\mathcal{S} \setminus \tilde{\mathcal{S}}) \geq \min(\kappa, \bar{f}_{s,a}(\tilde{\mathcal{S}}))$  **do**

    Let  $\tilde{s} \in \operatorname{argmax}_{s \notin \tilde{\mathcal{S}}} f(s)$

    Set  $\tilde{\mathcal{S}} = \tilde{\mathcal{S}} \cup \{\tilde{s}\}$

**end while**

**Output:**  $\tilde{\mathcal{S}}, \bar{p}_{s,a}$

---

## UCRL3: Revisiting EVI

- Recall that UCRL2 uses EVI as a subroutine, which involves computing

$$p_n^+ : s, a \mapsto \operatorname{argmax}\{P' u_n, p' \in \mathcal{C}_{t,\delta}(s, a)\}$$

at iteration  $n$  of EVI.

- Now, at iteration  $n$  of EVI, UCRL3 uses **Adaptive Support Selection** with  $f = u_n - \min_s u_n(s)$ .
- To optimize performance, we choose

$$\kappa = \kappa_{t,n}(s, a) = \frac{\mathbb{S}(u_n) |\operatorname{supp}(\hat{p}_t(\cdot | s, a))|}{\max_{s,a} N_t(s, a)^{2/3}}$$

## Theorem

*The regret under UCRL3 in any communicating MDP satisfies, uniformly over all  $T \geq 1$ ,*

$$\mathfrak{R}_T \leq 24D\sqrt{KSAT \log(\sqrt{T+1}/\delta)} + \tilde{O}(DS^{2/3}A^{2/3}T^{1/3})$$

*with probability at least  $1 - 2\delta$ .*

- Improves the regret of UCRL2 by a factor of  $\sqrt{S/K}$ .
- Holds uniformly over all  $T \geq 1$ .

## Theorem

*The regret under UCRL3 in any communicating MDP satisfies, uniformly over all  $T \geq 1$ ,*

$$\mathfrak{R}_T \leq 24D\sqrt{KSAT \log(\sqrt{T+1}/\delta)} + \tilde{O}(DS^{2/3}A^{2/3}T^{1/3})$$

*with probability at least  $1 - 2\delta$ .*

- Improves the regret of UCRL2 by a factor of  $\sqrt{S/K}$ .
- Holds uniformly over all  $T \geq 1$ .

# Outline

1 UCRL2

2 UCRL3

**3 KL-UCRL**

4 Numerical Experiments

5 Technical Tools

# Variants of UCRL2

There are variants of UCRL2, that mostly differ in the definition of models.

Two approaches to define the set of  $\mathcal{M}_{t,\delta}$  of models depending on how uncertainties in  $p$  and  $\mu$  are represented:

- **Polytopic** uncertainty sets
  - For example, models defined using Weissman's and Hoeffding's inequalities (as in UCRL2).
- **Non-polytopic** uncertainty sets
  - Smoother sets
  - For example, models defined using KL-divergence and Bernstein's inequality (as in (Burnetas & Katehakis, 1997), **KL-UCRL** (Filippi et al., 2010)).

**Polytopic uncertainty** models typically provide **poor representations** (cf. Robust control of MDPs (Nilim & El Ghaoui, 2005) and (Filippi et al., 2010)):

- (i) They could lead to **inconsistent** models by **excluding** an already observed element of kernel (i.e.,  $p'(x|s, a) = 0$  even though  $\hat{p}_t(x|s, a) \neq 0$  for some  $x$ ).
- (ii) The maximizer of a linear optimization over  $L_1$  ball could change **significantly** for a small change in the value function.

**Polytopic uncertainty** models typically provide **poor representations** (cf. Robust control of MDPs (Nilim & El Ghaoui, 2005) and (Filippi et al., 2010)):

