

# Problem 1:

P.1

(a). There are only three possible scenarios for  $\hat{\nabla}V$ :

Scenario #1:  $S_0 = S, a_1 = a, S_1 = \text{terminal}$  (this scenario happens with prob.  $\frac{1}{10}$ )

It is easy to verify that  $\frac{\partial \log \pi_\theta(a|s)}{\partial \theta_a} = 1 - \pi_\theta(a|s)$

$$\frac{\partial \log \pi_\theta(a|s)}{\partial \theta_b} = -\pi_\theta(b|s)$$

$$\frac{\partial \log \pi_\theta(a|s)}{\partial \theta_c} = -\pi_\theta(c|s)$$

$$\text{Therefore, we have } \hat{\nabla}V = \begin{bmatrix} 1 - \pi_\theta(a|s) \\ -\pi_\theta(b|s) \\ -\pi_\theta(c|s) \end{bmatrix} \times \gamma(s, a) = 100 \times \begin{bmatrix} \frac{9}{10} \\ -\frac{5}{10} \\ -\frac{4}{10} \end{bmatrix}$$

Scenario #2:  $S_0 = S, a_1 = b, S_1 = \text{terminal}$  (this scenario happens with prob.  $\frac{5}{10}$ )

Similar to scenario #1, we have

$$\hat{\nabla}V = \begin{bmatrix} -\pi_\theta(a|s) \\ 1 - \pi_\theta(b|s) \\ -\pi_\theta(c|s) \end{bmatrix} \times \gamma(s, b) = 98 \times \begin{bmatrix} -\frac{1}{10} \\ \frac{5}{10} \\ -\frac{4}{10} \end{bmatrix}$$

Scenario #3:  $S_0 = S, a_1 = c, S_1 = \text{terminal}$  (this scenario happens with prob.  $\frac{4}{10}$ )

$$\hat{\nabla}V = \begin{bmatrix} -\pi_\theta(a|s) \\ -\pi_\theta(b|s) \\ 1 - \pi_\theta(c|s) \end{bmatrix} \times \gamma(s, c) = 95 \times \begin{bmatrix} -\frac{1}{10} \\ -\frac{5}{10} \\ \frac{6}{10} \end{bmatrix}$$

(Cont.).

P.2

Therefore, we know

$$\begin{aligned} E[\hat{\nabla} V] &= \frac{1}{10} \times \left( 100 \times \begin{bmatrix} \frac{9}{10} \\ -\frac{5}{10} \\ -\frac{4}{10} \end{bmatrix} \right) + \frac{5}{10} \times \left( 98 \times \begin{bmatrix} -\frac{1}{10} \\ \frac{5}{10} \\ \frac{4}{10} \end{bmatrix} \right) + \frac{4}{10} \times \left( 95 \times \begin{bmatrix} -\frac{1}{10} \\ -\frac{5}{10} \\ \frac{6}{10} \end{bmatrix} \right) \\ &= \begin{bmatrix} \frac{30}{100} \\ \frac{50}{100} \\ -\frac{80}{100} \end{bmatrix} = \begin{bmatrix} \frac{3}{10} \\ \frac{5}{10} \\ -\frac{8}{10} \end{bmatrix} \leftarrow \text{This is actually the true policy gradient } \nabla_{\theta} V^{\pi_0} \end{aligned}$$

By definition, we have

$$\text{Cov}(\hat{\nabla} V, \hat{\nabla} V) = E[\hat{\nabla} V \hat{\nabla} V^T] - (E[\hat{\nabla} V]) \cdot (E[\hat{\nabla} V])^T = \begin{bmatrix} 894.03 & -509.75 & -384.28 \\ -509.75 & 2352.75 & -1843 \\ -384.28 & -1843 & 2227.28 \end{bmatrix}$$

$$\begin{aligned} \bullet E[\hat{\nabla} V \hat{\nabla} V^T] &= \frac{1}{10} \times \left( 100 \times \begin{bmatrix} \frac{9}{10} \\ -\frac{5}{10} \\ -\frac{4}{10} \end{bmatrix} \right) \cdot \left( 100 \times \begin{bmatrix} \frac{9}{10} \\ -\frac{5}{10} \\ -\frac{4}{10} \end{bmatrix} \right)^T \\ &\quad + \frac{5}{10} \times \left( 98 \times \begin{bmatrix} -\frac{1}{10} \\ \frac{5}{10} \\ \frac{4}{10} \end{bmatrix} \right) \cdot \left( 98 \times \begin{bmatrix} -\frac{1}{10} \\ \frac{5}{10} \\ \frac{4}{10} \end{bmatrix} \right)^T \\ &\quad + \frac{4}{10} \times \left( 95 \times \begin{bmatrix} -\frac{1}{10} \\ -\frac{5}{10} \\ \frac{6}{10} \end{bmatrix} \right) \cdot \left( 95 \times \begin{bmatrix} -\frac{1}{10} \\ -\frac{5}{10} \\ \frac{6}{10} \end{bmatrix} \right)^T \end{aligned}$$

