

EXERCISE 7.4 PROVE THAT THE N -STEP RETURN OF SARSA (7.4) CAN BE WRITTEN EXACTLY IN TERMS OF A NOVEL TD ERROR, AS

(7.4)

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n})$$

$$n \geq 1, 0 \leq t \leq T-n$$

(7.6)

$$G_{t:t+n} = Q_{t-1}(S_t, A_t) + \sum_{k=t}^{\min(t+n, T)-1} \gamma^{k-t} [R_{k+1} + \gamma Q_k(S_{k+1}, A_{k+1}) - Q_{k-1}(S_k, A_k)]$$

(start with (7.6))

$$G_{t:t+n} = Q_{t-1}(S_t, A_t) + \sum_{k=t}^{\min(t+n, T)-1} \gamma^{k-t} [R_{k+1} + \gamma Q_k(S_{k+1}, A_{k+1}) - Q_{k-1}(S_k, A_k)]$$

expand a lot to get

$$\begin{aligned} &= \cancel{Q_{t-1}(S_t, A_t)} + [R_{t+1} + \gamma \cancel{Q_t(S_{t+1}, A_{t+1})} - \cancel{Q_{t-1}(S_t, A_t)}] \\ &+ \gamma R_{t+2} + \gamma^2 \cancel{Q_{t+1}(S_{t+2}, A_{t+2})} - \gamma \cancel{Q_t(S_{t+1}, A_{t+1})} \\ &+ \gamma^2 R_{t+3} - \gamma^2 \cancel{Q_{t+2}(S_{t+3}, A_{t+3})} + \gamma^3 \cancel{Q_{t+2}(S_{t+3}, A_{t+3})} \end{aligned}$$

lots of terms will cancel out then

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n})$$

EXERCISE 7.6 PROVE THAT THE CONTROL VARIABLE IN THE ABOVE EQUATIONS DOES NOT CHANGE THE EXPECTED OF THE RETURN.

off policy

$$\begin{aligned} G_{t:h} &= R_{t+1} + \gamma (P_{t+1} b_{t+1:h} + \bar{V}_{h-1}(S_{t+1}) - P_{t+1} Q_{h-1}(S_{t+1}, A_{t+1})) \\ &= R_{t+1} + \gamma P_{t+1} (b_{t+1:h} - Q_{h-1}(S_{t+1}, A_{t+1})) + \gamma \bar{V}_{h-1}(S_{t+1}) \quad t < h \leq T. \end{aligned}$$

combine this w/ algorithm of n-step SARSA

$$\begin{aligned} E[G_{t:h}] &= E[R_{t+1} + \gamma P_{t+1} (b_{t+1:h} - Q_{h-1}(S_{t+1}, A_{t+1})) + \gamma \bar{V}_{h-1}(S_{t+1})] \\ &= E[R_{t+1}] + E[\gamma P_{t+1} (b_{t+1:h} - Q_{h-1}(S_{t+1}, A_{t+1}))] + E[\gamma \bar{V}_{h-1}(S_{t+1})] \end{aligned}$$

$$\underline{E[P] = 1}$$

$$\begin{aligned} &= R_{t+1} + E[\gamma b_{t+1:h} - \gamma Q_{h-1}(S_{t+1}, A_{t+1})] + \gamma \bar{V}_{h-1}(S_{t+1}) \\ &= R_{t+1} + E[G_{t:h} - R_{t+1} - \gamma Q_{h-1}(S_{t+1}, A_{t+1})] + \gamma \bar{V}_{h-1}(S_{t+1}) \\ &= E[G_{t:h}] + \gamma [E[-Q_{h-1}(S_{t+1}, A_{t+1})] + \bar{V}_{h-1}(S_{t+1})] \end{aligned}$$

$$= E[G_{t:h}]$$