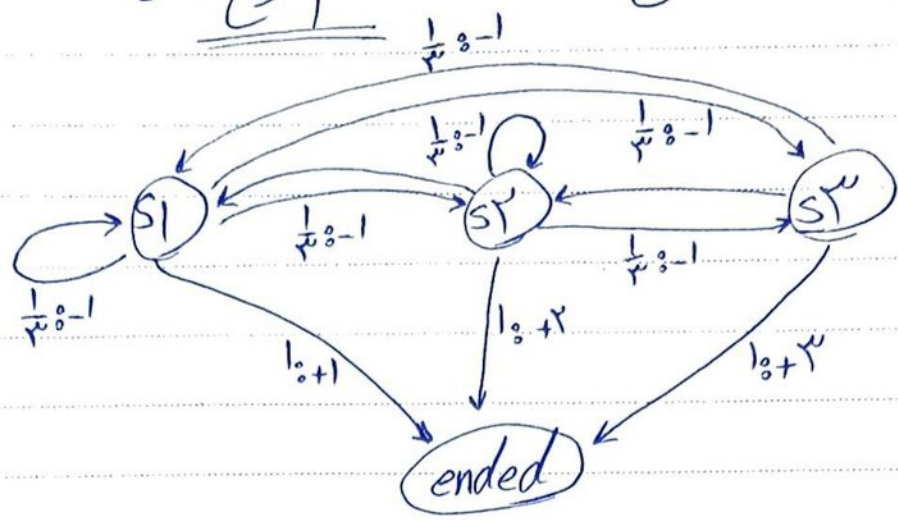


۱. آ، درست، در این صورت گران بالای برای $utility\ function$ وجود خواهد داشت و می توان دید که $value$ ها و $policy$ نهایی $converge$ می کنند.
- ب، نادرست، هر دو این روش ها سیاست بهینه را پیدا می کنند.
- ج، درست، پس از گذشت زمانی کافی، فارغ از $action$ های انتخاب شده، $Q-value$ علی صحیح به دست می آید و می توان از روی آنها سیاست بهینه را متوجه شد.
- د، نادرست، صرفاً به زمان بسیار بیشتری برای همگرا شدن Q ها نیاز دارد.
- ه، درست، چون نمودار Q بر اساس Q ها ابتدا صعودی است.
- و، انتقال سریع تر یا داشتن در محاسبات یادگیری سریع تر، وابستگی کمتر به مقدار موجود Q

الف) در هر دو قسمت ب و ج سیاست بهینه، اتمام مربع بازی است.

یادداشت: احتمال



	S_1	S_2	S_3	ended
U_0	0	0	0	0
U_1	+1	2	3	0
U_2	1	2	3	0
U_3	1	2	3	0

ب)

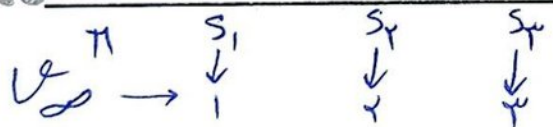
l	S_1	S_2	S_3	ended
$U^{\pi_0}_0$	0	0	0	0
$U^{\pi_0}_1$	1	2	-1	0
$U^{\pi_0}_2$	1	2	-0.4	0
$U^{\pi_0}_3$	1	2	-0.22	0
$U^{\pi_0}_\infty$	1	2	-0.143	0

ج) $U^{\pi_0}_{k+1}(S_3) = 0.3 \times U^{\pi_0}_k(S_3) - 0.1$

$U_3 = -0.144, U_\infty = -0.1491$

$U_4 = -0.14494$

$\pi_1(S_1) = \pi_1(S_2) = \text{اتمام}$, $\pi_1(S_3) = \text{اتمام}$ $\left\{ \begin{array}{l} \text{مانند} \rightarrow -1 + \frac{0.9}{3} \times 2.157 \\ \text{اتمام} \rightarrow 3 \end{array} \right.$
 $= -0.1429$



سیاست بهر در حالت quit است.

(۳) آ (افزاین زمانی) $\pi^*(s_1, s_2, s_3, s_4, s_5) = \text{climb}$, $\pi^*(s_v) = \text{quit}$

	s_1	s_2	s_3	s_4	s_5	s_v	(ب)
Q_0	%	%	%	%	%	%	climb/quit
Q_1	%	%	%	%	%	%	
Q_2	$\frac{3,5}{1}$	$\frac{4,5}{2}$	$\frac{3,5}{4}$	$\frac{4,5}{5}$	$\frac{1}{V}$		
Q_3	$\frac{4,5}{1}$	$\frac{4,5}{2}$	$\frac{4,5}{4}$	$\frac{5,1,5}{5}$	$\frac{3,5}{V}$		
Q_4	$\frac{5,1,2,5}{1}$	$\frac{5,2,5}{2}$	$\frac{5,1,2,5}{4}$	$\frac{5,1,5}{5}$	$\frac{4,5}{V}$		
Q_5	$\frac{5,5}{1}$	$\frac{4,4,3,5}{2}$	$\frac{5,5}{4}$	$\frac{4,1,5}{5}$	$\frac{5,1,5}{V}$		
Q_6	$\frac{V}{1}$	$\frac{V}{2}$	$\frac{V}{4}$	$\frac{V}{5}$	$\frac{V}{V}$		

