

Exploration Bonus for Regret Minimization in Discrete and Continuous Average Reward MDPs

J. Qian (INRIA Lille), R. Fruit (INRIA Lille), M. Pirotta (FAIR), A. Lazaric (FAIR)

Introduction

WHAT DO WE DO?

- Exploration in continuous MDPs
- With theoretical guarantees

How?

- Using exploration bonus and discretization

Online Learning in MDPs

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

- Optimality criterion: average reward

For any policy π starting from $s \in S$:

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

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

- Learning problem: cumulative regret minimization

$$\pi^* \text{ and } g^* \text{ are unknown} \quad \Delta(\mathcal{A}, T) = Tg^* - \sum_{t=1}^T r_t(s_t, a_t)$$

- Diameter: [Jaksch et al. 2010]

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

expected time $s \rightarrow s'$

Prior Knowledge on the Bias Span

- Provides a sense of what is realizable in the true MDP
- Avoids over-optimism
- Necessary to define the exploration bonus

Implicit in other settings
(infinite-horizon discounted,
finite-horizon)

$$|\tilde{p}(\cdot | s, a) - \hat{p}(\cdot | s, a)| \leq \|\tilde{p}(\cdot | s, a) - \hat{p}(\cdot | s, a)\|_1 \|h^*\|_\infty$$

Setting	MDP parameter	Horizon	Knowledge	Exploration Bonus
infinite-horizon discounted	γ	$\frac{1}{1-\gamma}$	$\ v^*\ _\infty \leq \frac{r_{\max}}{1-\gamma}$	$\tilde{\Theta} \left(\frac{r_{\max}}{1-\gamma} \sqrt{\frac{1}{N_k(s, a)}} \right)$
finite-horizon	H	H	$\ v^*\ _\infty \leq r_{\max} H$	$r_{\max} H \sqrt{\frac{1}{N_k(s, a)}}$
average reward	?	$+\infty$	$\text{rng}(h^*) \leq c$ ▲ assumption	$\tilde{\Theta} \left(c \sqrt{\frac{1}{N_k(s, a)}} \right)$

UCRL2-like Exploration

For $k = 1, 2, \dots$

1. Estimation of model and uncertainty

$$\mathcal{M}_k = \{M = (S, \mathcal{A}, \tilde{p}, \tilde{r}) : \tilde{p}(\cdot | s, a) \in B_k^p(s, a), \tilde{r}(s, a) \in B_k^r(s, a)\}$$

2. Planning for optimistic policy

$$(M_k, \pi_k) = \arg \max_{M \in \mathcal{M}_k} \max_{\pi} \{g^\pi(M)\}$$

3. Execution of policy π_k

- execute action $a_t \sim \pi_k$
- observe reward r_t and next state s_{t+1}

- Estimation: \mathcal{M}_k is the set of plausible MDP such that

$$\|\tilde{p}(\cdot | s, a) - p(\cdot | s, a)\|_1 \leq \beta_k^p(s, a) \approx \sqrt{\frac{SL}{N_k(s, a)}}$$

$$|\tilde{r}(s, a) - r(s, a)| \leq \beta_k^r(s, a) \approx r_{\max} \sqrt{\frac{L}{N_k(s, a)}}$$

- Planning: use Extended Value Iteration

$$v_{n+1}(s) = \tilde{L}v_n = \max_a \left\{ \max_{\tilde{r} \in B_k^r(s, a)} \tilde{r} + \max_{\tilde{p} \in B_k^p(s, a)} \tilde{p}^\top v_n \right\}$$

SCAL⁺: tabular MDP

- Exploration bonus: Used in deep RL [Bellemare et al. 2016, Tang et al. 2017] and/or when the intrinsic horizon is known [Azar et al. 2017, Jin et al. 2018]

For $k = 1, 2, \dots$

1. Estimation of empirical model

$$M_k = (S, \mathcal{A}, \hat{p}, \hat{r} + b_k)$$

average rewards \hat{r} + exploration bonus b_k

$$\hat{p}_k(s' | s, a) = \frac{N_k(s, a, s')}{\sum_{s'} N_k(s, a, s')}$$

average transitions

2. Planning for optimistic policy

$$\pi_k = \arg \max_{\pi} \{g^\pi(M_k)\}$$

3. Execution of policy π_k

- Plan using the empirical MDP

- Bonus is used to recover optimism

$$b_k(s, a) \approx (r_{\max} + c) \sqrt{\frac{L}{N_k(s, a)}} < \beta_k^r(s, a) + c\beta_k^p(s, a)$$

