# Exploration-Exploitation in Reinforcement Learning

**Ronan Fruit[*], Alessandro Lazaric[†] and Matteo Pirotta[†]**

**Facebook AI Research[†] and INRIA Lille[*]**

# Reinforcement Learning



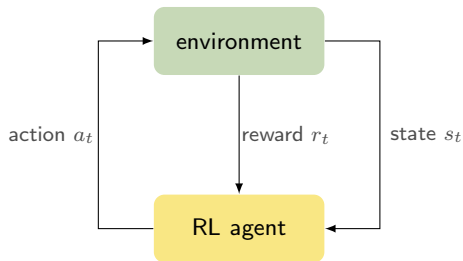action $a_t$     reward $r_t$     state $s_t$

environment

RL agent

"**Reinforcement learning** is learning how to map states to actions so as to **maximize** a numerical **reward** signal in an unknown and **uncertain** environment.

In the most interesting and challenging cases, **actions** affect not only the immediate reward but also the **next situation** and all subsequent rewards (**delayed reward**).

The agent is not told which actions to take but it must discover which actions yield the most reward by trying them (**trial-and-error**)."

— Sutton and Barto [1998]

# Reinforcement Learning



action $a_t$

reward $r_t$

state $s_t$

environment

RL agent

*Exploration*

"**Reinforcement learning** is learning how to map states to actions so as to **maximize** a numerical **reward** signal in an unknown and **uncertain** environment.

In the most interesting and challenging cases, **actions** affect not only the immediate reward but also the **next situation** and all subsequent rewards (**delayed reward**).

The agent is not told which actions to take but it must discover which actions yield the most reward by trying them (**trial-and-error**)."

— Sutton and Barto [1998]

# Reinforcement Learning



*Exploitation*

"**Reinforcement learning** is learning how to map states to actions so as to **maximize** a numerical **reward** signal in an unknown and **uncertain** environment.

In the most interesting and challenging cases, **actions** affect not only the immediate reward but also the **next situation** and all subsequent rewards (**delayed reward**).
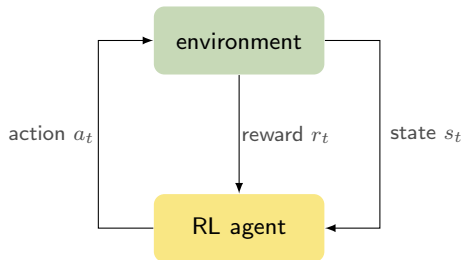
The agent is not told which actions to take but it must discover which actions yield the most reward by trying them (**trial-and-error**)."

*Exploration*

— Sutton and Barto [1998]

*Fruit, Lazaric and Pirotta*

# Disclaimer: the Real Title

Regret Minimization in

Infinite-Horizon

Finite Markov Decision Processes

# Organization

Website
`https://rlgammazero.github.io`

# Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

- State space $\mathcal{S}$, $|\mathcal{S}| = S < \infty$

- Action space $\mathcal{A}$, $|\mathcal{A}| = A < \infty$

- Transition distribution $p(\cdot|s,a) \in \Delta(\mathcal{S})$

- Reward distribution with expectation $r(s,a) \in [0, r_{\max}]$

# Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

- State space $\mathcal{S}, \; |\mathcal{S}| = S < \infty$

- Action space $\mathcal{A}, \; |\mathcal{A}| = A < \infty$

  finite

- Transition distribution $p(\cdot|s,a) \in \Delta(\mathcal{S})$

- Reward distribution with expectation $r(s,a) \in [0, r_{\max}]$

# Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

- State space $\mathcal{S}$, $|\mathcal{S}| = S < \infty$ ⎫
- Action space $\mathcal{A}$, $|\mathcal{A}| = A < \infty$ ⎭ finite

- Transition distribution $p(\cdot|s,a) \in \Delta(\mathcal{S})$ } Markov

- Reward distribution with expectation $r(s,a) \in [0, r_{\max}]$

# Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

- State space $\mathcal{S}$, $|\mathcal{S}| = S < \infty$ ⎫
  
- Action space $\mathcal{A}$, $|\mathcal{A}| = A < \infty$ ⎭ finite

- Transition distribution $p(\cdot|s, a) \in \Delta(\mathcal{S})$ } Markov

- Reward distribution with expectation $r(s, a) \in [0, r_{\max}]$

👍 The process generates history $H_t = (s_1, a_1, \ldots, s_{t-1}, a_{t-1}, s_t)$, with $s_{t+1} \sim p(\cdot|s_t, a_t)$

# Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

- State space $\mathcal{S}, \ |\mathcal{S}| = S < \infty$ ⎫
- Action space $\mathcal{A}, \ |\mathcal{A}| = A < \infty$ ⎭ finite

- Transition distribution $p(\cdot|s,a) \in \Delta(\mathcal{S})$ } Markov

- Reward distribution with expectation $r(s,a) \in [0, r_{\max}]$

☞ The process generates history $H_t = (s_1, a_1, \ldots, s_{t-1}, a_{t-1}, s_t)$, with $s_{t+1} \sim p(\cdot|s_t, a_t)$

📓 In (contextual) bandit, actions do not influence the evolution of states

# Policies

An agent acts according to a *policy*

|  | stationary | history-dependent |
|---|---|---|
| deterministic | $\pi : \mathcal{S} \to \mathcal{A}$ | $\pi_t : \mathcal{H}_t \to \mathcal{A}$ |
| stochastic | $\pi : \mathcal{S} \to \Delta(\mathcal{A})$ | $\pi_t : \mathcal{H}_t \to \Delta(\mathcal{A})$ |

# Classification

An MDP $M$ is

- *ergodic* if it is possible to go from any state to any other state under *any* deterministic stationary policy

$$\forall s, s', \; \forall \pi : \mathcal{S} \to \mathcal{A}, \; \exists t < \infty, \; \text{s.t. } \mathbb{P}_\pi^M \big( s_t = s' | s_0 = s \big) > 0$$

- *communicating* if it is possible to go from any state to any other state under *a specific* deterministic stationary policy

$$\forall s, s', \; \exists \pi : \mathcal{S} \to \mathcal{A}, \; \exists t < \infty, \; \text{s.t. } \mathbb{P}_\pi^M \big( s_t = s' | s_0 = s \big) > 0$$

☞ A communicating MDP has *finite diameter*

$$D_M = \max_{s, s' \in \mathcal{S}} \; \min_{\pi : \mathcal{S} \to \mathcal{A}} \; \mathbb{E} \big[ T_\pi^M (s, s') \big]$$

# Classification

An MDP $M$ is

- *ergodic* if it is possible to go from any state to any other state under *any* deterministic stationary policy

$$\forall s, s', \; \forall \pi : \mathcal{S} \to \mathcal{A}, \; \exists t < \infty, \; \text{s.t.} \; \mathbb{P}_\pi^M\big(s_t = s' | s_0 = s\big) > 0$$

- *communicating* if it is possible to go from any state to any other state under *a specific* deterministic stationary policy

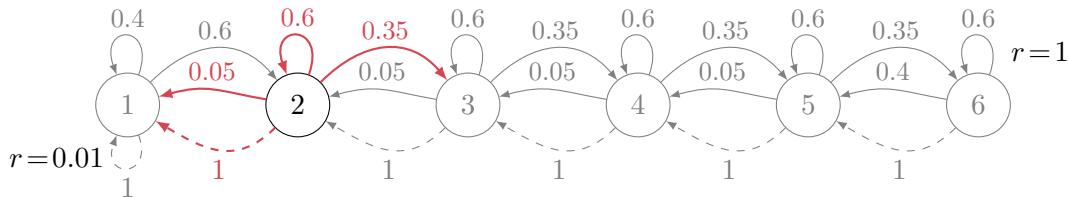$$\forall s, s', \; \exists \pi : \mathcal{S} \to \mathcal{A}, \; \exists t < \infty, \; \text{s.t.} \; \mathbb{P}_\pi^M\big(s_t = s' | s_0 = s\big) > 0$$

👍 A communicating MDP has *finite diameter*

$$D_M = \max_{s,s' \in \mathcal{S}} \underbrace{\min_{\pi:\mathcal{S} \to \mathcal{A}} \mathbb{E}\big[T_\pi^M(s, s')\big]}_{\text{shortest path}}$$

# River Swim: Markov Decision Processes

Strehl and Littman [2008]

- $\mathcal{S} = \{1, 2, 3, 4, 5, 6\}$, $\mathcal{A} = \{L, R\}$
- $\pi_L(s) = L$, $\pi_R(s) = R$
- $M \oplus \pi_R$ is *ergodic* but $M \oplus \pi_L$ is *not ergodic*
- $T_{\pi_L}^M(6, 1) = 5$,   $D_M = \mathbb{E}\big[T_{\pi_R}^M(1, 6)\big] \approx 14.7$

# Gain and Bias

*Gain* of a deterministic stationary policy $\pi$

$$g_M^\pi(s) = \lim_{T \to \infty} \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^{T} r(s_t, a_t) \Big| s_0 = s, a_t = \pi(s_t)\right]$$

*Bias* of a deterministic stationary policy $\pi$

$$h_M^\pi(s) := C\text{-}\lim_{T \to \infty} \mathbb{E}\left[\sum_{t=1}^{T} \big(r(s_t, a_t) - g_M^\pi(s_t)\big) \Big| s_0 = s, a_t = \pi(s_t)\right]$$

*Span* of the bias function

$$\mathsf{sp}\big(h_M^\pi\big) = \max_s h_M^\pi(s) - \min_s h_M^\pi(s)$$

# Bellman operators

*Bellman* operator $L_M^a : \mathbb{R}^S \to \mathbb{R}^S$

$$= \sum_{s'} p(s'|s,a)h(s')$$

$$L_M^a h(s) = r(s,a) + p(\cdot|s,a)^\mathsf{T} h$$

*Optimal Bellman* operator $L_M^\star : \mathbb{R}^S \to \mathbb{R}^S$

$$L_M^\star h(s) = \max_{a \in \mathcal{A}} \left\{ r(s,a) + p(\cdot|s,a)^\mathsf{T} h \right\}$$

*Optimality gap* of action $a$ at $s$

$$\delta_M^\star(s,a) = L_M^\star h_M^\star(s) - L_M^a h_M^\star(s)$$

a.k.a. advantage function

# Optimality

*Optimal policy* and *optimal gain*

$$\pi_M^\star \in \arg\max_\pi g_M^\pi(s) \qquad g_M^\star = g_M^{\pi^\star}(s) \;\; \forall s \in \mathcal{S}$$

*Optimality equation*

$$h_M^\star(s) + g_M^\star = L_M^\star h_M^\star(s)$$

*Greedy policy* w.r.t. $h_M^\star$ is optimal

$$\pi_M^\star(s) \in \arg\max_{a \in \mathcal{A}} \left\{ r(s,a) + p(\cdot|s,a)^\mathsf{T} h_M^\star \right\}$$

*Set of optimal actions* in state $s$

$$\Pi_M^\star(s) = \arg\max_{a \in \mathcal{A}} \left\{ r(s,a) + p(\cdot|s,a)^\mathsf{T} h_M^\star \right\}$$

# Optimality

deterministic stationary

*Optimal policy* and *optimal gain*

$$\pi_M^\star \in \arg\max_\pi g_M^\pi(s) \qquad g_M^\star = g_M^{\pi^\star}(s) \ \ \forall s \in \mathcal{S}$$

*Optimality equation*

$$h_M^\star(s) + g_M^\star = L_M^\star h_M^\star(s)$$

*Greedy policy* w.r.t. $h_M^\star$ is optimal

$$\pi_M^\star(s) \in \arg\max_{a\in\mathcal{A}}\left\{ r(s,a) + p(\cdot|s,a)^\mathsf{T} h_M^\star \right\}$$

*Set of optimal actions* in state $s$

$$\Pi_M^\star(s) = \arg\max_{a\in\mathcal{A}}\left\{ r(s,a) + p(\cdot|s,a)^\mathsf{T} h_M^\star \right\}$$

# Optimality

deterministic stationary

constant gain*

*Optimal policy* and *optimal gain*

$$\pi_M^\star \in \arg\max_\pi g_M^\pi(s) \qquad g_M^\star = g_M^{\pi^\star}(s) \;\; \forall s \in \mathcal{S}$$

*Optimality equation*

$$h_M^\star(s) + g_M^\star = L_M^\star h_M^\star(s)$$

*Greedy policy* w.r.t. $h_M^\star$ is optimal

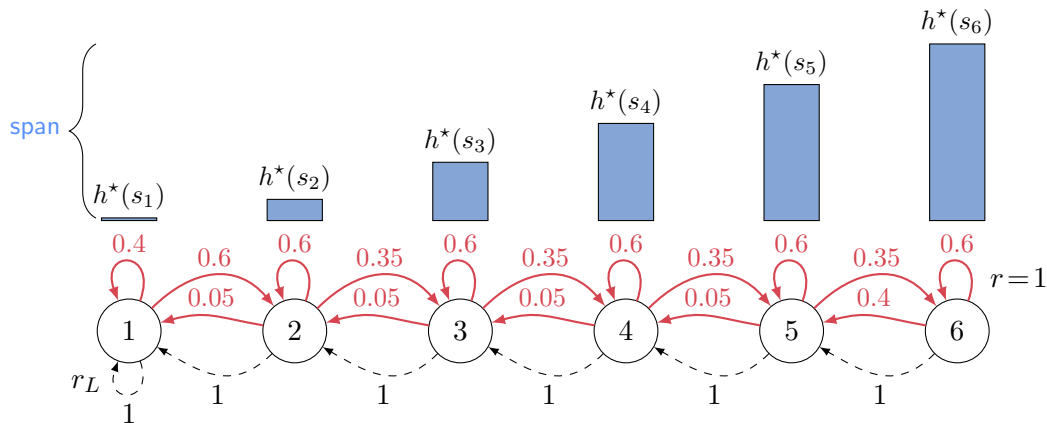$$\pi_M^\star(s) \in \arg\max_{a \in \mathcal{A}} \left\{ r(s,a) + p(\cdot|s,a)^\mathsf{T} h_M^\star \right\}$$

*Set of optimal actions* in state $s$

$$\Pi_M^\star(s) = \arg\max_{a \in \mathcal{A}} \left\{ r(s,a) + p(\cdot|s,a)^\mathsf{T} h_M^\star \right\}$$

*In communicating MDPs

# River Swim: Optimality



- $\pi^\star = \pi_R$
- If $r_L = 0.01$, $g^\star \approx 0.43$, $\mathsf{sp}(h^\star) \approx 6.4$

# River Swim: Optimality



- $\pi^\star = \pi_R$
- If $r_L = 0.01$, $g^\star \approx 0.43$, $\mathsf{sp}(h^\star) \approx 6.4$
- If $r_L = 0.4$, $g^\star \approx 0.43$, $\mathsf{sp}(h^\star) \approx 5.5$

# River Swim: Optimality



- $\pi^\star = \pi_R$
- If $r_L = 0.01$, $g^\star \approx 0.43$, $\mathsf{sp}(h^\star) \approx 6.4$
- If $r_L = 0.4$, $g^\star \approx 0.43$, $\mathsf{sp}(h^\star) \approx 5.5$

$D$ is constant

# Value Iteration

---

**initialize** $v_0(s) = 0 \ \forall s \in \mathcal{S}, \ n = 0, \ \varepsilon$

**repeat**

    **for** $s \in \mathcal{S}$ **do**

$$v_{n+1}(s) = L_M^\star v_n(s) = \max_{a \in \mathcal{A}} \left\{ r(s,a) + p(\cdot|s,a)^\mathsf{T} v_n \right\}$$

    **end**

    $n = n + 1$

**until** $sp(v_{n+1} - v_n) < \varepsilon$

**return** greedy policy

$$\pi_\varepsilon(s) = \arg\max_{a \in \mathcal{A}} L_M^a v_n(s) = \arg\max_{a \in \mathcal{A}} \left\{ r(s,a) + p(\cdot|s,a)^\mathsf{T} v_n \right\}$$

