

HW 1

Problem 1

(a)

score function :

$$\frac{\frac{\partial \log \left(\frac{e^{\theta_a}}{e^{\theta_a} + e^{\theta_b} + e^{\theta_c}} \right)}{\partial \theta_a}}{\frac{\partial \log \left(\frac{e^{\theta_a}}{e^{\theta_a} + e^{\theta_b} + e^{\theta_c}} \right)}{\partial \theta_a}} = \frac{\frac{\partial \log e^{\theta_a}}{\partial \theta_a}}{\frac{\partial \log (e^{\theta_a} + e^{\theta_b} + e^{\theta_c})}{\partial \theta_a}} = 1 - \bar{\pi}(a|s)$$

$$\frac{\frac{\partial \log \left(\frac{e^{\theta_a}}{e^{\theta_a} + e^{\theta_b} + e^{\theta_c}} \right)}{\partial \theta_b}}{\frac{\partial \log \left(\frac{e^{\theta_a}}{e^{\theta_a} + e^{\theta_b} + e^{\theta_c}} \right)}{\partial \theta_b}} = \frac{\frac{\partial \log e^{\theta_a}}{\partial \theta_b}}{\frac{\partial \log (e^{\theta_a} + e^{\theta_b} + e^{\theta_c})}{\partial \theta_b}} = -\bar{\pi}(b|s)$$

↪ action c is similar to b

$$\therefore \nabla_{\theta} \log \bar{\pi}_{\theta}(a|s) = \begin{bmatrix} 1 - \bar{\pi}(a|s) \\ -\bar{\pi}(b|s) \\ -\bar{\pi}(c|s) \end{bmatrix}$$

Mean vector :

$$\hat{V} = \sum_{t=0}^{\infty} \gamma^t \cdot r_t(s_t) \cdot \nabla_{\theta} \log \bar{\pi}_{\theta}(a_t|s_t)$$

$$= r(s_0, a_0) \cdot \nabla_{\theta} \log \bar{\pi}_{\theta}(a_0|s_0)$$

$$E[\hat{V}] = \sum \pi(-|s) \cdot r(s_0, a_0) \cdot \nabla_{\theta} \log \bar{\pi}_{\theta}(-|s)$$

$$= \frac{1}{e^0 + e^{2.5} + e^{2.4}} \cdot \left(\overset{x_1}{e^0 \cdot 100 \cdot \begin{bmatrix} 1 - \frac{1}{10} \\ -\frac{5}{10} \\ -\frac{4}{10} \end{bmatrix}} + \overset{x_2}{e^{2.5} \cdot 98 \cdot \begin{bmatrix} \frac{11}{10} \\ 1 - \frac{5}{10} \\ -\frac{4}{10} \end{bmatrix}} \right. \\ \left. + \overset{x_3}{e^{2.4} \cdot 95 \cdot \begin{bmatrix} -\frac{1}{10} \\ -\frac{5}{10} \\ 1 - \frac{4}{10} \end{bmatrix}} \right) = \overset{\mu}{\frac{1}{10} \cdot \begin{bmatrix} 3 \\ 5 \\ -8 \end{bmatrix}} \quad \#$$

Covariance matrix :

$$\begin{aligned}
 & E[(\hat{\theta} - E[\hat{\theta}])(\hat{\theta} - E[\hat{\theta}])^T] \\
 &= \frac{1}{10} \left(100 \cdot \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right) \cdot \left(100 \cdot \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right)^T + \\
 & \quad \frac{5}{10} \cdot \left(98 \cdot \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right) \cdot \left(98 \cdot \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right)^T + \\
 & \quad \frac{4}{10} \cdot \left(95 \cdot \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right) \cdot \left(95 \cdot \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right)^T
 \end{aligned}$$

Where $\bar{\pi}(a_n | s_0) = \bar{\pi}(b_n | s_0) = \bar{\pi}(a_n | s_0) \cdot (X_1 - \mu) \cdot (X_2 - \mu)^T = 0$

$$= \begin{bmatrix} \frac{89405}{100} & \frac{-1039}{4} & \frac{-9607}{25} \\ \frac{-1039}{4} & \frac{9411}{4} & -1843 \\ \frac{-9607}{25} & -1843 & \frac{55682}{25} \end{bmatrix} \#$$

(b)

Baseline :

$$\begin{aligned}
 \sqrt{\pi_{\theta_1}} &= \sum_{a_0} \bar{\pi}(a_0 | s_0) \cdot Q(s_0, a_0) = \sum_{a_0} \bar{\pi}(a_0 | s_0) \cdot V(s_0, a_0) \\
 &= (0.1) \cdot (100) + (0.5) \cdot (98) + (0.4) \cdot (95) = 97
 \end{aligned}$$

Mean vector :

$$\begin{aligned}
 E[\hat{\theta}] &= \frac{1}{10} \cdot \left(1 \cdot (100 - 97) \cdot \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} + 5 \cdot (98 - 97) \cdot \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} \right. \\
 & \quad \left. + 4 \cdot (95 - 97) \cdot \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} \right) = \frac{1}{10} \cdot \begin{bmatrix} 3 \\ 5 \\ -8 \end{bmatrix} \#
 \end{aligned}$$

