

Online Markov Decision Processes under Bandit Feedback



Gergely Neu^{*,†}

gergely.neu@gmail.com

^{*}Department of Computer Science and Information Theory
Budapest University of Technology and Economics, Hungary

Csaba Szepesvári

szepesva@ualberta.ca

Department of Computing Science
University of Alberta, Canada

András György

gya@szit.bme.hu

[†]Machine Learning Research Group

MTA SZTAKI Institute for Computer Science and Control, Hungary

András Antos

antos@szit.bme.hu

Machine Learning Research Group

MTA SZTAKI Institute for Computer Science and Control, Hungary



Abstract

We consider online learning in finite, stochastic Markovian environments where in each time step a new reward function is chosen by an oblivious adversary. The goal of the learning agent is to compete with the best stationary policy in terms of the total reward received. In each time step the agent observes the current state and the reward associated with the last transition, however, the agent does not observe the rewards associated with other state-action pairs. The agent is assumed to know the transition probabilities. The state of the art result for this setting is a no-regret algorithm. In this paper we propose a new learning algorithm and, assuming that stationary policies mix uniformly fast, we show that the expected regret of the new algorithm in T time steps is $\mathcal{O}(T^{2/3}(\ln T)^{1/3})$, giving the first rigorously proved regret bound for the problem.

Assumptions

- **Assumption A1** Every policy π has a well-defined unique stationary distribution μ^π .
- **Assumption A2** The stationary distributions are uniformly bounded away from zero: $\inf_{\pi, x} \mu^\pi(x) \geq \beta > 0$.
- **Assumption A3** There exists some fixed positive **mixing time** τ such that for any two arbitrary μ and μ' over \mathcal{X} ,

$$\sup_{\pi} \|(\mu - \mu')P^\pi\|_1 \leq e^{-1/\tau} \|\mu - \mu'\|_1.$$

Definitions

- Value functions and average rewards:

$$\rho_t^\pi = \lim_{S \rightarrow \infty} \frac{1}{S} \sum_{s=1}^S \mathbb{E}[r_t(\mathbf{x}'_s, \mathbf{a}'_s)],$$

$$q_t^\pi(x, a) = \mathbb{E} \left[\sum_{s=1}^{\infty} (r_t(\mathbf{x}'_s, \mathbf{a}'_s) - \rho_t^\pi) \middle| \mathbf{x}'_1 = x, \mathbf{a}'_1 = a \right],$$

$$v_t^\pi(x, a) = \mathbb{E} \left[\sum_{s=1}^{\infty} (r_t(\mathbf{x}'_s, \mathbf{a}'_s) - \rho_t^\pi) \middle| \mathbf{x}'_1 = x \right].$$

- At time t , use only experience gathered up to time step $t - N$ and define

$$\mu_{t,x,a}^N(x') = \mathbb{P}[\mathbf{x}'_t = x' | \mathbf{x}_{t-N} = x, \mathbf{a}_{t-N} = a, \pi_{t-N+1}, \dots, \pi_{t-1}],$$

so that $\mu_{t,x,a}^N$ is positive.

- Estimate reward as

$$\hat{r}_t(x, a) = \begin{cases} \frac{r_t(x, a)}{\pi_t(a|x) \mu_{t,x,a}^N(x|\mathbf{x}_{t-N}, \mathbf{a}_{t-N})}, & \text{if } (x, a) = (\mathbf{x}_t, \mathbf{a}_t); \\ 0, & \text{otherwise.} \end{cases}$$

Mixing ensures that the probability of visiting state x at time t is positive for all x and t , that is,

$$\pi_t(a|x) \mu_{t,x,a}^N(x|\mathbf{x}_{t-N}, \mathbf{a}_{t-N}) > 0.$$

- Let $\hat{\rho}_t = \sum_{x,a} \mu_{t,x,a}^N(x|\mathbf{x}_{t-N}, \mathbf{a}_{t-N}) \hat{r}_t(x, a)$ and solve, for all x, a , the Bellman equations

$$\hat{q}_t(x, a) = \hat{r}_t(x, a) - \hat{\rho}_t + \sum_{x', a'} P(x'|x, a) \pi_t(a'|x') \hat{q}_t(x', a').$$

Algorithm

Set $N \geq 1$, $\mathbf{w}_1(x, a) = \mathbf{w}_2(x, a) = \dots = \mathbf{w}_N(x, a) = 1$, $\gamma \in (0, 1)$, $\eta \in (0, \gamma]$.