---

# Value Iteration

## Theorem (Thm. 8.5.5 [Puterman, 1994])

*In any communicating MDP $M$, value iteration is such that*

- *convergence: for any $\varepsilon$, there exists $n_\epsilon$ s.t. the stopping condition is met*

- *optimality: policy $\pi_\varepsilon$ is $\epsilon$-optimal*

$$g_M^{\pi_\varepsilon}(s) \geq g_M^\star - \varepsilon$$

# Regret Minimization

# Regret Minimization

# Regret Minimization



$$R(T, \ M^\star, \ \mathfrak{A}) = T g_M^\star - \sum_{t=1}^{T} r_t$$

unknown true MDP $M^\star = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

algorithm $\mathfrak{A} = \{\pi_t\}$

reward obtained by $\mathfrak{A}$

# Regret Minimization



reward

$g_M^\star$

reward $r_t$

regret $R$

$T$

unknown true MDP $M^\star = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$

algorithm $\mathfrak{A} = \{\pi_t\}$

reward obtained by $\mathfrak{A}$

$$R(T,\ M^\star,\ \mathfrak{A}) = T g_M^\star - \sum_{t=1}^{T} r_t$$

Expected regret w.r.t. randomness of $s_t$, $r_t$, and (possibly) $\mathfrak{A}$

$$\overline{R}(T, M^\star, \mathfrak{A}) = \mathbb{E}\big[R(T, M^\star, \mathfrak{A})\big]$$

# Problem-Dependent Lower Bound

Let $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$ and $M' = \langle \mathcal{S}, \mathcal{A}, r, p' \rangle$

- *Difference* between $M$ and $M'$ at $s, a$ (w.l.o.g. assuming reward known)

$$\mathsf{KL}_{M,M'}(s, a) = \mathsf{KL}\big(p(\cdot|s, a) \| p'(\cdot|s, a)\big)$$

- *Set of alternative* (confusing) models w.r.t. $M$

  same everywhere but in $(s, a)$

$$\mathcal{M}_M^{\mathsf{alt}}(s, a) = \Big\{ M' \colon p'(\cdot|s', a') = p(\cdot|s', a'), \text{ for all } (s', a') \neq (s, a),$$

$$a \notin \Pi_M^\star(s) \ , \ a \in \Pi_{M'}^\star(s) \Big\}$$

sub-optimal in $M$

optimal in $M'$

# Problem-Dependent Lower Bound

**Theorem** (Thm. 1 Burnetas and Katehakis [1997], Thm. 2 Ok et al. [2018])

*Let $\mathfrak{A}$ be s.t. $\overline{R}(T, M, \mathfrak{A}) = o(T^\alpha)$ for all $\alpha > 0$ and ergodic MDP $M$. For any ergodic MDP $M^\star$ with $r_{\max} = 1$, the expected regret is lower bounded as*

$$\liminf_{T \to \infty} \frac{\overline{R}(T, M^\star, \mathfrak{A})}{\log T} \geq K_{M^\star}$$

cumulative regret

*where*

$$K_{M^\star} = \inf_{\eta \geq 0} \sum_{s,a} \eta(s,a) \delta_{M^\star}^\star(s,a)$$

$$\textit{s.t. } \sum_{s,a} \eta(s,a) KL_{M^\star, M}(s,a) \geq 1 \quad \forall M \in \mathcal{M}_{M^\star}^{alt}(s,a)$$

"evidence" of difference between $M^\star$ and $M$

# Problem-Dependent Lower Bound

**Theorem** (Thm. 1 Burnetas and Katehakis [1997], Thm. 2 Ok et al. [2018])

*Let $\mathfrak{A}$ be s.t. $\overline{R}(T, M, \mathfrak{A}) = o(T^\alpha)$ for all $\alpha > 0$ and ergodic MDP $M$. For any ergodic MDP $M^\star$ with $r_{\max} = 1$, the expected regret is lower bounded as*

$$\liminf_{T \to \infty} \frac{\overline{R}(T, M^\star, \mathfrak{A})}{\log T} \geq K_{M^\star}$$

cumulative regret

*where*

$$K_{M^\star} = \inf_{\eta \geq 0} \sum_{s,a} \eta(s,a) \delta^\star_{M^\star}(s,a)$$

$$\textit{s.t.} \sum_{s,a} \eta(s,a) KL_{M^\star, M}(s,a) \geq 1 \quad \forall M \in \mathcal{M}^{alt}_{M^\star}(s,a)$$

"evidence" of difference between $M^\star$ and $M$

📋 Similar to [Lai and Robbins, 1985] for MAB but alternative models and regret are different.

# Problem-Dependent Lower Bound

**Theorem** (Thm. 1 Burnetas and Katehakis [1997], Thm. 2 Ok et al. [2018])

*Let $\mathfrak{A}$ be s.t. $\overline{R}(T, M, \mathfrak{A}) = o(T^{\alpha})$ for all $\alpha > 0$ and ergodic MDP $M$. For any ergodic MDP $M^{\star}$ with $r_{\max} = 1$, the expected regret is lower bounded as*

$$\liminf_{T \to \infty} \frac{\overline{R}(T, M^{\star}, \mathfrak{A})}{\log T} \geq K_{M^{\star}}$$

*where*

$$K_{M^{\star}} \leq 2 \frac{(C+1)^2}{\min_{s,a} \delta_{M^{\star}}(s,a)} SA \qquad C = sp(h_{M^{\star}}^{\star})$$

# Minimax Lower Bound

**Theorem** (Thm. 5 Jaksch et al. [2010])

*For any communicating MDP $M^\star$ with $r_{\max} = 1$, $S, A \geq 10$, $D \geq 20 \log_A S$, any algorithm $\mathfrak{A}$ at any time $T \geq DSA$ suffers a regret*

$$\sup_{M^\star} \overline{R}(T, M^\star, \mathfrak{A}) \geq 0.015 \sqrt{DSAT}$$

# Minimax Lower Bound

> **Theorem** (Thm. 5 Jaksch et al. [2010])
>
> *For any* communicating *MDP* $M^\star$ *with* $r_{\max} = 1$, $S, A \geq 10$, $D \geq 20 \log_A S$, *any algorithm* $\mathfrak{A}$ *at any time* $T \geq DSA$ *suffers a regret*
>
> $$\sup_{M^\star} \overline{R}(T, M^\star, \mathfrak{A}) \geq 0.015\sqrt{DSAT}$$

📓 In MAB $\Omega(\sqrt{AT})$ since $D = 1$ and $S = 1$.

# Open Questions

$C$ could be arbitrarily large
($C = \infty$ for non ergodic)

**1** *Asymptotic* regime and *ergodicity* assumption

$$\mathbb{P}_M^\pi \big[ N_T(s) \geq \rho T \big] \geq 1 - C \ \exp(-\rho T/2) \qquad \text{[Prop.2 Burnetas and Katehakis [1997]]}$$

**2** *Span vs. diameter*

$D = 2\mathsf{sp}(h^\star)$ in the proof

$$\overline{R}(T, M^\star, \mathfrak{A}) \geq 0.015 \sqrt{D \ SAT}$$

**3** *Number of states vs branching factor* $\Gamma = \max\limits_{s,a} |\mathsf{supp}(p(\cdot|s,a))|$

$$\overline{R}(T, M^\star, \mathfrak{A}) \geq 0.015 \sqrt{D \ S \ AT}$$

$\Gamma = 2$ in the proof

# The Optimism Principle: Intuition

# The Optimism Principle: Intuition

Exploration vs. Exploitation

Fruit, Lazaric and Pirotta

# The Optimism Principle: Intuition

Exploration vs. Exploitation

*Optimism in Face of Uncertainty*

When you are uncertain, consider the **best possible world (reward-wise)**

# The Optimism Principle: Intuition

Exploration vs. Exploitation

> *Optimism in Face of Uncertainty*
>
> When you are uncertain, consider the **best possible world (reward-wise)**

If the best possible world is **correct**

$\implies$ **no regret**

**Exploitation**

If the best possible world is **wrong**

$\implies$ **learn useful information**

**Exploration**

# The Optimism Principle: Intuition

Exploration vs. Exploitation

Optimism in **gain**

Optimism in Face of Uncertainty

When you are uncertain, consider the **best possible world (reward-wise)**

If the best possible world is **correct**

$\implies$ **no regret**

**Exploitation**

If the best possible world is **wrong**

$\implies$ **learn useful information**

**Exploration**

Fruit, Lazaric and Pirotta

# History: OFU for Regret Minimization in RL

FH: finite-horizon
AR: average reward



Agrawal [1990]

Auer and Ortner [2006] (AR)
Bartlett and Tewari [2009] (AR)
Filippi et al. [2010] (AR)
Jaksch et al. [2010] (AR)

Talebi and Maillard [2018] (AR)
Fruit et al. [2018b] (AR)
Fruit et al. [2018a] (AR)
Qian et al. [2018] (AR)

Azar et al. [2017] (FH)
Zanette and Brunskill [2018] (FH)
Kakade et al. [2018] (FH)
Jin et al. [2018] (FH)

Fruit, Lazaric and Pirotta

# Gain Optimism:  Example



$a_0,\ r(s, a_0)?$

$s$

$a_1,\ r(s, a_1)?$

$a_2,\ r(s, a_2)?$

- Deterministic *policies*:
  - $\pi_0(s) = a_0$
  - $\pi_1(s) = a_1$
  - $\pi_2(s) = a_2$

# Gain Optimism: Example

$a_0, \ r(s, a_0)?$

$s$

$a_1, \ r(s, a_1)?$

$a_2, \ r(s, a_2)?$

- Deterministic *policies*:
  - $\pi_0(s) = a_0$
  - $\pi_1(s) = a_1$
  - $\pi_2(s) = a_2$

- Optimism

$$\widetilde{\pi} = \arg\max_{\pi_i} \mathsf{UCB}(g^{\pi_i})$$

# Gain Optimism: Example



$a_0$, $r(s, a_0) = g^{\pi_0}$?

$a_1$, $r(s, a_1) = g^{\pi_1}$?

$a_2$, $r(s, a_2) = g^{\pi_2}$?

$s$

- Deterministic *policies*:
  - $\pi_0(s) = a_0$
  - $\pi_1(s) = a_1$
  - $\pi_2(s) = a_2$

- Reward $r(s, a_i) = $ *gain* $g^{\pi_i}$

- Optimism
$$\widetilde{\pi} = \arg\max_{\pi_i} \mathsf{UCB}(g^{\pi_i})$$

# Gain Optimism: Example

- Deterministic *policies*:
  - $\pi_0(s) = a_0$
  - $\pi_1(s) = a_1$
  - $\pi_2(s) = a_2$

- Reward $r(s, a_i) = $ *gain* $g^{\pi_i}$

- Upper confidence bound
  $$\text{UCB}(g^{\pi_i}) = \text{UCB}(r(s, a_i))$$

- Optimism
  $$\widetilde{\pi} = \arg\max_{\pi_i} \text{UCB}(g^{\pi_i})$$

# Gain Optimism: Example

Deterministic *policies*:
- $\pi_0(s) = a_0$
- $\pi_1(s) = a_1$
- $\pi_2(s) = a_2$

Reward $r(s, a_i) = $ *gain* $g^{\pi_i}$

Upper confidence bound
$$\mathsf{UCB}(g^{\pi_i}) = \mathsf{UCB}(r(s, a_i))$$

Optimism
$$\widetilde{\pi} = \arg\max_{\pi_i} \mathsf{UCB}(g^{\pi_i})$$

In the diagram:

$a_0, \widehat{r}(s, a_0) + r_{\max}\sqrt{\dfrac{\log(1/\delta)}{N(s, a_0)}}$

$a_1, \widehat{r}(s, a_1) + r_{\max}\sqrt{\dfrac{\log(1/\delta)}{N(s, a_1)}}$

$a_2, \widehat{r}(s, a_2) + r_{\max}\sqrt{\dfrac{\log(1/\delta)}{N(s, a_2)}}$

confidence

num visits

estimated reward

Fruit, Lazaric and Pirotta

# Gain Optimism: Example



Deterministic *policies*:
- $\pi_0(s) = a_0$
- $\pi_1(s) = a_1$
- $\pi_2(s) = a_2$

Reward $r(s, a_i) = gain\ g^{\pi_i}$

Upper confidence bound
$$\mathsf{UCB}(g^{\pi_i}) = \mathsf{UCB}(r(s, a_i))$$

Optimism
$$\widetilde{\pi} = \arg\max_{\pi_i} \mathsf{UCB}(g^{\pi_i})$$

☝ UCB algorithm (Bandit)

From the diagram:

$a_0,\ \widehat{r}(s, a_0) + r_{\max}\sqrt{\dfrac{\log(1/\delta)}{N(s, a_0)}}$

$a_1,\ \widehat{r}(s, a_1) + r_{\max}\sqrt{\dfrac{\log(1/\delta)}{N(s, a_1)}}$

$a_2,\ \widehat{r}(s, a_2) + r_{\max}\sqrt{\dfrac{\log(1/\delta)}{N(s, a_2)}}$

confidence

num visits

estimated reward

$s$

# Gain Optimism: Implementation

---

### Tentative algorithm

---

Observe $s_1$

**for** $t = 1, 2, \ldots$ **do**

  *Compute* $\pi_t \leftarrow \arg\max_{\pi} UCB_t(g^{\pi})$

  Take action $a_t = \pi_t(s_t)$

  Observe reward $r_t$ and next state $s_{t+1}$

  Compute $\mathrm{UCB}_{t+1}(g^{\pi})$ for all $\pi$ based on $\mathrm{UCB}_t(g^{\pi})$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

**end**

---

# Gain Optimism: Implementation

---

**Tentative algorithm**

---

Observe $s_1$
**for** $t = 1, 2, \ldots$ **do**
    *Compute* $\pi_t \leftarrow \arg\max_{\pi} UCB_t(g^\pi)$
    Take action $a_t = \pi_t(s_t)$
    Observe reward $r_t$ and next state $s_{t+1}$
    Compute $\mathrm{UCB}_{t+1}(g^\pi)$ for all $\pi$ based on $\mathrm{UCB}_t(g^\pi)$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$
**end**

---

⚠ *3 major issues:*

■ *Upper confidence bounds*: construct $\mathrm{UCB}_t(g^\pi)$ with unknown dynamics

■ *Computational complexity*: exponential number of policies

■ *Frequent policy update*: inefficient exploration

# Gain Optimism: Implementation

**Tentative algorithm**

Observe $s_1$
for $t = 1, 2, \ldots$ do

    *Compute* $\pi_t \leftarrow \arg\max_{\pi} UCB_t(g^\pi)$

    Take action $a_t = \pi_t(s_t)$
    Observe reward $r_t$ and next state $s_{t+1}$
    Compute $UCB_{t+1}(g^\pi)$ for all $\pi$ based on $UCB_t(g^\pi)$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

end

⚠ *3 major issues:*

■ *Upper confidence bounds*: construct $UCB_t(g^\pi)$ with unknown dynamics

■ *Computational complexity*: exponential number of policies

■ *Frequent policy update*: inefficient exploration

# Bounded Parameter MDP: Definition

*Bounded parameter MDP* [Strehl and Littman, 2008]

$$\mathcal{M}_t = \left\{ \langle \mathcal{S}, \mathcal{A}, r, p \rangle : \ r(s,a) \in B_t^r(s,a), \ p(\cdot|s,a) \in B_t^p(s,a), \forall (s,a) \in \mathcal{S} \times \mathcal{A} \right\}$$

