

# REINFORCEMENT LEARNING

STS/H – Harmonic step sizes. Show that the step sizes

$$\alpha_n := \frac{1}{an+b}, \quad a, b \in \mathbb{R},$$

(where  $a, b \in \mathbb{R}$  are chosen such that  $an+b \neq 0$ ) satisfy the convergence conditions

$$\sum_{n=1}^{\infty} \alpha_n = \infty, \quad \sum_{n=1}^{\infty} \alpha_n^2 < \infty.$$

$$\alpha_n = \frac{1}{an+b} = \frac{1}{a} \cdot \underbrace{\frac{1}{n+\frac{b}{a}}}_{=: \tilde{\alpha}_n}$$

da  $\frac{1}{a}$  aus den Summen herausgezogen werden kann, reicht es  $a=1$  zu betrachten.

$\Rightarrow$  OBdA  $a=1, b \in \mathbb{R}$

$$\begin{aligned} \sum_{n=0}^{\infty} \alpha_n &= \sum_{n=0}^{\infty} \frac{1}{n+b} \\ &= \sum_{n=0}^{\infty} \frac{1}{n+|b|+(b-|b|)} \\ \underline{b \geq 0} \quad &= \sum_{n=|b|}^{\infty} \frac{1}{n+(b-|b|)} \\ &\geq \sum_{n=|b|}^{\infty} \frac{1}{n+1} \\ &= \sum_{n=|b|}^{\infty} \frac{1}{n} = \infty \end{aligned} \quad \begin{aligned} \underline{b < 0} \quad &= \sum_{n=0}^{|b|} (\dots) + \sum_{n=|b|}^{\infty} (\dots) \\ &\geq \sum_{n=0}^{\infty} \frac{1}{n} \end{aligned}$$

D

$$\sum_{n=0}^{\infty} \alpha_n^2 = \sum_{n=|b|}^{\infty} \frac{1}{(n+(b-|b|))^2} \quad \leftarrow \text{if } b \in -\mathbb{N} \text{ division by 0 happens} \\ \Rightarrow b \in \mathbb{R}/(-\mathbb{N})$$

$$\underline{b < 0} \quad = \sum_{n=|b|}^0 (\dots) + \sum_{n=0}^{\infty} (\dots) \\ \underbrace{\leq \infty}_{\leq \sum \frac{1}{n^2} < \infty}$$

B

$$\underline{b \geq 0} \quad \leq \sum_{n=|b|}^{\infty} \frac{1}{n^2} < \infty$$

sts/U – Unbiased step sizes. We use the iteration

$$Q_1 \in \mathbb{R},$$

$$Q_{n+1} := Q_n + \alpha_n(R_n - Q_n), \quad n \geq 1,$$

to estimate  $Q_n$  using  $R_n$ , where

$$\alpha_n := \frac{\alpha}{\beta_n}, \quad \alpha \in (0, 1), \quad n \geq 1,$$

and

$$\beta_0 := 0,$$

$$\beta_n := \beta_{n-1} + \alpha(1 - \beta_{n-1}), \quad n \geq 1.$$

Show that the iteration for  $Q_n$  above yields an exponential recency-weighted average *without initial bias* (i.e., the  $Q_n$  do not depend on the initial value  $Q_1$ ).

$$Q_2 = Q_1 + \alpha_1(R_1 - Q_1) \quad // \alpha_1 = \frac{\alpha}{\beta_1} = 1$$

$$= Q_1 + 1 \cdot (R_1 - Q_1)$$

$$= R_1 \Rightarrow Q_n \text{ is independent of } Q_1 \forall n \geq 2$$

$$n=1 \quad Q_1 = Q_1$$

~~$$n=2 \quad Q_2 = Q_1 + \alpha_1(R_1 - Q_1) = Q_1(1 - \alpha_1) + R_1\alpha_1$$~~

~~$$n=3 \quad Q_3 = Q_2 + \alpha_2(R_2 - Q_2)$$~~

$$= Q_1(1 - \alpha_1) + R_1\alpha_1 + R_2\alpha_2 - Q_1\alpha_1(1 - \alpha_1) - R_1\alpha_1\alpha_2$$

$$= Q_1(1 - \alpha_1)(1 - \alpha_2) + R_1\alpha_1(1 - \alpha_2) + R_2\alpha_2$$

~~$$n=N \quad Q_N = Q_1 \prod_{i=1}^{N-1} (1 - \alpha_i) + \sum_{i=1}^{N-1} R_i \alpha_i \prod_{k=i+1}^{N-1} (1 - \alpha_k)$$~~

