

سوال ۱:

الف) در مسائل RL، با توجه به اینکه ابتدا باید صدای sample از

محیط جمع کنیم، تعامل با محیط به صورت online انجام می شود.

ب) نادرست. در مسائل MDP، در نهایت به دنبال سیاست

حتمی به بیشترین reward را داشته باشد. در نهایت اگر دو حالتی

$$Q(s, a) = Q(s', a) \text{ باشد، با دو سیاست مختلف به یک مقدار}$$

رسیده ایم، در نتیجه سیاست بهینه لزوماً یکتا نیست.

ج) نادرست. الگوریتم REINFORCE در واقع یک الگوریتم Model-free

است. زیرا به دنبال پیدا کردن یک Policy بهینه است و \hat{A} و \hat{R} را به ندرت

د) درست. به صورت کلی $V(s)^*$ ، در واقع جواب بهینه برای

$V(s)$ است. در نتیجه باید حاکم $V(s)$ باشد و در نتیجه $V(s)^* \geq V(s)$

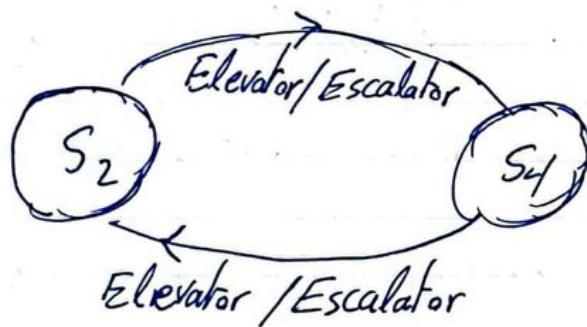
سؤال (2):

(آ) با توجه به Infinite Horizon بین، آینده مقدار Q -Value را دارم
می‌خواهم با استفاده از الگوریتم Policy Extracting، سیاست بین را بیابم

$$\forall s \in S : \pi^*(s) = \underset{a}{\operatorname{argmax}} Q^*(s, a)$$

$$\pi^*(s_1) = \text{Elevator} \quad \pi^*(s_2) = \text{Elevator}$$

$$\pi^*(s_3) = \text{Escalator}$$



(ج) مانند قسمت الف برای به دست آوردن سیاست بین از Policy Extracting استفاده می‌کنیم:

$$\pi^*(s_2) = \text{Escalator} \quad \pi^*(s_4) = \text{Escalator}$$

(د) فرض ماکد به این صورت است که احتمال رسیدن به s_1 تنها در زمان $t+1$ تنها به این s و اکشن a در زمان t بستگی دارد. (بی حافظگی)

در اینجا با توجه به اینکه از Q -learning استفاده کردیم ممکن است که داده‌هایی که جمع کردیم به اندازه کافی نباشد. در نتیجه بی حافظگی در فرض ماکد ممکن است که رعایت نشده باشد. در نتیجه سیاست بین عوض می‌شود.

سوال (3) :

$$E_{S \sim P(s), a \sim \pi_0(s,a)} \frac{\pi_1(s,a)}{\hat{\pi}_0(s,a)} R(s,a) = \quad (1)$$

$$= E R(s,a) = E \frac{\pi_1(s,a)}{\hat{\pi}}$$

$$= \sum \frac{\pi_1(s,a)}{\hat{\pi}_0(s,a)} R(s,a) P(s) \cancel{P(a|s)} \xrightarrow{\pi_0(s,a)} \rightarrow$$

$$\xrightarrow{\hat{\pi}_0(s,a) = \pi_0(s,a)} = \sum \pi_1(s,a) R(s,a) P(s) \Rightarrow$$

$$\Rightarrow E_{S \sim P(s), a \sim \pi_0(s,a)} \frac{\pi_1(s,a)}{\hat{\pi}_0(s,a)} R(s,a) = E_{S \sim P(s), a \sim \pi_1(s,a)} R(s,a)$$

$$P(\|\hat{\theta} - \theta\| > \epsilon) = 0 : \forall \epsilon > 0 : \text{سازگاری تخمین}$$

یعنی با توجه به اینکه تخمین که داده شده تابع R برای π_1 را دقیق تخمین می زند،
این تخمین سازگار است.

