

(۱) الف) در صورتی که RL آنلاین باشد و عامل با انجام عمل به درستی یا غلطی آن پیوسته

ب) یادگیری در صورتی که $Q(s, a) = Q(s, a')$ هر دو آن a و a' می‌توانند در policy بجهت باشند در شبکه عصبی یادگیری

ج) یادگیری در الگوریتم reinforce باید دنبال به درستی آوردن $reward$ function و $transition$ function

نیستیم در مفاهیم policy بجهت را به صورت مستقیم به درستی آوردیم در شبکه عصبی این الگوریتم $model$ است

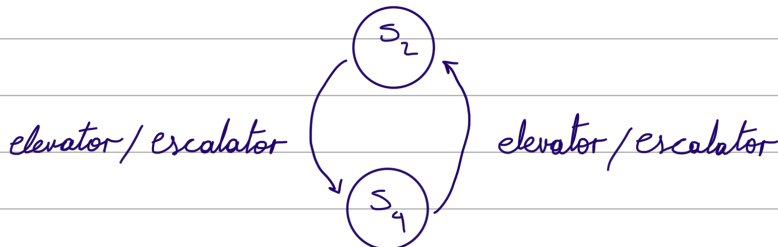
د) در صورتی که V^* هست بجهت است و $value$ را در حالتی که $agent$ بجهت عمل می‌کند اندازه می‌گیرد و مقدار آن از $value$ هر بجهت دیگری بیشتر است

Q-values	Elevator	Escalator
S_1	$+1/5$	$-0/5$
S_2	$+2/0$	$-0/3$
S_3	$+0/8$	$+0/9$

(۲) الف) به کمک الگوریتم $policy$ extraction می‌دانیم که باید این را انتخاب کنیم که Q -state ای با بیشترین Q -value

$\pi^*(S_1) = \text{elevator}$ $\pi^*(S_2) = \text{elevator}$ باید در شبکه π^* به صورتی ادغام شود:

$\pi^*(S_3) = \text{escalator}$



(ج) محاسبه انتظاری داریم:

Q^*	Elevator	Escalator
S_2	+۰/۰۸	+۲/۰۸
S_4	+۱/۱۸	+۱/۳۸

$$\pi^*(S_2) = \text{escalator}$$

$$\pi^*(S_4) = \text{escalator}$$

(د) فرض مارتینال بودن به صرفه‌های گذشته و آینده از هم مستقل هستند و هر استیج اطلاعاتی را می‌تواند برای پیش‌بینی صرفه بعدی خود دارد و نیازی به اطلاعات استیج‌ها در آن‌ها نیست. در اینجا ما فضای حالت را کاهش دادیم و معادله دینامیک موجود است. ما از دست داده ایم در نتیجه استقلال صرفه‌های گذشته و آینده در اینجا محقق است. به‌طور کلی باید دید با این تغییر چه اتفاقی می‌افتد.

(3) انتظاری برای اثبات سازه نظری تغییر باید ثابت کنیم:

$$\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| > \epsilon) = 0$$

$$E_{s \sim p(s), a \sim \pi_0(s, a)} \frac{\pi_1(s, a)}{\pi_0(s, a)} R(s, a) = \sum \frac{\pi_1(s, a)}{\pi_0(s, a)} \times R(s, a) \times p(s) \times \pi_0(s, a)$$

که برابر صورت سوال است:

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

(ب)

$$E_{s \sim p(s), a \sim \pi_0(s, a)} \frac{\pi_1(s, a)}{\pi_0(s, a)} = \sum \frac{\pi_1(s, a)}{\pi_0(s, a)} \times p(s) \times \pi_0(s, a) = \sum \pi_1(s, a) \times p(s)$$

حذف کنیم برابر است. صورت دوم مورد الف است. اثبات کردیم در نتیجه این تغییر سازه نظری است.

$$\sum p(a|s) p(s) = 1 \rightarrow$$

(ج) در هر مقطعی نمونه داشته باشیم تغییر برابر $R_{\pi_0}(s, a)$ می‌شود که از توزیع π_0 پیروی می‌کند. در نتیجه امید ریاضی $R_{\pi_0}(s, a)$ برابر $E[R_{\pi_0}(s, a)]$ می‌شود و تغییر سازه نظری داریم.

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

(4) انفع

$$U^{\pi_0}(M) = 1 + \gamma (0.5 \times U^{\pi_0}(R) + 0.5 \times U^{\pi_0}(D))$$

$$2 \rightarrow U^{\pi_0}(M) = \frac{2-\gamma}{2-2\gamma}$$

$$U^{\pi_0}(R) = 2 + \gamma (1 \times U^{\pi_0}(R)) \rightarrow U^{\pi_0}(R) = \frac{2}{1-\gamma}$$

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

$\pi_0 = \text{peace}$

$$U^{\pi_0} = \begin{matrix} M & R & D \\ (5.5, 20, -10) \end{matrix}$$

(←)

$$\pi_1(M) = \arg\max \left(\begin{matrix} \text{peace: } 1 + 0.5(0.5 \times 20 - 0.5 \times 10) \\ \text{war: } 1 + 0.7(0.1 \times 5.5 - 0.2 \times 10 + 0.7 \times 20) \end{matrix} \right)$$

$\pi_1(M) = \text{War}$

$$\pi_1(R) = \arg\max \left(\begin{matrix} \text{peace: } 2 + 0.7 \times 20 \\ \text{war: } 2 + 0.7(0.2 \times 20 + 0.8 \times 5.5) \end{matrix} \right)$$

$\pi_1(R) = \text{peace}$

$$\pi_1(D) = \arg\max \left(\begin{matrix} \text{peace: } -1 + 0.7 \times -10 \\ \text{war: } -1 + 0.7 \times 5.5 \end{matrix} \right)$$

$\pi_1(D) = \text{war}$

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

$$Q(s, a) \leftarrow (1-\alpha)Q(s, a) + \alpha \cdot \text{sample} \quad (\text{ج.})$$

(فرض کنید $\gamma = 1, \alpha = 0.5$)

Start State	Action	State End	Reward
Desert	War	Desert	-2/0
Desert	War	Riverside	3/0
Riverside	Peace	Mountain	1/0
Mountain	Peace	Riverside	1/0

Mountain-peace	Riverside-peace	Desert-war
.	.	.
0	0	-1
0	0	1
0	0.5	1

0.75

0.5

1