Compact *confidence sets*

$$B_t^r(s,a) := \left[ \widehat{r}_t(s,a) - \beta_t^r(s,a), \ \widehat{r}_t(s,a) + \beta_t^r(s,a) \right]$$

$$B_t^p(s,a) := \left\{ p(\cdot|s,a) \in \Delta(\mathcal{S}) : \ \|p(\cdot|s,a) - \widehat{p}_t(\cdot|s,a)\|_1 \leq \ \beta_t^p(s,a) \right\}$$

# Bounded Parameter MDP: Definition

*Bounded parameter MDP* [Strehl and Littman, 2008]

$$\mathcal{M}_t = \left\{ \langle \mathcal{S}, \mathcal{A}, r, p \rangle : \; r(s,a) \in B_t^r(s,a), \; p(\cdot|s,a) \in B_t^p(s,a), \forall (s,a) \in \mathcal{S} \times \mathcal{A} \right\}$$

Compact *confidence sets*

$$B_t^r(s,a) := \left[ \widehat{r}_t(s,a) - \beta_t^r(s,a), \; \widehat{r}_t(s,a) + \beta_t^r(s,a) \right]$$

$$B_t^p(s,a) := \left\{ p(\cdot|s,a) \in \Delta(\mathcal{S}) : \; \|p(\cdot|s,a) - \widehat{p}_t(\cdot|s,a)\|_1 \leq \; \beta_t^p(s,a) \right\}$$

*Confidence bounds* based on [Hoeffding, 1963] and [Weissman et al., 2003]

$$\beta_t^r(s,a) \propto \sqrt{\frac{\log(N_t(s,a)/\delta)}{N_t(s,a)}}$$

$$\beta_t^p(s,a) \propto \sqrt{\frac{S \log(N_t(s,a)/\delta)}{N_t(s,a)}}$$

# Bounded Parameter MDP: Optimism

$g_M^\pi$      Fix a *policy* $\pi$

# Bounded Parameter MDP: Optimism

Fix a *policy* $\pi$

# Bounded Parameter MDP: Optimism



Fix a *policy* $\pi$

# Bounded Parameter MDP: Optimism

Optimism: $\mathsf{UCB}_t(g^\pi) = \max_{M \in \mathcal{M}_t} g^\pi_M \geq g^\pi_{M^\star}$ $\quad \boxed{\mathsf{UCB}_t(g^\pi)}$

Fix a *policy* $\pi$

# Gain Optimism:  Implementation

---

Tentative algorithm

---

Observe state $s_1$
**for** $t = 1, 2, \dots$ **do**

  *Compute* $\pi_t \leftarrow \arg \max_{\pi} UCB_t(g^\pi)$

  Take action $a_t = \pi_t(s_t)$
  Observe reward $r_t$ and next state $s_{t+1}$
  Compute $\text{UCB}_{t+1}(g^\pi)$ for all $\pi$ based on $\text{UCB}_t(g^\pi)$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

**end**

---

⚠ *3 major issues:*

■ *Upper confidence bounds*: construct $\text{UCB}_t(g^\pi)$ with unknown dynamics?  ✔

■ *Computational complexity*: exponential number of policies

■ *Frequent policy update*: inefficient exploration

# Gain Optimism: Implementation

---

**Tentative algorithm**

---

Observe state $s_1$
**for** $t = 1, 2, \ldots$ **do**
  $\quad$ *Compute* $\pi_t \leftarrow \arg \max_{\pi} UCB_t(g^\pi)$
  $\quad$ Take action $a_t = \pi_t(s_t)$
  $\quad$ Observe reward $r_t$ and next state $s_{t+1}$
  $\quad$ Compute $\text{UCB}_{t+1}(g^\pi)$ for all $\pi$ based on $\text{UCB}_t(g^\pi)$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$
**end**

---

⚠️ *3 major issues:*

- 🟧 *Upper confidence bounds*: construct $\text{UCB}_t(g^\pi)$ with unknown dynamics? ✔
- 🟧 *Computational complexity*: exponential number of policies
- 🟧 *Frequent policy update*: inefficient exploration

# Gain Optimism: Implementation

---

**Tentative algorithm**

---

Observe state $s_1$

**for** $t = 1, 2, \ldots$ **do**

$\quad$ *Compute* $\pi_t \leftarrow \arg\max_\pi \left\{ \max_{M \in \mathcal{M}_t} \boldsymbol{g}_M^{\boldsymbol{\pi}} \right\}$

$\quad$ Take action $a_t = \pi_t(s_t)$

$\quad$ Observe reward $r_t$ and next state $s_{t+1}$

$\quad$ Compute $\mathsf{UCB}_{t+1}(g^\pi)$ for all $\pi$ based on $\mathsf{UCB}_t(g^\pi)$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

**end**

---

⚠ *3 major issues:*

■ *Upper confidence bounds*: construct $\mathsf{UCB}_t(g^\pi)$ with unknown dynamics? ✔

■ *Computational complexity*: exponential number of policies

■ *Frequent policy update*: inefficient exploration

# Gain Optimism: Implementation

---

**Tentative algorithm**

---

Observe state $s_1$

**for** $t = 1, 2, \ldots$ **do**

> *Compute* $\pi_t \leftarrow \arg\max\limits_{\pi} \left\{ \max\limits_{M \in \mathcal{M}_t} \boldsymbol{g}_M^{\boldsymbol{\pi}} \right\}$
>
> Take action $a_t = \pi_t(s_t)$
>
> Observe reward $r_t$ and next state $s_{t+1}$
>
> Compute $\mathsf{UCB}_{t+1}(g^\pi)$ for all $\pi$ based on $\mathsf{UCB}_t(g^\pi)$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

**end**

---

⚠ *3 major issues:*

■ *Upper confidence bounds*: construct $\mathsf{UCB}_t(g^\pi)$ with unknown dynamics? ✔

    ■ How to efficiently *compute* $\max\limits_{M \in \mathcal{M}_t} g_M^\pi$ for every $\pi$?

■ *Computational complexity*: exponential number of policies

■ *Frequent policy update*: inefficient exploration

# Gain Optimism: Implementation

---

**Tentative algorithm**

---

Observe state $s_1$

**for** $t = 1, 2, \ldots$ **do**

$\quad$ *Compute* $\pi_t \leftarrow \arg\max\limits_{\pi} \left\{ \max\limits_{M \in \mathcal{M}_t} \boldsymbol{g_M^\pi} \right\}$

$\quad$ Take action $a_t = \pi_t(s_t)$

$\quad$ Observe reward $r_t$ and next state $s_{t+1}$

$\quad$ Compute $\mathsf{UCB}_{t+1}(g^\pi)$ for all $\pi$ based on $\mathsf{UCB}_t(g^\pi)$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

**end**

---

⚠️ *3 major issues:*

■ *Upper confidence bounds*: construct $\mathsf{UCB}_t(g^\pi)$ with unknown dynamics? ✔

$\quad$ ■ How to efficiently *compute* $\max\limits_{M \in \mathcal{M}_t} g_M^\pi$ for every $\pi$?

■ *Computational complexity*: exponential number of policies

■ *Frequent policy update*: inefficient exploration

# Extended MDP
[Strehl and Littman, 2008, Jaksch et al., 2010]

---

## Theorem (Bounded parameter MDP $\iff$ Extended MDP)

Let $\mathcal{M}_t^+ := \langle \mathcal{S}, \mathcal{A}_t^+, r^+, p^+ \rangle$ be an *extended* MDP such that

$$\mathcal{A}_t^+(s) = \mathcal{A}(s) \times B_t^r(s, a) \times B_t^p(s, a)$$

with $a^+ = (a, r, p) \in \mathcal{A}_t^+(s)$, $r^+(s, a^+) = r$, $p^+(\cdot|s, a^+) = p$.

> Continuous **compact** action space

Then the optimal gain of $\mathcal{M}_t^+$ satisfies

$$g_{\mathcal{M}_t^+}^* := \max_\pi \left\{ \max_{M \in \mathcal{M}_t} g_M^\pi \right\}$$

Let $\pi_t^+ = \arg \max_\pi g_{\mathcal{M}_t^+}^\pi$, then

$$\pi_t = \arg \max_\pi \left\{ \max_{M \in \mathcal{M}_t} g_M^\pi \right\} \text{ s.t. } \pi_t(s) = \pi_t^+(s)[a]$$

# Extended MDP
[Strehl and Littman, 2008, Jaksch et al., 2010]

---

**Theorem (Bounded parameter MDP $\iff$ Extended MDP)**

Let $\mathcal{M}_t^+ := \langle \mathcal{S}, \mathcal{A}_t^+, r^+, p^+ \rangle$ be an *extended* MDP such that

$$\mathcal{A}_t^+(s) = \mathcal{A}(s) \times B_t^r(s,a) \times B_t^p(s,a)$$

with $a^+ = $ _Abuse of notation_: $\mathcal{M}_t$ denotes the extended MDP **ompact** pace

Then the optimal gain of $\mathcal{M}_t^+$ satisfies

$$g_{\mathcal{M}_t^+}^* := \max_\pi \left\{ \max_{M \in \mathcal{M}_t} g_M^\pi \right\}$$

Let $\pi_t^+ = \arg\max_\pi g_{\mathcal{M}_t^+}^\pi$, then

$$\pi_t = \arg\max_\pi \left\{ \max_{M \in \mathcal{M}_t} g_M^\pi \right\} \; \text{s.t.} \; \pi_t(s) = \pi_t^+(s)[a]$$

---

# Extended Value Iteration

Value iteration on $\mathcal{M}_t$

$$v_{n+1}(s) = \mathcal{L}_t v_n(s) = \max_{(a,r,p)\in\mathcal{A}(s)\times B_t^r(s,a)\times B_t^p(s,a)} \left\{ r + p^\mathsf{T} v_n \right\}$$

$$= \max_{a\in\mathcal{A}(s)} \left\{ \max_{r\in B_t^r(s,a)} r + \max_{p\in B_t^p(s,a)} p^\mathsf{T} v_n \right\}$$

$$= \max_{a\in\mathcal{A}(s)} \left\{ \widehat{r}_t(s,a) + \beta_t^r(s,a) + \max_{p\in B_t^p(s,a)} p^\mathsf{T} v_n \right\}$$

$\pi_t = $ *Greedy policy* w.r.t. $v_n$

# Gain Optimism: Implementation

---

**Tentative algorithm**

Observe state $s_1$

**for** $t = 1, 2, \ldots$ **do**

    *Compute* $\pi_t \leftarrow \arg\max_{\pi} \left\{ \max_{M \in \mathcal{M}_t} g_M^\pi \right\}$

    Take action $a_t = \pi_t(s_t)$

    Observe reward $r_t$ and next state $s_{t+1}$

    Compute $\mathsf{UCB}_{t+1}(g^\pi)$ for all $\pi$ based on $\mathsf{UCB}_t(g^\pi)$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

**end**

---

⚠ *3 major issues:*

■ *Upper confidence bounds*: construct $\mathsf{UCB}_t(g^\pi)$ with unknown dynamics ✔

    ■ How to efficiently *compute* $\max_{M \in \mathcal{M}_t} g_M^\pi$ for every $\pi$? ✔

■ *Computational complexity*: exponential number of policies ✔

■ *Frequent policy update*: inefficient exploration

# Gain Optimism: Implementation

---

Tentative algorithm

---

Observe state $s_1$

**for** $t = 1, 2, \ldots$ **do**

    *Compute* $\pi_t \leftarrow \arg\max_\pi \left\{ \max_{M \in \mathcal{M}_t} g_M^\pi \right\}$

    Take action $a_t = \pi_t(s_t)$

    Observe reward $r_t$ and next state $s_{t+1}$

    Compute $\mathsf{UCB}_{t+1}(g^\pi)$ for all $\pi$ based on $\mathsf{UCB}_t(g^\pi)$ and $\langle s_t, a_t, r_t, s_{t+1} \rangle$

**end**

---

⚠ *3 major issues:*

■ *Upper confidence bounds*: construct $\mathsf{UCB}_t(g^\pi)$ with unknown dynamics ✔

    ■ How to efficiently *compute* $\max_{M \in \mathcal{M}_t} g_M^\pi$ for every $\pi$? ✔

■ *Computational complexity*: exponential number of policies ✔

■ *Frequent policy update*: inefficient exploration

# Optimism: Frequency of Policy Updates

> **Proposition** [Ortner, 2010]
>
> There exists an MDP s.t.
>
> $$\Omega(T) \text{ number of policy updates} \implies \textit{linear regret}.$$

$$\implies \quad o(T) \text{ number of policy updates}$$

# Final Algorithm: UCRL2

Initialize $t \leftarrow 1$
Observe state $s_1$
Initialize empirical means $\widehat{r}_1 = r_{\max}$ and $\widehat{p}_1 = (1/S, \ldots, 1/S)^\mathsf{T}$
Initialize visit counts $N_1 = 0$
**for** *episodes* $k = 1, 2, \ldots$ **do**

    Set $t_k \leftarrow t$
    Build extended MDP $\mathcal{M}_k := \mathcal{M}_{t_k}$
    Using EVI, compute *optimistic policy* $\pi_k$ and $(h_k, g_k) \in \mathbb{R}^S \times [0, r_{\max}]$ such that

$$\boxed{\mathcal{L}_{\mathcal{M}_k} h_k = \mathcal{L}_{\mathcal{M}_k}^{\pi_k} h_k = h_k + g_k e} \quad \text{with} \quad g_k = g_{\mathcal{M}_k}^\star \geq g_{M^\star}^\star$$

    **while** $N_t(s_t, a_t) < \max\{1, \ N_{t_k}(s_t, a_t)\}$ **do**

        Take action $a_t = \pi_k(s_t)$
        Observe reward $r_t$ and next state $s_{t+1}$
        Compute new empirical means $\widehat{r}_{t+1}(s_t, a_t)$ and $\widehat{p}_{t+1}(\cdot|s_t, a_t)$
        Compute new visit count $N_{t+1}(s_t, a_t)$
        $t \leftarrow t + 1$
    **end**
**end**

# Final Algorithm: UCRL2

Initialize $t \leftarrow 1$
Observe state $s_1$
Initialize empirical means $\widehat{r}_1 = r_{\max}$ and $\widehat{p}_1 = (1/S, \ldots, 1/S)^{\mathsf{T}}$
Initialize visit counts $N_1 = 0$
**for** *episodes* $k = 1, 2, \ldots$ **do**

 Set $t_k \leftarrow t$
 Build extended MDP $\mathcal{M}_k := \mathcal{M}_{t_k}$
 Using EVI, compute *optimistic policy* $\pi_k$ and $(h_k, g_k) \in \mathbb{R}^S \times [0, r_{\max}]$ such that

$$\boxed{\mathcal{L}_{\mathcal{M}_k} h_k = \mathcal{L}_{\mathcal{M}_k}^{\pi_k} h_k = h_k + g_k e} \quad \text{with} \quad g_k = g^\star_{\mathcal{M}_k} \geq g^\star_{M^\star}$$

*Optimism*

 **while** $N_t(s_t, a_t) < \max\{1, N_{t_k}(s_t, a_t)\}$ **do**

  Take action $a_t = \pi_k(s_t)$
  Observe reward $r_t$ and next state $s_{t+1}$
  Compute new empirical means $\widehat{r}_{t+1}(s_t, a_t)$ and $\widehat{p}_{t+1}(\cdot | s_t, a_t)$
  Compute new visit count $N_{t+1}(s_t, a_t)$
  $t \leftarrow t + 1$
 **end**
**end**

# Final Algorithm: UCRL2

Initialize $t \leftarrow 1$

Observe state $s_1$

Initialize empirical means $\widehat{r}_1 = r_{\max}$ and $\widehat{p}_1 = (1/S, \ldots, 1/S)^{\mathsf{T}}$

