

Efficient Bias-Span-Constrained Exploration-Exploitation in Reinforcement Learning

Ronan Fruit (INRIA), Matteo Pirodda (INRIA), Alessandro Lazaric (FAIR), Ronald Ortner (Montanuniversität Leoben)



Motivations

- Learning in an unknown environment means to balance
 - Exploration
 - Exploitation
- Optimistic (OFU) methods
 - Construct a set of plausible MDPs
 - Execute the optimal policy of the “best” MDP in the set
- OFU may lead to over-optimism in some MDPs
 - even fail in learning
- Regularization has proved to be effective in ML \Rightarrow in Exp-Exp?

Online Learning in MDPs

Markov Decision Process $M = \{S, A, r, p\}$

- S : states
- $A = (A_s)_{s \in S}$: actions
- $r(s, a)$: mean rewards
- $p(s'|s, a)$: transition probabilities, $\Gamma = \max_{s,a} \|p(\cdot|s, a)\|_0 \leq S$

Optimality criterion: long-term average reward

For any policy $\pi \in \Pi^{\text{SR}}(M)$ starting from $s \in S$:

$$\text{GAIN: } g_M^\pi(s) := \lim_{T \rightarrow +\infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T r(s_t, a_t) \right]$$

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

In weakly communicating MDPs, any optimal policy

$$\pi^* \in \arg \max_{\pi} \{g^\pi(s)\}$$

has constant gain, i.e., $g^{\pi^*}(s) = g^*$ for all $s \in S$.

Moreover, $(g^*, h^*) = (g^{\pi^*}, h^{\pi^*})$ satisfy the optimality equation

$$h^* = [L]h^* - g^*e, \quad \text{where } e = (1, \dots, 1)^\top.$$

optimal Bellman operator: $Lv(s) = \max_{a \in A_s} \{r(s, a) + p(\cdot|s, a)^\top v\}$

Cumulative regret minimization $\Delta(\mathfrak{A}, T) = \sum_{t=1}^T (g^* - r_t(s_t, a_t))$

Diameter and Span: [Jaksch et al. 2010; Bartlett and Tewari, 2009]

$$D = \max_{s, s' \in S} \left\{ \min_{\pi: S \rightarrow \mathcal{P}(A)} \left[\mathbb{E}_\pi [T(s')|s] \right] \right\}$$

$$sp\{h^*\} = \max_{s \in S} \{h^*(s)\} - \min_{s \in S} \{h^*(s)\}$$

- $sp\{h^*\} \leq D$ (always)
- $D = \infty$ but $sp\{h^*\} < \infty$

Regret Results

UCRL [Jaksch et al. 2010] & Optimistic PSRL [Agrawal and Jia, 2017]	$\tilde{O}(D\sqrt{\Gamma SAT})$
REGAL.C [Bartlett and Tewari, 2009]	$\tilde{O}(C\sqrt{\Gamma SAT})$ with $sp\{h^*\} \leq C$
Lower bound: $\Omega(\sqrt{DSAT})$ or $\Omega(\sqrt{sp\{h^*\}SAT})$? [Jaksch et al. 2010]	

References

- Jaksch, Ortner, and Auer. Near-optimal regret bounds for reinforcement learning. Journal of Machine Learning Research, 2010.
- Bartlett and Tewari. REGAL: A regularization based algorithm for reinforcement learning in weakly communicating MDPs. In UAI 2009.
- Agrawal and Jia. Optimistic posterior sampling for reinforcement learning: worst-case regret bounds. In NIPS 2017, 2017.