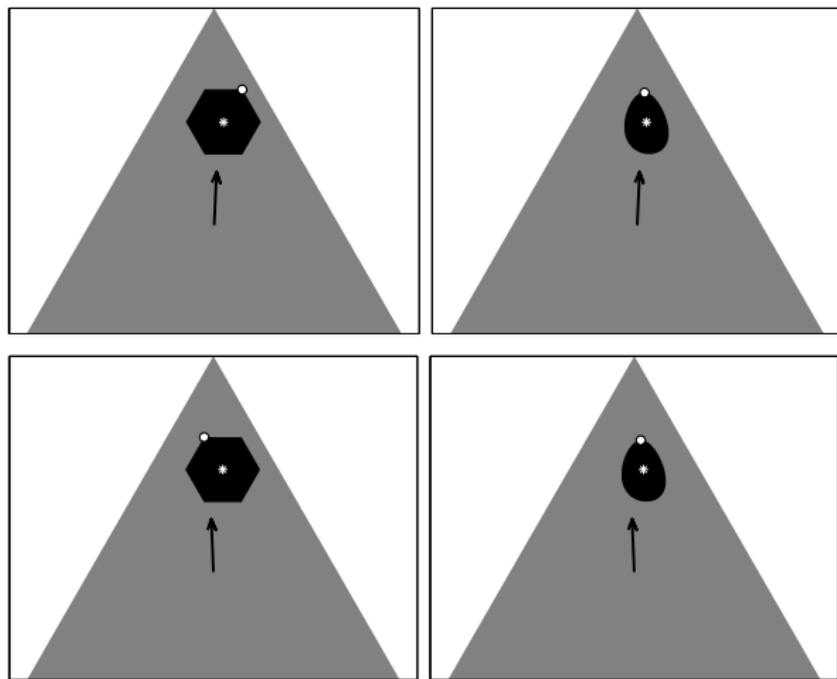
- (i) They could lead to **inconsistent** models by **excluding** an already observed element of kernel (i.e.,  $p'(x|s, a) = 0$  even though  $\hat{p}_t(x|s, a) \neq 0$  for some  $x$ ).
- (ii) The maximizer of a linear optimization over  $L_1$  ball could change **significantly** for a small change in the value function.

**Polytopic uncertainty** models typically provide **poor representations** (cf. Robust control of MDPs (Nilim & El Ghaoui, 2005) and (Filippi et al., 2010)):

- (i) They could lead to **inconsistent** models by **excluding** an already observed element of kernel (i.e.,  $p'(x|s, a) = 0$  even though  $\hat{p}_t(x|s, a) \neq 0$  for some  $x$ ).
- (ii) The maximizer of a linear optimization over  $L_1$  ball could change **significantly** for a small change in the value function.

# $L_1$ -Norm vs. KL

Linear optimization over  $L_1$ -ball (left) vs. KL-ball (right): The vector represents a value function (e.g., in EVI).



(Filippi et al., 2010)

These shortcomings are avoided by resorting to KL-based confidence bounds (as in KL-UCRL):

$$\text{KL}(\hat{p}_t(\cdot|s, a), p'(\cdot|s, a)) \leq \frac{\square S \log(\log(T)/\delta)}{N_t(s, a)}$$
$$|\hat{\mu}_t(s, a) - \mu'(s, a)| \leq \sqrt{\frac{\square \log(\log(T)/\delta)}{N_t(s, a)}}$$

- Numerically, KL-UCRL **outperforms** UCRL2 (uniformly in all environment).
- Yet the best known regret bound for KL-UCRL:  $\tilde{O}(DS\sqrt{AT})$

**Our contribution** is to investigate the benefit of using KL theoretically.

These shortcomings are avoided by resorting to KL-based confidence bounds (as in KL-UCRL):

$$\text{KL}(\hat{p}_t(\cdot|s, a), p'(\cdot|s, a)) \leq \frac{\square S \log(\log(T)/\delta)}{N_t(s, a)}$$
$$|\hat{\mu}_t(s, a) - \mu'(s, a)| \leq \sqrt{\frac{\square \log(\log(T)/\delta)}{N_t(s, a)}}$$

- Numerically, KL-UCRL **outperforms** UCRL2 (uniformly in all environment).
- Yet the best known regret bound for KL-UCRL:  $\tilde{O}(DS\sqrt{AT})$

**Our contribution** is to investigate the benefit of using KL theoretically.

# Variance-Aware Regret Bounds for KL-UCRL

Variance of bias function w.r.t. transition law  $p(\cdot|s, a)$ :

$$\mathbb{V}_{p(\cdot|s,a)}(b^*) := \sum_{x \in \mathcal{S}} p(x|s, a) \left( b^*(x) - \mathbb{E}_{p(\cdot|s,a)}[b^*] \right)^2$$

with  $\mathbb{E}_{p(\cdot|s,a)}[b^*] = \sum_x p(\cdot|s, a) b^*(x)$ .

## Theorem

*The regret under KL-UCRL in any ergodic MDP satisfies*

$$\mathfrak{R}_T \leq \left( 31 \sqrt{S \sum_{s,a} \mathbb{V}_{p(\cdot|s,a)}(b^*)} + 35S\sqrt{A} + 2D \right) \sqrt{T \log(\log(T)/\delta)} + \tilde{\mathcal{O}}(\text{polylog}(T))$$

*and with probability at least  $1 - \delta$ .*

- Improves over the previous bound of  $\tilde{\mathcal{O}}(DS\sqrt{AT})$  for KL-UCRL (since  $\mathbb{V}_{p(\cdot|s,a)}(b^*) \leq D^2$ ).
- Proof: Uses novel concentration inequalities

# Variance-Aware Regret Bounds for KL-UCRL

Variance of bias function w.r.t. transition law  $p(\cdot|s, a)$ :

$$\mathbb{V}_{p(\cdot|s,a)}(b^*) := \sum_{x \in \mathcal{S}} p(x|s, a) \left( b^*(x) - \mathbb{E}_{p(\cdot|s,a)}[b^*] \right)^2$$

with  $\mathbb{E}_{p(\cdot|s,a)}[b^*] = \sum_x p(\cdot|s, a) b^*(x)$ .

