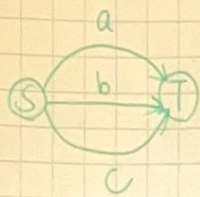


Problem 1:

(a) mean vector of $\hat{\nabla} V$ ($E[\hat{\nabla} V]$)



$$r_a = 100$$

$$r_b = 98$$

$$r_c = 95$$

$$\theta_a = 0$$

$$\theta_b = \ln 5$$

$$\theta_c = \ln 4$$

$$\pi_\theta(a|s) = \frac{e^0}{e^0 + e^{\ln 5} + e^{\ln 4}} = \frac{1}{1+5+4} = 0.1$$

$$\pi_\theta(b|s) = 0.5$$

$$\pi_\theta(c|s) = 0.4$$

$$\log \pi_\theta(a|s) = \log e^{\theta_a} - \log (e^{\theta_a} + e^{\theta_b} + e^{\theta_c})$$

$$\frac{\partial}{\partial \theta_a} [\log \pi_\theta(a|s)] = 1 - \pi_\theta(a|s) = 0.9, \quad \frac{\partial}{\partial \theta_b} = -\pi_\theta(b|s) = -0.5, \quad \frac{\partial}{\partial \theta_c} = -\pi_\theta(c|s) = -0.4$$

$$\log \pi_\theta(b|s) = \log e^{\theta_b} - \log (e^{\theta_a} + e^{\theta_b} + e^{\theta_c})$$

$$\frac{\partial}{\partial \theta_b} [\log \pi_\theta(b|s)] = 1 - \pi_\theta(b|s) = 0.5, \quad \frac{\partial}{\partial \theta_a} = -\pi_\theta(a|s) = -0.1, \quad \frac{\partial}{\partial \theta_c} = -\pi_\theta(c|s) = -0.4$$

$$\log \pi_\theta(c|s) = \log e^{\theta_c} - \log (e^{\theta_a} + e^{\theta_b} + e^{\theta_c})$$

$$\frac{\partial}{\partial \theta_c} [\log \pi_\theta(c|s)] = 1 - \pi_\theta(c|s) = 0.6, \quad \frac{\partial}{\partial \theta_a} = -\pi_\theta(a|s) = -0.1, \quad \frac{\partial}{\partial \theta_b} = -\pi_\theta(b|s) = -0.5$$

$$\Rightarrow \nabla_\theta \log \pi_\theta(a|s) = \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix}, \quad \nabla_\theta \log \pi_\theta(b|s) = \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix}, \quad \nabla_\theta \log \pi_\theta(c|s) = \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix}$$

$$E[\hat{\nabla} V] = \pi_\theta(a|s) [r_a \cdot \nabla_\theta \log \pi_\theta(a|s)] + \pi_\theta(b|s) [r_b \cdot \nabla_\theta \log \pi_\theta(b|s)] + \pi_\theta(c|s) [r_c \cdot \nabla_\theta \log \pi_\theta(c|s)]$$

$$= 0.1 \times 100 \times \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} + 0.5 \times 98 \times \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} + 0.4 \times 95 \times \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix}$$

$$= \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix}$$

② covariance matrix of \hat{V}

$$\begin{aligned}
 & E \left[(\hat{V} - E[\hat{V}]) (\hat{V} - E[\hat{V}])^T \right] \\
 &= 0.1 \times \left(\left(100 \times \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} \right) - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right) \cdot \left(\left(100 \times \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} \right) - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right)^T + \\
 & \quad 0.5 \times \left(\left(98 \times \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} \right) - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right) \left(\left(98 \times \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} \right) - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right)^T + \\
 & \quad 0.4 \times \left(\left(95 \times \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} \right) - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right) \left(\left(95 \times \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} \right) - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right)^T \\
 &= 0.1 \times \begin{bmatrix} 89.7 \\ -50.5 \\ -39.2 \end{bmatrix} \begin{bmatrix} 89.7 & -50.5 & -39.2 \end{bmatrix} + 0.5 \times \begin{bmatrix} -10.1 \\ 48.5 \\ -38.4 \end{bmatrix} \begin{bmatrix} -10.1 & 48.5 & -38.4 \end{bmatrix} \\
 & \quad + 0.4 \times \begin{bmatrix} -9.8 \\ -48 \\ 57.8 \end{bmatrix} \begin{bmatrix} -9.8 & -48 & 57.8 \end{bmatrix} = \begin{bmatrix} 894.03 & -509.75 & -384.28 \\ -509.75 & 2352.75 & -1843 \\ -384.28 & -1843 & 2227.28 \end{bmatrix}
 \end{aligned}$$