Initialize visit counts $N_1 = 0$

**for** *episodes* $k = 1, 2, \ldots$ **do**

  Set $t_k \leftarrow t$

  Build extended MDP $\mathcal{M}_k := \mathcal{M}_{t_k}$

  Using EVI, compute *optimistic policy* $\pi_k$ and $(h_k, g_k) \in \mathbb{R}^S \times [0, r_{\max}]$ such that

$$\boxed{\mathcal{L}_{\mathcal{M}_k} h_k = \mathcal{L}_{\mathcal{M}_k}^{\pi_k} h_k = h_k + g_k e} \quad \text{with} \quad g_k = g_{\mathcal{M}_k}^\star \geq g_{M^\star}^\star$$

> Bellman equation in $\mathcal{M}_k$

> Optimism

  **while** $N_t(s_t, a_t) < \max\{1, \ N_{t_k}(s_t, a_t)\}$ **do**

    Take action $a_t = \pi_k(s_t)$

    Observe reward $r_t$ and next state $s_{t+1}$

    Compute new empirical means $\widehat{r}_{t+1}(s_t, a_t)$ and $\widehat{p}_{t+1}(\cdot | s_t, a_t)$

    Compute new visit count $N_{t+1}(s_t, a_t)$

    $t \leftarrow t + 1$

  **end**

**end**

# Final Algorithm: UCRL2

Initialize $t \leftarrow 1$

Observe state $s_1$

Initialize empirical means $\widehat{r}_1 = r_{\max}$ and $\widehat{p}_1 = (1/S, \ldots, 1/S)^{\mathsf{T}}$

Initialize visit counts $N_1 = 0$

**for** *episodes* $k = 1, 2, \ldots$ **do**

    Set $t_k \leftarrow t$

    Build extended MDP $\mathcal{M}_k := \mathcal{M}_{t_k}$

    Using EVI, compute *optimistic policy* $\pi_k$ and $(h_k, g_k) \in \mathbb{R}^S \times [0, r_{\max}]$ such that

$$\boxed{\mathcal{L}_{\mathcal{M}_k} h_k = \mathcal{L}_{\mathcal{M}_k}^{\pi_k} h_k = h_k + g_k e} \quad \text{with} \quad g_k = g_{\mathcal{M}_k}^{\star} \geq g_{M^\star}^{\star}$$

*Bellman equation in $\mathcal{M}_k$*

*Optimism*

    **while** $N_t(s_t, a_t) < \max\{1, \ N_{t_k}(s_t, a_t)\}$ **do**

        Take action $a_t = \pi_k(s_t)$

        Observe reward $r_t$ and next state $s_{t+1}$

        Compute new empirical means $\widehat{r}_{t+1}(s_t, a_t)$ and $\widehat{p}_{t+1}(\cdot | s_t, a_t)$

        Compute new visit count $N_{t+1}(s_t, a_t)$

        $t \leftarrow t + 1$

    **end**

**end**

*Stopping condition of an episode*

# UCRL2: Regret Guarantees

## Theorem (Thm.2 of [Jaksch et al., 2010])

*There exists a numerical constant $\beta > 0$ such that in any communicating MDP $M^\star = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, with probability at least $1 - \delta$, UCRL2 suffers a regret bounded as*

$$\forall T \geq 1, \ R(T, M^\star, \text{UCRL2}) \leq \beta \cdot r_{\max} DS \sqrt{AT \log\left(\frac{T}{\delta}\right)}$$

Fruit, Lazaric and Pirotta

# UCRL2: Regret Guarantees

> **Theorem** (Thm.2 of [Jaksch et al., 2010])
>
> *There exists a numerical constant $\beta > 0$ such that in any communicating MDP $M^\star = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, with probability at least $1 - \delta$, UCRL2 suffers a regret bounded as*
>
> $$\forall T \geq 1, \; R(T, M^\star, \text{UCRL2}) \leq \beta \cdot r_{\max} DS \sqrt{AT \log\left(\frac{T}{\delta}\right)}$$

Comparison to lower bound

$$\overline{R}(T, M^\star, \text{UCRL}) \geq 0.015 \sqrt{DSAT}$$

# UCRL2: Regret Guarantees

> **Theorem** (Thm.2 of [Jaksch et al., 2010])
>
> *There exists a numerical constant $\beta > 0$ such that in any communicating MDP $M^\star = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, with probability at least $1 - \delta$, UCRL2 suffers a regret bounded as*
>
> $$\forall T \geq 1, \ R(T, M^\star, \mathsf{UCRL2}) \leq \beta \cdot r_{\max} DS \sqrt{AT \log \left( \frac{T}{\delta} \right)}$$

Comparison to lower bound

$$\overline{R}(T, M^\star, \mathsf{UCRL}) \geq 0.015 \sqrt{DSAT}$$

- Can the gap between upper and lower bound be closed?  👍 More on this later

# UCRL2: Regret Guarantees (cont'd.)

## Theorem (Thm.4 of [Jaksch et al., 2010])

*There exists a numerical constant $\beta > 0$ such that in any ergodic MDP*
*$M^\star = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, for all $T \geq 1$, UCRL2 (with $\delta = 1/T$) suffers a regret bounded as*

$$\overline{R}(T, M^\star, \text{UCRL2}) \leq \beta \cdot r_{\max} \frac{D^2 S^2 A \log(T)}{\delta_g^\star} + \text{Big constant independent of } T$$

*with*

- $\delta_g^\star := g_{M^\star}^\star - \max\limits_{s \in \mathcal{S}, \pi} \left\{ g_{M^\star}^\pi(s) < g_M^\star \right\} \quad \sim$ *"gap in gain"*

# UCRL2: Regret Guarantees (cont'd.)

**Theorem** (Thm.4 of [Jaksch et al., 2010])

*There exists a numerical constant $\beta > 0$ such that in any ergodic MDP $M^\star = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, for all $T \geq 1$, UCRL2 (with $\delta = 1/T$) suffers a regret bounded as*

$$\overline{R}(T, M^\star, \text{UCRL2}) \leq \beta \cdot r_{\max} \frac{D^2 S^2 A \log(T)}{\delta_g^\star} + \text{Big constant independent of } T$$

*with*

- $\delta_g^\star := g_{M^\star}^\star - \max\limits_{s \in \mathcal{S}, \pi} \left\{ g_{M^\star}^\pi(s) < g_M^\star \right\} \quad \sim \text{"gap in gain"}$

Comparison to lower bound

$$\liminf_{T \to \infty} \frac{\overline{R}(T, M^\star, \mathfrak{A})}{\log T} \geq K_{M^\star}, \text{ with } K_{M^\star} \lesssim \frac{D^2 S A}{\min\limits_{s,a} \delta_{M^\star}^\star(s, a)}$$

# UCRL2: Regret Guarantees (cont'd.)

**Theorem** (Thm.4 of [Jaksch et al., 2010])

*There exists a numerical constant $\beta > 0$ such that in any ergodic MDP $M^\star = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$, for all $T \geq 1$, UCRL2 (with $\delta = 1/T$) suffers a regret bounded as*

$$\overline{R}(T, M^\star, \text{UCRL2}) \leq \beta \cdot r_{\max} \frac{D^2 S^2 A \log(T)}{\delta_g^\star} + \text{Big constant independent of } T$$

*with*

- $\delta_g^\star := g_{M^\star}^\star - \max_{s \in \mathcal{S}, \pi} \left\{ g_{M^\star}^\pi(s) < g_M^\star \right\}$   $\sim$ *"gap in gain"*

how do they compare?

Comparison to lower bound

$$\liminf_{T \to \infty} \frac{\overline{R}(T, M^\star, \mathfrak{A})}{\log T} \geq K_{M^\star}, \text{ with } K_{M^\star} \lesssim \frac{D^2 S A}{\min_{s,a} \delta_{M^\star}^\star(s, a)}$$

# Qualitative Regret Shape



$R(T, M^{\star}, \text{UCRL2})$

$\mathcal{O}\left(\dfrac{D^2 S^2 A}{\delta_g^{\star}} \log(T)\right)$

$T$

$\mathcal{O}\left(DS\sqrt{AT \log(T)}\right)$

Regret upper-bound

$0$    $T_1$    $T_2$    $T$

*illustrative plot

# Regret Bound of UCRL2: Proof Sketch

1 $$R(T, M^\star, \texttt{UCRL2}) = \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} g_{M^\star}^\star - r(s_t, a_t) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} g_k - r(s_t, a_t)$$

Split in episodes

Optimism: $g_k \geq g_{M^\star}^\star$

# Regret Bound of UCRL2: Proof Sketch

1 $R(T, M^\star, \mathsf{UCRL2}) = \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} g_{M^\star}^\star - r(s_t, a_t) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} g_k - r(s_t, a_t)$

2 $\sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} g_k = \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} r_k(s_t, a_t) + p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

Bellman equation $(a_t = \pi_k(s_t))$:
$$L_{\mathcal{M}_k}^{\pi_k} h_k(s_t) = h_k(s_t) + g_k$$

# Regret Bound of UCRL2: Proof Sketch

1-2 $$R(T, M^\star, \mathsf{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} r_k(s_t, a_t) - r(s_t, a_t) + p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$$

# Regret Bound of UCRL2: Proof Sketch

$$1\text{-}2 \quad R(T, M^\star, \texttt{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} r_k(s_t, a_t) - r(s_t, a_t) + p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$$

Assumption: true reward is known $r = r_k$

# Regret Bound of UCRL2: Proof Sketch

$$1\text{-}2 \quad R(T, M^\star, \mathsf{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$$

# Regret Bound of UCRL2: Proof Sketch

$$\boxed{\text{1-2}} \quad R(T, M^\star, \mathtt{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^{\mathsf{T}} h_k - h_k(s_t)$$

$$\boxed{\text{3}} \quad p_k(\cdot|s_t, a_t)^{\mathsf{T}} h_k - h_k(s_t) = \Big(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t)\Big)^{\mathsf{T}} h_k + p(\cdot|s_t, a_t)^{\mathsf{T}} h_k - h_k(s_t)$$

# Regret Bound of UCRL2: Proof Sketch

**1-2** $R(T, M^\star, \text{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**3** $p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \left( p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k + p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**4** $\sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_{t+1})$

$$+ \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} h_k(s_{t+1}) - h_k(s_t)$$

# Regret Bound of UCRL2: Proof Sketch

**1-2** $R(T, M^\star, \texttt{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**3** $p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t)\right)^\mathsf{T} h_k + p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**4** $\sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) \implies \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_{t+1})$

Martingale Difference Sequence
(Azuma's inequality)

$+ \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} h_k(s_{t+1}) - h_k(s_t)$

# Regret Bound of UCRL2: Proof Sketch

$$1\text{-}2 \quad R(T, M^\star, \texttt{UCRL2}) \le \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$$

$$3 \quad p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \left( p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k + p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$$

$$4 \quad \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) \lesssim \sup_k \{\mathsf{sp}(h_k)\} \sqrt{T \log(T/\delta)}$$

$$+ \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} h_k(s_{t+1}) - h_k(s_t)$$

# Regret Bound of UCRL2: Proof Sketch

**1-2** $R(T, M^\star, \mathsf{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**3** $p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \left(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t)\right)^\mathsf{T} h_k + p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**4** $\sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) \lesssim \sup_{k} \{\mathsf{sp}(h_k)\} \sqrt{T \log(T/\delta)}$

$$+ \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} h_k(s_{t+1}) - h_k(s_t) \longleftarrow \text{Telescopic sum}$$

# Regret Bound of UCRL2: Proof Sketch

$$\boxed{\text{1-2}} \quad R(T, M^\star, \mathsf{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$$

$$\boxed{3} \quad p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \Big(p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t)\Big)^\mathsf{T} h_k + p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$$

$$\boxed{4} \quad \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) \lesssim \sup_k \{\mathsf{sp}(h_k)\} \sqrt{T \log(T/\delta)}$$

$$+ \; m \sup_k \{\mathsf{sp}(h_k)\}$$

Number of episodes
(stopping condition)

# Regret Bound of UCRL2: Proof Sketch

**1-2** $R(T, M^\star, \mathsf{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**3** $p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \left( p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k + p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**4** $\sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) \lesssim \sup_k \{ \mathsf{sp}(h_k) \} \sqrt{T \log (T/\delta)}$

$$+ SA \log(T) \sup_k \{ \mathsf{sp}(h_k) \}$$

# Regret Bound of UCRL2: Proof Sketch

1-2 $R(T, M^\star, \text{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

3 $p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \left( p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k + p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

4 $\sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) \lesssim \sup_k \{\text{sp}(h_k)\} \sqrt{T \log(T/\delta)}$

$$+ SA \log(T) \sup_k \{\text{sp}(h_k)\}$$

$$\boxed{\text{sp}(h_k) \leq r_{\max} D}$$ [Bartlett and Tewari, 2009, Jaksch et al., 2010]

# Regret Bound of UCRL2: Proof Sketch

1-2 $\quad R(T, M^\star, \mathsf{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

3 $\quad p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \left( p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k + p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

4 $\quad \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) \lesssim r_{\max} D \sqrt{T \log (T/\delta)} + r_{\max} D S A \log(T)$

$\mathsf{sp}\big(h_k\big) \leq r_{\max} D$ [Bartlett and Tewari, 2009, Jaksch et al., 2010]

# Regret Bound of UCRL2: Proof Sketch

1-2 $\quad R(T, M^\star, \mathsf{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

3 $\quad p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \left( p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k + p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

4 $\quad \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) \lesssim r_{\max} D \sqrt{T \log(T/\delta)} + r_{\max} D S A \log(T)$

5 $\quad \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} \left( p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k = \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} \underbrace{\left( p_k(\cdot|s_t, a_t) - \widehat{p}_k(\cdot|s_t, a_t) \right)^\mathsf{T} h_k}_{\leq \mathsf{sp}(h_k) \beta_k^p(s,a)}$

$$+ \underbrace{\left( \widehat{p}_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k}_{\leq \mathsf{sp}(h_k) \beta_k^p(s,a)}$$

# Regret Bound of UCRL2: Proof Sketch

**1-2** $R(T, M^\star, \text{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**3** $p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \left( p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k + p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**4** $\sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) \lesssim r_{\max} D \sqrt{T \log(T/\delta)} + r_{\max} D S A \log(T)$

**5** $\sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} \left( p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k \lesssim r_{\max} D \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} \sqrt{\frac{S \log(T/\delta)}{N_{t_k}(s_t, a_t)}}$

# Regret Bound of UCRL2: Proof Sketch

**1-2** $R(T, M^\star, \text{UCRL2}) \leq \sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**3** $p_k(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) = \left( p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k + p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t)$

**4** $\sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} p(\cdot|s_t, a_t)^\mathsf{T} h_k - h_k(s_t) \lesssim r_{\max} D \sqrt{T \log(T/\delta)} + r_{\max} DSA \log(T)$

**5** $\sum_{k=1}^{m} \sum_{t=t_k}^{t_{k+1}-1} \left( p_k(\cdot|s_t, a_t) - p(\cdot|s_t, a_t) \right)^\mathsf{T} h_k \lesssim r_{\max} DS \sqrt{AT \log(T/\delta)}$

# Refined Confidence Bounds

- UCRL2 with *Bernstein bounds* (instead of Hoeffding/Weissman):
  - i see tutorial website

$$R(T, M^\star, \mathsf{UCRL2B}) = \mathcal{O}\left(\sqrt{D\Gamma SAT \log\left(\frac{T}{\delta}\right) \log(T)}\right)$$

- 👎 Still not matching the lower bound!
- 👍 For most MPDs: $\Gamma \ll S$

# Refined Confidence Bounds

- UCRL2 with *Bernstein bounds* (instead of Hoeffding/Weissman):
  - **i** see tutorial website

$$R(T, M^{\star}, \text{UCRL2B}) = \mathcal{O}\left(\sqrt{D\Gamma SAT \log\left(\frac{T}{\delta}\right) \log(T)}\right)$$

  - 👎 Still not matching the lower bound!
  - 👍 For most MPDs: $\Gamma \ll S$

- *Kullback-Leibler* UCRL [Filippi et al., 2010, Talebi and Maillard, 2018]:

$$R(T, M^{\star}, \text{UCRL-KL}) = \mathcal{O}\left(\sqrt{\underbrace{\sum_{s,a} \mathbb{V}_{X \sim p^{\star}(\cdot|s,a)}\left(h_{M^{\star}}^{\star}(X)\right)}_{\leq D^2 SA} ST \log\left(\frac{T}{\delta}\right)} + D\sqrt{T}\right)$$

  - 👎 Only for ergodic MDPs!