## Theorem

*The regret under KL-UCRL in any ergodic MDP satisfies*

$$\mathfrak{R}_T \leq \left( 31 \sqrt{S \sum_{s,a} \mathbb{V}_{p(\cdot|s,a)}(b^*)} + 35S\sqrt{A} + 2D \right) \sqrt{T \log(\log(T)/\delta)} + \tilde{\mathcal{O}}(\text{polylog}(T))$$

*and with probability at least  $1 - \delta$ .*

- Improves over the previous bound of  $\tilde{\mathcal{O}}(DS\sqrt{AT})$  for KL-UCRL (since  $\mathbb{V}_{p(\cdot|s,a)}(b^*) \leq D^2$ ).
- Proof: Uses novel concentration inequalities

# Variance-Aware Regret Bounds for KL-UCRL

Variance of bias function w.r.t. transition law  $p(\cdot|s, a)$ :

$$\mathbb{V}_{p(\cdot|s,a)}(b^*) := \sum_{x \in \mathcal{S}} p(x|s, a) \left( b^*(x) - \mathbb{E}_{p(\cdot|s,a)}[b^*] \right)^2$$

with  $\mathbb{E}_{p(\cdot|s,a)}[b^*] = \sum_x p(\cdot|s, a) b^*(x)$ .

## Theorem

*The regret under KL-UCRL in any ergodic MDP satisfies*

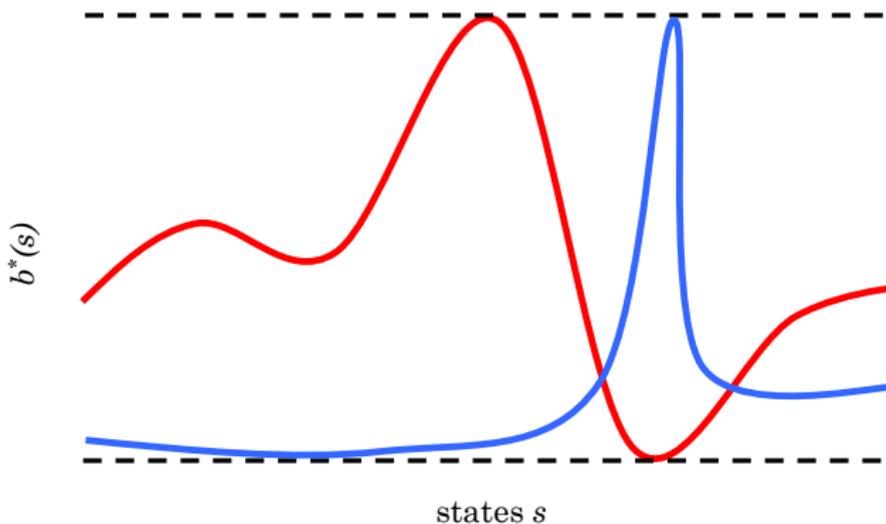
$$\mathfrak{R}_T \leq \left( 31 \sqrt{S \sum_{s,a} \mathbb{V}_{p(\cdot|s,a)}(b^*)} + 35S\sqrt{A} + 2D \right) \sqrt{T \log(\log(T)/\delta)} + \tilde{\mathcal{O}}(\text{polylog}(T))$$

*and with probability at least  $1 - \delta$ .*

- Improves over the previous bound of  $\tilde{\mathcal{O}}(DS\sqrt{AT})$  for KL-UCRL (since  $\mathbb{V}_{p(\cdot|s,a)}(b^*) \leq D^2$ ).
- Proof: Uses novel concentration inequalities

