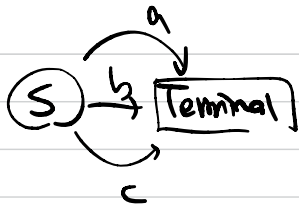


Problem 1) a. Find the mean Vector of $E[\hat{V}]$



$$\begin{matrix} r_a = 100 & \theta_a = 0 \\ r_b = 98 & \theta_b = 2.5 \\ r_c = 95 & \theta_c = 2.4 \end{matrix} \Rightarrow \begin{cases} \pi(a|s) = 0.1 \\ \pi(b|s) = 0.5 \\ \pi(c|s) = 0.4 \end{cases}$$

$$\therefore \pi(a|s) = \frac{\exp(0)}{\exp(0) + \exp(2.5) + \exp(2.4)} \Rightarrow \log \pi(a|s) = \log e^0 - \log(e^0 + e^{2.5} + e^{2.4})$$

$$\begin{aligned} \frac{\partial}{\partial \theta_a} [\log \pi(a|s)] &= 1 - \pi(a|s) = 1 - 0.1 = 0.9 \\ \frac{\partial}{\partial \theta_b} [\log \pi(a|s)] &= -\pi(b|s) = -0.5 \\ \frac{\partial}{\partial \theta_c} [\log \pi(a|s)] &= -\pi(c|s) = -0.4 \end{aligned} \Rightarrow \nabla_{\theta} \log \pi(a|s) = \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix}$$

同理可得

$$\nabla_{\theta} \log \pi(b|s) = \begin{bmatrix} 0.1 \\ 0.5 \\ -0.4 \end{bmatrix}$$

$$\nabla_{\theta} \log \pi(c|s) = \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix}$$

$$\text{可得 } E[\hat{V}] = E \left[\sum_{t=0}^{T-1} \gamma^t \nabla_{\theta} (s_t, a_t) \nabla_{\theta} \log \pi(a_t | s_t) \right]$$

$$\gamma \nabla_{\theta} (s_t, a_t) = r(s, a) + \gamma \sum_{s'} P_{ss'}^a V(s') - V(s) = r(s, a)$$

$$\begin{aligned} E[\hat{V}] &= \pi(a|s) \left[r(s, a) \times \nabla_{\theta} \log \pi(a|s) \right] + \pi(b|s) \left[r(s, b) \times \nabla_{\theta} \log \pi(b|s) \right] + \pi(c|s) \left[r(s, c) \times \nabla_{\theta} \log \pi(c|s) \right] \\ &= 0.1 \times \begin{bmatrix} 100 \times 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} + 0.5 \times \begin{bmatrix} 98 \times -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} + 0.4 \times \begin{bmatrix} 95 \times -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} \\ &= \begin{bmatrix} 9.7 \\ -5 \\ -0.4 \end{bmatrix} + \begin{bmatrix} -4.9 \\ 24.5 \\ -19.6 \end{bmatrix} + \begin{bmatrix} -3.8 \\ -19 \\ 22.8 \end{bmatrix} = \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \end{aligned}$$

Problem 1) - a Find the Covariance Matrix of \hat{V}

$$\text{Since } E[\hat{V}] = \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix}, \text{ and } (\hat{V}_a, \hat{V}_b, \hat{V}_c) = \left(100 \times \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix}, 98 \times \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix}, 95 \times \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} \right)$$

$$\text{可得 } (\hat{V}_a - E[\hat{V}], \hat{V}_b - E[\hat{V}], \hat{V}_c - E[\hat{V}]) = \left(\begin{bmatrix} 89.7 \\ -50.5 \\ -39.2 \end{bmatrix}, \begin{bmatrix} -10.1 \\ 48.5 \\ -38.4 \end{bmatrix}, \begin{bmatrix} -9.8 \\ -48 \\ 57.8 \end{bmatrix} \right)$$

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

$$\begin{aligned} &= 0.1 \times \begin{bmatrix} 89.7 \\ -50.5 \\ -39.2 \end{bmatrix} \begin{bmatrix} 89.7 & -50.5 & -39.2 \end{bmatrix} + 0.5 \times \begin{bmatrix} -10.1 \\ 48.5 \\ -38.4 \end{bmatrix} \begin{bmatrix} -10.1 & 48.5 & -38.4 \end{bmatrix} + 0.4 \times \begin{bmatrix} -9.8 \\ -48 \\ 57.8 \end{bmatrix} \begin{bmatrix} -9.8 & -48 & 57.8 \end{bmatrix} \\ &= \begin{bmatrix} 897.03 & -597.75 & -384.28 \\ -597.75 & 2252.75 & -1843 \\ -384.28 & -1843 & 2227.28 \end{bmatrix} \end{aligned}$$

Problem (2)

$$RHS = \frac{1}{1-r} E_{S_0 \sim \pi_0} E_{a \sim \pi_0(\cdot|s)} [f(s, a)]$$

$$= \frac{1}{1-r} \sum_s \pi_0(s) E_{a \sim \pi_0(\cdot|s)} [f(s, a)]$$

$$= \frac{1}{1-r} \sum_s \sum_a \pi_0(a|s) \times f(s, a) \times \boxed{\pi_0(s)} \quad \xrightarrow{S \sim M} E[\pi_0(s)] = E\left[(1-r) \sum_{t=0}^{\infty} r^t P(S_t = s | S_0, \pi_0)\right]$$

$$= \frac{1}{1-r} E_{S \sim M} \left[(1-r) \sum_s \sum_a \pi_0(a|s) \times f(s, a) \sum_{t=0}^{\infty} r^t P(S_t = s | S_0, \pi_0) \right]$$

$$= E_{S \sim M} \left[\sum_s \sum_a \pi_0(a|s) \sum_{t=0}^{\infty} P(S_t = s | S_0, \pi_0) \cdot r^t \times f(s, a) \right]$$

$$= \sum_z \sum_{t \in \mathbb{N}} P(z | M, \pi_0) r^t \cdot f(s_t, a_t)$$

$$= E_{z \sim p_M} \left[\sum_{t=0}^{\infty} r^t f(s_t, a_t) \right]$$

Problem (3)

Since only have one state S and Terminal T . So we only visit once

Property 1: $V(s) = P_S(R_S + V(s)) + P_T(R_T + V_T)$

$$2(1-P_S)V(s) = P_S R_S + P_T R_T$$

$$\therefore V(s) = \frac{P_S}{2P_T} R_S + R_T$$

Property 2: $E[\hat{V}_{MC}(s; \tau)] = \sum_k P(k) \times E_z[x(x) | k]$

$$= \sum_k P_T P_S^k \left(\frac{R_S + 2R_S + 3R_S + \dots + kR_S + (k+1)R_T}{k+1} \right)$$

$$= \sum_k P_T P_S^k \left(\frac{k}{2} R_S + R_T \right)$$

$$= \frac{P_S}{2P_T} R_S + R_T$$