

# Reinforcement Learning in a Birth and Death Process: Breaking the Dependence on the State Space

**Louis-Sébastien Rebuffi** Jonatha Anselmi, Bruno Gaujal

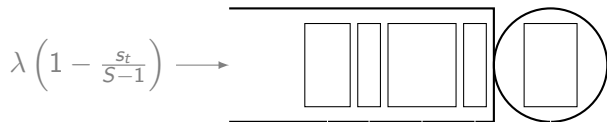
December 2nd, 2022

- 1 Introducing an Example of a basic MDP
- 2 Classic Reinforcement Learning Algorithms on MDPs
- 3 Our Contribution

## Example of a Simple Queue

Consider a processor with:

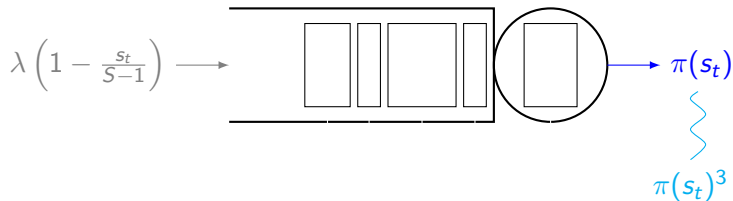
- Poisson arrivals with rate  $\lambda \left(1 - \frac{s_t}{S-1}\right)$ .



## Example of a Simple Queue

Consider a processor with:

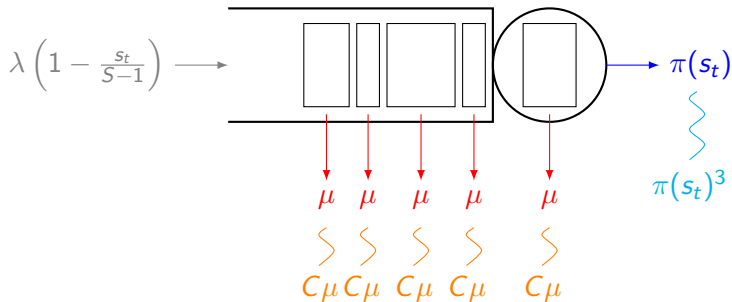
- Poisson arrivals with rate  $\lambda \left(1 - \frac{s_t}{S-1}\right)$ ,
- Service rate  $\pi(s_t)$  and power dissipation  $\pi(s_t)^3$ .



# Example of a Simple Queue

Consider a processor with:

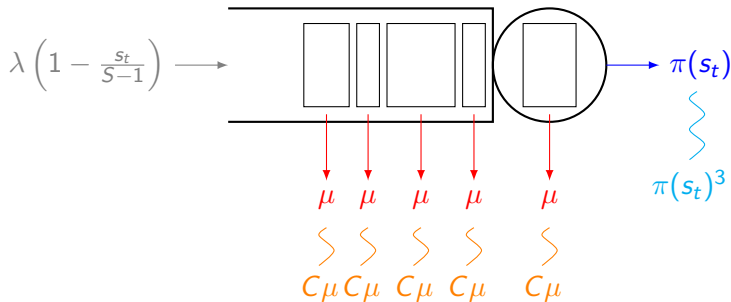
- Poisson arrivals with rate  $\lambda \left(1 - \frac{s_t}{S-1}\right)$ ,
- Service rate  $\pi(s_t)$  and power dissipation  $\pi(s_t)^3$ ,
- Jobs with Markovian deadlines that induce a cost  $C$  when dropped from the queue.



# Example of a Simple Queue

Consider a processor with:

- Poisson arrivals with rate  $\lambda \left(1 - \frac{s_t}{S-1}\right)$ ,
- Service rate  $\pi(s_t)$  and power dissipation  $\pi(s_t)^3$ ,
- Jobs with Markovian deadlines that induce a cost  $C$  when dropped from the queue.



Objective: Find the optimal speeds that minimize the long term energy spent by the processor and the cost of missed deadlines.

# Model as an MDP

Define the MDP  $\mathcal{M} := (\mathcal{S} = \{0, \dots, S-1\}, \mathcal{A} = \{0, \dots, A-1\}, Q, c)$ .

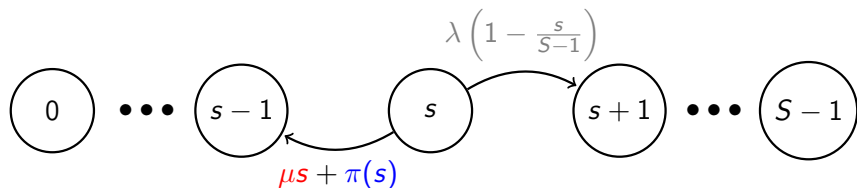


Figure: Transition diagram of the Markov chain induced by a policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ .

With the expected instant cost:

$$c(s, \pi) := C\mu s + \pi(s)^3.$$

## Definition of Rewards

Classically, different types of reward are considered, like the total discounted reward, for a discount  $\gamma < 1$ :

$$\rho_\gamma(\pi) := \lim_{T \rightarrow \infty} \sum_{t=1}^T \gamma^t \mathbb{E}[r(s_t, \pi(s_t))].$$

or the **long-run average reward**:

$$\rho(\pi) := \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[r(s_t, \pi(s_t))].$$

While using discounted rewards would guarantee the stability of the system, the average reward is more suitable for queueing systems.



# Definition of the Regret

Let  $\rho^* := \sup_{\pi} \rho(\pi)$  be the optimal average reward.

## Definition (Regret)

For an MDP  $M$ , the *regret* at time  $T$  of a learning algorithm  $\mathbb{L}$  is:

$$\text{Reg}(M, \mathbb{L}, T) := T\rho^* - \sum_{t=1}^T r_t.$$

We place ourself in the context of tabular MDPs:

- Exploration: Learn the transitions probabilities and rewards for each state-action pair.
- Exploitation: Use the best policies to minimize the regret.

## Definition of the Bias

Define the bias of any policy  $\pi$ :

$$h_{\pi}(s) := \mathbb{E}_{\pi} \left[ \sum_{t=1}^{\infty} (r(s_t^{\pi}) - \rho(\pi)) \mid s_1^{\pi} = s \right], \quad \forall 0 \leq s \leq S - 1,$$

In the computations of the regret, we use the Bellman equation:

$$\rho(\pi) + h_{\pi}(s) = r_{\pi}(s) + \sum_{s'} P_{\pi}(s' \mid s) h_{\pi}(s').$$

## Relating the Diameter to the Bias

Usually, to control the bias, we need to introduce the diameter:

### Definition

Letting  $\tau(s'|\pi, s)$  be the time to go from  $s$  to  $s'$ , the the diameter  $D$  of the MDP is:

$$D := \max_{s \neq s'} \min_{\pi: \mathcal{S} \rightarrow \mathcal{A}} \mathbb{E} [\tau(s'|\pi, s)].$$

Note that the bias and the diameter are related:

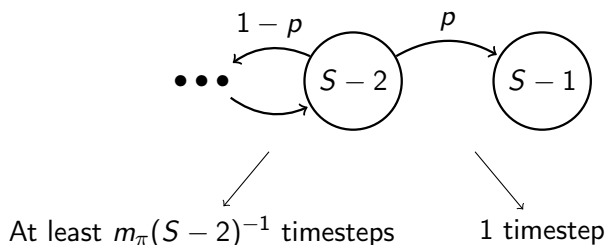
$$h_{\pi}(s) - h_{\pi}(s') \leq r_{\max} D_{\pi}$$

The diameter is an important quantity in the average reward case.

# Diameter in our Model

Queues have a **large diameter**.

→ We indeed have that:  $D_\pi \geq \mathbb{E}[\tau(S-1 \mid \pi, S-2)]$ . Letting  $p := \mathbb{P}(S-1 \mid S-2, \pi(S-2))$ .



In our example, the stationary measure  $m_\pi$  decreases exponentially, the diameter itself is therefore exponential in  $S$ .

In the average reward case:

- UCRL2-Jaksch et al. (2010), upper bound of the regret in  $\tilde{O}(r_{\max}DS\sqrt{AT})$ .
- UCRL2B Fruit et al. (2019), regret bounded in  $\tilde{O}(r_{\max}\sqrt{D\Gamma SAT})$  with  $\Gamma$  the highest number of neighbours of any state.
- Using additional information in the algorithm, such as an upper bound on the bias  $H$ , in Zhang et al. (2019) the regret is bounded in  $\tilde{O}(r_{\max})\sqrt{HSAT}$ .

In these algorithms, the parameters  $D$  and  $H$  still depend on  $S$ .

# Using the MDP Structure in the Literature

Examples of algorithms using the structure of MDPs:

- We could think to use linear mixture models with  $d$  parameters. With discount  $\gamma$ , the regret is upper bounded by  $r_{\max} d \sqrt{T} / (1 - \gamma)^2$  [Zhou et al. 2021], but it does not deal with the long-run average reward case.
- Model free algorithms with Q learning: Wei et al. (2020), regret bound close to  $\mathcal{O}\left(r_{\max} \sqrt{t_{\text{mix}}^3 SAT}\right)$  for ergodic MDPs, where  $t_{\text{mix}}$  is the mixing time of the MDP, which depends on  $S$ . This algorithm also requires additional information, such as a bound on  $t_{\text{mix}}$  and also on the worst hitting time.

## Theorem (Universal Lower Bound - Jaksch et al. (2010))

*For any learning algorithm  $\mathbb{L}$ , any  $S, A \geq 10$ ,  $D \geq 20 \log_A S$ , and  $T \geq DSA$ , there is an MDP  $M$ , such that:*

$$\mathbb{E}[\text{Reg}(M, \mathbb{L}, T)] \geq 0.015 r_{\max} \sqrt{DSAT}.$$

Existing learning algorithms have upper bounds in the regret that almost match this lower bound.

