

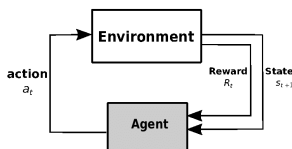
On the complexity of learning good policies with and without rewards

Emilie Kaufmann, Pierre Ménard, Omar Darwiche Domingues,
Anders Jonsson, Edouard Leurent and Michal Valko



AI Seminar, University of Alberta
Novembre 20th, 2020

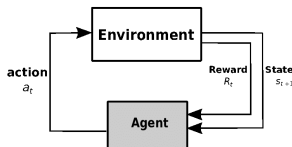
RL setup: an agent interacts with an environment (MDP)



Several Performance measures:

- 1 the agent should *adopt* a good behavior
 - maximize the total rewards (*regret minimization*)
 - use as much as possible an ϵ -optimal policy (*PAC-MDP*)
- 2 the agent should *learn* a good behavior
 - learn an optimal policy for a **given reward function**
 - learn the dynamics so that to be robust to find the optimal policy for **any reward function**

RL setup: an agent interacts with an environment (MDP)



Several Performance measures:

- 1 the agent should *adopt* a good behavior
 - maximize the total rewards (*regret minimization*)
 - use as much as possible an ϵ -optimal policy (*PAC-MDP*)
- 2 the agent should *learn* a good behavior
 - learn an optimal policy for a **given reward function**
 - learn the dynamics so that to be robust to find the optimal policy for **any reward function**

two Pure Exploration problems

Episodic MDP: horizon H and MDP $(\mathcal{S}, \mathcal{A}, P, r)$ for

- a state space \mathcal{S} of size $S < \infty$
- an action space \mathcal{A} of size $A < \infty$
- a transition kernel $P = (p_h(s'|s, a))_{\substack{(s,a,s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \\ h \in [H]}}$
- a reward function $r = (r_h(s, a))_{\substack{(s,a) \in \mathcal{S} \times \mathcal{A} \\ h \in [H]}}$

Value of a policy $\pi = (\pi_h)_{h=1}^H$, $\pi_h : \mathcal{S} \rightarrow \mathcal{A}$:

$$V_h^\pi(s; r) \triangleq \mathbb{E}^\pi \left[\sum_{\ell=h}^H r_\ell(s_\ell, \pi_\ell(s_\ell)) \middle| s_{\ell+1} \sim p_\ell(\cdot | s_\ell, \pi_\ell(s_\ell)), s_h = s \right]$$

Optimal policy: π_r^* such that $V_h^{\pi_r^*}(s; r) \geq V_h^\pi(s; r)$ for all π, s, h .

Episodic MDP: horizon H and MDP $(\mathcal{S}, \mathcal{A}, P, r)$ for

- a state space \mathcal{S} of size $S < \infty$
- an action space \mathcal{A} of size $A < \infty$
- a transition kernel $P = (p_h(s'|s, a))_{\substack{(s,a,s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \\ h \in [H]}}$
- a reward function $r = (r_h(s, a))_{\substack{(s,a) \in \mathcal{S} \times \mathcal{A} \\ h \in [H]}}$

Value of a policy $\pi = (\pi_h)_{h=1}^H$, $\pi_h : \mathcal{S} \rightarrow \mathcal{A}$:

$$V_h^\pi(s; r) \triangleq \mathbb{E}^\pi \left[\sum_{\ell=h}^H r_\ell(s_\ell, \pi_\ell(s_\ell)) \middle| s_{\ell+1} \sim p_\ell(\cdot | s_\ell, \pi_\ell(s_\ell)), s_h = s \right]$$

Optimal policy: π_r^* such that $V_h^{\pi_r^*}(s; r) \geq V_h^\pi(s; r)$ for all π, s, h .

Setting: Episodic Markov Decision Process

Episodic MDP: horizon H and MDP $(\mathcal{S}, \mathcal{A}, P, r)$ for

- a state space \mathcal{S} of size $S < \infty$
- an action space \mathcal{A} of size $A < \infty$
- a transition kernel $P = (p_h(s'|s, a))_{\substack{(s,a,s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \\ h \in [H]}}$
- a reward function $r = (r_h(s, a))_{\substack{(s,a) \in \mathcal{S} \times \mathcal{A} \\ h \in [H]}}$

Q-value of a policy $\pi = (\pi_h)_{h=1}^H$, $\pi_h : \mathcal{S} \rightarrow \mathcal{A}$:

$$Q_h^\pi(s, a; r) \triangleq \mathbb{E}^\pi \left[r_h(s, a) + \sum_{\ell=h+1}^H r_\ell(s_\ell, \pi_\ell(s_\ell)) \mid \substack{s_h=s, a_h=a \\ s_{\ell+1} \sim p_\ell(\cdot | s_\ell, \pi_\ell(s_\ell))} \right]$$

Optimal policy: π_r^* such that $V_h^{\pi_r^*}(s; r) \geq V_h^\pi(s; r)$ for all π, s, h .