~~$$\text{IA: } n=1 \quad Q_1 = Q_1 + 0 = Q_1 \quad \checkmark$$~~

~~$$\text{IS: } n \rightarrow n+1 \quad // \text{IV: } Q_n = Q_1 \prod_{i=1}^{n-1} (1 - \alpha_i) + \sum_{i=1}^{n-1} R_i \alpha_i \prod_{k=i+1}^{n-1} (1 - \alpha_k)$$~~

~~$$Q_{n+1} = Q_n + \alpha_n(R_n - Q_n)$$~~

~~$$= Q_n(1 - \alpha_n) + R_n\alpha_n$$~~

$$\begin{aligned}
 &= Q_1 \underbrace{\left( \prod_{i=1}^{n-1} (1-\alpha_i) \right) (1-\alpha_n)}_{= \prod_{i=1}^n (1-\alpha_i)} + \underbrace{\left[ \sum_{i=1}^{n-1} R_i \alpha_i \prod_{k=i+1}^{n-1} (1-\alpha_k) \right] (1-\alpha_n)}_{\substack{n < 0 \\ = \sum_{i=1}^{n-1} R_i \alpha_i \prod_{k=i+1}^n (1-\alpha_k)}} + R_n \alpha_n \\
 &= \sum_{i=1}^{n-1} R_i \alpha_i \prod_{k=i+1}^n (1-\alpha_k)
 \end{aligned}$$

$$= Q_1 \prod_{i=1}^n (1-\alpha_i) + \sum_{i=1}^n R_i \alpha_i \prod_{k=i+1}^n (1-\alpha_k)$$

We need to show that  $\prod_{i=1}^n (1-\alpha_i) \xrightarrow{n \rightarrow \infty} 0$

$$\beta_n = 1 - \sum \binom{n}{k} (-1)^{n-k} \alpha^k = 1 - (1-\alpha)^n$$

$$\boxed{IA: n=1} \quad \beta_1 = 1 - (1-\alpha) = \alpha \quad \checkmark$$

$$\boxed{IS: n \rightarrow n+1} \quad \beta_{n+1} = \beta_n + \alpha(1-\beta_n)$$

$$= \beta_n(1-\alpha) + \alpha$$

$$= (1-\alpha) - (1-\alpha)^{n+1} \quad \times$$

$$= 1 - (1-\alpha)^{n+1} \quad \checkmark$$

$$\begin{aligned}
 \Rightarrow (1-\alpha_n) &= 1 - \frac{\alpha}{1-(1-\alpha)^n} = \frac{1 - (1-\alpha)^n - \alpha}{1 - (1-\alpha)^n} \\
 &= (1-\alpha) \frac{1 - (1-\alpha)^{n-1}}{1 - (1-\alpha)^n} \\
 &= (1-\alpha) \frac{\beta_{n-1}}{\beta_n}
 \end{aligned}$$

$$\|\beta_n = 1 - ((1-\alpha)^n) \xrightarrow{n \rightarrow \infty} 1 \Rightarrow \frac{\beta_{n-1}}{\beta_n} \xrightarrow{n \rightarrow \infty} 1$$

$$\rightarrow (1-\alpha)$$

$$\Rightarrow \prod_{i=1}^n (1-\alpha_i) \xrightarrow{n \rightarrow \infty} 0$$

D

**AS** -  $\epsilon$ -greedy action selection. Suppose  $|\mathcal{A}| = 2$  and  $\epsilon = 1/2$ . When using  $\epsilon$ -greedy action selection, what is the probability that the greedy action is selected?

$$\begin{aligned} P(\text{greedy action picked}) &= P(\text{greedy choice}) \\ &\quad + P(\text{non-greedy choice AND greedy action picked}) \\ &= 0,5 + 0,5 \cdot \frac{1}{|\mathcal{A}|} = 0,75 \end{aligned}$$

**ENV/EX** - Think of application. Think of a (preferably creative) application of reinforcement learning. Specify the states, actions, and rewards as well as what is needed to satisfy the Markov property.