Optimism in Face of Uncertainty

For episodes $k = 1, 2, \dots$

- Compute a set of plausible MDPs \mathcal{M}_k
- Planning: Compute an optimistic policy $\tilde{\pi}_k$
- Execute policy $\tilde{\pi}_k$

Confidence set:

$$\mathcal{M}_k = \{M = \langle S, A, \tilde{r}, \tilde{p} \rangle : \tilde{r}(s, a) \in B_r^k(s, a), \tilde{p}(s'|s, a) \in B_p^k(s, a, s')\}$$

$$B_r^k, B_p^k = \text{high-probability confidence intervals around empirical estimates}$$

$$|\hat{r}(s, a) - \tilde{r}(s, a)| \leq \beta_{r,k}^{sa}, \quad |\hat{p}(s'|s, a) - \tilde{p}(s'|s, a)| \leq \beta_{p,k}^{sas'}$$

Planning:

$$\text{UCRL [Jaksch et al., 2010]} \quad (\tilde{M}_k^*, \tilde{\pi}_k^*) \in \arg \max_{M \in \mathcal{M}_k, \pi \in \Pi^{\text{SD}}(M)} g_M^\pi$$

- No prior knowledge \Rightarrow might be over-optimistic

$$\text{REGAL.C [Bartlett and Tewari, 2009]} \quad (\tilde{M}_R^*, \tilde{\pi}_R^*) = \arg \max_{M \in \mathcal{M}_{\text{RC}}, \pi \in \Pi^{\text{SD}}(M)} g_M^\pi$$

where $\mathcal{M}_{\text{RC}} := \{M \in \mathcal{M}_k : sp\{h_M^*\} \leq c\}$

- Constrains the set of possible MDPs \Rightarrow intractable
- Overall problem is ill-posed

$$\text{SCAL} \quad (\tilde{M}_c^*, \tilde{\pi}_c^*) \in \arg \max_{M \in \mathcal{M}_k, \pi \in \Pi_c(M)} g_M^\pi$$

where $\Pi_c(M) := \{\pi \in \Pi^{\text{SR}} : sp\{h_M^*\} \leq c \wedge sp\{g_M^*\} = 0\}$

- Relaxation of REGAL.C: see PROP. 1
- Overall problem is well-posed

Use prior knowledge about $sp\{h_{M^*}^*\} \leq c$

Planning under bias-span constraint (SCOPT)

Optimization problem: $g_c^* := \sup_{\pi \in \Pi_c(M)} \{g_M^\pi\}$

where $\Pi_c(M) := \{\pi \in \Pi^{\text{SR}} : sp\{h_M^*\} \leq C \wedge sp\{g_M^*\} = 0\}$

Why $\pi \in \Pi^{\text{SR}}?$

the maximizer $\pi_c^*(M)$ may be a randomize policy

Why $sp\{g_M^*\} = 0?$

there may be no dominating policy $\pi \in \Pi^{\text{SR}}$ with constant bias span

$\pi \in \Pi^{\text{SR}}$ is dominating $\Rightarrow \forall \pi' \in \Pi^{\text{SR}}, \forall s \in S, g^{\pi'}(s) \geq g^{\pi}(s)$

The supremum always exists, the set $\Pi_c^*(M)$ of maximizers?

Lemma. If M is unichain then $\Pi_c^*(M) \neq \emptyset$.

Value Operator: Given $v \in \mathbb{R}^S$ and $C \geq 0$, we define $\tilde{\mathcal{S}}(C, v) = \{s \in S | Lv(s) \leq \min_s \{Lv(s)\} + C\}$ and

$$T_c v = \Gamma_c L v = \begin{cases} Lv(s) & \forall s \in \tilde{\mathcal{S}}(C, v), \\ C + \min_s \{Lv(s)\} & \forall s \in S \setminus \tilde{\mathcal{S}}(C, v), \end{cases}$$

T_c is feasible at (v, s) i.i.f.

$$s \in \tilde{\mathcal{S}}(C, v) = \left\{ \min_{a \in A_s} \{r(s, a) + p(\cdot|s, a)^\top v\} \leq \min_{s'} \{Lv(s')\} + C \right\}$$

$$\Rightarrow \exists \delta_v^+ T_c v(s) = \sum_{a \in A_s} \delta_v^+(s, a) [r(s, a) + p(\cdot|s, a)^\top v]$$