# Infinite Diameter (weakly communicating MDPs)

- *Known* bound on the optimal bias span $C \geq \mathsf{sp}(h^{\star}_{M^{\star}})$
  [Bartlett and Tewari, 2009, Fruit et al., 2018b]

$$R(T, M^{\star}, \mathsf{SCAL}) = \mathcal{O}\left(\sqrt{C\Gamma SAT \log\left(\frac{T}{\delta}\right)} \log(T)\right)$$

👎 Requires prior knowledge!

# Infinite Diameter (weakly communicating MDPs)

- *Known* bound on the optimal bias span $C \geq \mathsf{sp}(h^\star_{M^\star})$
  [Bartlett and Tewari, 2009, Fruit et al., 2018b]

$$R(T, M^\star, \mathsf{SCAL}) = \mathcal{O}\left(\sqrt{C\Gamma SAT \log\left(\frac{T}{\delta}\right)} \log(T)\right)$$

  👎 Requires prior knowledge!

- No prior knowledge: TUCRL [Fruit et al., 2018a]:

$$R(T, M^\star, \mathsf{SCAL}) = \mathcal{O}\left(\sqrt{D_{\mathsf{com}} S_{\mathsf{com}} \Gamma AT \log\left(\frac{T}{\delta}\right)} \log(T)\right)$$

  👎 Never achieves *logarithmic* regret! Intrinsic limitation of the setting!

# Open Questions

1. *Tightness of minimax $\mathcal{O}(\sqrt{T})$ regret bounds for infinite horizon problems*
   - Dependency on $\Gamma$: regret + sample complexity bounds?
   - Analysis not tight *vs.* change in the algorithm?
   - Lower bound not tight?

2. *Finite time logarithmic upper and lower regret bounds*
   - Non-asymptotic lower bounds
   - Tighter analysis of UCRL-like algorithms? New algorithms?

Fruit, Lazaric and Pirotta

# Posterior Sampling
a.k.a. Thompson Sampling [Thompson, 1933]

Keep Bayesian posterior for the *unknown* MDP

👍 A sample from the posterior is used as an
estimate of the unknown MDP

Set of MDPs

Exploration

Few samples $\implies$ uncertainty in the
estimate

More samples $\implies$ posterior concentrates
on the true MDP

Exploitation

Posterior
distribution $\mu_t$

# History: PS for Regret Minimization in RL

FH: finite-horizon
AR: average reward

# Posterior Sampling

$t \leftarrow 1$

**for** *episode* $k = 1, 2, \ldots$ **do**

    $t_k \leftarrow t$

    $M_k \sim \mu_{t_k}$

    $\pi_k \in \arg \max_{\pi} \{g_{M_k}^\pi\}$

    **while** *not enough knowledge* **do**

        Take action $a_t \sim \pi_k(\cdot | s_t)$

        Observe reward $r_t$ and next state $s_{t+1}$

        Compute $\mu_{t+1}$ based on $\mu_t$ and

        $(s_t, a_t, r_t, s_{t+1})$

        $t \leftarrow t + 1$

    **end**

**end**

# Posterior Sampling

$t \leftarrow 1$
**for** *episode* $k = 1, 2, \dots$ **do**
$\quad t_k \leftarrow t$

$\quad M_k \sim \mu_{t_k}$
$\quad \pi_k \in \arg \max_\pi \{g^\pi_{M_k}\}$

$\quad$ **while** *not enough knowledge* **do**
$\quad\quad$ Take action $a_t \sim \pi_k(\cdot|s_t)$
$\quad\quad$ Observe reward $r_t$ and next state $s_{t+1}$
$\quad\quad$ Compute $\mu_{t+1}$ based on $\mu_t$ and
$\quad\quad$ $(s_t, a_t, r_t, s_{t+1})$
$\quad\quad$ $t \leftarrow t + 1$
$\quad$ **end**
**end**

Prior distribution:
$$\forall \Theta, \quad \mathbb{P}(M^* \in \Theta) = \mu_1(\Theta)$$

Posterior distribution:
$$\forall \Theta, \quad \mathbb{P}(M^* \in \Theta | H_t, \mu_1) = \mu_t(\Theta)$$

Priors
- Dirichlet (transitions)
- Beta, Normal-Gamma, etc. (rewards)

# Bayesian Regret

$$R^B(T, \mu_1, \mathfrak{A}) = \mathbb{E}_{M^\star \sim \mu_1}\Big[\ \underbrace{\overline{R}(T, M^\star, \mathfrak{A})}_{:=\mathbb{E}\big[R(T, M^\star, \mathfrak{A})\big]}\ \Big] = \mathbb{E}\left[\sum_{t=1}^{T} g_{M^\star}^\star - r(s_t, a_t)\right]$$

# TSDE: Thompson Sampling with Dynamic Episodes
[Ouyang et al., 2017b]

*Episode length $l_k = t_{k+1} - t_k$ is dynamically determined* by

**1** Doubling of visits (stochastic)

**2** Increasing length of previous episode by one (deterministic)

$$t_{k+1} = \min\left\{ t > t_k \ : \ \underbrace{\exists(s,a), N_t(s,a) > 2N_{t_k}(s,a)}_{(ST1)} \ \text{or} \ \underbrace{t > t_k + l_{k-1}}_{(ST2)} \right\}$$

☞ (ST2) is $\sigma(H_{t_k})$-measurable

$l_k \leq l_{k-1} + 1$

# TSDE: Regret Guarantees

**Theorem** ([Ouyang et al., 2017b])

*There exists a numerical constant $\beta > 0$ such that for any prior $\mu_1$ whose support is a subset of communicating MDPs,* TSDE *suffers a regret bounded as*

$$\forall T \geq 1, \quad R^B(T, \mu_1, \text{TSDE}) \leq \beta \cdot \left( CS\sqrt{AT \log(AT)} \right)$$

*where*

$$\mu_1 \quad \text{is such that} \quad \sup_{M^\star \sim \mu_1} \left\{ sp(h^\star_{M^\star}) \right\} \leq C < +\infty \qquad \text{(ASM-SP)}$$

# Proof Step 1: Regret Decomposition

☞ The support of the prior $\mu_1$ is a subset of communicating MDPs
   $M_k$ is communicating and optimality equation (i.e., constant gain)

$$R^B(T, \mu_1, \mathsf{TSDE}) \leq \underbrace{T\mathbb{E}\left[g^\star_{M^\star}\right] - \mathbb{E}\left[\sum_{k=1}^{k_T} l_k \; g^\star_{M_k}\right]}_{R_g}$$

$$+ \; \mathbb{E}\left[\sum_{k=1}^{k_T} \sum_{t=t_k}^{t_{k+1}-1} (h_k(s_t) - h_k(s_{t+1}))\right]$$

$$+ \; \mathbb{E}\left[\sum_{k=1}^{k_T} \sum_{t=t_k}^{t_{k+1}-1} (h_k(s_{t+1}) - p_k(\cdot|s_t, a_t)^\mathsf{T} h_k) + r_k(s_t, a_t) - r(s_t, a_t)\right]$$

# Proof Step 1: Regret Decomposition

👉 The support of the prior $\mu_1$ is a subset of communicating MDPs
$M_k$ is communicating and optimality equation (i.e., constant gain)

$$R^B(T, \mu_1, \mathsf{TSDE}) \leq \underbrace{T\mathbb{E}\left[g^\star_{M^\star}\right] - \mathbb{E}\left[\sum_{k=1}^{k_T} l_k \; g^\star_{M_k}\right]}_{R_g}$$

$$+ \; \mathbb{E}\left[\sum_{k=1}^{k_T} \sum_{t=t_k}^{t_{k+1}-1} (h_k(s_t) - h_k(s_{t+1}))\right]$$

Telescopic sum
+ span bound (ASM-SP)[†]

$$+ \; \mathbb{E}\left[\sum_{k=1}^{k_T} \sum_{t=t_k}^{t_{k+1}-1} (h_k(s_{t+1}) - p_k(\cdot|s_t, a_t)^\mathsf{T} h_k) + r_k(s_t, a_t) - r(s_t, a_t)\right]$$

[†] as in UCRL2

Fruit, Lazaric and Pirotta

# Proof Step 1: Regret Decomposition

👉 The support of the prior $\mu_1$ is a subset of communicating MDPs
   $M_k$ is communicating and optimality equation (i.e., constant gain)

$$R^B(T, \mu_1, \mathsf{TSDE}) \leq \underbrace{T\mathbb{E}\left[g_{M^\star}^\star\right] - \mathbb{E}\left[\sum_{k=1}^{k_T} l_k \ g_{M_k}^\star\right]}_{R_g}$$

$$+ \ \mathbb{E}\left[\sum_{k=1}^{k_T}\sum_{t=t_k}^{t_{k+1}-1} (h_k(s_t) - h_k(s_{t+1}))\right]$$

Telescopic sum
+ span bound (ASM-SP)[†]

Confidence sets[†]

$$+ \ \mathbb{E}\left[\sum_{k=1}^{k_T}\sum_{t=t_k}^{t_{k+1}-1} (h_k(s_{t+1}) - p_k(\cdot|s_t, a_t)^\mathsf{T} h_k) + r_k(s_t, a_t) - r(s_t, a_t)\right]$$

[†] as in UCRL2

# Proof Step 2: Bounding $R_g$

**Thompson Sampling Lemma** [Osband et al., 2013, Ouyang et al., 2017b]

Let $t_k$ be an almost surely finite $\sigma(H_{t_k})$-stopping time. For any measurable function $f$ and $\sigma(H_{t_k})$-measurable variable $X$

$$\mathbb{E}[f(M_k, X)|H_{t_k}] = \mathbb{E}[f(M^\star, X)|H_{t_k}]$$

# Proof Step 2: Bounding $R_g$

$$R_g = \mathbb{E}\left[\sum_{t=1}^{k_T} l_k \ g_{M^\star}^\star\right] - \mathbb{E}\left[\sum_{k=1}^{k_T} l_k \ g_{M_k}^\star\right]$$

# Proof Step 2: Bounding $R_g$

random duration of episode $k$
not $\sigma(H_{t_k})$-measurable

$$R_g = \mathbb{E}\left[\sum_{t=1}^{k_T} l_k \; g_{M^\star}^\star\right] - \mathbb{E}\left[\sum_{k=1}^{k_T} l_k \; g_{M_k}^\star\right]$$

# Proof Step 2: Bounding $R_g$

> random duration of episode $k$
> not $\sigma(H_{t_k})$-measurable

$$R_g = \mathbb{E}\left[\sum_{t=1}^{k_T} l_k \ g_{M^\star}^\star\right] - \mathbb{E}\left[\sum_{k=1}^{k_T} l_k \ g_{M_k}^\star\right]$$

$$\leq \mathbb{E}\left[\sum_{k=1}^{k_T}(l_{k-1}+1)\left(g_{M^\star}^\star - g_{M_k}^\star\right)\right] + \mathbb{E}\left[\sum_{k=1}^{k_T}(l_{k-1}+1-l_k)g_{M_k}^\star\right] \quad \left(\begin{array}{c}\text{by (ST2)} \\ l_k \leq l_{k+1}+1\end{array}\right)$$

# Proof Step 2: Bounding $R_g$

random duration of episode $k$
not $\sigma(H_{t_k})$-measurable

$$R_g = \mathbb{E}\left[\sum_{t=1}^{k_T} l_k \ g_{M^\star}^\star\right] - \mathbb{E}\left[\sum_{k=1}^{k_T} l_k \ g_{M_k}^\star\right]$$

$$\leq \mathbb{E}\left[\sum_{k=1}^{k_T}(l_{k-1}+1)\left(g_{M^\star}^\star - g_{M_k}^\star\right)\right] + \mathbb{E}\left[\sum_{k=1}^{k_T}(l_{k-1}+1-l_k)g_{M_k}^\star\right] \quad \begin{pmatrix} \text{by (ST2)} \\ l_k \leq l_{k+1}+1 \end{pmatrix}$$

$$\leq \qquad\qquad 0 \qquad\qquad + \ r_{\max}\mathbb{E}[k_T]$$

$t_k$ is a stopping time
$(l_{k-1} + 1)$ is $\sigma(H_{t_k})$-measurable
$\implies$ use TS lemma

Fruit, Lazaric and Pirotta

# Proof Step 2: Bounding $R_g$

random duration of episode $k$
not $\sigma(H_{t_k})$-measurable

$$R_g = \mathbb{E}\left[\sum_{t=1}^{k_T} l_k \ g_{M^\star}^\star\right] - \mathbb{E}\left[\sum_{k=1}^{k_T} l_k \ g_{M_k}^\star\right]$$

$$\leq \mathbb{E}\left[\sum_{k=1}^{k_T}(l_{k-1}+1)\left(g_{M^\star}^\star - g_{M_k}^\star\right)\right] + \mathbb{E}\left[\sum_{k=1}^{k_T}(l_{k-1}+1-l_k)g_{M_k}^\star\right] \quad \left(\begin{array}{c} \text{by (ST2)} \\ l_k \leq l_{k+1}+1 \end{array}\right)$$

$$\leq \qquad\qquad\qquad 0 \qquad\qquad + \ r_{\max}\mathbb{E}[k_T]$$

$t_k$ is a stopping time
$(l_{k-1}+1)$ is $\sigma(H_{t_k})$-measurable
$\implies$ use TS lemma

$$\sum_{k=1}^{k_T} l_{k-1} = l_0 + \sum_{k=1}^{k_T-1} l_k \leq T$$
$g_{M_k}^{\pi_k} \in [0, r_{\max}], \forall k$

# Proof Step 2: Bounding $R_g$

random duration of episode $k$
not $\sigma(H_{t_k})$-measurable

$$R_g = \mathbb{E}\left[\sum_{t=1}^{k_T} l_k \; g^{\star}_{M^{\star}}\right] - \mathbb{E}\left[\sum_{k=1}^{k_T} l_k \; g^{\star}_{M_k}\right]$$

$$\leq \mathbb{E}\left[\sum_{k=1}^{k_T}(l_{k-1}+1)\left(g^{\star}_{M^{\star}} - g^{\star}_{M_k}\right)\right] + \mathbb{E}\left[\sum_{k=1}^{k_T}(l_{k-1}+1-l_k)g^{\star}_{M_k}\right] \quad \left(\begin{array}{c} \text{by (ST2)} \\ l_k \leq l_{k+1}+1 \end{array}\right)$$

$$\leq \qquad\qquad 0 \qquad\qquad + \; r_{\max}\mathbb{E}[k_T]$$

$$\leq r_{\max}\sqrt{2SAT\log(T)} \qquad\qquad\qquad\qquad \text{([Ouyang et al., 2017b] Lem. 1)}$$

$t_k$ is a stopping time
$(l_{k-1}+1)$ is $\sigma(H_{t_k})$-measurable
$\implies$ use TS lemma

$$\sum_{k=1}^{k_T} l_{k-1} = l_0 + \sum_{k=1}^{k_T-1} l_k \leq T$$
$g^{\pi_k}_{M_k} \in [0, r_{\max}], \forall k$