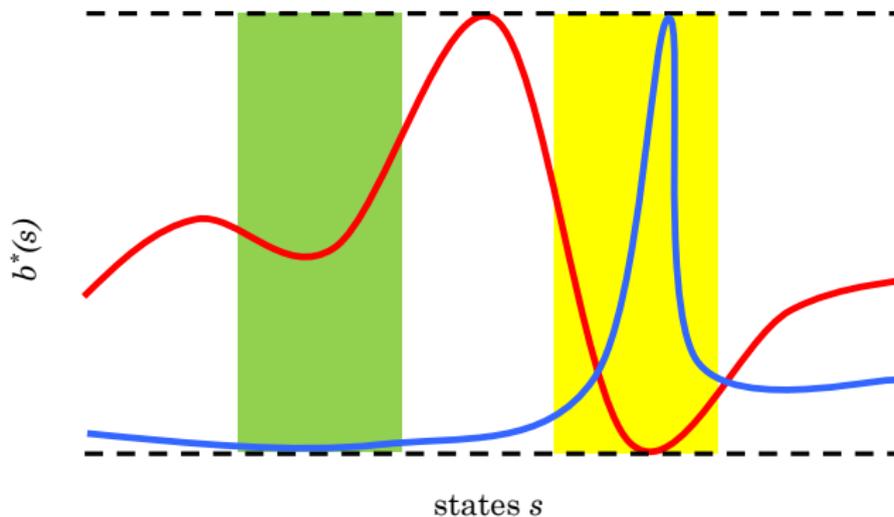
# Variance vs. Diameter

In contrast to diameter  $D$  (global measures), variance  $\mathbb{V}_{p(\cdot|s,a)}(b^*)$  is a **local measure**, which is aware of variations of  $b^*$  over state-space.



# Variance vs. Diameter

In contrast to diameter  $D$  (global measures), variance  $\mathbb{V}_{p(\cdot|s,a)}(b^*)$  is a **local measure**, which is aware of variations of  $b^*$  over state-space.



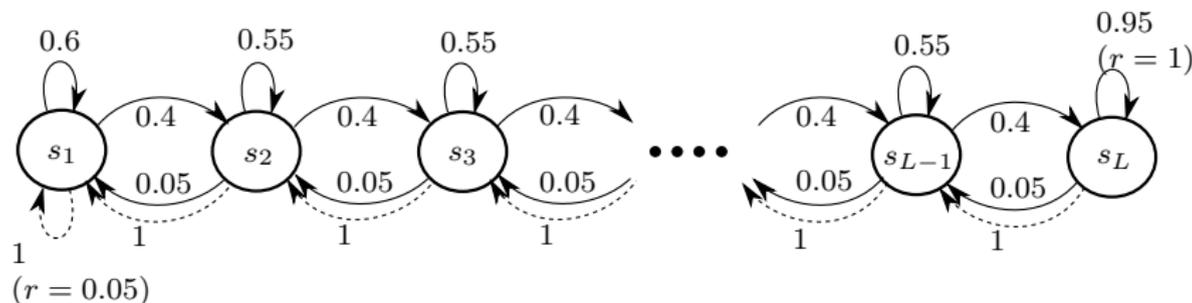
# Outline

- 1 UCRL2
- 2 UCRL3
- 3 KL-UCRL
- 4 Numerical Experiments**
- 5 Technical Tools

# Numerical Experiments

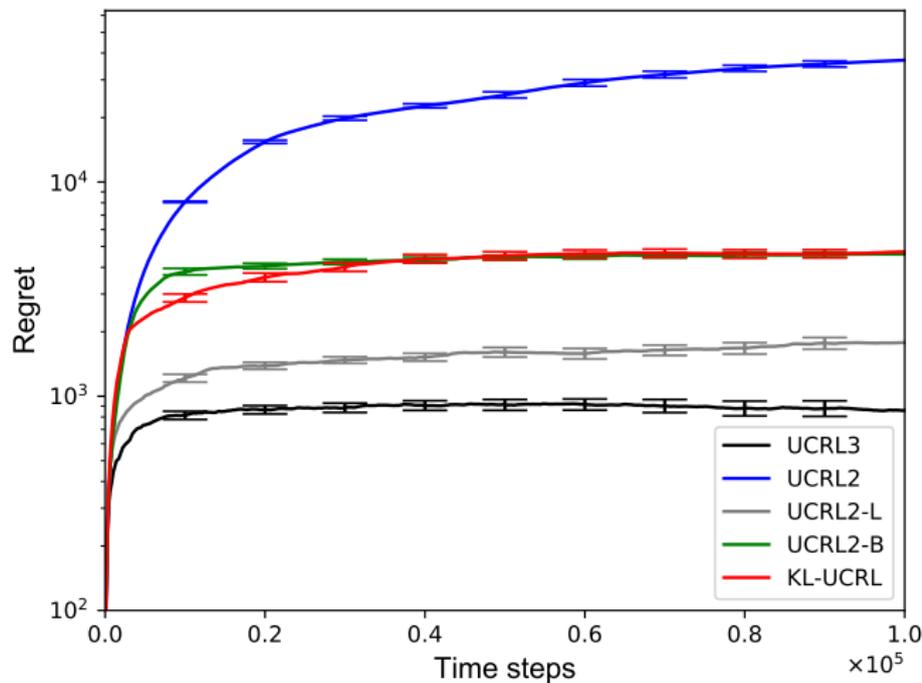
We examine UCRL2, KL-UCRL, UCRL3, **UCRL-L**, and **UCRL-B** on the *RiverSwim* environment (shown below).