Greedy operator:

$$G_c(v) = \begin{cases} \delta_v^+(s) & s \in \tilde{\mathcal{S}}(C, v), \\ \arg \min_{a \in A_s} \{r(s, a) + p(\cdot|s, a)^\top v\} & s \in S \setminus \tilde{\mathcal{S}}(C, v). \end{cases}$$

Planning algorithm: SCOPT \rightarrow

relative value iteration
c-span truncation (from above)
 $\forall n, sp\{v_n\} \leq c$

- Initialize $n = 0$ and $v_1 = T_c v_0 - (T_c v_0)(\bar{s})e$
- While $sp\{v_{n+1} - v_n\} + \frac{2\gamma^n}{1-\gamma} sp\{v_1 - v_0\} > \varepsilon$
 - $n \leftarrow n + 1$
 - $v_{n+1} = T_c v_n - (T_c v_n)(\bar{s})e$.
- return v_n and $\pi^n = G_c(v_n)$

Asm. 1. L is a γ -span contraction

$\forall u, v \in \mathbb{R}^S, sp\{Lu - Lv\} \leq \gamma sp\{u - v\}$

Asm. 2. T_c is globally feasible:

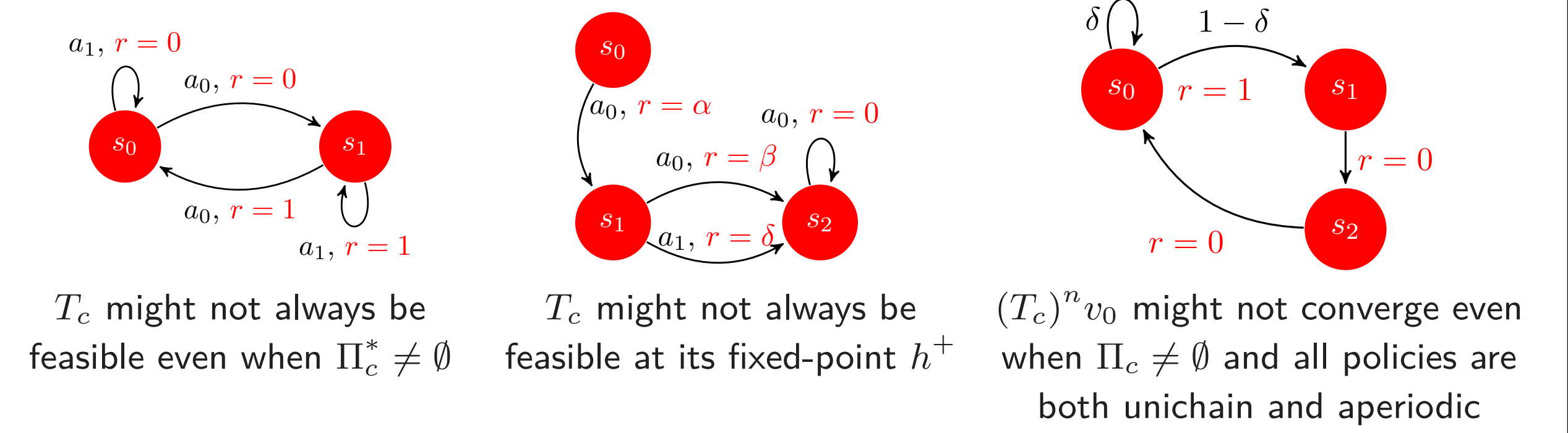
$\exists \delta_v^+, T_c v = L_{\delta_v^+} v$ (feasible at each s)

SCOPT PROPERTIES

- Span contraction: T_c is a γ -span contraction
- Optimality equation (OE): $\exists (g^+, h^+) \in \mathbb{R} \times \mathbb{R}^S$ s.t. $T_c h^+ = h^+ + g^+ e$.
- Uniqueness: if (g, h) solution of OE, then $g = g^+$ and $h = h^+ + \lambda e$
- Convergence: $\forall v_0 \in \mathbb{R}^S$, the sequence (v_n) converges to h^+ and $\lim_{n \rightarrow +\infty} T_c^{n+1} v_0 - T_c^n v_0 = g^+ e$.
- Dominance: the gain g^+ is an upper-bound on the supremum, i.e., $g^+ \geq g_c^*$.