$$= \begin{bmatrix} 894.12 & -509.6 & -384.52 \\ -509.6 & 2353 & -1843.4 \\ -384.52 & -1843.4 & 2227.92 \end{bmatrix}.$$

$$\bullet (E[\hat{\nabla} V]) (E[\hat{\nabla} V])^T = \begin{bmatrix} 0.09 & 0.15 & -0.24 \\ 0.15 & 0.25 & -0.40 \\ -0.24 & -0.40 & 0.64 \end{bmatrix}$$

□

$$(b). \quad V^{\pi_{\theta}}(s) = \pi_{\theta}(a|s) \cdot Y(s,a) + \pi_{\theta}(b|s) \cdot Y(s,b) + \pi_{\theta}(c|s) \cdot Y(s,c)$$

P.3

$$= \frac{1}{10} \times 100 + \frac{5}{10} \times 98 + \frac{4}{10} \times 95$$

$$= 97$$

Then, we can get the estimated policy gradient with baseline for the three scenarios:

$$\text{Scenario \#1:} \quad \tilde{\nabla} V = \begin{bmatrix} 1 - \pi_{\theta}(a|s) \\ -\pi_{\theta}(b|s) \\ -\pi_{\theta}(c|s) \end{bmatrix} \times (Y(s,a) - V^{\pi_{\theta}}(s)) = 3 \times \begin{bmatrix} \frac{9}{10} \\ \frac{5}{10} \\ -\frac{4}{10} \end{bmatrix}$$

$$\text{Scenario \#2:} \quad \tilde{\nabla} V = \begin{bmatrix} -\pi_{\theta}(a|s) \\ 1 - \pi_{\theta}(b|s) \\ -\pi_{\theta}(c|s) \end{bmatrix} \times (Y(s,b) - V^{\pi_{\theta}}(s)) = 1 \times \begin{bmatrix} -\frac{1}{10} \\ \frac{5}{10} \\ -\frac{4}{10} \end{bmatrix}$$

$$\text{Scenario \#3:} \quad \tilde{\nabla} V = \begin{bmatrix} -\pi_{\theta}(a|s) \\ -\pi_{\theta}(b|s) \\ 1 - \pi_{\theta}(c|s) \end{bmatrix} \times (Y(s,c) - V^{\pi_{\theta}}(s)) = -2 \times \begin{bmatrix} -\frac{1}{10} \\ \frac{5}{10} \\ \frac{6}{10} \end{bmatrix}$$

Then, similar to Problem (a), we have

$$\begin{aligned} E[\tilde{\nabla} V] &= \frac{1}{10} \times \left( 3 \times \begin{bmatrix} \frac{9}{10} \\ \frac{5}{10} \\ -\frac{4}{10} \end{bmatrix} \right) + \frac{5}{10} \times \left( 1 \times \begin{bmatrix} -\frac{1}{10} \\ \frac{5}{10} \\ -\frac{4}{10} \end{bmatrix} \right) + \frac{4}{10} \times \left( -2 \times \begin{bmatrix} -\frac{1}{10} \\ \frac{5}{10} \\ \frac{6}{10} \end{bmatrix} \right) \\ &= \begin{bmatrix} \frac{3}{10} \\ \frac{5}{10} \\ -\frac{8}{10} \end{bmatrix} \end{aligned}$$

(Conti.)

$$Cov(\tilde{V}, \tilde{V}) = E[\tilde{V}\tilde{V}^T] - (E[\tilde{V}]) \cdot (E[\tilde{V}])^T = \begin{bmatrix} 0.66 & -0.50 & -0.16 \\ -0.50 & 0.50 & 0 \\ -0.16 & 0 & 0.16 \end{bmatrix}$$

$$\begin{aligned} \bullet E[\tilde{V}\tilde{V}^T] &= \frac{1}{10} \times \left( 3 \times \begin{bmatrix} \frac{9}{10} \\ \frac{5}{10} \\ -\frac{4}{10} \end{bmatrix} \right) \cdot \left( 3 \times \begin{bmatrix} \frac{9}{10} \\ \frac{5}{10} \\ -\frac{4}{10} \end{bmatrix} \right)^T \\ &\quad + \frac{5}{10} \times \left( 1 \times \begin{bmatrix} \frac{1}{10} \\ \frac{5}{10} \\ -\frac{4}{10} \end{bmatrix} \right) \cdot \left( 1 \times \begin{bmatrix} \frac{1}{10} \\ \frac{5}{10} \\ -\frac{4}{10} \end{bmatrix} \right)^T \\ &\quad + \frac{4}{10} \times \left( -2 \times \begin{bmatrix} \frac{1}{10} \\ \frac{5}{10} \\ \frac{6}{10} \end{bmatrix} \right) \cdot \left( -2 \times \begin{bmatrix} \frac{1}{10} \\ \frac{5}{10} \\ \frac{6}{10} \end{bmatrix} \right)^T \end{aligned}$$