- **UCRL-L**: Uses  $L_1$  confidence bounds (as UCRL2) combined with the Laplace method.
- **UCRL-B**: Uses element-wise empirical Bernstein confidence bounds combined with peeling.



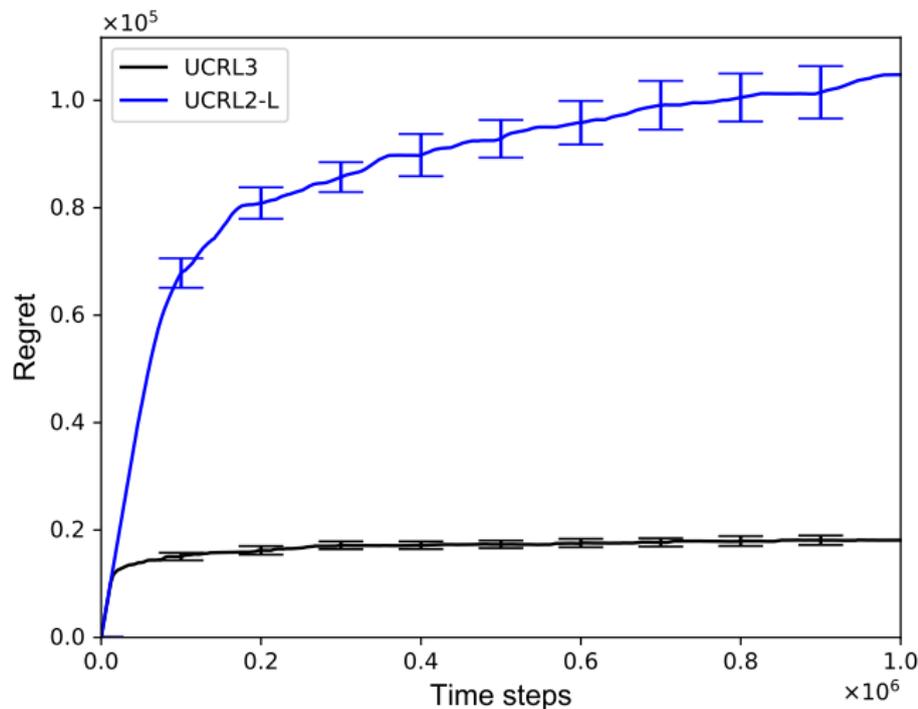
# Numerical Experiments

Regret of various algorithms in 6-state RiverSwim:



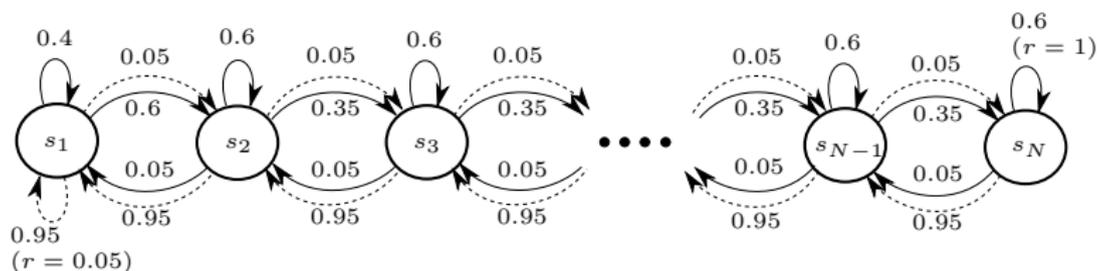
# Numerical Experiments

Comparison between UCRL2-L and UCRL3 in 25-state RiverSwim:



# Numerical Experiments

Examining the main terms in the regret bounds of KL-UCRL in  $N$ -state ergodic *RiverSwim* MDP:



$S$	$D$	$\max_{s,a} \mathbb{V}_{p(\cdot s,a)}(b^*)$	$D\sqrt{SA}$	$\sqrt{\sum_{s,a} \mathbb{V}_{p(\cdot s,a)}(b^*)}$
6	6.3	0.6322	21.9	1.8
12	14.9	0.6327	72.9	2.8
20	26.3	0.6327	166.4	3.7
40	54.9	0.6327	490.9	5.3
100	140.6	0.6327	1988.3	8.5

# Outline

- 1 UCRL2
- 2 UCRL3
- 3 KL-UCRL
- 4 Numerical Experiments
- 5 Technical Tools**

# Old Proof(s)

The regret is decomposed into **per-episode** regret terms.

Consider episode  $k$  with **optimistic model**  $\tilde{M}_k$  (with kernel  $\tilde{p}_k$  and bias function  $\tilde{b}_k$ ), and assume  $M \in \mathcal{M}_t$ .