Setting: Episodic Markov Decision Process

Episodic MDP: horizon H and MDP $(\mathcal{S}, \mathcal{A}, P, r)$ for

- a state space \mathcal{S} of size $S < \infty$
- an action space \mathcal{A} of size $A < \infty$
- a transition kernel $P = (p_h(s'|s, a))_{\substack{(s,a,s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \\ h \in [H]}}$
- a reward function $r = (r_h(s, a))_{\substack{(s,a) \in \mathcal{S} \times \mathcal{A} \\ h \in [H]}}$ (**step-dependent**)

Q-value of a **policy** $\pi = (\pi_h)_{h=1}^H$, $\pi_h : \mathcal{S} \rightarrow \mathcal{A}$:

$$Q_h^\pi(s, a; r) \triangleq \mathbb{E}^\pi \left[r_h(s, a) + \sum_{\ell=h+1}^H r_\ell(s_\ell, \pi_\ell(s_\ell)) \mid \substack{s_h=s, a_h=a \\ s_{\ell+1} \sim p_\ell(\cdot | s_\ell, \pi_\ell(s_\ell))} \right]$$

Optimal policy: π_r^* such that $V_h^{\pi_r^*}(s; r) \geq V_h^\pi(s; r)$ for all π, s, h .

- 1 The BPI and RFE objectives
- 2 Reward-Free UCRL
- 3 BPI Algorithms

Online episodic algorithm

Collect data from the MDP by generating **trajectories** (episodes)
 \neq generative model

In each episode $t = 1, 2, \dots$, the agent

- selects an **exploration policy** π^t
- generates an episode under this policy

$$(s_1^t, a_1^t, s_2^t, a_2^t, \dots, s_H^t, a_H^t)$$

where $s_1^t \sim \rho$, $a_h^t = \pi_h^t(s_h^t)$ and $s_{h+1}^t \sim p_h(\cdot | s_h^t, a_h^t)$

- can decide to **stop exploration**
- if decides to stop, **outputs a prediction**

→ three **data-dependent** components

Online episodic algorithm

Collect data from the MDP by generating **trajectories** (episodes)
 \neq generative model

In each episode $t = 1, 2, \dots$, the agent

- selects an **exploration policy** π^t
- generates an episode under this policy

$$(s_1^t, a_1^t, s_2^t, a_2^t, \dots, s_H^t, a_H^t)$$

where $\mathbf{s}_1^t = \mathbf{s}_1$, $a_h^t = \pi_h^t(s_h^t)$ and $s_{h+1}^t \sim p_h(\cdot | s_h^t, a_h^t)$

- can decide to **stop exploration**
- if decides to stop, **outputs a prediction**

→ three **data-dependent** components

Best Policy Identification (BPI)

→ Learn the optimal policy for a **known reward function** r

[Fiechter, 1994]

BPI algorithm

- **exploration policy** π^t : may depend on past data \mathcal{D}_{t-1} and r

$$\mathcal{D}_t = \mathcal{D}_{t-1} \cup \{(s_1^t, a_1^t, s_2^t, a_2^t, \dots, s_H^t, a_H^t)\}$$

- **stopping rule** τ : stopping time w.r.t. $(\mathcal{D}_t)_{t \in \mathbb{N}}$
(can depend on r)
- **prediction** $\hat{\pi}$: a **policy** that may depend on \mathcal{D}_τ and r

(ε, δ) -PAC algorithm for Best Policy Identification

$$\mathbb{P} \left(V_1^*(\mathbf{s}_1; r) - V_1^{\hat{\pi}}(\mathbf{s}_1; r) \leq \varepsilon \right) \geq 1 - \delta$$

Wanted: (ε, δ) -PAC algorithm with a small sample complexity τ

Reward-Free Exploration (RFE)

→ Learn the optimal policy for **any** reward function r

[Jin et al., 2020]