(b) ① $B(s) = V^{\pi_0}(s) = \sum_a \pi(a|s) Q^{\pi_0}(s, a) = 0.1 \times 100 + 0.5 \times 98 + 0.4 \times 95 = 97$

$$\begin{aligned}
 E[\tilde{V}] &= 0.1 \times (100 - 97) \times \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} + 0.5 \times (98 - 97) \times \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} + 0.4 \times (95 - 97) \times \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} \\
 &= \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix}
 \end{aligned}$$

② $E \left[(\tilde{V} - E[\tilde{V}]) (\tilde{V} - E[\tilde{V}])^T \right]$

$$= 0.1 \times \begin{bmatrix} 2.4 \\ -2 \\ -0.4 \end{bmatrix} \begin{bmatrix} 2.4 & -2 & -0.4 \end{bmatrix} + 0.5 \times \begin{bmatrix} -0.4 \\ 0 \\ 0.4 \end{bmatrix} \begin{bmatrix} -0.4 & 0 & 0.4 \end{bmatrix}$$

$$+ 0.4 \times \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} \begin{bmatrix} -0.1 & 0.5 & -0.4 \end{bmatrix} = \begin{bmatrix} 0.66 & -0.5 & -0.16 \\ -0.5 & 0.5 & 0 \\ -0.16 & 0 & 0.16 \end{bmatrix}$$

(c) find optimal $B(s)$

$$\tilde{V} - E[\tilde{V}]$$

$$\text{for } a: \begin{bmatrix} (100-b) \times 0.9 - 0.3 \\ (100-b) \times (-0.5) - 0.5 \\ (100-b) \times (-0.4) + 0.8 \end{bmatrix}, \quad \text{for } b: \begin{bmatrix} (98-b) \times (-0.1) - 0.3 \\ (98-b) \times 0.5 - 0.5 \\ (98-b) \times (-0.4) + 0.8 \end{bmatrix}, \quad \text{for } c: \begin{bmatrix} (95-b) \times (-0.1) - 0.3 \\ (95-b) \times (-0.5) - 0.5 \\ (95-b) \times 0.6 + 0.8 \end{bmatrix}$$

$$E[(\tilde{V} - E[\tilde{V}])(\tilde{V} - E[\tilde{V}])^T]$$

$$\begin{aligned} \text{the trace} = & 0.1 \times [(89.7 - 0.9b)^2 + (-50.5 + 0.5b)^2 + (-39.2 + 0.4b)^2] + \\ & 0.5 \times [(-10.1 + 0.1b)^2 + (48.5 - 0.5b)^2 + (-38.4 + 0.4b)^2] + \\ & 0.4 \times [(-9.8 + 0.1b)^2 + (-48 + 0.5b)^2 + (57.8 - 0.6b)^2] \end{aligned}$$

has minimum when $b = 97.138$ (set $\frac{d}{db} = 0$ to get).