Covariance matrix :

$$\begin{aligned}
 & E[(\hat{\theta}V - E[\hat{\theta}V]) \cdot (\hat{\theta}V - E[\hat{\theta}V])^T] \\
 &= \frac{1}{10} \cdot \left(3 \cdot \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right) \cdot \left(3 \cdot \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right)^T + \\
 & \quad \frac{5}{10} \cdot \left(1 \cdot \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right) \cdot \left(1 \cdot \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right)^T + \\
 & \quad \frac{4}{10} \cdot \left((-2) \cdot \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right) \cdot \left((-2) \cdot \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \right)^T
 \end{aligned}$$

where $\pi(a \wedge b | S_0) = \pi(b \wedge c | S_0) = \pi(a \wedge c | S_0) \cdot (X_1 - \mu) \cdot (X_2 - \mu)^T = 0$

$$= \begin{bmatrix} 0.66 & -0.5 & -0.16 \\ -0.5 & 0.5 & 0 \\ -0.16 & 0 & 0.16 \end{bmatrix} \#$$

12)

Suppose baseline $B(s) = b + 97$

Mean vector :

$$\begin{aligned}
 E[\hat{\theta}V] &= \frac{1}{10} \cdot \left(1 \cdot (3 - b) \cdot \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} + 5 \cdot (1 - b) \cdot \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} \right. \\
 & \quad \left. + 4 \cdot (-2 - b) \cdot \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} \right) = \frac{1}{10} \cdot \begin{bmatrix} 3 \\ 5 \\ -8 \end{bmatrix}
 \end{aligned}$$

Trace of covariance matrix :

$$x_1 - E[\hat{\theta}V] = (3-b) \cdot \begin{bmatrix} 0.9 \\ -0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} = \begin{bmatrix} 2.4 - 0.9b \\ -2.0 + 0.5b \\ -0.4 + 0.4b \end{bmatrix}$$

$$x_2 - E[\hat{\theta}V] = (1-b) \cdot \begin{bmatrix} -0.1 \\ 0.5 \\ -0.4 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} = \begin{bmatrix} -0.4 + 0.1b \\ -0.5b \\ 0.4 + 0.4b \end{bmatrix}$$

$$x_3 - E[\hat{\theta}V] = (-1-b) \cdot \begin{bmatrix} -0.1 \\ -0.5 \\ 0.6 \end{bmatrix} - \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} = \begin{bmatrix} -0.1 + 0.1b \\ 0.5 + 0.5b \\ -0.4 - 0.6b \end{bmatrix}$$

$$\begin{aligned} \text{trace} &= \frac{1}{1000} \cdot \left[1 \cdot (2.4 - 0.9b)^2 + 1 \cdot (-2.0 + 0.5b)^2 + 1 \cdot (-0.4 + 0.4b)^2 \right. \\ &\quad + 5 \cdot (-0.4 + 0.1b)^2 + 5 \cdot (-0.5b)^2 + 5 \cdot (0.4 + 0.4b)^2 \\ &\quad \left. + 4 \cdot (-0.1 + 0.1b)^2 + 4 \cdot (0.5 + 0.5b)^2 + 4 \cdot (-0.4 - 0.6b)^2 \right] \\ &= \frac{1}{1000} \left[580 \cdot \left(b - \frac{80}{596} \right)^2 + 1320 - \frac{80^2}{596} \right] \end{aligned}$$

$$\therefore \text{the optimal baseline} = 97 + \frac{80}{596} \doteq 97.138 \quad \#$$

Problem 2

(a)

$$P_4 \text{ under softmax policy: } \frac{d\sqrt{\pi_{\theta_{i,\mu}}}}{d\theta_{i,a}} = \frac{1}{1-\gamma} \cdot d_{\mu}^{\pi_{\theta_{i,\mu}}}(s) \cdot \pi_{\theta}(a|s) \cdot A^{\pi_{\theta_{i,\mu}}}(s,a)$$

$$\begin{aligned} \left\| \frac{d\sqrt{\pi_{\theta_{i,\mu}}}}{d\theta} \right\|_2 &= \sqrt{\sum_s \sum_a \left(\frac{d\sqrt{\pi_{\theta_{i,\mu}}}}{d\theta_{i,a}} \right)^2} \quad (\ell_2 \text{ norm}) \\ &= \sqrt{\sum_s \sum_a \left(\frac{1}{1-\gamma} \cdot d_{\mu}^{\pi_{\theta_{i,\mu}}}(s) \cdot \pi_{\theta}(a|s) \cdot A^{\pi_{\theta_{i,\mu}}}(s,a) \right)^2} \\ &\geq \sqrt{\sum_s \left(\frac{1}{1-\gamma} \cdot d_{\mu}^{\pi_{\theta_{i,\mu}}}(s) \cdot \pi_{\theta}(a^*(s)|s) \cdot A^{\pi_{\theta_{i,\mu}}}(s, a^*(s)) \right)^2} \cdot \sqrt{S} \cdot \frac{1}{\sqrt{S}} \\ &\quad \downarrow \text{by Cauchy-Schwarz inequality} \\ &\geq \sqrt{\left(\sum_s \frac{1}{1-\gamma} \cdot d_{\mu}^{\pi_{\theta_{i,\mu}}}(s) \cdot \pi_{\theta}(a^*(s)|s) \cdot A^{\pi_{\theta_{i,\mu}}}(s, a^*(s)) \right)^2} \cdot \frac{1}{\sqrt{S}} \\ &\geq \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{S}} \cdot \sum_s d_{\mu}^{\pi_{\theta_{i,\mu}}}(s) \cdot \pi_{\theta}(a^*(s)|s) \cdot |A^{\pi_{\theta_{i,\mu}}}(s, a^*(s))| \quad \# \end{aligned}$$