RFE algorithm

- exploration policy π^t : may depend on past data \mathcal{D}_{t-1}

$$\mathcal{D}_t = \mathcal{D}_{t-1} \cup \{(s_1^t, a_1^t, s_2^t, a_2^t, \dots, s_H^t, a_H^t)\}$$

- stopping rule τ : stopping time w.r.t. $(\mathcal{D}_t)_{t \in \mathbb{N}}$
- prediction $\hat{P} = (\hat{p}_h(\cdot|s, a))_{h,s,a}$: a **transition kernel** that may depend on \mathcal{D}_τ

$\hat{\pi}_r^*$: optimal policy in the MDP (\hat{P}, r)

(ε, δ) -PAC algorithm for Reward-Free Exploration

$$\mathbb{P} \left(\text{for all reward function } r, V_1^*(\mathbf{s}_1; r) - V_1^{\hat{\pi}_r^*}(\mathbf{s}_1; r) \leq \varepsilon \right) \geq 1 - \delta$$

Wanted: (ε, δ) -PAC algorithm with a small sample complexity τ

- 1 The BPI and RFE objectives
- 2 Reward-Free UCRL
- 3 BPI Algorithms

A model-based algorithm

Based on the available data \mathcal{D}_t , builds estimates of the transition probabilities $p_h(s, a)$

→ estimates of the Q-values $Q_h^\pi(s, a; r)$

Number of visits:

$$n_h^t(s, a) = \sum_{k=1}^t \mathbb{1}_{\{(s_h^k, a_h^k) = (s, a)\}} \quad n_h^t(s, a, s') = \sum_{k=1}^t \mathbb{1}_{\{(s_h^k, a_h^k, s_{h+1}^k) = (s, a, s')\}}$$

Empirical transitions: $\hat{P}^t = (\hat{p}_h^t(s'|s, a))_{h,s,a,s'}$

$$\hat{p}_h^t(s'|s, a) = \begin{cases} \frac{n_h^t(s, a, s')}{n_h^t(s, a)} & \text{if } n_h^t(s, a) > 0 \\ \frac{1}{S} & \text{else} \end{cases}$$

Empirical values:

- $\hat{V}_h^{t,\pi}(s; r)$ values in the empirical MDP $(\mathcal{S}, \mathcal{A}, \hat{P}^t, r)$
- $\hat{Q}_h^{t,\pi}(s; r)$ Q-values in the empirical MDP $(\mathcal{S}, \mathcal{A}, \hat{P}^t, r)$

Central observation

A sufficient condition to be (ϵ, δ) -PAC is to have **accurate estimates of the value function for all π and r** :

$$\mathbb{P} \left(\forall \pi, \forall r, |\hat{V}_1^{t,\pi}(s_1; r) - V_1^\pi(s_1; r)| \leq \epsilon/2 \right) \geq 1 - \delta.$$

RF-UCRL:

- builds upper bounds on the errors

$$\hat{e}_h^{t,\pi}(s, a; r) := |\hat{Q}_h^{t,\pi}(s, a; r) - Q_h^\pi(s, a; r)|$$

... that are **independent of π and r** !

- greedily reduces the upper bounds

$$\hat{e}_h^{t,\pi}(s, a; r) := |\hat{Q}_h^{t,\pi}(s, a; r) - Q_h^\pi(s, a; r)|$$

We define inductively $\bar{E}_{H+1}^t(s, a) = 0$ and

$$\bar{E}_h^t(s, a) = \min \left[(H-h); (H-h) \sqrt{\frac{2\beta(n_h^t(s, a), \delta)}{n_h^t(s, a)}} + \sum_{s'} \hat{p}_h^t(s'|s, a) \max_b \bar{E}_{h+1}^t(s', b) \right]$$

for some **threshold function** $\beta(n, \delta)$.