So the optimal $B(s)$ is 97.138

Problem 2:

Notations : Discount state visitation distribution :

$$d_{s_0}^{\pi}(s) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t P(S_t = s | s_0, \pi)$$

$$d_{\mu}^{\pi}(s) = E_{s_0 \sim \mu} [d_{s_0}^{\pi}(s)]$$

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

$$= \frac{1}{1-\gamma} \sum_s d_{\mu}^{\pi_0}(s) E_{a \sim \pi_{\theta}(\cdot|s)} [f(s, a)]$$

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

$$= \frac{1}{1-\gamma} E_{s_0 \sim \mu} \left[\cancel{(1-\gamma)} \sum_s \sum_a \pi_{\theta}(a|s) f(s, a) \sum_{t=0}^{\infty} \gamma^t P(S_t = s | s_0, \pi_{\theta}) \right]$$

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

$$= \sum_{\tau} \sum_{t=0}^{\infty} P_{\mu}^{\pi_{\theta}}(\tau) \gamma^t f(s_t, a_t)$$

$$= E_{\tau \sim P_{\mu}^{\pi_{\theta}}} \left[\sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] = \text{LHS}$$

Problem 3:

Property 1:

consider all possible trajectories:

$$\begin{cases} S \rightarrow T & : R_T P_T \\ S \rightarrow S \rightarrow T & : (R_S + R_T) P_S P_T \\ S \rightarrow S \rightarrow S \rightarrow T & : (R_S + R_S + R_T) P_S^2 P_T \\ \vdots & \vdots \end{cases}$$

$$\begin{aligned} V(s) &= R_T P_T + (R_S + R_T) P_S P_T + (2R_S + R_T) P_S^2 P_T + (3R_S + R_T) P_S^3 P_T + \dots \\ &= \underline{R_T P_T} + \underline{R_S P_S P_T} + \underline{R_T P_S P_T} + \underline{2R_S P_S^2 P_T} + \underline{R_T P_S^2 P_T} + \underline{3R_S P_S^3 P_T} + \underline{R_T P_S^3 P_T} + \dots \\ &= \underline{\frac{R_T P_T}{1 - P_S}} + \underline{R_S P_T (P_S + 2P_S^2 + 3P_S^3 + \dots)} \\ &= \frac{R_T P_T}{1 - P_S} + R_S P_T \cdot \sum_{n=1}^{\infty} n P_S^n \\ &= \frac{R_T P_T}{1 - P_S} + \frac{R_S P_T P_S}{(1 - P_S)^2} \stackrel{(1 - P_S = P_T)}{=} \frac{R_T P_T}{P_T} + \frac{R_S P_T P_S}{P_T^2} = \frac{P_S}{P_T} R_S + R_T \end{aligned}$$

Property 2:

consider all possible trajectories:

$$\begin{cases} S \rightarrow T & : R_T P_T \\ S \rightarrow S \rightarrow T & : \left(\frac{R_S + 2R_T}{2}\right) P_T P_S \\ S \rightarrow S \rightarrow S \rightarrow T & : \left(\frac{R_S + 2R_S + 3R_T}{3}\right) P_T P_S^2 \\ S \rightarrow S \rightarrow S \rightarrow S \rightarrow T & : \left(\frac{R_S + 2R_S + 3R_S + 4R_T}{4}\right) P_T P_S^3 \\ & \vdots R_S + R_T \\ & \vdots 2R_S + R_T \\ & \vdots 3R_S + R_T \end{cases}$$

$$\begin{aligned} V(s) &= \sum_{k=0}^{\infty} P_T P_S^k \left(\frac{(1+k)R_S + (k+1)R_T}{k+1} \right) \\ &= \sum_{k=0}^{\infty} P_T P_S^k \left(\frac{\frac{k(k+1)}{2} R_S + (k+1)R_T}{k+1} \right) \\ &= \sum_{k=0}^{\infty} P_T P_S^k \cdot \left(\frac{k}{2} R_S + R_T \right) \\ &= \frac{1}{2} P_T R_S \sum_{k=0}^{\infty} k P_S^k + P_T R_T \sum_{k=0}^{\infty} P_S^k \\ &= \frac{1}{2} P_T R_S \cdot \frac{P_S}{(1 - P_S)^2} + P_T R_T \cdot \frac{1}{(1 - P_S)} \\ &= \frac{1}{2} P_T R_S \cdot \frac{P_S}{P_T^2} + P_T R_T \cdot \frac{1}{P_T} \\ &= \frac{P_S}{2P_T} R_S + R_T \end{aligned}$$