A program that designs new proteins. The state would be the current sequence of amino acids; The actions would correspond to the possible amino acids that can be added in the next step; The reward would correspond to the chemical properties of the protein (for example, one could look for antibacterial properties; these could be determined by experiment, from known proteins or by simulations). // Alternative football manager: states position and state of all players' action for every player next action rewards goals - enemy goals

**ENV/COUNTEREX** - Goal-directed learning task that is not an MDP. Try to find a goal-directed learning task that cannot be represented by a Markov decision process.

A MDP requires, that a description exists such that the current state contains all information of a system. This is not fulfilled for systems where:

- 1) The state cannot be completely measured because it is infeasible  
(For example a weather prediction method would require precise measurements of temp, humidity etc. at every point in space)
- 2) The state cannot be measured because it is physically not possible  
(e.g. a cosmic ray prediction system. EM travels with c as does information)

$$\begin{aligned} v_{\pi}(s) &= \mathbb{E}_{\pi}(G_t | S_t = s) \\ &= \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} | S_t = s\right], \quad \forall s \in S \end{aligned}$$

$$\begin{aligned} q_{\pi}(s, a) &= \mathbb{E}_{\pi}(G_t | S_t = s, A_t = a) \\ &= \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} | S_t = s, A_t = a\right] \end{aligned}$$

MDP/V - Equation for  $v_{\pi}$ . Give an equation for  $v_{\pi}$  in terms of  $q_{\pi}$  and  $\pi$ .

$$\boxed{\sum_a \pi(a|s) q_{\pi}(s, a)}$$

$$S \xrightarrow{\pi} \underbrace{v_{\pi}}_{q_{\pi}(s, a)}$$

$$\mathbb{E}(X|Y, Z) = \sum_x x \mathbb{P}(X|Y, Z) \quad \text{a.}$$

$$( \mathbb{P}(X, Y) = \mathbb{P}(X|Y) \mathbb{P}(Y) )$$

$$\Rightarrow \mathbb{P}(X, Y|Z) = \mathbb{P}(X|Y, Z) \mathbb{P}(Y|Z) \Rightarrow \mathbb{P}(X|Y, Z) = \frac{\mathbb{P}(X, Y|Z)}{\mathbb{P}(Y|Z)} \\ = \frac{1}{\mathbb{P}(Y)} \sum_x x \mathbb{P}(X, Y|Z)$$

$$\Rightarrow \sum_Y \mathbb{E}(X|Y, Z) \mathbb{P}(Y|Z) = \mathbb{E}(X|Z)$$

In our case  $X \triangleq G_t$ ;  $Y \triangleq A_t$ ;  $Z \triangleq S_t$  and  
the expectation is with respect to  $\pi$ .

$$\Rightarrow v_{\pi}(s) = \mathbb{E}_{\pi}(G_t | S_t = s)$$

$$\begin{aligned} &= \sum_a \pi(a|s) \underbrace{\mathbb{E}_{\pi}(G_t | S_t = s, A_t = a)}_{\mathbb{P}(A_t = a | S_t = s) = q(s, a)} \\ &= q(s, a) \end{aligned}$$

**MDP/Q – Equation for  $q_\pi$ .** Give an equation for  $q_\pi$  in terms of  $v_\pi$  and the four-argument  $p$ .

$$\begin{aligned}
 p(s', r | s, a) &= P(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a) = p_{sa}(s', a) \\
 q_\pi(s, a) &= \mathbb{E}_\pi(G_t | S_t = s, A_t = a) \\
 &= \sum_{k=0}^{\infty} \gamma^k \mathbb{E}_\pi(R_{t+k+1} | S_t = s, A_t = a) \\
 &= \mathbb{E}_\pi(R_{t+1} | S_t = s, A_t = a) + \gamma \mathbb{E}_\pi(G_{t+1} | S_t = s, A_t = a) \\
 &= \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a) + \gamma \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} p(s', r | s, a) \mathbb{E}_\pi(G_{t+1} | S_t = s') \\
 &= \sum_{r, s'} p(s', r | s, a) [r + \gamma v_\pi(s')]
 \end{aligned}$$

**MDP/RET – Change of return.** In episodic tasks and in continuing tasks, how does the return  $G_t$  change if a constant  $c$  is added to all rewards  $R_t$ ?