(→) صورت کسر را در سمت الف) بدست آورده ایم برای مخرج نیز
به صورت زیر عمل می کنیم :

$$E_{S \sim P(s), a \sim \pi_0(s,a)} \frac{\pi_1(s,a)}{\hat{\pi}_0(s,a)} = \sum \frac{\pi_1(s,a)}{\hat{\pi}_0(s,a)} P(s) \pi_0(s,a) =$$

$$= \sum P(s) \pi_1(s,a) = \sum P(s) P(a|s) = 1 \Rightarrow$$

$$\frac{E_{SNP(s), a \sim \pi_0(s,a)} \frac{\pi_1(s,a)}{\hat{\pi}_0(s,a)} R(s,a)}{E_{SNP(s), a \sim \pi_0(s,a)} \frac{\pi_1(s,a)}{\hat{\pi}_0(s,a)}} = E_{SNP(s), a \sim \pi_1(s,a)} R(s,a)$$

در نتیجه این محاسبه تخمین کرد، تخمین گوی سازگار است.

(ج) فرض می کنیم تنها یک داده داشته باشیم، در این صورت

مقدار تخمینگرها برابر با $R_{\pi_0}(s,a)$ که از تقاضای احتمالی π_0

سروی می کند در نتیجه $E(R_{\pi_1}(s,a)) \neq E(R_{\pi_0}(s,a))$ که بین محاسبه که این تخمینگر می تواند بایس باشد.

سوال 9 :

$$V_{(s)}^{\pi} = R_{(s)} + \gamma \sum_{s'} T_{(s, \pi(s), s')} V_{(s')}^{\pi} \quad (1)$$

↓ با توجه به این مقدار Reward نسبتی به s دارد. می توان $V_{(s)}^{\pi}$

M: Mountain R: Riverside D: Desert π : Peace

$$V_{(M)}^{\pi} = 1 + \gamma \left\{ \frac{1}{2} \times V_{(D)}^{\pi} + \frac{1}{2} \times V_{(R)}^{\pi} \right\}$$

$$V_{(R)}^{\pi} = 2 + \gamma \{ 1 \times V_{(R)}^{\pi} \} \Rightarrow V_{(R)}^{\pi} = \frac{2}{1-\gamma}$$

$$V_{(D)}^{\pi} = -1 + \gamma \{ 1 \times V_{(D)}^{\pi} \} \Rightarrow V_{(D)}^{\pi} = -\frac{1}{1-\gamma}$$

$$V_{(M)}^{\pi} = 1 + \frac{\gamma}{2} \times \left\{ \frac{2}{1-\gamma} + \left(-\frac{1}{1-\gamma} \right) \right\} = 1 + \frac{\gamma}{2(1-\gamma)} = \frac{2-\gamma}{2(1-\gamma)}$$

←

$$\pi_0: \text{Peace} \Rightarrow V_{(M)}^{\pi_0} = 5.5 \quad V_{(R)}^{\pi_0} = 20 \quad V_{(D)}^{\pi_0} = -10$$

$$\pi_1(s) = \arg \max_a \sum_{s'} T_{(s,a,s')} \{ R_{(s,a,s')} + \gamma V_{(s')}^{\pi_0} \}$$

↳ Policy Iteration

$$\pi_{1(M)} : \arg \max \left(\text{peace} : 1 + 0,9 \times \left(\frac{1}{2} \times 20 - \frac{1}{2} \times 10 \right), \right. \\ \left. , \text{war} : 1 + 0,9 \times (0,1 \times 5,5 + 0,7 \times 20 - 0,2 \times 10) \right)$$

$$\Rightarrow \pi_{1(M)} : \text{War}$$

$$\pi_{1(D)} : \arg \max \left(\text{peace} : -1 + 0,9 \times (1 \times -1), \text{war} : -1 + 0,9 \times (1 \times 5,5) \right)$$

$$\pi_{1(D)} : \text{war}$$

$$\pi_{1(R)} : \arg \max \left(\text{peace} : 2 + 0,9 \times (1 \times 20), \right. \\ \left. \text{war} : 2 + 0,9 \times (0,2 \times 20 - 0,8 \times 10) \right)$$

$$\Rightarrow \pi_{1(R)} : \text{Peace}$$

$$\text{Sample} : R(s, a, s') + \gamma \max_{a'} (Q(s', a'))$$

(2)

$$Q(s, a) \leftarrow (1 - \alpha) Q(s, a) + \alpha \cdot \text{Sample}$$

| Mountain - peace | Riverside - peace | Desert - war |
|------------------|-------------------|--------------|
| 0 | 0 | 0 |
| 0 | 0 | -1 |
| 0 | 0 | 1 |
| 0 | 0,5 | 1 |
| 0,75 | 0,5 | 1 |