Asm. 1 + Asm. 2

- ϵ -optimality: the policy returned by SCOPT is ϵ -optimal: $\|g^+ e - g^{\pi^n}\|_\infty \leq \varepsilon$
- Convergence to g_c^* : if $\pi^+ = G_c(h^+)$ is unichain, then $g^+ = g_c^*$ and $\pi^+ \in \Pi_c^*$



Learning under bias-span constraint (SCAL)

Equivalent planning problem:

$$\tilde{\mu}^* \in \arg \max_{\mu \in \Pi_c(\tilde{\mathcal{M}}_k)} \{g_{\tilde{\mathcal{M}}_k}^\mu\}$$

$\tilde{\mathcal{M}}_k$ is an extended MDP

$\Rightarrow \mathcal{A}$ is “extended” to a compact space $\tilde{\mathcal{A}}$ by considering every possible value in B_r^k and B_p^k

Define

$$\tilde{L}v = \max_{a \in A_s} \left\{ \max_{\tilde{r} \in B_r^k(s, a)} \tilde{r} + \max_{\tilde{p} \in B_p^k(s, a, \cdot)} \{\tilde{p}^\top v\} \right\}$$

and $\tilde{T}_c = \Gamma_c \tilde{L}v$

\Rightarrow use SCOPT for planning

No convergence for SCOPT in $\tilde{\mathcal{M}}_k$

CAN WE STILL USE SCOPT?

We can alter $\tilde{\mathcal{M}}_k \Rightarrow \tilde{\mathcal{M}}_k^\dagger$

- η -perturbation of the transition model
we enforce that the “attractive” state \bar{s} is reached with non-zero probability from any state-action pair

$$\gamma = 1 - \min_{s, u \in \tilde{S}, a, b \in \mathcal{A}} \left\{ \sum_{j \in S} \min \{\tilde{p}(j|s, a), \tilde{q}(j|u, b)\} \right\}$$

$$\leq 1 - \eta < 1 \quad \Rightarrow \tilde{L}^\dagger \text{ is } \gamma\text{-contractive}$$

SCOPT in $\tilde{\mathcal{M}}_k^\dagger$ converges to $g_c^* \geq \max_{\pi \in \Pi_c(\tilde{\mathcal{M}}_k)} g_{\tilde{\mathcal{M}}_k}^\pi - \eta c$

There might not be any policy associated to g_c^*

Augment the reward in $\tilde{\mathcal{M}}_k^\eta \Rightarrow \tilde{\mathcal{M}}_k^\dagger$

- \mathcal{A}_k is expanded by duplicating every action
i.e., $B_{r,k}^\dagger(s, a) = [0, \max\{B_r^k(s, a)\}]$

$$\circ \tilde{L}^\eta v = \tilde{L}^\dagger v \text{ and } \tilde{T}^\eta v = \tilde{T}_c^\dagger v$$

$$g_c^*(\tilde{\mathcal{M}}_k^\eta) = g_c^*(\tilde{\mathcal{M}}_k^\dagger)$$

- \tilde{T}_c is globally feasible

$$\forall v \text{ s.t. } sp\{v\} \leq c, \exists \delta \text{ s.t. } sp\{\tilde{L}_\delta^\dagger v\} \leq c$$

- $\tilde{\mathcal{M}}_k^\dagger$ is unichain

$$\Rightarrow g_c^* \mapsto \tilde{\pi}_k \text{ (a policy exists)}$$

SCAL computes

$$\max_{\pi \in \Pi_c(\tilde{\mathcal{M}}_k^\dagger)} g_M^\pi$$

well defined problem and admits a maximizer $\pi_c^*(\tilde{\mathcal{M}}_k^\dagger)$
 \Rightarrow efficiently computed using SCOPT

