

EVALUATION  
ONLINE LEARNING  
LINKS WITH OPTIMIZATION AND GAMES  
UNIVERSITÉ PARIS–SACLAY



ADAPTIVE DIAGONAL SCALINGS FOR Q-LEARNING

*This project requires familiarity with reinforcement learning<sup>1</sup>*

Let  $\lambda \in (0, 1)$ . Consider a Markov decision process  $(\mathcal{S}, \mathcal{A}, \mathcal{R}, p)$  where  $\mathcal{S}$  is the set of states,  $\mathcal{A}$  the set of actions,  $\mathcal{R} \subset \mathbb{R}$  the set of possible rewards and  $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{R}$  the transition function where

$$p(r, s' | s, a) := p(s, a, r, s'), \quad (s, a, r, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{R} \times \mathcal{S},$$

corresponds to the probability of obtaining reward  $r$  and moving to state  $s'$  when action  $a$  is chosen at state  $s$ . All sets are finite.

An action-value function is a vector  $q = (q(s, a))_{(s,a) \in \mathcal{S} \times \mathcal{A}} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ . The Bellman optimality operator (for action-value functions)  $B_* : \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \rightarrow \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$  is defined as

$$(B_* q)(s, a) = \sum_{(r,s') \in \mathcal{R} \times \mathcal{S}} p(r, s' | s, a) \left( r + \lambda \max_{a' \in \mathcal{A}} q(s', a') \right), \quad (s, a) \in \mathcal{S} \times \mathcal{A},$$

---

<sup>1</sup>see e.g. <https://joon-kwon.github.io/rl-ups/reinforcement-learning-lecture-notes.pdf>

where we simply denote  $B_*q$  instead of  $B_*(q)$ .  $B_*$  is known to be a contraction: it thus admits a unique fixed point  $q_*$ , which is the optimal action-value function.

Without knowledge of  $p$ , evaluations of the map  $B_*$  cannot be computed, but a stochastic estimator can be obtained as follows. For  $q \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$  and  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , if  $(R, S') \sim p(\cdot | s, a)$ , in other words if  $R, S'$  are the actual (random) reward and new state obtained by choosing action  $a$  at state  $s$ , then

$$(\hat{B}_*q)(R, S') = R + \lambda \max_{a \in \mathcal{A}} q(S', a)$$

is an unbiased estimator of  $(B_*q)(s, a)$ .

Traditional Q-learning is defined as follows. Let  $q_0 \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$  be an initial action-value function. For all  $t \geq 0$ , let  $(S_t, A_t, R_t, S'_t)$  be such that  $(R_t, S'_t) | S_t, A_t \sim p(\cdot | S_t, A_t)$  (often,  $S'_t = S_{t+1}$ , unless the episode terminates), and

$$q_{t+1}(s, a) = \begin{cases} (1 - \gamma_t)q_t(s, a) + \gamma_t ((\hat{B}_*q_t)(R_t, S'_t)) & \text{if } (s, a) = (S_t, A_t) \\ q_t(s, a) & \text{otherwise,} \end{cases}$$

where  $\gamma_t \in (0, 1)$ . Q-learning is therefore an asynchronous<sup>2</sup> stochastic fixed point iteration.

- 1) Similarly to the way AdaGrad-Norm is used to solve fixed point problems, define AdaGrad-Diagonal in the context of Q-learning.

Numerous variants of AdaGrad have achieved great success in deep learning optimization. We here consider RMSprop and Adam. Let  $d \geq 1$  and  $x_0 \in \mathbb{R}^d$ . For a sequence  $(u_t)_{t \geq 0}$  in  $\mathbb{R}^d$ ,  $\gamma > 0$ , the associated RMSprop (resp. Adam) iterates are defined component-wise for  $t \geq 0$ , and  $0 \leq i \leq d$  as

$$x_{t+1,i} = x_{t,i} + \frac{\gamma}{\sqrt{\sum_{\tau=0}^t \beta^{t-\tau} u_{\tau,i}^2}} u_{t,i},$$

$$\left( \text{resp. } x_{t+1,i} = x_{t,i} + \frac{\gamma}{\sqrt{\sum_{\tau=0}^t \beta_2^{t-\tau} u_{\tau,i}^2}} \sum_{\tau=0}^t \beta_1^{t-\tau} u_{\tau,i} \right),$$

where  $\beta = .99$ ,  $\beta_1 = .9$  and  $\beta_2 = .999$  are common default values.

- 2) Adapt RMSprop and Adam to the context of Q-learning.

---

<sup>2</sup>meaning that not all components are updated at each iteration

- 3) Perform numerical experiments to compare the performance of the above algorithms with traditional Q-learning. Consider environments with finite number of states and actions from e.g. the Gymnasium package<sup>3</sup>.



---

<sup>3</sup><https://gymnasium.farama.org/>