For $t = 1, 2, \dots, T$, repeat

1. Set

$$\pi_t(a|x) = (1 - \gamma) \frac{\mathbf{w}_t(x, a)}{\sum_b \mathbf{w}_t(x, b)} + \frac{\gamma}{|\mathcal{A}|}$$

for all $(x, a) \in \mathcal{X} \times \mathcal{A}$.

2. Draw an action \mathbf{a}_t randomly, according to the policy $\pi_t(\cdot|\mathbf{x}_t)$.

3. Receive reward $r_t(\mathbf{x}_t, \mathbf{a}_t)$ and observe \mathbf{x}_{t+1} .

4. If $t \geq N + 1$

(a) Compute $\mu_{t,x,a}^N(x|\mathbf{x}_{t-N}, \mathbf{a}_{t-N})$ for all $x \in \mathcal{X}$.

(b) Construct estimates \hat{r}_t and compute \hat{q}_t using the Bellman equations for π_t .

(c) Set $\mathbf{w}_{t+N}(x, a) = \mathbf{w}_{t+N-1}(x, a) e^{\eta \hat{q}_t(x, a)}$ for all $(x, a) \in \mathcal{X} \times \mathcal{A}$.

Main result

Theorem 1. Let $N = \lceil \tau \ln T \rceil$,

$$C = (2\tau + 4)\tau |\mathcal{A}| \ln T + (3\tau + 1)^2,$$

$$\eta = T^{-2/3} \cdot (\ln |\mathcal{A}|)^{2/3} \cdot \left(\frac{4C(\tau + 2)}{\beta} \right)^{-1/3},$$

$$\gamma = T^{-1/3} \cdot (2\tau + 4)^{-2/3} \cdot \left(\frac{2C \ln |\mathcal{A}|}{\beta} \right)^{1/3}.$$

Then

$$\hat{L}_T \leq 3T^{2/3} \cdot \left(\frac{(4\tau + 8) \ln |\mathcal{A}|}{\beta} C \right)^{1/3} + \mathcal{O}(T^{1/3}).$$

Proof

The proof is based on ideas from Even-Dar et al. [2009]. The complication is of course that rewards are estimated. The regret can be decomposed as

$$R_T^\pi - \hat{R}_T = \underbrace{\left(R_T^\pi - \sum_{t=1}^T \rho_t^\pi \right)}_{(i)} + \underbrace{\left(\sum_{t=1}^T \rho_t^\pi - \sum_{t=1}^T \rho_t^{\pi_t} \right)}_{(ii)} + \underbrace{\left(\sum_{t=1}^T \rho_t^{\pi_t} - \hat{R}_T \right)}_{(iii)}.$$

Bounding (i) After Even-Dar et al. [2009]:

$$R_T^\pi - \sum_{t=1}^T \rho_t^\pi \leq 2\tau + 2.$$

The policies π_t change slowly

Lemma 1. Let $c = \frac{2\eta}{\beta} \left(\frac{1}{\gamma} + 4\tau + 6 \right)$. Assume that $c(3\tau + 1)^2 < \beta/2$ and $N \geq \lceil \tau \ln \left(\frac{4}{\beta - 2c(3\tau + 1)^2} \right) \rceil$. Then, for all $N < t \leq T$, $x, x' \in \mathcal{X}$ and $a \in \mathcal{A}$, we have

$$\mu_{t,x,a}^N(x') \geq \beta/2$$

and

$$\max_{x'} \sum_a |\pi_{t+1}(a'|x') - \pi_t(a'|x')| \leq c.$$

Bounding (ii)