Regret

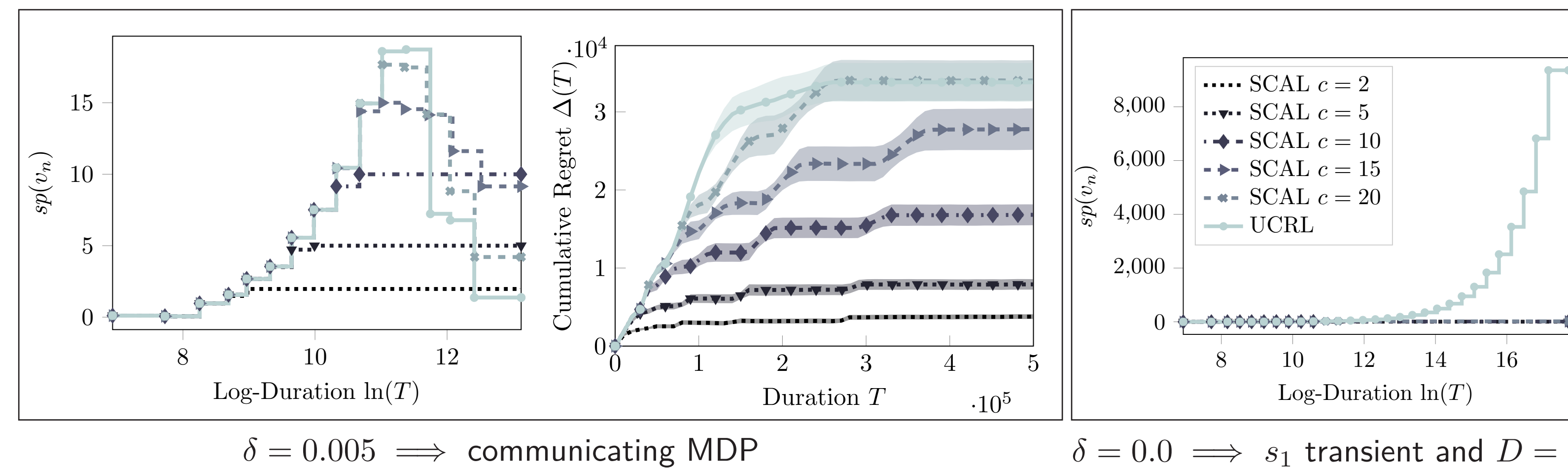
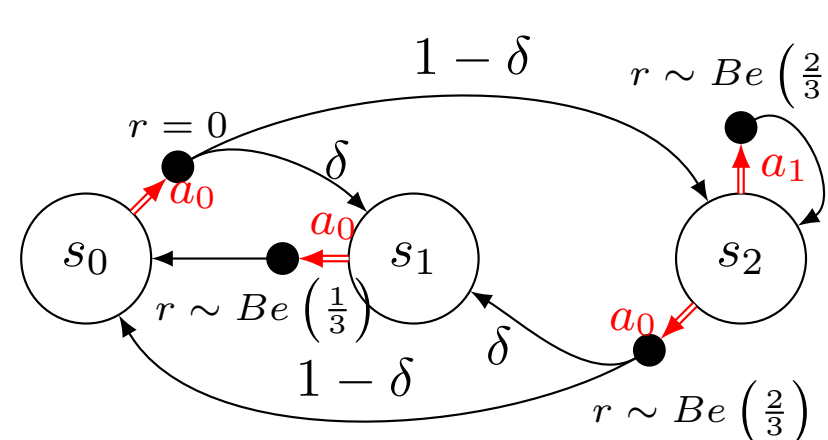
$$\Delta(\text{SCAL}, T) = \tilde{O}(\min\{D, c\} \sqrt{\Gamma SAT})$$