→ like in UCRL [Jaksch et al. 10], this construction relies on **confidence regions on the transitions probabilities**

Upper Bound Property

On the event

$$\mathcal{E} = \left\{ \forall t \in \mathbb{N}, \forall h \in [H], \forall (s, a), \text{KL}(\hat{p}_h^t(\cdot|s, a), p_h(\cdot|s, a)) \leq \frac{\beta(n_h^t(s, a), \delta)}{n_h^t(s, a)} \right\},$$

for all π and r , for all h, s, a , $\hat{e}_h^{t,\pi}(s, a; r) \leq \bar{E}_h^t(s, a)$.

$$\hat{e}_h^{t,\pi}(s, a; r) := |\hat{Q}_h^{t,\pi}(s, a; r) - Q_h^\pi(s, a; r)|$$

A simple consequence of **Bellman equations**:

$$\hat{Q}_h^{t,\pi}(s, a; r) = r_h(s, a) + \sum_{s'} \hat{p}_h^t(s'|s, a) \hat{Q}_{h+1}^{t,\pi}(s', \pi(s'); r)$$

$$\text{and } Q_h^\pi(s, a; r) = r_h(s, a) + \sum_{s'} p_h(s'|s, a) Q_{h+1}^\pi(s', \pi(s'); r).$$

Error decomposition:

$$\begin{aligned} \hat{e}_h^{t,\pi}(s, a; r) &\leq \sum_{s'} |\hat{p}_h^t(s'|s, a) - p_h(s'|s, a)| Q_{h+1}^\pi(s', \pi(s'); r) \\ &\quad + \sum_{s'} \hat{p}_h^t(s'|s, a) \left| \hat{Q}_{h+1}^{t,\pi}(s', \pi(s'); r) - Q_{h+1}^\pi(s', \pi(s'); r) \right| \\ &\leq (H-h) \underbrace{\|\hat{p}_h^t(\cdot|s, a) - p_h(\cdot|s, a)\|_1}_{\leq \sqrt{\frac{2\beta(n_h^t(s,a), \delta)}{n_h^t(s,a)}} \text{ (Pinsker} + \mathcal{E})} + \sum_{s'} \hat{p}_h^t(s'|s, a) \underbrace{\hat{e}_{h+1}^{t,\pi}(s', \pi(s'); r)}_{\leq \max_b \bar{E}_{h+1}^t(s', b) \text{ (induction)}}. \end{aligned}$$

$$\bar{E}_h^t(s, a) = \min \left[(H-h); (H-h) \sqrt{\frac{2\beta(n_h^t(s, a), \delta)}{n_h^t(s, a)}} + \sum_{s'} \hat{p}_h^t(s'|s, a) \max_b \bar{E}_{h+1}^t(s', b) \right]$$

Reward-Free UCRL

- **exploration policy:** π^{t+1} is the greedy policy wrt $\bar{E}^t(s, a)$:

$$\forall s \in \mathcal{S}, \forall h \in [H], \quad \pi_h^{t+1}(s) = \arg \max_{a \in \mathcal{A}} \bar{E}_h^t(s, a).$$

- **stopping rule:** $\tau = \inf \{ t \in \mathbb{N} : \bar{E}_1^t(s_1, \pi_1^{t+1}(s_1)) \leq \varepsilon/2 \}$
- **prediction:** transition kernel \hat{P}^τ

→ very close to an old algorithm by [Fiechter, 1994]
... originally proposed for Best Policy Identification!

Theorem [Kaufmann et al. 2020]

With $\beta(n, \delta) \simeq \log\left(\frac{1}{\delta}\right) + (S - 1) \log(n)$, RF-UCRL is (ε, δ) -PAC for Reward-Free Exploration and satisfies, w.p. $1 - \delta$,

$$\tau^{\text{RF-UCRL}} = \tilde{O} \left(\frac{H^4 SA}{\varepsilon^2} \left[\log \left(\frac{1}{\delta} \right) + S \right] \right)$$

→ improves over the state-of-the art bound of [Jin et al. 20]

$$\tau^{\text{RF-RL-Explore}} = \tilde{O} \left(\frac{S^2 AH^5}{\varepsilon^2} \log \left(\frac{1}{\delta} \right) + \frac{S^4 AH^7}{\varepsilon} \log^3 \left(\frac{1}{\delta} \right) \right)$$

with a very **different approach**

→ RF-UCRL is a natural **adaptive** approach to RFE
... with a simple sample complexity analysis

Sample complexity: Sketch of proof

$p_h^t(s, a)$: probability to visit state (s, a) at step h under policy π^t

If the algorithm does not stop after $t + 1$ episodes,

$$\begin{aligned}\epsilon/2 &\leq \bar{E}_1^t(s_1, \pi_1^{t+1}(s_1)) \\ &\lesssim \sum_{h,s,a} \hat{p}_h^t(s, a) \sqrt{H^2 \frac{\beta(n_h^t(s, a), \delta)}{n_h^t(s, a)}} \quad (\text{inductive definition of } \bar{E}_h) \\ &\lesssim \sum_{h,s,a} p_h^t(s, a) \sqrt{H^2 \frac{\beta(n_h^t(s, a), \delta)}{n_h^t(s, a)}} \quad (\text{concentration: on the event } \mathcal{E})\end{aligned}$$