After Even-Dar et al. [2009]:

$$\rho_t^\pi - \rho_t^{\pi_t} = \sum_{x,a} \mu^\pi(x) \pi(a|x) \left[q_t^{\pi_t}(x, a) - v_t^{\pi_t}(x) \right].$$

A simple modification of the proof of the regret bound of **Exp3** yields the following:

Proposition 1. Let $N \geq \lceil \tau \ln T \rceil$. For any policy π and for all T large enough, we have

$$\sum_{t=1}^T \mathbb{E}[\rho_t^\pi - \rho_t] \leq (4\tau + 10)N + \frac{\ln |\mathcal{A}|}{\eta} + (2\tau + 4)T \left(\gamma + \frac{2\eta}{\beta} |\mathcal{A}| (Nc + (e - 2)(2\tau + 4)) \right).$$

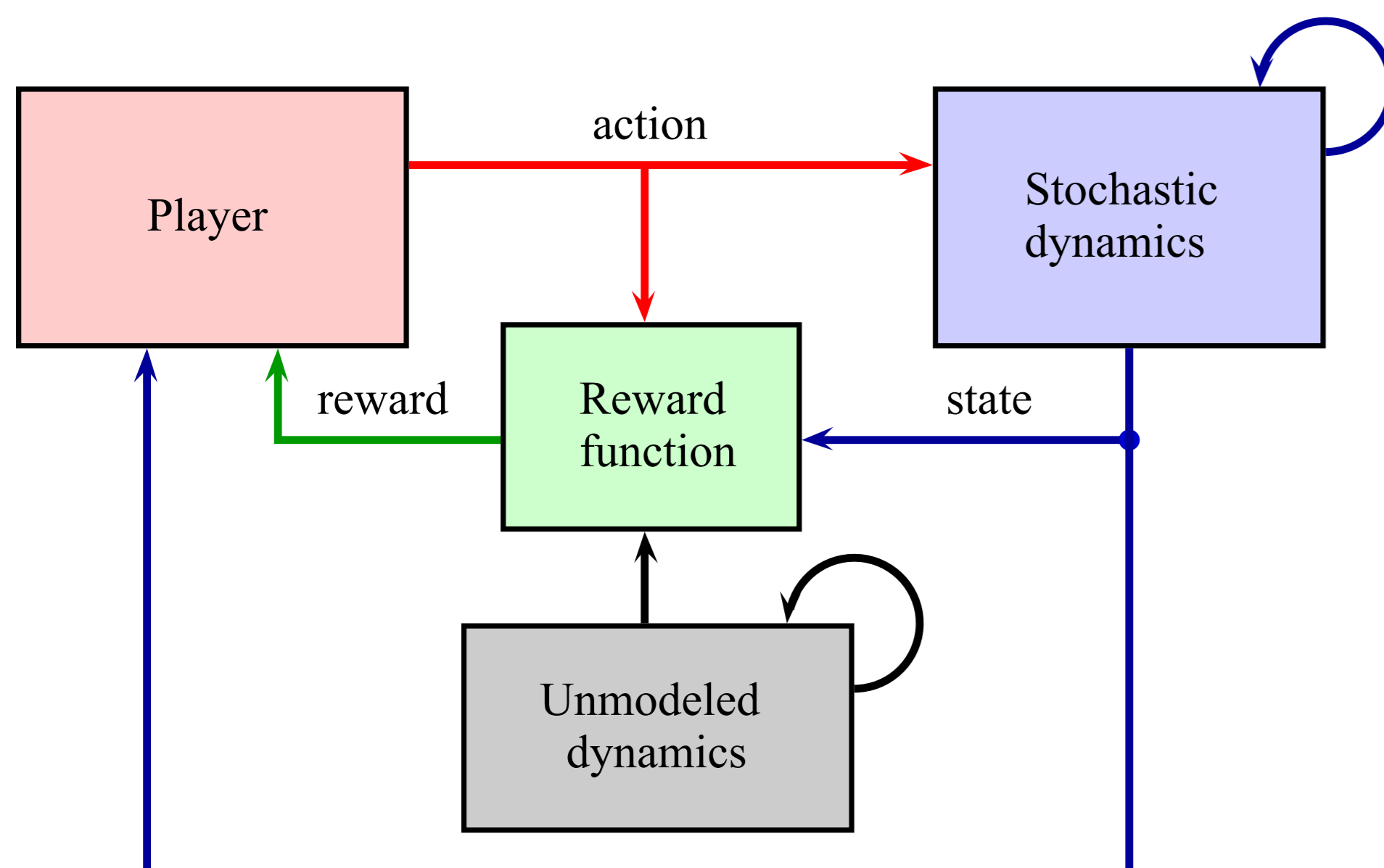
Bounding (iii)