$$= \begin{bmatrix} 0.75 & -0.35 & -0.40 \\ -0.35 & 0.75 & -0.40 \\ -0.40 & -0.40 & 0.80 \end{bmatrix}$$

$$\bullet (E[\tilde{V}]) \cdot (E[\tilde{V}])^T = \begin{bmatrix} 0.09 & 0.15 & -0.24 \\ 0.15 & 0.25 & -0.40 \\ -0.24 & -0.40 & 0.64 \end{bmatrix}$$

P.4

D

(c) Let  $B(s)$  be the baseline function.

P.5

Then, we can write down the estimated policy gradient (denoted by  $\nabla V_B$ ) for the three scenarios:

Scenario #1:  $S_0 = S$ ,  $a_1 = a$ ,  $S_1 = \text{terminal}$ .

$$\nabla V_B = \begin{bmatrix} 1 - \pi_\theta(a|s) \\ -\pi_\theta(b|s) \\ -\pi_\theta(c|s) \end{bmatrix} \times (r(s,a) - B(s)) = (100 - B(s)) \cdot \begin{bmatrix} \frac{9}{10} \\ -\frac{5}{10} \\ -\frac{4}{10} \end{bmatrix}$$

Scenario #2:  $S_0 = S$ ,  $a_1 = b$ ,  $S_1 = \text{terminal}$

$$\nabla V_B = \begin{bmatrix} -\pi_\theta(a|s) \\ 1 - \pi_\theta(b|s) \\ -\pi_\theta(c|s) \end{bmatrix} \times (r(s,b) - B(s)) = (98 - B(s)) \cdot \begin{bmatrix} -\frac{1}{10} \\ \frac{5}{10} \\ -\frac{4}{10} \end{bmatrix}$$

Scenario #3:  $S_0 = S$ ,  $a_1 = c$ ,  $S_1 = \text{terminal}$

$$\nabla V_B = \begin{bmatrix} -\pi_\theta(a|s) \\ -\pi_\theta(b|s) \\ 1 - \pi_\theta(c|s) \end{bmatrix} \times (r(s,c) - B(s)) = (95 - B(s)) \cdot \begin{bmatrix} -\frac{1}{10} \\ -\frac{5}{10} \\ \frac{6}{10} \end{bmatrix}$$

It is easy to verify that  $E[\nabla V_B] = \begin{bmatrix} \frac{3}{10} \\ \frac{5}{10} \\ -\frac{8}{10} \end{bmatrix}$ , which does not depend on  $B(s)$ .

Moreover,  $E[(\nabla V_B) \cdot (\nabla V_B)^T] = \begin{bmatrix} C_{11} & C_{12} & C_{13} \\ C_{21} & C_{22} & C_{23} \\ C_{31} & C_{32} & C_{33} \end{bmatrix}$

Let's evaluate  $C_{11}$ ,  $C_{22}$ ,  $C_{33}$  separately in the next page.



(Cont.).

$$C_{11} = \frac{1}{10} \times \left[ (100 - B(s)) \cdot \left(\frac{9}{10}\right) \right]^2 + \frac{5}{10} \times \left[ (98 - B(s)) \cdot \left(-\frac{1}{10}\right) \right]^2 + \frac{4}{10} \times \left[ (95 - B(s)) \cdot \left(-\frac{1}{10}\right) \right]^2$$

P.6

$$= \frac{1}{1000} \left[ 90 \cdot B(s)^2 - 17940 \cdot B(s) + \text{some constant} \right]$$

$$C_{22} = \frac{1}{10} \times \left[ (100 - B(s)) \cdot \left(-\frac{5}{10}\right) \right]^2 + \frac{5}{10} \times \left[ (98 - B(s)) \cdot \left(\frac{5}{10}\right) \right]^2 + \frac{4}{10} \times \left[ (95 - B(s)) \cdot \left(-\frac{5}{10}\right) \right]^2$$
$$= \frac{1}{1000} \left[ 250 \cdot B(s)^2 - 48500 \cdot B(s) + \text{some constant} \right]$$