Summing for all $t < \tau$,

$$\begin{aligned}\tau \times (\epsilon/2) &\lesssim \sum_{h,s,a} \sum_{t=0}^{\tau-1} p_h^{t+1}(s, a) \sqrt{H^2 \frac{\beta(n_h^t(s, a), \delta)}{n_h^t(s, a)}} \\ &\lesssim \sum_h \sqrt{H^2 SA \tau \beta(\tau, \delta)} = C_0 \sqrt{H^4 SA \tau \beta(\tau, \delta)}\end{aligned}$$

Therefore,

$$\tau \leq \inf \left\{ t \in \mathbb{N} : t(\epsilon/2)^2 > C_0^2 H^4 SA \beta(t, \delta) \right\} .$$

- 1 The BPI and RFE objectives
- 2 Reward-Free UCRL
- 3 BPI Algorithms**

First observation: RF-UCRL is also (ε, δ) -PAC for Best Policy Identification with the updated

- **prediction rule:** $\hat{\pi}$, the optimal policy if the MDP $(\mathcal{S}, \mathcal{A}, \hat{P}^\tau, r)$

$$\tau^{\text{RF-UCRL}} = \tilde{O} \left(\frac{H^4 SA}{\varepsilon^2} \left[\log \left(\frac{1}{\delta} \right) + S \right] \right) \text{ w.h.p.}$$

Lower bound for BPI [Domingues et al. 2020]

For every (ε, δ) -PAC BPI algorithm, there exists an MDP (with stage-dependent transitions) such that

$$\mathbb{E}[\tau] \geq c_1 \frac{H^3 SA}{\varepsilon^2} \log \left(\frac{1}{\delta} \right),$$

where c_1 is an absolute constant.

→ some room for improvement...

Building on Regret Minimizing algorithm

The UCB-VI algorithm of [Azar et al. 17] satisfies

$$\mathbb{E} \left[\sum_{t=1}^T \left(V_1^*(s_1; r) - V_1^{\pi^t}(s_1; r) \right) \right] \leq C \left(\sqrt{H^3 SAT} \right)$$

(minimax optimal **cumulative regret**)

From UCB-VI to a BPI algorithm [Jin et al. 18]

- **exploration policy:** that of the UCB-VI algorithm
- **stopping rule:** $T = \frac{C^2 SAH^3}{\epsilon^2 \delta^2}$
($\tau = T$ is fixed in advance)
- **prediction rule:** $\hat{\pi}$ is **one of the policies used by UCB-VI**, chosen uniformly at random: $\hat{\pi} = \hat{\pi}^n \quad n \sim \mathcal{U}(\{1, \dots, T\})$

Building on Regret Minimizing algorithm

The UCB-VI algorithm of [Azar et al. 17] satisfies

$$\mathbb{E} \left[\sum_{t=1}^T \left(V_1^*(s_1; r) - V_1^{\pi^t}(s_1; r) \right) \right] \leq C \left(\sqrt{H^3 SAT} \right)$$

(minimax optimal **cumulative regret**)

From UCB-VI to a BPI algorithm [Jin et al. 18]

- **exploration policy:** that of the UCB-VI algorithm
- **stopping rule:** $T = \frac{C^2 SAH^3}{\varepsilon^2 \delta^2}$
($\tau = T$ is fixed in advance)
- **prediction rule:** $\hat{\pi}$ is **one of the policies used by UCB-VI**, chosen uniformly at random: $\hat{\pi} = \hat{\pi}^n \quad n \sim \mathcal{U}(\{1, \dots, T\})$

→ optimal dependency in ε , **sub-optimal dependency in δ**

A more *adaptive* conversion from a regret minimizer:

→ associate a **data-dependent stopping rule** to a UCRL algorithm

BPI-UCRL

- **exploration policy:** $\pi^{t+1}(s) = \arg \max_{a \in \mathcal{A}} \bar{Q}_h^t(s, a; r)$
- **stopping rule:** $\tau = \inf \{t \in \mathbb{N} : \bar{V}_1^t(s_1; r) - \underline{V}_1^t(s_1; r) \leq \epsilon\}$
- **prediction rule:** $\hat{\pi}_h(s) = \arg \max_{a \in \mathcal{A}} \underline{Q}_h^r(s, a; r)$

where we have built upper and lower confidence bounds

$$\begin{aligned} \underline{Q}_h^t(s, a; r) &\leq Q_h^*(s, a; r) \leq \bar{Q}_h^t(s, a; r) \\ \underline{V}_h^t(s; r) &\leq V_h^*(s; r) \leq \bar{V}_h^t(s; r). \end{aligned}$$

A more *adaptive* conversion from a regret minimizer:

→ associate a **data-dependent stopping rule** to a UCRL algorithm

BPI-UCRL

- **exploration policy:** $\pi^{t+1}(s) = \arg \max_{a \in \mathcal{A}} \overline{Q}_h^t(s, a; r)$
- **stopping rule:** $\tau = \inf \{ t \in \mathbb{N} : \overline{V}_1^t(s_1; r) - \underline{V}_1^t(s_1; r) \leq \epsilon \}$
- **prediction rule:** $\hat{\pi}_h(s) = \arg \max_{a \in \mathcal{A}} \underline{Q}_h^r(s, a; r)$

where we have built upper and lower confidence bounds

$$\begin{aligned} \underline{Q}_h^t(s, a; r) &\leq Q_h^*(s, a; r) \leq \overline{Q}_h^t(s, a; r) \\ \underline{V}_h^t(s; r) &\leq V_h^*(s; r) \leq \overline{V}_h^t(s; r). \end{aligned}$$

Theorem [Kaufmann et al. 2020]

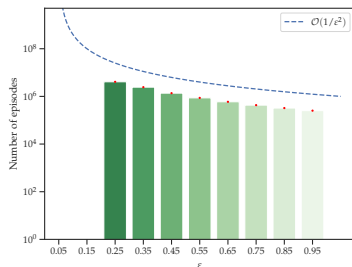
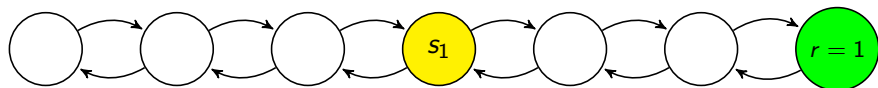
With $\beta(n, \delta) \simeq \log\left(\frac{1}{\delta}\right) + (S - 1) \log(n)$, BPI-UCRL is (ε, δ) -PAC for Best Policy Identification and satisfies, w.p. $\geq 1 - \delta$,

$$\tau^{\text{BPI-UCRL}} = \tilde{O}\left(\frac{H^4 SA}{\varepsilon^2} \left[\log\left(\frac{1}{\delta}\right) + S\right]\right)$$

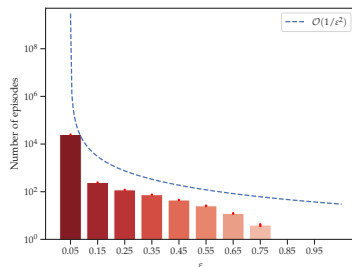
- similar sample complexity bound as RF-UCRL (obtained with a similar proof)
- yet the practical story is different...

RF-UCRL versus BPI-UCRL

Double Chain MDP with $L = 31, H = 20$:



$\mathbb{E}[\tau | \tau < 10^8]$ for RF-UCRL

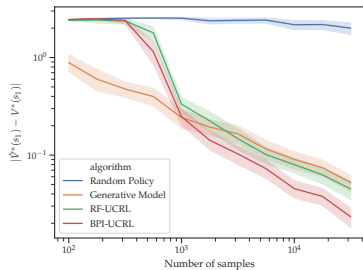
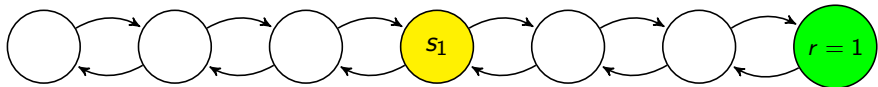


$\mathbb{E}[\tau | \tau < 10^6]$ for BPI-UCRL

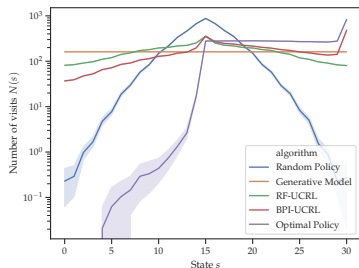
→ BPI-UCRL has a much smaller sample complexity!