Tighter uncertainty \Rightarrow better performance "equivalent" UCRL2 bonus

- ▲ prior knowledge:

$$\text{rng}(h^*) = \max_s h^*(s) - \min_s h^*(s) \leq c$$

SCAL⁺: regret

For any MDP such $\text{rng}(h^*) \leq c$, w.p. $1 - \delta$

$$R(\text{SCAL}^+, T) = \tilde{O} \left(c \sqrt{T \sum_{s, a} \Gamma(s, a)} \right)$$

- $\Gamma(s, a) = \|p(\cdot | s, a)\|_0$ (number of next states)

- Worst-case $\tilde{O}(cS\sqrt{AT})$
- as UCRL2 except that D is replaced by c

SCCAL⁺: continuous MDP

- S is continuous, \mathcal{A} is discrete
- MDP (reward and transitions) is Hölder continuous

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

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

- SCCAL⁺ combines SCAL⁺ with state aggregation

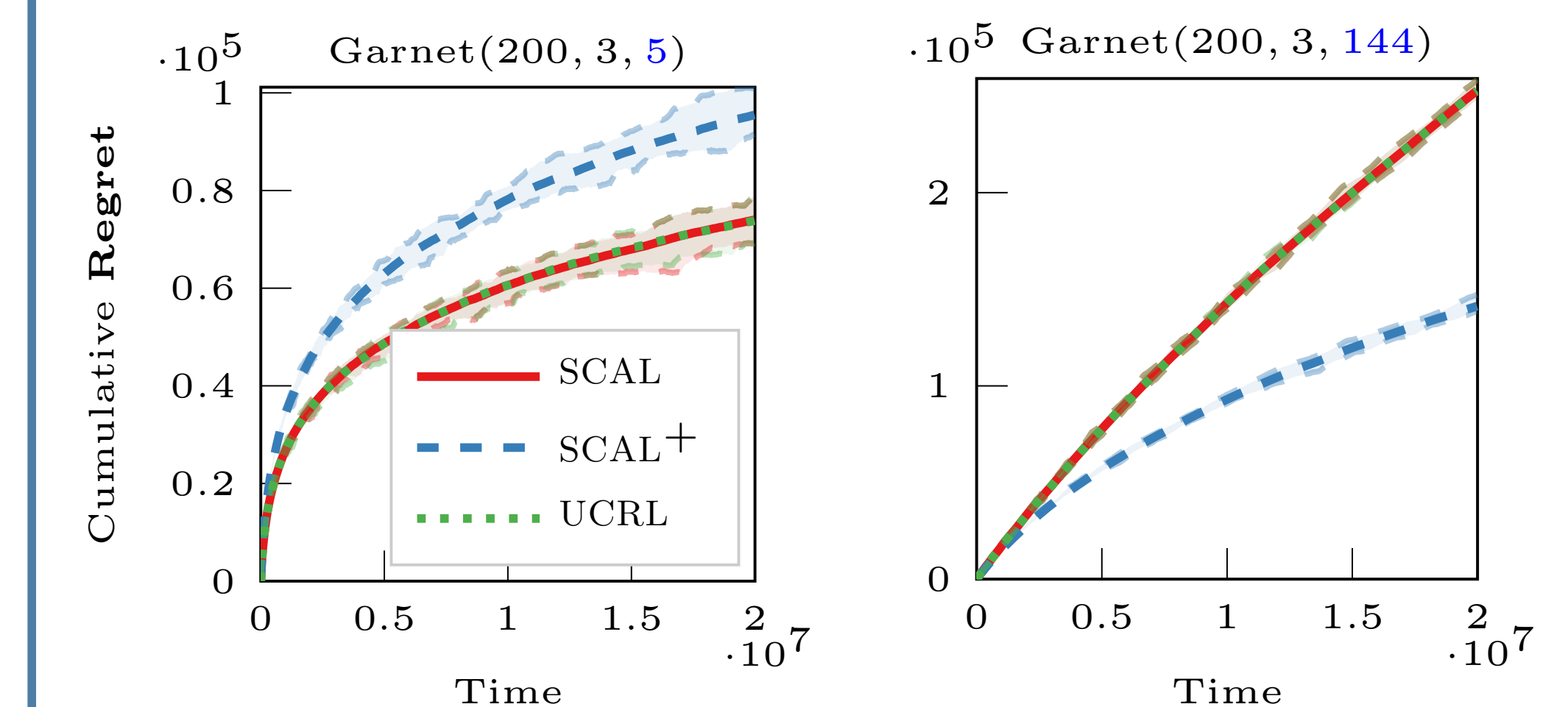
$$b(I, a) \approx (c + r_{\max}) \left(\frac{1}{\sqrt{N_k(I, a)}} + \underbrace{LMs^{-\alpha}}_{\text{bias}} \right)$$