Numerical Experiments

SIMPLE 2D DOMAIN

$\Rightarrow \delta$ regulates the complexity

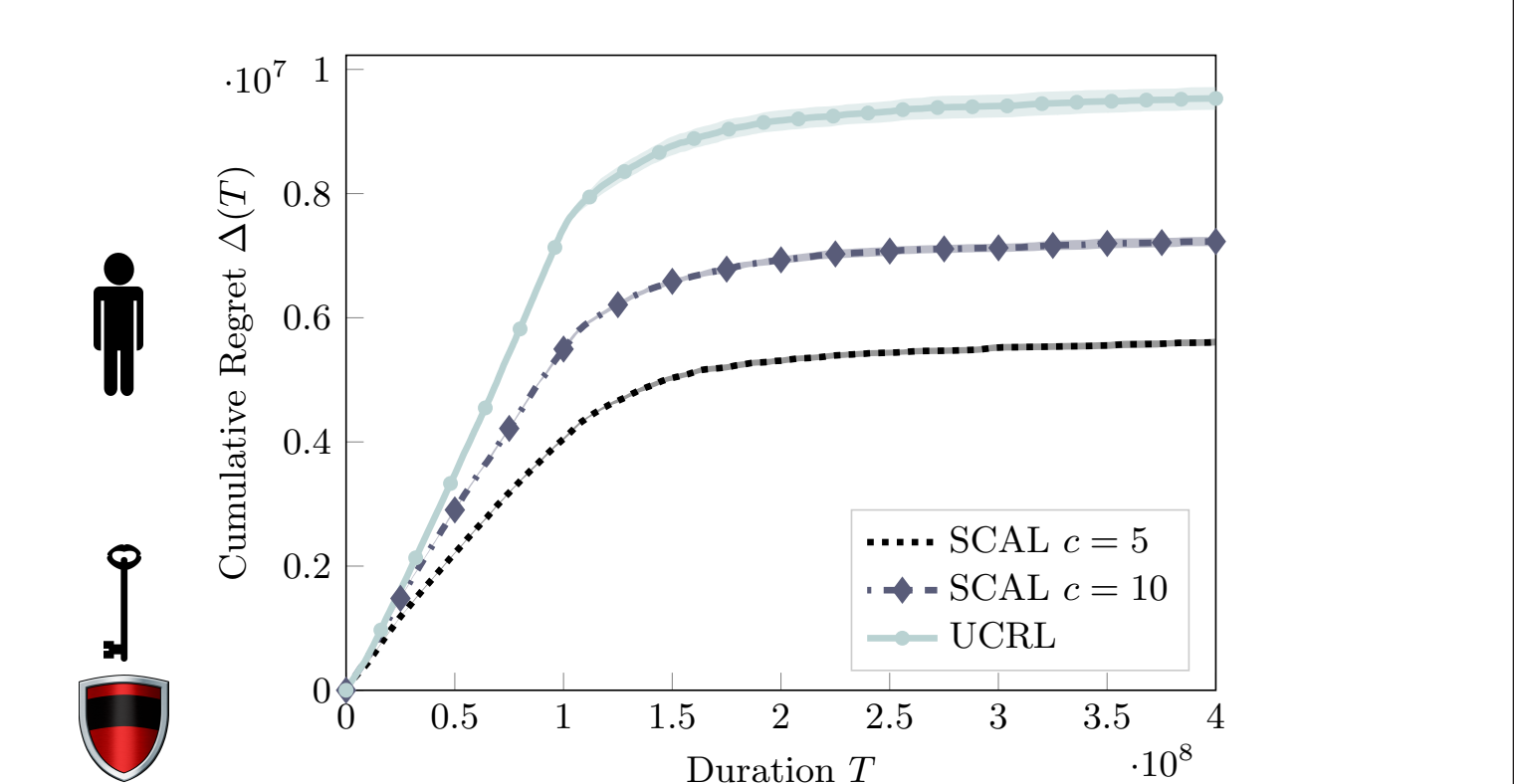
- $g^* = 2/3$ $\circ D \approx 1/\delta$
- $sp\{h^*\} = \frac{1}{1-\delta}$



$\delta = 0.005 \Rightarrow$ communicating MDP

$\delta = 0.0 \Rightarrow s_1$ transient and $D = \infty$

WAY TO FREEDOM



$S = 360, A = 8, D \approx 250$ (communicating) and $sp\{h^*\} \approx 3.2$