# OPT-PSRL: Optimistic Posterior Sampling
[Agrawal and Jia, 2017]



frequentist regret

gain optimism

OFU

PSRL

1. Sample posterior $\psi = \widetilde{O}(S)$ times

$$p_{sa}^i \sim \mu_{t_k}(s,a), \quad i = 1, \dots, \psi$$

2. Solve $\mathcal{M}_k$ for $\pi_k$

$\mathcal{M}_k$ is an *discrete extended* MDP

$$\widetilde{p}(\cdot, s, a^i) = p_{s,a}^i, \qquad a^i \in \mathcal{A} \times \{1, \dots, \psi\}$$

$$g_{\mathcal{M}_k}^\star \geq g_{M^\star}^\star - \widetilde{O}\left(D\sqrt{SA/T}\right)$$

# OPT-PSRL: Regret Guarantees

**Theorem** ([Agrawal and Jia, 2017])

*For any communicating MDP $M$, with probability $1 - \delta$, there exist two constant $\alpha, \beta > 0$ such that, for any $T \geq \alpha D A \log^2(T/\delta)$, the regret of* Opt-PSRL *is bound by*

$$R(T, M^\star, \text{Opt-PSRL}) \leq \beta r_{\max} \cdot \left( DS \sqrt{AT \log\left(\frac{T}{\delta}\right)} + poly(S, A) DT^{1/4} \log\left(\frac{T}{\delta}\right) \right)$$

# Open Questions

**1** *The nature of bounded bias span assumption (Asm. ASM-SP)*

- Used in [Ouyang et al., 2017b, Theocharous et al., 2018]
- $\mathrm{supp}(\mu_1)$ is continuous, then $\sup_{M^\star}\{\mathsf{sp}(h^\star_{M^\star})\} = +\infty$ [e.g., Fruit et al. [2018a]]

**2** *Statistical efficiency of* PSRL

- Claimed efficient Bayesian or frequentist $\widetilde{O}(D\sqrt{SAT})$ regret bound
- Not supported by proofs, incorrect Lem. C.1 [Osband and Roy, 2016a] and Lem. C.2 [Agrawal and Jia, 2017] [**i** see tutorial website]

# History:  Asymptotic Regret Minimization



Agrawal [1990] ($\infty$)

Graves and Lai [1997] ($\infty$)

Ok et al. [2018] ($\infty$)

Burnetas and Katehakis [1997] ($\infty$)

Tewari and Bartlett [2007] ($\infty$)

# Asymptotic Lower-Bound

**Theorem** (Thm. 2, [Burnetas and Katehakis, 1997])

*Any algorithm $\mathfrak{A}$ s.t. $\overline{R}(T, M, \mathfrak{A}) = o(T^\alpha)$ for all $\alpha > 0$ and ergodic MDP $M$ should satisfy*

$$\forall (s, a) : \mathcal{M}_{M^\star}^{alt}(s, a), \quad \lim_{T \to \infty} \inf \frac{\mathbb{E}[N_T(s, a)]}{\log T} \geq \frac{1}{\inf_{M \in \mathcal{M}_{M^\star}^{alt}(s, a)} KL_{M^\star, M}(s, a)}$$

☞ Should be satisfied by optimal algorithms

   *necessary* to be uniformly good on all the possible *alternative* models

# BKIA: Burnetas-Katehakis Index Algorithm

[Burnetas and Katehakis, 1997]

**for** $t = 1, \ldots, T$ **do**

$$D_t(s) \leftarrow \{a \in \mathcal{A}(s) : N_t(s,a) \geq \log^2(N_t(s))\}$$
$$(g_t, h_t) \leftarrow \text{solve } \widehat{M}_t = \langle \mathcal{S}, D_t, \widehat{p}_t, r \rangle$$

**A** Solve empirical MDP $\widehat{M}_t$ on a restricted action set

**if** $\exists a \in \Pi^\star_{\widehat{M}_t}(s_t), \ N_t(s_t, a) \geq \log^2(N_t(s_t) + 1)$ **then**

$$a_t \in \arg\max_{a \in \mathcal{A}(s_t)}\{b_t(s, a; h_t)\}$$

**B** Select maximum index action

**else**

$$a_t \in \arg\min_{a \in \Pi^\star_{\widehat{M}_t}(s_t)} \{N_t(s, a)\}$$

**C** Force exploration of "underestimated" actions

**end**

Observe reward $r_t$ and next state $s_{t+1}$

**end**

# BKIA: Interpretation

**B** *Exploration & Exploitation*

Optimistic greedy

$$a_t \in \arg\max_{a \in \mathcal{A}}\{b_t(s_t, a)\}$$

$$b_t(s, a) = \sup_{q \in \Delta(\mathcal{S})} \left\{ L_q^a h^\star_{\widehat{M}_t}(s) \; : \; N_t(s, a) \, \mathsf{KL}(\widehat{p}_t(\cdot|s_t, a)\|q) \leq \log(t) \right\}$$

$$\textit{related to} \; - \inf_{M \in \mathcal{M}^{\mathsf{alt}}_{\widehat{M}_t}(s,a)} \left\{ \delta^\star_{\widehat{M}_t}(s, a) \; : \; N_t(s, a) \, \mathsf{KL}_{\widehat{M}_t, M}(s, a) \leq \log(t) \right\}$$

⚠️ A not so explicit way of controlling the lower bound

# BKIA: Interpretation

**B** *Exploration & Exploitation*



$$a_t \in \arg\max_{a \in \mathcal{A}} \{b_t(s_t, a)\}$$

Optimistic greedy

$$b_t(s, a) = \sup_{q \in \Delta(\mathcal{S})} \left\{ L_q^a h_{\widehat{M}_t}^\star(s) \ : \ N_t(s, a) \, \mathsf{KL}(\widehat{p}_t(\cdot|s_t, a)\|q) \le \log(t) \right\}$$

$$\text{related to} \quad - \inf_{M \in \mathcal{M}_{\widehat{M}_t}^{\mathsf{alt}}(s, a)} \left\{ \delta_{\widehat{M}_t}^\star(s, a) \ : \ N_t(s, a) \, \mathsf{KL}_{\widehat{M}_t, M}(s, a) \le \log(t) \right\}$$

⚠ A not so explicit way of controlling the lower bound

📔 Computing $b_t$ is similar to KL-UCB [Garivier and Cappé, 2011] for MAB.

# BKIA: Interpretation

C *Forced Exploration*

$$\text{when} \quad \forall a \in \Pi^{\star}_{\widehat{M}_t}(s_t), \ N_t(s_t, a) < \log^2(N_t(s_t) + 1)$$

■ BKIA prevents that *all* optimal actions *will become* under-explored

$$\implies a_t \in \Pi^{\star}_{\widehat{M}_t}(s_t)$$

👍 Asymptotic monotonic property

$$\mathbb{P}\left(g^{\star}_{M^{\star}(D_{t+1})} \geq g^{\star}_{M^{\star}(D_t)}\right) = 1 - o\left(\frac{1}{t}\right) \quad \text{as } t \to \infty$$

# BKIA: Regret Guarantees

**Theorem** (Thm. 1, [Burnetas and Katehakis, 1997])

*For any ergodic MDP $M^\star$, the expected regret of BKIA is upper bounded as*

$$\limsup_{T \to \infty} \frac{\overline{R}(T, M^\star, BKIA)}{\log T} \leq K_{M^\star}^\star$$

# BKIA: Regret Guarantees

**Theorem** (Thm. 1, [Burnetas and Katehakis, 1997])

*For any ergodic MDP $M^\star$, the expected regret of BKIA is upper bounded as*

$$\limsup_{T \to \infty} \frac{\overline{R}(T, M^\star, BKIA)}{\log T} \leq K^\star_{M^\star}$$

👍 OLP [Tewari and Bartlett, 2007] replaces the KL constraint with an $L_1$

# BKIA: Regret Proof

By [Prop. 1, [Burnetas and Katehakis, 1997]]

$$\overline{R}(T, M^\star, \mathfrak{A}) = \sum_s \sum_{a \notin \Pi^\star_{M^\star}(s)} \mathbb{E}\left[N_T(s, a))\right] \delta^\star_{M^\star}(s, a) + O(1), \qquad \text{as } T \to +\infty$$

We define $W_T^1$ s.t.

$$\mathbb{E}[N_T(s, a)] \leq \mathbb{E}[W_T^1(s, a, \varepsilon)] + o(\log T)$$

Ergodicity of MDP ($g$ and $h$ continuity) about $h^\star_{\widehat{M}_t} \to h^\star_{M^\star}$

# BKIA: Regret Proof

Event

$$E_t^1 = \left\{ \|h_{\widehat{M}_t}^\star - h_{M^\star}^\star\|_\infty \leq \varepsilon \ \wedge \ \Pi_{\widehat{M}_t}^\star(s) \subseteq \Pi_{M^\star}^\star(s), \forall s \right\} \qquad \widehat{M}_t \approx M^\star$$

$$E_t^2 = \left\{ b_t(s,a) < L_{M^\star}^\star h_{M^\star}^\star(s) - 2\varepsilon \right\}$$

# BKIA: Regret Proof

Event

$$E_t^1 = \left\{ \|h_{\widehat{M}_t}^\star - h_{M^\star}^\star\|_\infty \leq \varepsilon \ \wedge \ \Pi_{\widehat{M}_t}^\star(s) \subseteq \Pi_{M^\star}^\star(s), \forall s \right\}$$

$\widehat{M}_t \approx M^\star$

$$E_t^2 = \left\{ b_t(s,a) < L_{M^\star}^\star h_{M^\star}^\star(s) - 2\varepsilon \right\}$$

$$W_T^1(s,a,\varepsilon) = \sum_{t=1}^T \mathbb{1}\left(s_t, a_t = s, a\right) \times \mathbb{1}\left(E_t^1 \wedge E_t^2\right)$$

**?** One Step Optimism

$$\forall(s,a): \mathcal{M}_{M^\star}^{\mathsf{alt}}(s,a) \neq \emptyset$$

$$\lim_{\varepsilon \to 0} \limsup_{T \to \infty} \frac{\mathbb{E}[W_T^1(s,a,\varepsilon)]}{\log T} \leq \frac{1}{\inf_{M \in \mathcal{M}_{M^\star}^{\mathsf{alt}}(s,a)} \mathsf{KL}_{M^\star,M}(s,a)}$$

# DEL: Directed Exploration Learning
[Ok et al., 2018]

- DEL exploits the same idea of BKIA

  *Explore suboptimal actions no more than what prescribed by the lower bound*

- Exploration rate of sub-optimal action is *directed by the lower bound*

$$\text{target} \quad \eta_t(s, a) \approx \mathbb{E}\big[N_T(s, a)\big]$$



BKIA + OSSB $\approx$ DEL

OSSB [Combes et al., 2017] asymptotic optimal algorithm for structured bandit

# DEL

**for** $t = 1, \ldots, T$ **do**

$D_t(s) \leftarrow \{a \in \mathcal{A}(s) : N_t(s,a) \geq \log^2(N_t(s))\}$
$(g_t, h_t) \leftarrow \text{solve } \widehat{M}_t = \langle \mathcal{S}, D_t, \widehat{p}_t, r \rangle$

**if** $\forall a \in \Pi^\star_{\widehat{M}_t}(s_t), \ N_t(s_t, a) < \log^2(N_t(s_t) + 1)$ **then**

$\quad a_t \in \underset{a \in \Pi^\star_{\widehat{M}_t}(s_t)}{\arg\min} \{N_t(s,a)\}$

**else if** $C^{xpt}(H_t)$ **then**

$\quad$ B1 $\ exploit \ (a_t \in \Pi^\star_{\widehat{M}_t}(s_t))$

**else**

$\quad$ B2 $\ explore$

**end**

Observe reward $r_t$ and next state $s_{t+1}$

**end**

A Solve empirical MDP $\widehat{M}_t$ on a restricted action set

C Force exploration of "underestimated" actions

⚠ BKIA automatically trade-off exploration and exploitation $B1 + B2 \approx B_{\mathsf{BKIA}}$

# DEL: Exploration

B2 Directly *optimize the lower bound* on the estimated MDP $\widehat{M_t}$

$$\eta_t = \arg\inf_{\eta \in \mathbb{R}^{S \times A}} \sum_{s,a} \eta(s,a)\delta^{\star}_{\widehat{M_t}}(s,a)$$

$$\text{s.t. } \sum_{s,a} \eta(s,a)\mathsf{KL}_{\widehat{M_t},M}(s,a) \geq 1 \quad \forall M \in \mathcal{M}^{\mathsf{alt}}_{\widehat{M_t}}(s,a)$$

$$a_t \in \arg\min_{\mathcal{A}:\, N_t(s_t,a) \leq \eta_t(s_t,a)\gamma_t} \{N_t(s_t,a)\} \qquad * \; \gamma_t = (1+\gamma)(1+\log t)$$

# DEL: Exploration

B2 Directly *optimize the lower bound* on the estimated MDP $\widehat{M}_t$

$$\eta_t = \arg\inf_{\eta \in \mathbb{R}^{S \times A}} \sum_{s,a} \eta(s,a) \delta^\star_{\widehat{M}_t}(s,a)$$

$$\text{s.t. } \sum_{s,a} \eta(s,a) \mathsf{KL}_{\widehat{M}_t,M}(s,a) \geq 1 \quad \forall M \in \mathcal{M}^{\mathsf{alt}}_{\widehat{M}_t}(s,a)$$

$$a_t \in \arg\min_{\mathcal{A}: N_t(s_t,a) \leq \eta_t(s_t,a)\gamma_t} \{N_t(s_t,a)\} \qquad * \; \gamma_t = (1+\gamma)(1 + \log t)$$

♀ Lower bound sets the desired number of visits

$$\eta_t(s_t,a) \approx \mathbb{E}_{\widehat{M}_t}\left[N_T(s_t,a)\right] \approx \mathbb{E}_{M^\star}\left[N_T(s_t,a)\right]$$

then track it (in one step)

🗩 $\eta_t$ computed on $\widehat{M}_t$ and not $M^\star$ (wrong target)

# DEL: Regret Guarantees

---

**Theorem** (Thm. 4, [Ok et al., 2018])

*For any ergodic MDP $M^\star$ and under some technical conditions, for any $\gamma > 0$, the expected regret of DEL($\gamma$) is upper bounded as*

$$\limsup_{T \to \infty} \frac{\overline{R}(T, M^\star, \mathsf{DEL}(\gamma))}{\log T} \leq (1 + \gamma) K_{M^\star}^\star$$

---

# DEL: Regret Guarantees

### Theorem (Thm. 4, [Ok et al., 2018])

*For any ergodic MDP $M^\star$ and under some technical conditions, for any $\gamma > 0$, the expected regret of DEL$(\gamma)$ is upper bounded as*

$$\limsup_{T \to \infty} \frac{\overline{R}(T, M^\star, \mathsf{DEL}(\gamma))}{\log T} \leq (1 + \gamma) K_{M^\star}^\star$$

👍 DEL works for MDPs with structure (e.g., Lipschitz continuity)

# Open Questions

- *The role of forced exploration*
  - Why do we need to force exploration?
  - Is it due to the lack of long-term optimism?
  - Is it really required at algorithmic level?