REGRET: for any Hölder MDP w.p. $1 - \delta$

$$R(\text{SCCAL}^+, T) = \tilde{O} \left(cL\sqrt{AT}^{(\alpha+2)/(2\alpha+2)} \right)$$

Numerical Results

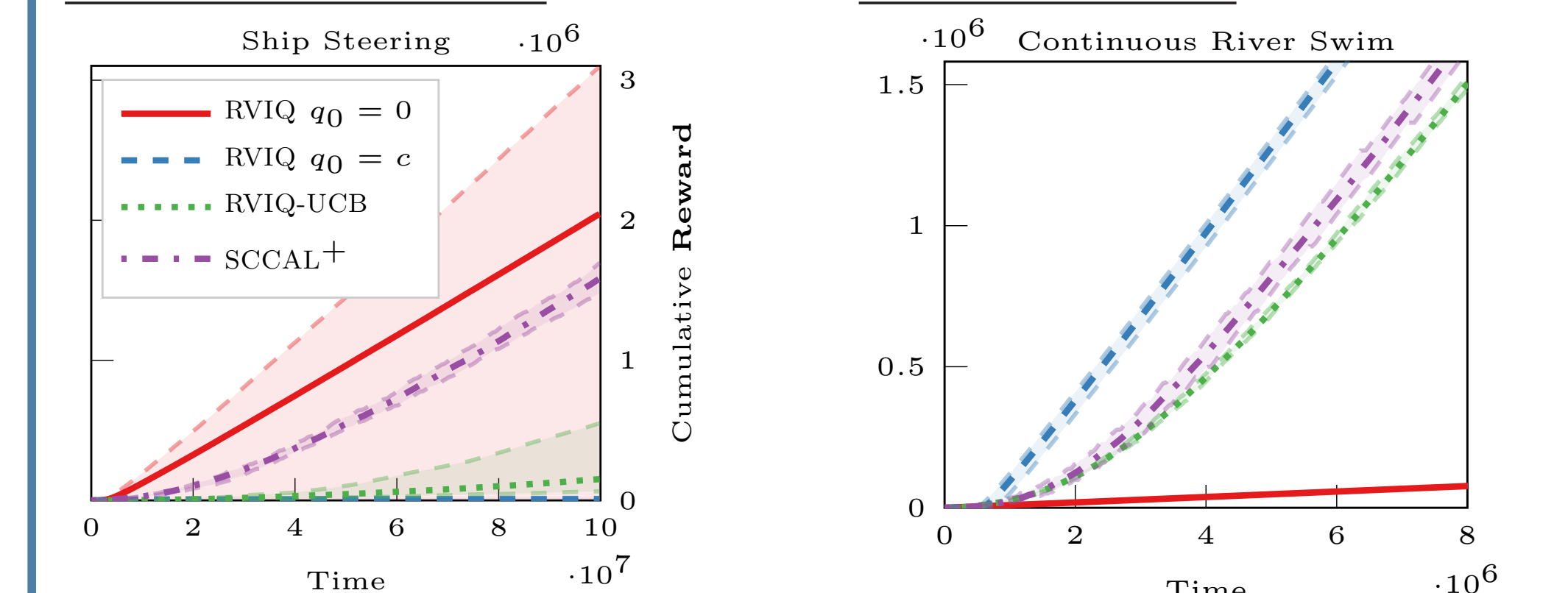
GARNET: tabular MDP



- SCAL⁺ can indeed outperform UCRL2

– Continuous MDPs –

SHIP STEERING:



- First implementable algorithm with guarantees in cont. MDPs (lots details are missing here, see paper)
- More stable than model-free version