

سوال ۱، الف، پیت است، و مسائل RL، agent ماه صورت online رفتاری کند، وابست خطا policy، را تقویت می کند

ب، تغییر کرد البته با یک $Q(s,a) = Q(s,b)$ ، فردر اکشن طوره می تواند جزر

policy بهینه باشد برای همین policy بهینه بهینه است

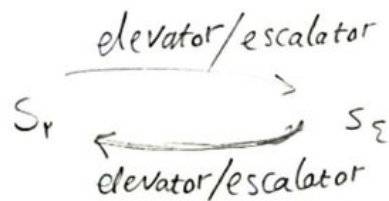
ج، تغییر در این الگوریتم های فراگیر policy بهینه، را مستقیم مطابق نیاز برای همین این

الگوریتم model free است.

(د) پیت، V^* بهینه است به همین دلیل V^* (یا بهینه/غیر بهینه) بهینه کمتر مساوی V^*

سوال ۲، مانند policy extraction اکشن ای که بهترین Q-value را دارد، را انتخاب می کنیم، پس:

$$\pi^*(s_1) = \text{elevator} \quad \pi^*(s_2) = \text{elevator} \quad \pi^*(s_3) = \text{escalator}$$



ج، مانند قسمت الف

$$\pi^*(s_r) = \text{escalator}$$

$$\pi^*(s_e) = \text{escalator}$$

(د) در فرض ما، ما می دانیم که حرکت می گذریم و آینده از هم مستقل اند و هر state تمام اطلاعات مورد نیاز برای تصمیم گیری را دارد و نیازی به اطلاعات action گذرته ندارد اما در اینجا ما فضا را گامی

داریم و مقداری دنیا را از دست داده ایم به همین دلیل ممکن است که انتقالی بین اکشن های

گذرته، و آینده موجود نداشته باشد، به همین دلیل ممکن است به policy متفاوتی برسیم.

$$u^{\pi}(s) = R(s) + \gamma \sum_{s'} T(s, \pi(s), s') u^{\pi}(s') \quad : \rho, b, (w), \epsilon$$

$$\rightarrow u^{\pi_0}(M) = 1 + \gamma (1/4 u^{\pi_0}(R) + 1/4 u^{\pi_0}(D))$$

$$\rightarrow u^{\pi_0}(R) = 1 + \gamma (1 \times u^{\pi_0}(R))$$

$$\rightarrow u^{\pi_0}(D) = -1 + \gamma (1 \times u^{\pi_0}(D))$$

$$u^{\pi_0}(R) = \frac{1}{1-\gamma}$$

$$u^{\pi_0}(D) = \frac{-1}{1-\gamma}$$

$$u^{\pi_0}(M) = \frac{1-\gamma}{1-\gamma}$$

$$\pi_0 = \rho \quad u^{\pi_0} = (\underbrace{0, 0, 1, -1}_{M, R, D})$$

$$\pi_0(M) = \arg \max (\rho: 1 + \gamma (1/4 \times 1 - 1/4 \times -1), \\ w: 1 + \gamma (1/4 \times 0 - 1/4 \times 0 + 1/2 \times 1))$$

$$\rightarrow \pi_0(M) = w$$

$$\pi_1(R) = \arg \max (\rho: 1 + \gamma \times 1, w: 1 + \gamma (1/4 \times 1 + 1/4 \times 0))$$

$$\rightarrow \pi_1(R) = \rho$$

$$\pi_1(D) = \arg \max (\rho: -1 + \gamma \times (-1), w: -1 + \gamma \times 0)$$

$$\rightarrow \pi_1(D) = w$$

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

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

$$\text{sample 1} = (D, w, D, -1)$$

$$\text{sample 2} = (D, w, R, 1)$$

$$\text{sample 3} = (R, p, M, 1)$$

$$\text{sample 4} = (M, p, R, 1)$$

$$Q(D, w) = 0 + 1/4(-1 + 0 + 0) = -1/4$$

$$Q(D, w) = -1/4 + 1/4(1 + 0 + 1) = 1/4$$

$$Q(R, p) = 0 + 1/4(1 + 0 + 0) = 1/4$$

$$Q(M, p) = 0 + 1/4(1 + 1/4 + 0) = 5/16$$

M-p	R-p	D-w
0	0	0
0	0	-1
0	0	+1
0	1/4	+1
1/4	1/4	+1