However this regret bound is unsatisfactory for queueing systems, where the diameter  $D$  is exponential in the number of states  $S$ .

MDPS for queueing systems have the following challenging characteristics:

- We consider the **long run average reward** rather than the total discounted reward.
- The considered MDPs have a **large diameter**  $D$ , *i.e.* a large expected time to cross the MDP.
- The transition matrices for queues are **sparse and structured**.



We have seen that the previous bounds depend on the diameter, meaning that they are inaccurate for birth and death processes.

### Question

When the underlying MDP has the structure of a queueing system, do the diameter  $D$  or the number of states  $S$  actually play a role in the regret?

To answer this question, we study the algorithm UCRL2 on our previous example.

# UCRL2 Algorithm

---

**Algorithm 1:** The UCRL2 algorithm.

---

Set  $s_1 = 0$

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

**Initialize** episode  $k$  with current reward and transition estimates  $\hat{r}_k$  and  $\hat{p}_k$ .

**Find** a policy  $\tilde{\pi}_k$  and an optimistic MDP  $\tilde{M}_k \in \mathcal{M}_k$ .

**Execute** policy  $\tilde{\pi}_k$  on the true MDP  $M$  until the end of the episode.

**end**

---

The optimistic MDP  $\tilde{M}_k$  with policy  $\tilde{\pi}_k$  is the queue of largest gain with rewards  $\tilde{r}$  and transitions  $\tilde{p}$  such that:

$$\forall(s, a), \quad |\tilde{r}(s, a) - \hat{r}_k(s, a)| \leq r_{\max} \sqrt{\frac{2 \log(At_k)}{\max\{1, N_{t_k}(s, a)\}}}$$

$$\forall(s, a), \quad \|\tilde{p}(\cdot|s, a) - \hat{p}_k(\cdot|s, a)\|_1 \leq \sqrt{\frac{8 \log(2At_k)}{\max\{1, N_{t_k}(s, a)\}}}$$

# Monotonicity

For the class of MDPs  $\mathcal{M}$ , we assume the following assumptions hold:

## Monotonicity

Denoting by  $\pi^0$  a reference policy and  $\pi$  any other policy, if  $s_0^\pi \leq_{st} s_0^{\pi^0}$ , then for all  $t$ ,  $s_t^\pi \leq_{st} s_t^{\pi^0}$ .

The reference policy controls the number of visits of a given state regardless of the chosen policy for the current episode:

## Lemma

For  $f : \mathcal{S} \rightarrow \mathbb{R}^+$  non-decreasing non-negative, we obtain

$$\mathbb{E} \left[ \sum_{s \geq 0} f(s) N_t(s) \right] \leq t \sum_{s \geq 0} f(s) m^{\pi^0}(s).$$

# Bound on the Bias

Reminder: define the optimal bias:

$$h_{\pi^*}(s) := \mathbb{E}_{\pi^*} \left[ \sum_{t=1}^{\infty} \left( r \left( s_t^{\pi^*} \right) - \rho(\pi) \right) \mid s_1^{\pi^*} = s \right], \quad \forall 0 \leq s \leq S - 1,$$

## Bias Bound

There is a positive, bounded function  $\Delta$  such that:

$$-\Delta(s) \leq H(s) - H(s - 1) \leq 0.$$

Here, in our example  $\Delta(s) = C$  is constant.

# Main Result

Independently of  $S$  and  $D$ , let  $E_2 := \left(\sum_{s \in \mathcal{S}} f(s)^{-1}\right) \mathbb{E}_{m^{\pi^0}} [(\Delta + r_{\max})^2 f]$ , where  $f : s \mapsto \frac{\max\{1, s(s-1)\}}{(\Delta(s) + r_{\max})^2}$ .

## Theorem

*The expected regret achieved by UCRL2 is upper bounded as follows:*

$$\mathbb{E} [\text{Reg}(M, \text{UCRL2}, T)] \leq 19\sqrt{E_2 AT \log(2AT)} + \mathcal{O}\left(T^{1/4}\right),$$

*where the lower order term contains terms polynomial in  $D$  and  $S$ .*

In the example,  $E_2 \leq 12r_{\max}^2 \left(1 + \frac{\lambda^2}{\mu^2}\right)$ , so that the regret satisfies

$$\mathbb{E} [\text{Reg}(M, \text{UCRL2}, T)] = \mathcal{O}\left(r_{\max} \sqrt{AT \left(1 + \frac{\lambda^2}{\mu^2}\right) \log(AT)}\right).$$

- Despite using a basic and non-specific reinforcement learning algorithm, the analysis of the regret can be greatly improved when studying queueing systems.
- The regret bounds should not involve  $D$  nor  $S$ . Our bound relies instead on the stationary measure of a reference policy.
- This type of regret bound could be generalized to other queueing systems, such as optimal routing or admission control for example.