The leading term in regret bound for episode  $k$  is due to:

$$\sum_x (\tilde{p}_k(x|s, a) - p(x|s, a)) \tilde{b}_k(x) \leq \underbrace{\|\tilde{p}_k(\cdot|s, a) - p(\cdot|s, a)\|_1}_{\mathcal{O}\left(\sqrt{\frac{S \log(T)}{N_k(s, a)}}\right)} \underbrace{\|\tilde{b}_k\|_\infty}_{\leq D}$$

Summing over episodes  $k$  and state-action pairs  $(s, a)$ , this leads to  $\tilde{\mathcal{O}}(DS\sqrt{AT})$ .

Using Cauchy-Schwarz in the above leads to a **too conservative** bound!

# Old Proof(s)

The regret is decomposed into **per-episode** regret terms.

Consider episode  $k$  with **optimistic model**  $\tilde{M}_k$  (with kernel  $\tilde{p}_k$  and bias function  $\tilde{b}_k$ ), and assume  $M \in \mathcal{M}_t$ .

The leading term in regret bound for episode  $k$  is due to:

$$\sum_x (\tilde{p}_k(x|s, a) - p(x|s, a)) \tilde{b}_k(x) \leq \underbrace{\|\tilde{p}_k(\cdot|s, a) - p(\cdot|s, a)\|_1}_{\mathcal{O}\left(\sqrt{\frac{S \log(T)}{N_k(s, a)}}\right)} \underbrace{\|\tilde{b}_k\|_\infty}_{\leq D}$$

Summing over episodes  $k$  and state-action pairs  $(s, a)$ , this leads to  $\tilde{\mathcal{O}}(DS\sqrt{AT})$ .

Using Cauchy-Schwarz in the above leads to a **too conservative** bound!

Decomposition:

$$\begin{aligned} \sum_x (\tilde{p}_k(x|s, a) - p(x|s, a)) \tilde{b}_k(x) &= \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[\tilde{b}_k] - \mathbb{E}_{p(\cdot|s, a)}[\tilde{b}_k]}_{\text{transportation cost of } \tilde{b}_k} \\ &= \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[b^*] - \mathbb{E}_{p(\cdot|s, a)}[b^*]}_{\text{transportation cost of } b^*} + \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[\tilde{b}_k - b^*] - \mathbb{E}_{p(\cdot|s, a)}[\tilde{b}_k - b^*]}_{\text{correction term}} \end{aligned}$$

⇒ **Transportation cost of  $b^*$** : using (novel) transportation inequalities

⇒ **Correction term**: using ergodic property of MDP + contraction of induced transition matrices. The total contribution of correction terms (over all  $(s, a)$  and  $k$ ):

$$\tilde{O}(S\sqrt{AT})$$

Decomposition:

$$\begin{aligned} \sum_x (\tilde{p}_k(x|s, a) - p(x|s, a)) \tilde{b}_k(x) &= \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[\tilde{b}_k] - \mathbb{E}_{p(\cdot|s, a)}[\tilde{b}_k]}_{\text{transportation cost of } \tilde{b}_k} \\ &= \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[b^*] - \mathbb{E}_{p(\cdot|s, a)}[b^*]}_{\text{transportation cost of } b^*} + \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[\tilde{b}_k - b^*] - \mathbb{E}_{p(\cdot|s, a)}[\tilde{b}_k - b^*]}_{\text{correction term}} \end{aligned}$$

⇒ **Transportation cost of  $b^*$** : using (novel) transportation inequalities

⇒ **Correction term**: using ergodic property of MDP + contraction of induced transition matrices. The total contribution of correction terms (over all  $(s, a)$  and  $k$ ):

$$\tilde{O}(S\sqrt{AT})$$

Decomposition:

$$\begin{aligned} \sum_x (\tilde{p}_k(x|s, a) - p(x|s, a)) \tilde{b}_k(x) &= \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[\tilde{b}_k] - \mathbb{E}_{p(\cdot|s, a)}[\tilde{b}_k]}_{\text{transportation cost of } \tilde{b}_k} \\ &= \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[b^*] - \mathbb{E}_{p(\cdot|s, a)}[b^*]}_{\text{transportation cost of } b^*} + \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[\tilde{b}_k - b^*] - \mathbb{E}_{p(\cdot|s, a)}[\tilde{b}_k - b^*]}_{\text{correction term}} \end{aligned}$$