RF-UCRL versus BPI-UCRL

Double Chain MDP with $L = 31, H = 20$:



approximation error



number of visits of each state

➔ RF-UCRL explores more uniformly

- The sample complexity of...

	Upper Bound	Lower Bound
BPI	$\frac{H^4 SA}{\epsilon^2} \left[\log \left(\frac{1}{\delta} \right) + S \right]$ BPI-UCRL / RF-UCRL	$\frac{H^3 SA}{\epsilon^2} \log \left(\frac{1}{\delta} \right)$ [Darwiche Domingues et al. 2020]
RFE	$\frac{H^4 SA}{\epsilon^2} \left[\log \left(\frac{1}{\delta} \right) + S \right]$ RF-UCRL	$\frac{H^3 SA}{\epsilon^2} \left[\log \left(\frac{1}{\delta} \right) + S \right]$ + [Jin et al. 2020]

Follow-up work: shaving the remaining H factor for BPI and RFE for more sophisticated algorithms using Bernstein bonuses

- BPI-UCBVI for Best Policy Identification
- RF-Express for Reward Free Exploration

Ménard et al. 2020, *Fast active learning for pure exploration in reinforcement learning*, arXiv:2007.13442

- The sample complexity of...

	Upper Bound	Lower Bound
BPI	$\frac{H^4 SA}{\epsilon^2} \left[\log \left(\frac{1}{\delta} \right) + S \right]$ BPI-UCRL / RF-UCRL	$\frac{H^3 SA}{\epsilon^2} \log \left(\frac{1}{\delta} \right)$ [Darwiche Domingues et al. 2020]
RFE	$\frac{H^4 SA}{\epsilon^2} \left[\log \left(\frac{1}{\delta} \right) + S \right]$ RF-UCRL	$\frac{H^3 SA}{\epsilon^2} \left[\log \left(\frac{1}{\delta} \right) + S \right]$ + [Jin et al. 2020]

Future work: beyond worst-case guarantees

- problem-dependent sample complexity for the simpler *planning* problem (= find the best first action)
[Jonsson et al., 2020]
- problem-dependent regret guarantees
[Simchowitz and Jamieson, 2019]

... how about BPI?

- Azar et al., *Minimax Regret Bounds for Reinforcement Learning*, ICML 2017
- Darwiche Domingues et al., *Episodic Reinforcement Learning in Finite MDPs: Minimax Lower Bounds Revisited*, arXiv:2010.03531, 2020
- Jaksch et al., *Near-optimal Regret Bounds for Reinforcement Learning*, JMLR 2010
- Jin et al., *Reward-Free Exploration for Reinforcement Learning*, ICML 2020
- Jin et al., *Is Q-Learning Provably Efficient?*, NeurIPS 2018
- Jonsson et al., *Planning in Markov Decision Processes with Gap-Dependent Sample Complexity*, NeurIPS 2020
- Fiechter, *Efficient Reinforcement Learning*, COLT 1994
- Kaufmann et al., *Adaptive Reward-Free Exploration*, arXiv:2006.06294, 2020
- Ménard et al., *Fast active learning for pure exploration in reinforcement learning*, arXiv:2007.13442, 2020
- Simchowitz and Jamieson, *Non-Asymptotic Gap-Dependent Regret Bounds for Tabular MDPs*, NeurIPS 2019

New bonuses, of order $\beta(\delta, n)/n$:

$$W_h^t(s, a) \approx H^2 \frac{\beta(n_h^t(s, a), \delta)}{n_h^t(s, a)} + \sum_{s'} \hat{p}_h^t(s'|s, a) \max_{a'} W_{h+1}^t(s', a')$$

$$\pi_h^{t+1}(s) = \arg \max_a W_h^t(s, a)$$

Upper bound on the error:

$$\hat{e}_1^{t, \pi}(s_1, \pi_1(s_1)) \lesssim \sqrt{\max_{a \in \mathcal{A}} W_1^t(s_1, a)} + \max_{a \in \mathcal{A}} W_1^t(s_1, a).$$

→ control of the error **only in the initial state** s_1 .

$$\tau = \inf \{ t \in \mathbb{N}^* : \sqrt{\max_{a \in \mathcal{A}} W_1^t(s_1, a)} + \max_{a \in \mathcal{A}} W_1^t(s_1, a) \leq \epsilon/2 \}$$