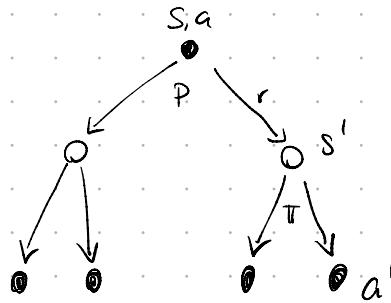
$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$\tilde{R}_t = R_t + c$$

$$\begin{aligned}
 \tilde{G}_t &= \sum_{k=0}^{\infty} \gamma^k \tilde{R}_{t+k+1} \\
 &= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} + \underbrace{\sum_{k=0}^{\infty} \gamma^k c}_{= \frac{c}{1-\gamma}}
 \end{aligned}$$

**MDP/BELLMAN/QPI – Bellman equation for  $q_\pi$ .** Analogous to the derivation of the Bellman equation for  $v_\pi$ , derive the Bellman equation for  $q_\pi$ .

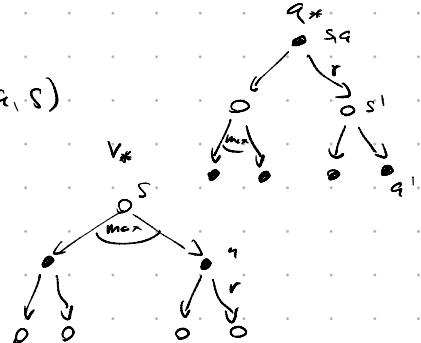
$$\begin{aligned}
 q_\pi(s, a) &= \mathbb{E}_\pi(G_t | S_t = s, A_t = a) \\
 &= \mathbb{E}_\pi(R_{t+1} + G_{t+1} | S_t = s, A_t = a) \\
 &= \sum_{r \in R} \sum_{s' \in S} p(s', r | s, a) r + \gamma \sum_{r, s'} p(s', r | s, a) \sum_{a'} \pi(a' | s') \mathbb{E}_\pi(G_{t+1} | S_{t+1} = s', A_{t+1} = a') \\
 &= \sum_{r, s'} p(s', r | s, a) [r + \gamma \sum_{a'} \pi(a' | s') q_\pi(s', a')]
 \end{aligned}$$



**MDP/VSTAR – Equation for  $v_*$ .** Give an equation for  $v_*$  in terms of  $q_*$ .

$$V_*(s) = \max_{a \in \mathcal{A}(s)} q_{\pi^*}(s, a)$$

$$\begin{aligned}
 V_*(s) &= \max_{\pi} V_\pi(s) = \max_{\pi} \sum_{a \in \mathcal{A}} \pi(a | s) q_\pi(a, s) \\
 &= \max_{a \in \mathcal{A}} q_*(a, s)
 \end{aligned}$$



**MDP/QSTAR – Equation for  $q_*$ .** Give an equation for  $q_*$  in terms of  $v_*$  and the four-argument  $p$ .

$$p(s', r | s, a) = P(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$$

$$\begin{aligned}
 q_*(s, a) &= \sum_{k=0}^{\infty} \gamma^k \mathbb{E}_{\pi^*}(R_{t+k+1} | S_t = s, A_t = a) \\
 &= \dots \quad // \text{see MDP/Q} \\
 &= \sum_{r, s'} p(s', r | s, a) [r + \gamma v_{\pi^*}(s')]
 \end{aligned}$$

MDP/PISTAR/VSTAR - Equation for  $\pi_*$ . Give an equation for  $\pi_*$  in terms of  $q_*$ .

$$\pi_*(a|s) = \begin{cases} 1, & a = \arg \max_a q_{\pi_*}(a, s) \\ 0, & \text{else} \end{cases}$$

↑ There might be multiple values where  $q_{\pi_*}(a, s)$  is max  
→ pick one

MDP/PISTAR/QSTAR - Equation for  $\pi_*$ . Give an equation for  $\pi_*$  in terms of  $v_*$  and the four-argument  $p$ .

Combine the last two answers.