⇒ **Transportation cost of  $b^*$** : using (novel) transportation inequalities

⇒ **Correction term**: using ergodic property of MDP + contraction of induced transition matrices. The total contribution of correction terms (over all  $(s, a)$  and  $k$ ):

$$\tilde{O}(S\sqrt{AT})$$

Decomposition:

$$\begin{aligned} \sum_x (\tilde{p}_k(x|s, a) - p(x|s, a)) \tilde{b}_k(x) &= \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[\tilde{b}_k] - \mathbb{E}_{p(\cdot|s, a)}[\tilde{b}_k]}_{\text{transportation cost of } \tilde{b}_k} \\ &= \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[b^*] - \mathbb{E}_{p(\cdot|s, a)}[b^*]}_{\text{transportation cost of } b^*} + \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s, a)}[\tilde{b}_k - b^*] - \mathbb{E}_{p(\cdot|s, a)}[\tilde{b}_k - b^*]}_{\text{correction term}} \end{aligned}$$

- ⇒ **Transportation cost of  $b^*$** : using (novel) transportation inequalities
- ⇒ **Correction term**: using ergodic property of MDP + contraction of induced transition matrices. The total contribution of correction terms (over all  $(s, a)$  and  $k$ ):

$$\tilde{O}(S\sqrt{AT})$$

# Transportation Inequalities

## Lemma (Transportation Lemma)

For any function  $f$ , introduce  $\varphi_f : \lambda \mapsto \log \mathbb{E}_P[\exp(\lambda(f(X) - \mathbb{E}_P[f]))]$ .

Then for all  $Q \ll P$ ,

$$\mathbb{E}_Q[f] - \mathbb{E}_P[f] \leq \inf\{x \geq 0 : \varphi_{*,f}(x) > KL(Q, P)\}$$

$$\mathbb{E}_Q[f] - \mathbb{E}_P[f] \geq \sup\{x \leq 0 : \varphi_{*,f}(x) > KL(Q, P)\}$$

where  $\varphi_{*,f}(x) = \sup_{\lambda} \lambda x - \varphi_f(\lambda)$ .

## Lemma (Transportation Inequality I)

For any function  $f$  and distribution  $P$ , such that  $\mathbb{V}_P(f)$  and  $\mathbb{S}(f)$  are finite

$$\mathbb{E}_Q[f] - \mathbb{E}_P[f] \leq \sqrt{2\mathbb{V}_P(f)KL(Q, P)} + \frac{2}{3}\mathbb{S}(f)KL(Q, P)$$

$$\mathbb{E}_P[f] - \mathbb{E}_Q[f] \leq \sqrt{2\mathbb{V}_P(f)KL(Q, P)}$$

# Transportation Inequalities

## Lemma (Transportation Lemma)

For any function  $f$ , introduce  $\varphi_f : \lambda \mapsto \log \mathbb{E}_P[\exp(\lambda(f(X) - \mathbb{E}_P[f]))]$ .  
Then for all  $Q \ll P$ ,

$$\mathbb{E}_Q[f] - \mathbb{E}_P[f] \leq \inf\{x \geq 0 : \varphi_{*,f}(x) > KL(Q, P)\}$$

$$\mathbb{E}_Q[f] - \mathbb{E}_P[f] \geq \sup\{x \leq 0 : \varphi_{*,f}(x) > KL(Q, P)\}$$

where  $\varphi_{*,f}(x) = \sup_{\lambda} \lambda x - \varphi_f(\lambda)$ .

## Lemma (Transportation Inequality I)

For any function  $f$  and distribution  $P$ , such that  $\mathbb{V}_P(f)$  and  $\mathbb{S}(f)$  are finite

$$\mathbb{E}_Q[f] - \mathbb{E}_P[f] \leq \sqrt{2\mathbb{V}_P(f)KL(Q, P)} + \frac{2}{3}\mathbb{S}(f)KL(Q, P)$$

$$\mathbb{E}_P[f] - \mathbb{E}_Q[f] \leq \sqrt{2\mathbb{V}_P(f)KL(Q, P)}$$

# Transportation Inequalities

A novel refinement of previous transportation inequality:

## Lemma (Transportation Inequality II)

For any function  $f$  and distributions  $P, Q$  defined on a finite alphabet  $\mathcal{X}$ ,

$$\mathbb{E}_Q[f] - \mathbb{E}_P[f] \leq \left( \sqrt{\mathcal{V}_{P,Q}(f)} + \sqrt{\mathcal{V}_{Q,P}(f)} \right) \sqrt{2\text{KL}(P, Q)} + \mathfrak{S}(f)\text{KL}(P, Q)$$

where  $\mathcal{V}_{P,Q}(f) := \sum_{x \in \mathcal{X}: P(x) \geq Q(x)} P(x)(f(x) - \mathbb{E}_P[f])^2$ .