$$C_{33} = \frac{1}{10} \times \left[ (100 - B(s)) \cdot \left(-\frac{4}{10}\right) \right]^2 + \frac{5}{10} \times \left[ (98 - B(s)) \cdot \left(-\frac{4}{10}\right) \right]^2 + \frac{4}{10} \times \left[ (95 - B(s)) \cdot \left(\frac{6}{10}\right) \right]^2$$
$$= \frac{1}{1000} \left[ 240 \cdot B(s)^2 - 46240 \cdot B(s) + \text{some constant} \right]$$

$$\text{Tr}(\text{Cov}(\nabla V_B, \nabla V_B)) = C_{11} + C_{22} + C_{33} + (\text{Some constant from } (E[\nabla V_B])(E[\nabla V_B])^T)$$
$$= \frac{1}{1000} \left[ 580 \cdot B(s)^2 - 112680 \cdot B(s) + \text{some constant} \right]$$

Then, it is easy to verify that  $\text{Tr}(\text{Cov}(\nabla V_B, \nabla V_B))$  is minimized at

$$B(s) = \frac{112680}{2 \times 580} \approx 97.14.$$

(Remark: One can see that  $V^{\pi_0}(s)$  is a near-optimal baseline in this problem) .

□ .

# Problem 2

LHS

RHS

P.17

(a) Show that 
$$E_{\tau \sim p_{\mu}^{\pi_0}} \left[ \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] = \frac{1}{1-\gamma} E_{s_0 \sim d_{\mu}^{\pi_0}} E_{a \sim \pi_{\theta}(\cdot|s)} [f(s, a)]$$

pf:

$$\text{RHS} = \frac{1}{1-\gamma} \sum_s d_{\mu}^{\pi_0}(s) \cdot \sum_a \pi_{\theta}(a|s) \cdot f(s, a)$$

by the definition  
of  $d_{\mu}^{\pi_0}(s)$

$$= \frac{1}{1-\gamma} \sum_s \left( \sum_{s_0} \mu(s_0) \cdot (1-\gamma) \cdot \sum_{t=0}^{\infty} \gamma^t P(s_t=s | s_0, \pi_{\theta}) \right) \cdot \sum_a \pi_{\theta}(a|s) \cdot f(s, a)$$

$$= \sum_{s_0} \mu(s_0) \cdot \left[ \sum_s \sum_a \pi_{\theta}(a|s) \cdot \sum_{t=0}^{\infty} \gamma^t P(s_t=s | s_0, \pi_{\theta}) \cdot f(s, a) \right]$$

$$= \sum_{s_0} \mu(s_0) \cdot \left[ \sum_s \sum_a \left( \sum_{t=0}^{\infty} \gamma^t \cdot P(s_t=s, a_t=a | s_0, \pi_{\theta}) \cdot f(s, a) \right) \right]$$

$\equiv$   
 $\sum_{\tau=\{s_0, a_0, \dots\}} P(\tau | s_0, \pi_{\theta})$   
 with  $s_t=s, a_t=a$

Let's reorganize this  
by figuring out how much  
each trajectory  $\tau$  contribute

$\Rightarrow$  Each  $\tau$  shall contribute  
 $P(\tau | s_0, \pi_{\theta}) \cdot \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t)$

$$= \sum_{s_0} \mu(s_0) \cdot \left[ \sum_{\tau=\{s_0, a_0, \dots\}} P(\tau | s_0, \pi_{\theta}) \cdot \left( \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right) \right]$$

$$= \sum_{\text{all possible } \tau} P(\tau | \pi_{\theta}) \cdot \left( \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right)$$

$$= E_{\tau \sim p_{\mu}^{\pi_0}} \left[ \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] \equiv \text{LHS.}$$

□

(b). For episodic environments, show that

$$\nabla_{\theta} V^{\pi_{\theta}}(\mu) = E_{\tau \sim p_{\mu}^{\pi_{\theta}}} \left[ \sum_{t=0}^{T-1} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

Pf. The major difference between episodic and continuing environments is the existence of a "terminal state".

Let  $S_*$  be the terminal state of an episodic environment.

Once the agent reaches state  $S_*$ , it will stay at  $S_*$  forever (and hence the episode ends).

Moreover,  $Q^{\pi}(S_*, a) = 0$ ,  $V^{\pi}(S_*) = 0$ ,  $A^{\pi}(S_*, a) = 0$ , for all  $a$  and all  $\pi$

Let  $T$  be the episode length of a trajectory  $\tau$ .

Then, we have by the policy gradient expression (P5) in Lec 8

$$\nabla_{\theta} V^{\pi_{\theta}}(\mu) = E_{\tau \sim p_{\mu}^{\pi_{\theta}}} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

$$= E_{\tau \sim p_{\mu}^{\pi_{\theta}}} \left[ \sum_{t=0}^{T-1} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

□