Proposition 2. Let $N \geq \lceil \tau \ln T \rceil$. For any T large enough,

$$\sum_{t=1}^T \mathbb{E}[\rho_t] - \hat{R}_T \leq Tc(3\tau + 1)^2 + 2Te^{-N/\tau} + 2N.$$

Follows from the slow change of policies π_t .

Online Markov Decision Processes



- Reward sequence: an **arbitrary sequence** r_1, r_2, \dots, r_T **fixed in advance**, where $r_t: \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$.
- History up to time t : $\mathbf{u}_t = (\mathbf{x}_1, \mathbf{a}_1, r_1(\mathbf{x}_1, \mathbf{a}_1), \mathbf{x}_2, \mathbf{a}_2, r_2(\mathbf{x}_2, \mathbf{a}_2), \dots, \mathbf{x}_t, \mathbf{a}_t, r_t(\mathbf{x}_t, \mathbf{a}_t))$
- At time step t :
 - agent selects action \mathbf{a}_t based on \mathbf{u}_{t-1} and \mathbf{x}_t ;
 - observes reward $r_t(\mathbf{x}_t, \mathbf{a}_t)$;
 - observes new state $\mathbf{x}_{t+1} \sim P(\cdot|\mathbf{x}_t, \mathbf{a}_t)$.

The learning problem

- Policy: $\pi_t(a|x) = \mathbb{P}[\mathbf{a}_t = a | \mathbf{x}_t = x, \mathbf{u}_{t-1}]$

- Expected total reward of the player:

$$\hat{R}_T = \mathbb{E} \left[\sum_{t=1}^T r_t(\mathbf{x}_t, \mathbf{a}_t) \right],$$

where $\mathbf{a}_t \sim \pi_t(\cdot|\mathbf{x}_t)$ and $\mathbf{x}_{t+1} \sim P(\cdot|\mathbf{x}_t, \mathbf{a}_t)$

- Expected total reward of a fixed policy π :

$$R_T^\pi = \mathbb{E} \left[\sum_{t=1}^T r_t(\mathbf{x}'_t, \mathbf{a}'_t) \right],$$

where $\mathbf{a}'_t \sim \pi(\cdot|\mathbf{x}'_t)$ and $\mathbf{x}'_{t+1} \sim P(\cdot|\mathbf{x}'_t, \mathbf{a}'_t)$

- Goal: minimize **regret**

$$\hat{L}_T = \sup_{\pi} R_T^\pi - \hat{R}_T.$$

Previous work

- Full information: $\hat{L}_T = \mathcal{O}(\sqrt{T})$: Even-Dar et al. [2009]
- Bandit information: $\hat{L}_T = \mathcal{O}(T)$: Yu et al. [2009]
- Bandit information for episodic loop-free MDPs: $\hat{L}_T = \mathcal{O}(\sqrt{T})$ Neu et al. [2010]

References

References

- E. Even-Dar, S. M. Kakade, and Y. Mansour. Online Markov decision processes. *Mathematics of Operations Research*, 34(3):726–736, 2009.
- G. Neu, A. György, and Cs. Szepesvári. The online loop-free stochastic shortest-path problem. In *COLT-10*, pages 231–243, 2010.
- J. Y. Yu, S. Mannor, and N. Shimkin. Markov decision processes with arbitrary reward processes. *Mathematics of Operations Research*, 34(3):737–757, 2009.