

سؤال ۱

(الف) نادرست؛ زیرا ممکن است داده‌ها از قبل به دست آمده باشند و دردی این داده‌ها یا دگیری تقویتی را انجام دهیم.

(ب) نادرست؛ اگر در حالتی بین از یک حالت باعث دست آوردن بیشینه یاداش شود در این حالت سیاست بهینه نیست.

~~ج) نادرست؛ این الگوریتم برای یادگیری در محیط‌های مکرر کارایی دارد و به عنوان یک الگوریتم یادگیری در محیط‌های مکرر کارایی دارد.~~

(ج) نادرست؛ این الگوریتم نمی‌تواند model based و model free باشد.

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') \{ R(s, a, s') + \gamma V^*(s') \}$$

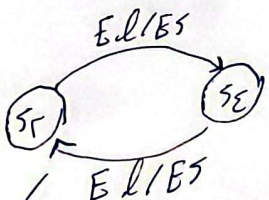
$$V^\pi(s) = \sum_{s'} T(s, \pi(s), s') \{ R(s, \pi(s), s') + \gamma V^\pi(s') \}$$

چون در حالت اول بیشینه یاداش را انتخاب می‌کنیم $V^*(s) \geq V^\pi(s)$

سؤال ۲

$$\pi(s) = \arg \max_a Q(s, a)$$

$$\Rightarrow \pi(s_1) = E, \pi(s_2) = E, \pi(s_3) = E$$



$$\pi(s_2) = E, \pi(s_3) = E$$

(د) در حالت دوم که مسئله ساده‌تر شده است، اطلاعاتی از بین رفته است، با توجه به اینکه در حالت باید بتوانیم با توجه به اطلاعات state فعلی به طور مستقل از state‌های قبلی برای state جدیدی تصمیم بگیریم ولی در اینجا اطلاعات نداریم که چگونه به state رسیده ایم پس سیاست بهینه نیز تغییر خواهد کرد.

9/11/04/44

مقدارهای آماری
سوال ۳:

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

$$= \sum \frac{\pi_1(s, a)}{\hat{\pi}_0(s, a)} R(s, a) p(s) p(a|s) =$$

$$= \sum \frac{\pi_1(s, a)}{\hat{\pi}_0(s, a)} R(s, a) p(s) \pi_0(a|s) = \sum \pi_1(s, a) R(s, a) p(s)$$

~~$$= E_{s \sim p(s), a \sim \pi_1(s, a)} R(s, a)$$~~

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

$$= \sum \frac{\pi_1(s, a)}{\hat{\pi}_0(s, a)} p(s) \pi_0(a|s) = \sum p(s) \pi_1(s, a) = 1$$

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

$$= E_{s \sim p(s), a \sim \pi_1(s, a)} R(s, a)$$

(2)

$$\pi_0 \rightarrow a=0 \quad R(s, a)=a$$

$$\pi_1 \rightarrow a=1$$

$$\Rightarrow E \pi_1 = 1 \quad \text{but مقدار خروجی تعیین نمی کند}$$

$$V^{\pi}(M) = \gamma \Delta (\tau + \gamma V^{\pi}(R)) + \gamma \Delta (-1 + \gamma V^{\pi}(D)) \quad \text{, } \varepsilon \text{ الـ } \text{الف}$$

$$V^{\pi}(R) = (\tau + \gamma V^{\pi}(R))$$

$$V^{\pi}(D) = (-1 + \gamma V^{\pi}(D))$$

$$\Rightarrow V^{\pi}(M) = \frac{1}{\gamma(1-\gamma)}, V^{\pi}(R) = \frac{\tau}{1-\gamma}, V^{\pi}(D) = \frac{-1}{1-\gamma}$$

$$\gamma = 0.9 \Rightarrow V^{\pi}(M) = 2, V^{\pi}(R) = 1, V^{\pi}(D) = -1$$

$$\Rightarrow Q^{\pi}(M, P) = 2, Q^{\pi}(R, P) = 1, Q^{\pi}(D, P) = -1$$

$$Q^{\pi}(M, W) = \gamma \Delta (\tau + \gamma V^{\pi}(R)) + \gamma \Delta (-1 + \gamma V^{\pi}(D)) + \gamma \Delta (\tau + \gamma V^{\pi}(M))$$

$$= 1.7, 2.2$$

$$\text{, } Q^{\pi}(R, W) = 1.7, Q^{\pi}(D, W) = 2.2$$

$$\Rightarrow V^{\pi_1}(M) = \text{war}, V^{\pi_1}(R) = \text{peace}, V^{\pi_1}(D) = \text{war}$$

(ج)

$$I \rightarrow 0 + \gamma \Delta (-1 + \gamma \Delta (-1)) = -1$$

$$\Gamma \rightarrow -1 + \gamma \Delta (\tau + \gamma \Delta (0)) = 1$$

$$\Upsilon \rightarrow 0 + \gamma \Delta (1 + \gamma \Delta (0 - 0)) = \gamma \Delta$$

$$\Sigma \rightarrow 0 + \gamma \Delta (1 + \gamma \Delta (0 - 0)) = \gamma \Delta$$

\Rightarrow

M-P	R-P	D-W
0	0	0
0	0	-1
0	0	1
0	$\gamma \Delta$	1
$\gamma \Delta$	$\gamma \Delta$	1