

(14).

$$(a). \quad \pi_{\theta}(a|s) = \frac{\theta^T x_{s,a}}{\sum_b \theta^T x_{s,b}}.$$

$$\log \pi_{\theta}(a|s) = \log \theta^T x_{s,a} - \log \sum_b \theta^T x_{s,b}.$$

$$\frac{1}{\pi_{\theta}(a|s)} \cdot \nabla \pi_{\theta}(a|s) = \frac{1}{\theta^T x_{s,a}} x_{s,a} - \frac{1}{\sum_b \theta^T x_{s,b}} \sum_b x_{s,b}.$$

$$\therefore \nabla \pi_{\theta}(a|s) = \pi_{\theta}(a|s) \left[\frac{x_{s,a}}{\theta^T x_{s,a}} - \frac{\sum_b x_{s,b}}{\sum_b \theta^T x_{s,b}} \right]$$

(b).

$$\ln \pi_{\theta}(a|s) = \ln \left(\frac{1}{\sigma(s, \theta) \sqrt{2\pi}} \right) - \left(\frac{(a - \mu(s, \theta))^2}{2\sigma(s, \theta)^2} \right).$$

$$(i). \quad \nabla_{\theta_{\mu}} \ln \pi_{\theta}(a|s) = \frac{(a - \mu(s, \theta))}{\sigma(s, \theta)^2} \cdot x_{\mu}(s).$$

$$\begin{aligned} (ii). \quad \nabla_{\theta_{\sigma}} \ln \pi_{\theta}(a|s) &= - \left(\frac{1}{\sigma(s, \theta)} \right) \cdot \sigma(s, \theta) \cdot x_{\sigma}(s) + \\ &\quad (a - \mu(s, \theta))^2 \left(\frac{1}{\sigma(s, \theta)^3} \right) (\sigma(s, \theta) \cdot x_{\sigma}(s)) \\ &= -x_{\sigma}(s) + \frac{(a - \mu(s, \theta))^2}{\sigma(s, \theta)^2} \cdot x_{\sigma}(s). \end{aligned}$$

(2A).

$$(a). \quad (H, w); (H, s); (L, Re); (L, w); (L, s).$$

$$\begin{aligned} (b). \quad Q(H, w) &\leftarrow (0.9)(0.5) + 0.1 [0.5 + 0.5] \\ &= 0.55 \\ Q(H, s) &= (0.9)(0.5) + 0.1 [1 + 0.5] \\ &= 0.6 \end{aligned}$$

$$Q(L, Re) = (0.9)(0.5) + 0.1 [0 + 0.6]$$

$$= 0.51$$

$$Q(H, S) = (0.9)(0.6) + 0.1 [1 + 0]$$

$$= 0.64.$$

↳ as given in question
terminal (state, action)
value taken to be 0.

$$(c). \quad Q(H, w) = (0.9)(0.5) + 0.1 [0.5 + (0.6 \times 0.5 + 0.4 \times 0.5)]$$

$$= 0.55$$

$$Q(H, S) = (0.9)(0.5) + 0.1 [1 + (0.3 \times 0.5 + 0.3 \times 0.5 + 0.4 \times 0.5)]$$

$$= 0.6$$

$$Q(L, Re) = (0.9)(0.5) + 0.1 [0 + [0.6 \times 0.6 + 0.4 \times 0.55]]$$

$$= 0.45 + 0.1 [0.58]$$

$$= 0.508$$

$$Q(H, S) = (0.9)(0.6) + 0.1 [1 + 0]$$

$$= 0.64.$$

$$(d). \quad Q(H, w) = (0.9)(0.5) + (0.1) [0.5 + 0.5]$$

$$= 0.55$$

$$Q(H, S) = (0.9)(0.5) + 0.1 [1 + 0.5]$$

$$= 0.6$$

$$Q(L, Re) = (0.9)(0.5) + 0.1 [0 + \max\{0.6, 0.55\}]$$

$$= 0.45 + 0.06$$

$$= 0.51.$$

$$Q(H, S) = (0.9)(0.6) + 0.1 [1 + 0]$$

$$= 0.64.$$

(3A). (a). From episode 1: $V(1) \leftarrow 0 + \gamma + 0 + \gamma^3 + \dots = \frac{\gamma}{1-\gamma^2}$
 episode 2: $1 + \gamma + \gamma^2 + \dots = \frac{1}{1-\gamma}$
 $\therefore V(1) \leftarrow \frac{\frac{\gamma}{1-\gamma^2} + \frac{1}{1-\gamma}}{2} = \frac{\gamma(1+\gamma) + 1}{2(1-\gamma^2)}.$

(b). From episode 1:- every-visit will always yield the infinite seq: 1, 2, 1, 2, ...
 \therefore estimate will still be $\frac{\gamma}{1-\gamma^2}$

Similarly for episode 2:- Same infinite seq. will be yielded

\therefore estimate is $\frac{1}{1-\gamma}.$

$$\therefore V(1) \leftarrow \frac{\gamma(1+\gamma) + 1}{2(1-\gamma^2)}.$$

(4A). $V(A) = 0.5 [V(B)] + 0.5 [V(C)].$

(a).
$$\left. \begin{aligned} V(B) &= 0.5 [V(D)] + 0.5 [V(E)] \\ V(C) &= 0.5 [V(D)] + 0.5 [V(E)] \end{aligned} \right\} V(B) = V(C).$$

$$V(D) = 0.5 [V(B)] + 0.5 (1+0). = 0.5 V(B) + 0.5$$

$$V(E) = 0.5 [V(C)] + 0.5 (0+0). = 0.5 V(C).$$

solving it, we get:

$$V(A) = V(B) = V(C) = \frac{1}{2}; \quad V(D) = \frac{3}{4}; \quad V(E) = \frac{1}{4}.$$

$$(b). \quad v_1(A) = (0.9)(0.5) + 0.1[0 + 0.5].$$

$$= 0.5$$

$$v_1(C) = (0.9)(0.5) + 0.1[0 + 0.5]$$

$$= 0.5$$

$$v_1(E) = (0.9)(0.5) + 0.1[0 + 0]$$

$$= 0.45$$

$$v_1(B) = v_1(D) = 0.5 \quad ; \quad v_1(F) = v_1(G) = 0.$$