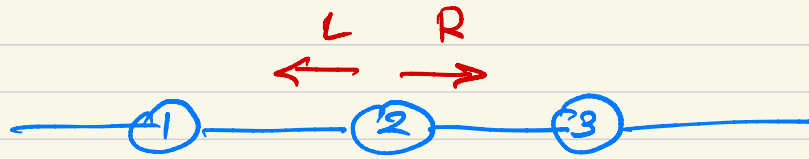




HW5 Solutions Q4

Multi-agent systems (Vrije Universiteit Amsterdam)

HW 5: Q5: Q-learning & SARSA



$\alpha = 0.9$, $\gamma = 2/3$: $2 \xrightarrow{R} 3$ (deterministic transitions)

Q-learning :

$$q(2, R) \leftarrow q(2, R) + \alpha \left[r(2, R) + \gamma \max_{a'} q(3, a') - q(2, R) \right]$$

$\begin{matrix} \text{"5"} & & & & & & \text{"5"} \\ \downarrow & & \downarrow & & \downarrow & & \downarrow \\ 0.9 & & -1 & & 2/3 & & \max\{6, 3\} = 6 \end{matrix}$

$$= 5 + 0.9 \left[-1 + \frac{2}{3} \cdot 6 - 5 \right] = 5 + 0.9 \left[-2 \right]$$

$$= \underline{\underline{3.2}}$$

Expected SARSA

$$\sum_a \pi(a|s_t) q_\pi(s_t, a) = \sum_a \pi(a|3) q_\pi(3, a)$$

$$= \underbrace{\pi(L|3)}_{1/2} \underbrace{q_\pi(3, L)}_3 + \underbrace{\pi(R|3)}_{1/2} \underbrace{q_\pi(3, R)}_6$$

$$= 9/2$$

Hence: $q_\pi(2, R) \leftarrow q_\pi(2, R) + 0.9 \left[-1 + \gamma \cdot \frac{9}{2} - 5 \right]$

$\begin{matrix} \text{"5"} & & & & \text{"4.5"} \\ \downarrow & & \downarrow & & \downarrow \\ 5 & & -1 & & 2/3 \cdot 9/2 = 4.5 \end{matrix}$

$$= 5 + 0.9 \left[-3.5 \right] = \underline{\underline{2.3}}$$