- *Finite Time Analysis*

- *Refined lower bound*
  - Current lower bound is derived from a bandit perspective

# Markov Decision Process

A discrete-time finite Markov decision process (MDP) is a tuple $M = \langle \mathcal{S}, \mathcal{A}, r, p \rangle$
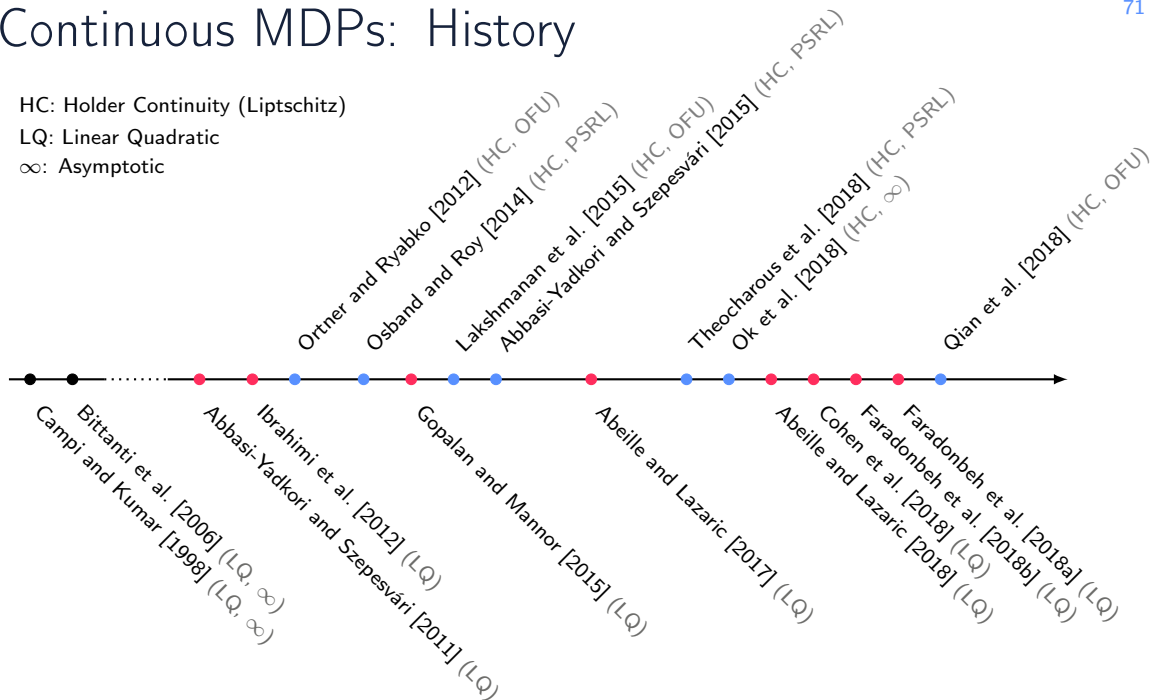
- State space $\mathcal{S}$, $|\mathcal{S}| = S < \infty$
- Action space $\mathcal{A}$, $|\mathcal{A}| = A < \infty$ } finite

- Transition distribution $p(\cdot|s,a) \in \Delta(\mathcal{S})$

- Reward distribution with expectation $r(s,a) \in [0, r_{\max}]$

👍 The process generates history $H_t = (s_1, a_1, \ldots, s_{t-1}, a_{t-1}, s_t)$, with $s_{t+1} \sim p(\cdot|s_t, a_t)$

# Continuous MDPs: History

HC: Holder Continuity (Liptschitz)
LQ: Linear Quadratic
∞: Asymptotic



Ortner and Ryabko [2012] (HC, OFU)
Osband and Roy [2014] (HC, PSRL)
Lakshmanan et al. [2015] (HC, OFU)
Abbasi-Yadkori and Szepesvári [2015] (HC, PSRL)
Theocharous et al. [2018] (HC, PSRL)
Ok et al. [2018] (HC, ∞)
Qian et al. [2018] (HC, OFU)

Campi and Kumar [1998] (LQ, ∞)
Bittanti et al. [2006] (LQ, ∞)
Abbasi-Yadkori and Szepesvári [2011] (LQ)
Ibrahimi et al. [2012] (LQ)
Gopalan and Mannor [2015] (LQ)
Abeille and Lazaric [2017] (LQ)
Abeille and Lazaric [2018] (LQ)
Cohen et al. [2018] (LQ)
Faradonbeh et al. [2018b] (LQ)
Faradonbeh et al. [2018a] (LQ)

Fruit, Lazaric and Pirotta

# Hölder Continuity

$\mathcal{S}$ continuous
$\mathcal{A}$ discrete

$L, \alpha > 0$  s.t.  $\forall s, s' \in \mathcal{S}, a \in \mathcal{A}$:

$$|r(s,a) - r(s',a)| \leq r_{\max} L |s - s'|^\alpha$$

$$\|p(\cdot|s,a) - p(\cdot|s',a)\|_1 \leq L |s - s'|^\alpha$$

HC1 Asm.

$$\mathsf{sp}(h^\star_{M^\star}) \leq C$$

HC2 Asm.

# Hölder Continuity

$\mathcal{S}$ continuous
$\mathcal{A}$ discrete

$L, \alpha > 0$  s.t.  $\forall s, s' \in \mathcal{S}, a \in \mathcal{A}$:
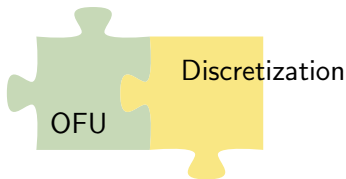
$$|r(s, a) - r(s', a)| \leq r_{\max} L |s - s'|^{\alpha}$$

$$\|p(\cdot | s, a) - p(\cdot | s', a)\|_1 \leq L |s - s'|^{\alpha}$$

HC1 Asm.

$$\mathsf{sp}(h^{\star}_{M^{\star}}) \leq C$$

HC2 Asm.

OFU

Discretization

[Ortner and Ryabko, 2012,
Lakshmanan et al., 2015,
Qian et al., 2018]

👉 $L, \alpha, C, T$ known in advance

# OFU: Hölder Continuity

**Theorem** (Ortner and Ryabko [2012], Lakshmanan et al. [2015], Qian et al. [2018])

*For any MDP $M$ satisfying Asm. (HC1) and (HC2), with probability at least $1 - \delta$ it holds that for any $T \geq 1$, the regret of UCCRL and SCCAL$^+$ is bounded as*

$$R(T, M^\star, \{\text{UCCRL}, \text{SCCAL}^+\}) \leq \beta \cdot C L \sqrt{A \log\left(\frac{T}{\delta}\right)} T^{(2+\alpha)/(2+2\alpha)}$$

*If the transition function is $\kappa$-times smoothly differentiable ($\gamma = \alpha + \kappa$)*

$$R(T, M^\star, \text{UCCRL-KD}) \leq \beta \cdot C L \sqrt{A \log\left(\frac{T}{\delta}\right)} T^{(\gamma+2\alpha+\alpha\gamma)/(\gamma+\alpha+2\alpha\gamma)}$$

# OFU: Liptschitz Continuity ($\alpha = 1$)

## Theorem (Ok et al. [2018])

*For any MDP $M$ satisfying Asm. (HC1) and (HC2) with $\alpha = 1$ the regret of DEL is bounded as*

$$\limsup_{T \to \infty} \frac{\overline{R}(T, M^\star, \mathsf{DEL})}{\log T} \leq S_L A \frac{(C+1)^3}{(\min_{s,a} \delta_{M^\star}(s,a))^2}$$

*with*

$$S_L = \min\{S, \frac{8L(C+1)}{\min_{s,a} \delta_{M^\star}(s,a)} + 1\}$$

*Comparison*

$$R(T, M^\star, \{\mathsf{UCCRL}, \mathsf{SCCAL}^+\}) = \widetilde{O}(T^{3/4})$$

$$R(T, M^\star, \mathsf{UCCRL\text{-}KD}) = \widetilde{O}(T^{2/3}) \text{ as } \kappa \to \infty$$

# Linear Quadratic Systems

$$\max_{\pi} \quad \lim_{T \to \infty} \frac{1}{T} \mathbb{E}\left[\sum_{t=1}^{T} r(s_t, a_t)\right]$$

$$\text{s.t.} \quad s_{t+1} = f(s_t, a_t, \epsilon_{t+1})$$

$$a_t \sim \pi(s_t)$$

# Linear Quadratic Systems

$$\max_{\pi} \quad \lim_{T \to \infty} \frac{1}{T} \mathbb{E}\left[ \sum_{t=1}^{T} -\left( s_t^\mathsf{T} Q s_t + a_t^\mathsf{T} R a_t \right) \right]$$

Quadratic Reward

$$\text{s.t.} \quad s_{t+1} = A s_t + B a_t + \epsilon_{t+1}$$

$$a_t \sim \pi(s_t)$$

Linear Dynamics

LQ system $M = \langle A, B, Q, R \rangle$

# Linear Quadratic Systems: Optimal Policy

- *Optimal policy*

solution of Discrete Algebraic Riccati Equation (DARE)

$$\pi^{\star}_M(s) = K^{\star}_M s$$
$$K^{\star}_M = -(R + B^{\mathsf{T}} P_M B)^{-1}(B^{\mathsf{T}} P_M A)$$

- *Optimal gain*

$$g^{\star}_M = Tr(P_M)$$

# Linear Quadratic Systems:  Optimal Policy

- *Optimal policy*

$$\pi_M^\star(s) = K_M^\star s$$
$$K_M^\star = -(R + B^\mathsf{T} P_M B)^{-1}(B^\mathsf{T} P_M A)$$

- *Optimal gain*

$$g_M^\star = Tr(P_M)$$

if $(A, B)$ are controllable, $K_M^\star$ makes the system *stable*

# OFU-LQ
[Abbasi-Yadkori and Szepesvári, 2011]

assume $Q$ and $R$ are known

*Optimism in LQ*

■ Estimation

Regularized Least Squares

$$\widehat{M}_t = \langle \widehat{A}_t, \widehat{B}_t, Q, R \rangle \qquad \text{where} \qquad (\widehat{A}_t, \widehat{B}_t) = \widehat{\theta}_t \leftarrow H_t$$

# OFU-LQ
[Abbasi-Yadkori and Szepesvári, 2011]

assume $Q$ and $R$ are known

*Optimism in LQ*

- Estimation

Regularized Least Squares

$$\widehat{M}_t = \langle \widehat{A}_t, \widehat{B}_t, Q, R \rangle \qquad \text{where} \qquad (\widehat{A}_t, \widehat{B}_t) = \widehat{\theta}_t \leftarrow H_t$$

Statistically admissible models

$$B_t^{\mathsf{RLS}} = \left\{ \theta \; : \; Tr\big((\theta - \widehat{\theta}_t)^{\mathsf{T}} V_t (\theta - \widehat{\theta}_t)\big) \leq \beta_t \right\}$$

# OFU-LQ

[Abbasi-Yadkori and Szepesvári, 2011]

assume $Q$ and $R$ are known

*Optimism in LQ*

- **Estimation**

Regularized Least Squares

<span style="color:blue">Statistically admissible models</span>

$$\widehat{M}_t = \langle \widehat{A}_t, \widehat{B}_t, Q, R \rangle \qquad \text{where} \qquad (\widehat{A}_t, \widehat{B}_t) = \widehat{\theta}_t \leftarrow H_t$$

$$B_t^{\mathsf{RLS}} = \left\{ \theta \ : \ Tr\left((\theta - \widehat{\theta}_t)^{\mathsf{T}} V_t (\theta - \widehat{\theta}_t)\right) \leq \beta_t \right\}$$

so that $\theta_t$ is controllable

- **Planning**

$$\theta_t = \underset{\theta \in \ \Theta \ \cap B_t^{\mathsf{RLS}}}{\arg\max} \ \{g_\theta^\star\}$$

# OFU-LQ
[Abbasi-Yadkori and Szepesvári, 2011]

assume $Q$ and $R$ are known

*Optimism in LQ*

- Estimation

$$\widehat{M}_t = \langle \widehat{A}_t, \widehat{B}_t, Q, R \rangle \qquad \text{where} \qquad (\widehat{A}_t, \widehat{B}_t) = \widehat{\theta}_t \leftarrow H_t$$

Regularized Least Squares

Statistically admissible models

$$B_t^{\mathsf{RLS}} = \{ \theta \ : \ Tr\big((\theta - \widehat{\theta}_t)^{\mathsf{T}} V_t (\theta - \widehat{\theta}_t)\big) \leq \beta_t \}$$

so that $\theta_t$ is controllable

- Planning

$$\theta_t = \underset{\theta \in \ \Theta \ \cap B_t^{\mathsf{RLS}}}{\arg\max} \ \{ g_\theta^\star \}$$

👎 Hard non-convex optimization problem

# OFU-LQ: Regret

### Theorem ([Abbasi-Yadkori and Szepesvári, 2011])

*For any $\delta \in ]0, 1[$, for any time $T$, with probability at least $1 - \delta$, the regret of* OFU-LQ *algorithm is bounded as*

$$R(T, M^\star, \text{OFU-LQ}) = \widetilde{O}\big(\sqrt{T \log(1/\delta)}\big)$$

Fruit, Lazaric and Pirotta

# OFU-LQ: Regret

## Theorem ([Abbasi-Yadkori and Szepesvári, 2011])

*For any $\delta \in ]0, 1[$, for any time $T$, with probability at least $1 - \delta$, the regret of OFU-LQ algorithm is bounded as*

$$R(T, M^\star, \text{OFU-LQ}) = \widetilde{O}\big(\sqrt{T \log(1/\delta)}\big)$$

💡 major challenge

$$K_{M^\star}^\star \to M^\star \quad \text{stable controller} \checkmark$$
$$K_t \to M^\star \quad \text{???}$$

central to the proof is how to control $\|s_t\|$

# Open Question

## Hölder continuity

**1** *Posterior Sampling*

- [Theocharous et al., 2018] proved $\widetilde{O}(C\sqrt{T})$
    - Under Asm. ASM-SP and Hölder continuity
    - Only for system parametrized by $1$-dimensional parameter

**2** Matching Lower Bound

## LQ Systems

**1** *Posterior Sampling*

- [Ouyang et al., 2017a] prove $\widetilde{O}(\sqrt{T})$ Bayesian regret under restrictive assumptions
- [Abeille and Lazaric, 2017, 2018] proved $\widetilde{O}(\sqrt{T})$ regret for PSRL with rejection sampling but only for $1$-dimensional systems

**2** *Efficient* OFU: many recent advances [Faradonbeh et al., 2018a, Cohen et al., 2019]

# Other Settings

- **Non-realizable approximated MDP** (e.g. [Jiang et al., 2017])

- **Non-stationary/adversarial environments** (e.g. [Even-Dar et al., 2009, Neu et al., 2014])

- **MDPs with arbitrary structure** (e.g. [Gopalan and Mannor, 2015])

- **Hierarchical exploration** (e.g. [Fruit and Lazaric, 2017, Fruit et al., 2017])

- **Low-exploration MDPs** (e.g. [Zanette and Brunskill, 2018])

- **Active/unsupervised exploration** (e.g. [Lim and Auer, 2012, Hazan et al., 2018, Tarbouriech and Lazaric, 2019])

- **Partially observable MDPs and beyond** (e.g. [Jiang et al., 2017, Azizzadenesheli et al., 2016])

Fruit, Lazaric and Pirotta

# Summary

| Alg. | Asymptotic (ergodic) | Finite-time (comm.) |
|------|---------------------|---------------------|
| Lower bound | $\dfrac{C^2 SA}{\min_{s,a} \delta^\star_{M^\star}(s,a)} \ln(T)$ | $\sqrt{DSAT}$ |
| UCRL2B | $\dfrac{D^2 S^2 A}{\delta^\star_g} \ln(T)$ | $\sqrt{DS\Gamma AT \ln(T)}$ |
| SCAL | $\dfrac{C^2 S^2 A}{\delta^\star_g} \ln(T)$ | $\sqrt{CS\Gamma AT \ln(T)}$ |
| TSDE | ? | $CS\sqrt{AT \ln(T)}$ |
| BKIA/DEL | $\dfrac{C^2 SA}{\min_{s,a} \delta^\star_{M^\star}(s,a)} \ln(T)$ | ? |

- $\Gamma = \max\limits_{s,a} |\mathsf{supp}(p(\cdot|s,a))|$
- $D_M = \max\limits_{s,s' \in \mathcal{S}} \min\limits_{\pi : \mathcal{S} \to \mathcal{A}} \mathbb{E}\big[T^M_\pi(s,s')\big]$
- $C \geq \mathsf{sp}(h^\star)$

- $\delta^\star_M(s,a) = L^\star_M h^\star_M(s) - L^a_M h^\star_M(s)$
- $\delta^\star_g := g^\star_M - \max\limits_{s \in \mathcal{S}, \pi} \big\{ g^\pi_{M^\star}(s) < g^\star_M \big\}$

# Open Question: Summary

| Alg. | Asymptotic (ergodic) | Finite-time (comm.) |
|------|----------------------|---------------------|
| Lower bound | $\dfrac{C^2 SA}{\min_{s,a} \delta^\star_{M^\star}(s,a)} \ln(T)$ | $\sqrt{DSAT}$ |
| UCRL2B | $\dfrac{D^2 S^2 A}{\delta^\star_g} \ln(T)$ | $\sqrt{DS\Gamma AT \ln(T)}$ |
| SCAL | $\dfrac{C^2 S^2 A}{\delta^\star_g} \ln(T)$ | $\sqrt{CS\Gamma AT \ln(T)}$ |
| TSDE | ? | $CS\sqrt{AT \ln(T)}$ (Bayes) |
| BKIA/DEL | $\dfrac{C^2 SA}{\min_{s,a} \delta^\star_{M^\star}(s,a)} \ln(T)$ | ? |

*Closing the gap* between upper and lower bounds and settings (ergodic/asymptotic vs communicating/worst-case)

▤ Many lessons learned from bandit but need to deal with dynamical nature of the problem.