The operator  $\mathcal{V}_{P,Q}(f)$  is closely related to the local variance of  $f$  (under  $P$  and  $Q$ ):

$$\begin{aligned} \mathcal{V}_{P,Q}(f) &\leq \mathbb{V}_P(f) \\ \sqrt{\mathcal{V}_{P,Q}(f)} &\leq \sqrt{2\mathbb{V}_Q(f) + 3\mathfrak{S}(f)\sqrt{|\mathcal{X}|\text{KL}(Q, P)}} \end{aligned}$$

Proof: Cauchy-Schwarz + local Pinsker's inequalities

# Transportation Inequalities

A novel refinement of previous transportation inequality:

## Lemma (Transportation Inequality II)

For any function  $f$  and distributions  $P, Q$  defined on a finite alphabet  $\mathcal{X}$ ,

$$\mathbb{E}_Q[f] - \mathbb{E}_P[f] \leq \left( \sqrt{\mathcal{V}_{P,Q}(f)} + \sqrt{\mathcal{V}_{Q,P}(f)} \right) \sqrt{2\text{KL}(P, Q)} + \mathbb{S}(f)\text{KL}(P, Q)$$

where  $\mathcal{V}_{P,Q}(f) := \sum_{x \in \mathcal{X}: P(x) \geq Q(x)} P(x)(f(x) - \mathbb{E}_P[f])^2$ .

The operator  $\mathcal{V}_{P,Q}(f)$  is closely related to the local variance of  $f$  (under  $P$  and  $Q$ ):

$$\begin{aligned} \mathcal{V}_{P,Q}(f) &\leq \mathbb{V}_P(f) \\ \sqrt{\mathcal{V}_{P,Q}(f)} &\leq \sqrt{2\mathbb{V}_Q(f) + 3\mathbb{S}(f)\sqrt{|\mathcal{X}|\text{KL}(Q, P)}} \end{aligned}$$

Proof: Cauchy-Schwarz + local Pinsker's inequalities

# Transportation Inequalities

A novel refinement of previous transportation inequality:

## Lemma (Transportation Inequality II)

For any function  $f$  and distributions  $P, Q$  defined on a finite alphabet  $\mathcal{X}$ ,

$$\mathbb{E}_Q[f] - \mathbb{E}_P[f] \leq \left( \sqrt{\mathcal{V}_{P,Q}(f)} + \sqrt{\mathcal{V}_{Q,P}(f)} \right) \sqrt{2\text{KL}(P, Q)} + \mathbb{S}(f)\text{KL}(P, Q)$$

where  $\mathcal{V}_{P,Q}(f) := \sum_{x \in \mathcal{X}: P(x) \geq Q(x)} P(x)(f(x) - \mathbb{E}_P[f])^2$ .

The operator  $\mathcal{V}_{P,Q}(f)$  is closely related to the local variance of  $f$  (under  $P$  and  $Q$ ):

$$\begin{aligned} \mathcal{V}_{P,Q}(f) &\leq \mathbb{V}_P(f) \\ \sqrt{\mathcal{V}_{P,Q}(f)} &\leq \sqrt{2\mathbb{V}_Q(f) + 3\mathbb{S}(f)\sqrt{|\mathcal{X}|\text{KL}(Q, P)}} \end{aligned}$$

Proof: Cauchy-Schwarz + local Pinsker's inequalities

$$\begin{aligned} \text{transportation cost of } b^* &= \underbrace{\mathbb{E}_{\tilde{p}_k(\cdot|s,a)}[b^*] - \mathbb{E}_{\hat{p}_k(\cdot|s,a)}[b^*]}_{T_1} \\ &\quad + \underbrace{\mathbb{E}_{\hat{p}_k(\cdot|s,a)}[b^*] - \mathbb{E}_{p(\cdot|s,a)}[b^*]}_{T_2} \end{aligned}$$

$\Rightarrow$  **Term  $T_1$ : Transportation Inequality II** with  $P = \tilde{p}_k(\cdot|s, a)$  and  $Q = \hat{p}_k(\cdot|s, a)$

$\Rightarrow$  **Term  $T_2$ : Transportation Inequality I** with  $Q = \hat{p}_k(\cdot|s, a)$  and  $P = p(\cdot|s, a)$

Combining, and summing over  $(s, a)$  and episodes  $k$ , the contribution of  $T_2$  terms become

$$\tilde{\mathcal{O}}\left(\sqrt{S \sum_{s,a} \mathbb{V}_{p(\cdot|s,a)}(b^*) T}\right)$$

Two variants of UCRL2: UCRL3 and KL-UCRL

## **UCRL3:**

- A novel variant of UCRL2 using (i) improved confidence sets, and (ii) novel efficient approach for computing an optimistic policy.
- Beats all existing variants of UCRL2 in practice yet enjoying the same regret guarantees.

## **KL-UCRL:**

- A variant of UCRL2, which uses KL-divergence to define confidence sets.
- We provided improved regret analysis for it in ergodic MDPs, thanks to novel variants of transportation concentration inequalities.

- Optimal stopping criterion for UCRL2-style algorithms
- Problem-dependent regret **lower** and **upper** bounds for average-reward RL

Thank you!