(b)

$$\begin{aligned} \left\| \frac{d\sqrt{\pi_{\theta_{i,\mu}}}}{d\theta} \right\|_2 &\geq \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{S}} \cdot \sum_s d_{\mu}^{\pi_{\theta_{i,\mu}}}(s) \cdot \pi_{\theta}(a^*(s)|s) \cdot |A^{\pi_{\theta_{i,\mu}}}(s, a^*(s))| \\ &\geq \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{S}} \cdot \sum_{s'} d_{\mu}^{\pi^*}(s') \cdot \frac{d_{\mu}^{\pi_{\theta_{i,\mu}}}(s')}{d_{\mu}^{\pi^*}(s')} \cdot \pi_{\theta}(a^*(s')|s') \cdot |A^{\pi_{\theta_{i,\mu}}}(s', a^*(s'))| \\ &\geq \frac{1}{1-\gamma} \cdot \frac{1}{\sqrt{S}} \cdot \sum_{s'} d_{\mu}^{\pi^*}(s') \cdot \frac{1}{\max_{s'' \in S} \left(\frac{d_{\mu}^{\pi^*}(s'')}{d_{\mu}^{\pi_{\theta_{i,\mu}}}(s'')} \right)} \cdot \min_{s \in S} \pi_{\theta}(a^*(s)|s) \cdot |A^{\pi_{\theta_{i,\mu}}}(s', a^*(s'))| \\ &\geq \frac{1}{\sqrt{S}} \cdot \left\| \frac{d_{\mu}^{\pi^*}}{d_{\mu}^{\pi_{\theta_{i,\mu}}}} \right\|_{\infty}^{-1} \cdot \min \pi_{\theta}(a^*(s)|s) \cdot \frac{1}{1-\gamma} \cdot \mathbb{E}_{s' \sim d_{\mu}^{\pi^*}} \left[\mathbb{E}_{a' \sim \pi^*} [A^{\pi_{\theta_{i,\mu}}}(s', a')] \right] \\ &\geq \frac{1}{\sqrt{S}} \cdot \left\| \frac{d_{\mu}^{\pi^*}}{d_{\mu}^{\pi_{\theta_{i,\mu}}}} \right\|_{\infty}^{-1} \cdot \min \pi_{\theta}(a^*(s)|s) \cdot (\sqrt{\gamma_{i,\mu}} - \sqrt{\gamma_{i,\mu}^*}) \quad \# \end{aligned}$$

Problem 3

Property 1 :

$$\begin{aligned} V(S) &= \sum_{k=0}^{\infty} p_S^k \cdot p_T \cdot (K \cdot R_S + R_T) \\ &= \frac{p_S}{p_T} \cdot R_S + \frac{R_T}{1-p_S} \cdot \cancel{p_T} = \frac{p_S}{p_T} \cdot R_S + R_T \quad (\text{by lemma 1}) \\ &\quad \# \end{aligned}$$

Property 2 :

$$\begin{aligned} \tilde{V}_T(\hat{V}_{MC}(S; T)) &= \sum_{k=0}^{\infty} p_S^k \cdot p_T \cdot \left(\frac{R_S + \dots + K \cdot R_S + (K+1) \cdot R_T}{K+1} \right) \\ &= \sum_{k=0}^{\infty} p_S^k \cdot p_T \cdot \left(\frac{\frac{K}{2} (K+1) \cdot R_S}{K+1} + R_T \right) \\ &= \frac{p_S}{2 \cdot p_T} \cdot R_S + R_T \quad (\text{by lemma 1}) \\ &\quad \# \end{aligned}$$

Lemma 1.

$$\begin{aligned} f &: \sum_{k=0}^{\infty} k \cdot p_S^k \cdot R_S \cdot p_T = \sum_{k=0}^{\infty} (k+1) \cdot p_S^{k+1} \cdot R_S \cdot p_T \\ p_S \cdot f &= \sum_{k=0}^{\infty} k \cdot p_S^{k+1} \cdot R_S \cdot p_T \\ -) \quad (1-p_S) \cdot f &= \sum_{k=0}^{\infty} p_S^{k+1} \cdot R_S \cdot p_T = \frac{p_S}{1-p_S} \cdot R_S \cdot \cancel{p_T} \\ \therefore p_S + p_T &= 1 \quad \therefore f = \frac{p_S}{p_T} \cdot R_S \end{aligned}$$