# Open Questions

- Unifying finite-horizon, infinite-horizon regret and discounted PAC-MDP guarantees (e.g. [Dann et al., 2017])

- Model-based vs model-free (e.g. [Jin et al., 2018, Szepesvari et al., 2019])

- Scalable exp-exp (e.g. Bellemare et al. [2016], Tang and Agrawal [2018], Fortunato et al. [2017])

# Open Questions: Model-free vs Model-based

Model-based exploration

👍 sample efficient (regret $O(\sqrt{T})$)

👎 solves an MDP at each episode ($O(S^2 A)$)

👎 difficult to extend to function approximation

Model-free exploration

👎 sample inefficient (regret $O(T^{2/3})$?)

👍 simple update at each step ($O(1)$)

👍 easy to extend to function approximation

# Open Questions: Model-free vs Model-based

Model-based exploration

👍 sample efficient (regret $O(\sqrt{T})$)

👎 solves an MDP at each episode ($O(S^2 A)$)

👎 difficult to extend to function approximation

Model-free exploration

👎 sample inefficient (regret $O(T^{2/3})$?)

👍 simple update at each step ($O(1)$)

👍 easy to extend to function approximation

*Sample and computationally efficient* exploration algorithm? (see Jin et al. [2018])

# Resources

## Reinforcement Learning

- Books

  - Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*.
    John Wiley & Sons, Inc., New York, NY, USA, 1994

  - Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1.
    MIT press Cambridge, 1998

  - Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control, Vol. II*.
    Athena Scientific, 3rd edition, 2007

  - Csaba Szepesvari. *Algorithms for Reinforcement Learning*.
    Morgan and Claypool Publishers, 2010

- Courses (with good references for exploration)

  - Nan Jiang. Cs598 statistical reinforcement learning.
    http://nanjiang.cs.illinois.edu/cs598/

  - Emma Brunskill. Cs234 reinforcement learning winter 2019.
    http://web.stanford.edu/class/cs234/index.html

  - Alessandro Lazaric. Mva reinforcement learning.
    http://chercheurs.lille.inria.fr/~lazaric/Webpage/Teaching.html

  - Alexandre Proutiere. Reinforcement learning: A graduate course.
    http://www.it.uu.se/research/systems_and_control/education/2017/relearn/

# Resources

**Exploration-Exploitation and Regret Minimization**

- Books

  - Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems.
    *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012

  - Tor Lattimore and Csaba Szepesvári. Bandit algorithms.
    Pre-publication version, 2018.
    URL `http://downloads.tor-lattimore.com/banditbook/book.pdf`

# Thank you!

**facebook**
Artificial Intelligence Research

Yasin Abbasi-Yadkori and Csaba Szepesvári. Regret bounds for the adaptive control of linear quadratic systems. In *COLT*, volume 19 of *JMLR Proceedings*, pages 1–26. JMLR.org, 2011.

Yasin Abbasi-Yadkori and Csaba Szepesvári. Bayesian optimal control of smoothly parameterized systems. In *UAI*, pages 1–11. AUAI Press, 2015.

Marc Abeille and Alessandro Lazaric. Thompson sampling for linear-quadratic control problems. In *AISTATS*, volume 54 of *Proceedings of Machine Learning Research*, pages 1246–1254. PMLR, 2017.

Marc Abeille and Alessandro Lazaric. Improved regret bounds for thompson sampling in linear quadratic control problems. In *ICML*, volume 80 of *JMLR Workshop and Conference Proceedings*, pages 1–9. JMLR.org, 2018.

Rajeev Agrawal. Adaptive control of markov chains under the weak accessibility. In *29th IEEE Conference on Decision and Control*, pages 1426–1431. IEEE, 1990.

Shipra Agrawal and Randy Jia. Optimistic posterior sampling for reinforcement learning: worst-case regret bounds. In *NIPS*, pages 1184–1194, 2017.

Peter Auer and Ronald Ortner. Logarithmic online regret bounds for undiscounted reinforcement learning. In *NIPS*, pages 49–56. MIT Press, 2006.

Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 263–272. PMLR, 2017.

Kamyar Azizzadenesheli, Alessandro Lazaric, and Animashree Anandkumar. Reinforcement learning of pomdps using spectral methods. In *COLT*, volume 49 of *JMLR Workshop and Conference Proceedings*, pages 193–256. JMLR.org, 2016.

Peter L. Bartlett and Ambuj Tewari. REGAL: A regularization based algorithm for reinforcement learning in weakly communicating MDPs. In *UAI*, pages 35–42. AUAI Press, 2009.

Marc G. Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Rémi Munos. Unifying count-based exploration and intrinsic motivation. In *NIPS*, pages 1471–1479, 2016.

Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control, Vol. II*. Athena Scientific, 3rd edition, 2007.

Sergio Bittanti, Marco C Campi, et al. Adaptive control of linear time invariant systems: the "bet on the best" principle. *Communications in Information & Systems*, 6(4):299–320, 2006.

Emma Brunskill. Cs234 reinforcement learning winter 2019. http://web.stanford.edu/class/cs234/index.html.

Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.

Apostolos N. Burnetas and Michael N. Katehakis. Optimal adaptive policies for markov decision processes. *Mathematics of Operations Research*, 22(1):222–255, 1997.

Marco C Campi and PR Kumar. Adaptive linear quadratic gaussian control: the cost-biased approach revisited. *SIAM Journal on Control and Optimization*, 36(6):1890–1907, 1998.

Alon Cohen, Avinatan Hassidim, Tomer Koren, Nevena Lazic, Yishay Mansour, and Kunal Talwar. Online linear quadratic control. In *ICML*, volume 80 of *JMLR Workshop and Conference Proceedings*, pages 1028–1037. JMLR.org, 2018.

Alon Cohen, Tomer Koren, and Yishay Mansour. Learning Linear-Quadratic Regulators Efficiently with only $O(\sqrt{T})$ Regret. *arXiv e-prints*, art. arXiv:1902.06223, Feb 2019.

Richard Combes, Stefan Magureanu, and Alexandre Proutière. Minimal exploration in structured stochastic bandits. In *NIPS*, pages 1761–1769, 2017.

Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying PAC and regret: Uniform PAC bounds for episodic reinforcement learning. In *NIPS*, pages 5717–5727, 2017.

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

Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. Input perturbations for adaptive regulation and learning. *CoRR*, abs/1811.04258, 2018a.

Mohamad Kazem Shirani Faradonbeh, Ambuj Tewari, and George Michailidis. On optimality of adaptive linear-quadratic regulators. *CoRR*, abs/1806.10749, 2018b.

Sarah Filippi, Olivier Cappé, and Aurélien Garivier. Optimism in reinforcement learning and kullback-leibler divergence. *2010 48th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 115–122, 2010.

Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Ian Osband, Alex Graves, Vlad Mnih, Rémi Munos, Demis Hassabis, Olivier Pietquin, Charles Blundell, and Shane Legg. Noisy networks for exploration. *CoRR*, abs/1706.10295, 2017.

Ronan Fruit and Alessandro Lazaric. Exploration-exploitation in mdps with options. In *AISTATS*, volume 54 of *Proceedings of Machine Learning Research*, pages 576–584. PMLR, 2017.

Ronan Fruit, Matteo Pirotta, Alessandro Lazaric, and Emma Brunskill. Regret minimization in mdps with options without prior knowledge. In *NIPS*, pages 3169–3179, 2017.

Ronan Fruit, Matteo Pirotta, and Alessandro Lazaric. Near optimal exploration-exploitation in non-communicating markov decision processes. In *NeurIPS*, pages 2998–3008, 2018a.

Ronan Fruit, Matteo Pirotta, Alessandro Lazaric, and Ronald Ortner. Efficient bias-span-constrained exploration-exploitation in reinforcement learning. In *ICML*, Proceedings of Machine Learning Research. PMLR, 2018b.

Aurélien Garivier and Olivier Cappé. The KL-UCB algorithm for bounded stochastic bandits and beyond. In *COLT*, volume 19 of *JMLR Proceedings*, pages 359–376. JMLR.org, 2011.

Aditya Gopalan and Shie Mannor. Thompson sampling for learning parameterized markov decision processes. In *COLT*, volume 40 of *JMLR Workshop and Conference Proceedings*, pages 861–898. JMLR.org, 2015.

Todd L Graves and Tze Leung Lai. Asymptotically efficient adaptive choice of control laws incontrolled markov chains. *SIAM journal on control and optimization*, 35(3):715–743, 1997.

Elad Hazan, Sham M. Kakade, Karan Singh, and Abby Van Soest. Provably Efficient Maximum Entropy Exploration. *arXiv e-prints*, art. arXiv:1812.02690, Dec 2018.

Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963. URL http://www.jstor.org/stable/2282952.

Morteza Ibrahimi, Adel Javanmard, and Benjamin Van Roy. Efficient reinforcement learning for high dimensional linear quadratic systems. In *NIPS*, pages 2645–2653, 2012.

Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11:1563–1600, 2010.

Nan Jiang. Cs598 statistical reinforcement learning. http://nanjiang.cs.illinois.edu/cs598/.

Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E. Schapire. Contextual decision processes with low bellman rank are pac-learnable. In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 1704–1713. PMLR, 2017.

Chi Jin, Zeyuan Allen-Zhu, Sébastien Bubeck, and Michael I. Jordan. Is q-learning provably efficient? In *NeurIPS*, pages 4868–4878, 2018.

Sham Kakade, Mengdi Wang, and Lin F. Yang. Variance reduction methods for sublinear reinforcement learning. *CoRR*, abs/1802.09184, 2018.

T.L Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4 – 22, 1985. ISSN 0196-8858. doi: https://doi.org/10.1016/0196-8858(85)90002-8. URL http://www.sciencedirect.com/science/article/pii/0196885885900028.

K. Lakshmanan, Ronald Ortner, and Daniil Ryabko. Improved regret bounds for undiscounted continuous reinforcement learning. In *ICML*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 524–532. JMLR.org, 2015.

Tor Lattimore and Csaba Szepesvári. Bandit algorithms. Pre-publication version, 2018. URL http://downloads.tor-lattimore.com/banditbook/book.pdf.

Alessandro Lazaric. Mva reinforcement learning.
http://chercheurs.lille.inria.fr/~lazaric/Webpage/Teaching.html.

Shiau Hong Lim and Peter Auer. Autonomous exploration for navigating in mdps. In *COLT*, volume 23 of *JMLR Proceedings*, pages 40.1–40.24. JMLR.org, 2012.

Gergely Neu, András György, Csaba Szepesvári, and András Antos. Online markov decision processes under bandit feedback. *IEEE Trans. Automat. Contr.*, 59(3):676–691, 2014.

Jungseul Ok, Alexandre Proutière, and Damianos Tranos. Exploration in structured reinforcement learning. In *NeurIPS*, pages 8888–8896. 2018.

Ronald Ortner. Online regret bounds for markov decision processes with deterministic transitions. *Theor. Comput. Sci.*, 411(29-30):2684–2695, 2010.

Ronald Ortner and Daniil Ryabko. Online regret bounds for undiscounted continuous reinforcement learning. In *NIPS*, pages 1772–1780, 2012.

Ian Osband and Benjamin Van Roy. Model-based reinforcement learning and the eluder dimension. In *NIPS*, pages 1466–1474, 2014.

Ian Osband and Benjamin Van Roy. On lower bounds for regret in reinforcement learning. *CoRR*, abs/1608.02732, 2016a.

Ian Osband and Benjamin Van Roy. Posterior sampling for reinforcement learning without episodes. *CoRR*, abs/1608.02731, 2016b.

Ian Osband and Benjamin Van Roy. Why is posterior sampling better than optimism for reinforcement learning? In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 2701–2710. PMLR, 2017.

Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via posterior sampling. In *NIPS*, pages 3003–3011, 2013.

Yi Ouyang, Mukul Gagrani, and Rahul Jain. Learning-based control of unknown linear systems with thompson sampling. *CoRR*, abs/1709.04047, 2017a. URL http://arxiv.org/abs/1709.04047.

Yi Ouyang, Mukul Gagrani, Ashutosh Nayyar, and Rahul Jain. Learning unknown markov decision processes: A thompson sampling approach. In *NIPS*, pages 1333–1342, 2017b.

Alexandre Proutiere. Reinforcement learning: A graduate course. http://www.it.uu.se/research/systems_and_control/education/2017/relearn/.

Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1994.

Jian Qian, Ronan Fruit, Matteo Pirotta, and Alessandro Lazaric. Exploration bonus for regret minimization in undiscounted discrete and continuous markov decision processes. *CoRR*, 2018.

Alexander L Strehl and Michael L Littman. An analysis of model-based interval estimation for markov decision processes. *Journal of Computer and System Sciences*, 74(8):1309–1331, 2008.

Malcolm Strens. A bayesian framework for reinforcement learning. In *In Proceedings of the Seventeenth International Conference on Machine Learning*, pages 943–950. ICML, 2000.

Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.

Csaba Szepesvari. *Algorithms for Reinforcement Learning*. Morgan and Claypool Publishers, 2010.

Csaba Szepesvari, Nevena Lazic, and Yasin Abbasi-Yadkori. Model-free linear quadratic control via reduction to expert prediction. In *AISTATS*, 2019.

Mohammad Sadegh Talebi and Odalric-Ambrym Maillard. Variance-aware regret bounds for undiscounted reinforcement learning in mdps. In *ALT*, volume 83 of *Proceedings of Machine Learning Research*, pages 770–805. PMLR, 2018.

Yunhao Tang and Shipra Agrawal. Exploration by distributional reinforcement learning. *CoRR*, abs/1805.01907, 2018.

Jean Tarbouriech and Alessandro Lazaric. Active Exploration in Markov Decision Processes. *arXiv e-prints*, art. arXiv:1902.11199, Feb 2019.

Ambuj Tewari and Peter L. Bartlett. Optimistic linear programming gives logarithmic regret for irreducible mdps. In *NIPS*, pages 1505–1512. Curran Associates, Inc., 2007.

Georgios Theocharous, Zheng Wen, Yasin Abbasi, and Nikos Vlassis. Scalar posterior sampling with applications. In *NeurIPS*, pages 7696–7704, 2018.

William R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.

Tsachy Weissman, Erik Ordentlich, Gadiel Seroussi, Sergio Verdú, and Marcelo J. Weinberger. Inequalities for the L1 deviation of the empirical distribution. 2003.

Andrea Zanette and Emma Brunskill. Problem dependent reinforcement learning bounds which can identify bandit structure in mdps. In *ICML*, volume 80 of *JMLR Workshop and Conference Proceedings*, pages 5732–5740. JMLR.org, 2018.