(ج) $\pi^*(s_1, s_2, s_3, s_4) = \text{climb}$, $\pi^*(s_v) = \text{quit}$

$\frac{\delta}{r} \times (\max(2, Q(s_2, \text{climb})) + \max(1, Q(s_v, \text{climb}))) > 5$? climb: quit

$$V_{stay}(1) = \sum_{i=0}^{\infty} \gamma^i = 1 + \frac{1}{2} + \frac{1}{4} + \dots = 2 \quad (4) \quad (آ)$$

$$V_{(1)}^* = V_{stay}(1) = 2 \quad (ب)$$

ج، از آنجا که $\gamma = 1$ ، باید به شهری برویم که بیشترین r_i را دارد و سپس در آن شهر بمانیم و تجارت کنیم. از آنجا که p ها مثبت اند، انجام این کار ناممکن نیست.

د، i را کوچکترین مقداری در نظر بگیرید که به ازای آن $r_i > 0$ ، آنگاه بیشترین مقدار K که در آن ممکن است $V_K(1) = 0$ برابر است با $i-1$. اگر فرض شود به ازای هر i

$r_i > 0$ آنگاه بیشترین مقدار ممکن برای K صفر است.

ه، پس از یک مرحله $iter$ $V_K(i) = r_i > 0$ ، بنابراین بیشترین مقدار ممکن برای K صفر است.

sarsa	$Q(1, S)$	$Q(1, E)$	$Q(2, W)$	$Q(2, S)$	(و)
initial	0	0	0	0	
$1, S, 1$	2	0	0	0	
$1, E, 2$	2	0	0	0	
$2, S, 2$	2	0	0	3	
$2, W, 1$	2	0	1	3	
$1, S, 1$	4	0	1	3	

$$\gamma = 1, \alpha = LR = 0.5$$

(۶) (آ)

$\frac{1}{4}$	$\frac{1}{4}$	$\boxed{+1}$
$\frac{1}{8}$	$\frac{1}{4}$	$\frac{1}{4}$
$\frac{1}{14}$	$\frac{1}{8}$	$\frac{1}{4}$

(ب)

↓		X
↓		↑
→	→	↑

$$w_1 = w_2 = 1$$

$$V(s) = f_1(s) + f_2(s) = x + y$$

۲	۳	۴
۱	۲	۳
۰	۱	۲

(ج)

(د) بله، می‌توانیم. اگر فاصله منتهی خانه s با خانه $+1$ را با $d(s)$ نشان دهیم:

$$V^*(s) = 2^{-d(s)} \quad , \quad d(s) = 2 - x + 2 - y = 4 - (f_1(s) + f_2(s))$$

(ه) ا. در قسمت ج نشان دادیم که برای A قابل حل است اما برای B و C به دلیل

ساختار شبکه آنها قابل حل نیست.

ii. با این وزن‌ها برای A قابل حل است:

مشابه قسمت ج می‌شود:

برای شبکه B نیز قابل حل است:

$$w_{2,0} = 4, w_{2,1} = 5, w_{2,2} = 6$$

برای شبکه C نیز قابل حل است:

$$w_{0,0} = w_{0,2} = w_{1,0} = w_{1,2} = 0, w_{1,1} = 2$$

$$w_1 = w_2 = -1, w_3 = 2$$

۶) iii. برای C قابل حل است:

برای A, B قابل حل نیست.

$$W^*(y) = \sum_{s'} P(s'|y) V^*(s') \quad (۵) \text{ آ}$$

$$Q^*(s, a) = R(s, a) + \sum_{s'} T(s, a, s') \delta V^*(s') \quad (ب)$$

$$V^*(s) = \max_a (R(s, a) + \delta W^*(f(s, a))) \quad \leftarrow \text{به گدازت الف}$$

$$W^*(y) = \sum_{s'} p(s'|y) \left(\max_a (R(s', a) + \delta W^*(f(s', a))) \right) \quad (ج)$$

$$\pi^{(l+1)}(s) = \arg \max_a (R(s, a) + \delta W^{\pi^i}(f(s, a))) \quad (د) \text{ جا خالی اول}$$

$$W(y_t) \leftarrow (1-\alpha)W(y_t) + \alpha (R(s_{t+1}, a_{t+1}) + \delta W(y_{t+1})) \quad (ه)$$

$$W^{\pi^{l+1}}(s) \quad \text{اگر مقدار } W^{\pi^i}(s) \text{ همگرا نشده بود و اختلاف آن با } W^{\pi^{l+1}}(s) \text{ جا خالی دوم}$$

ناجیز و قابل جستجو نیستی بود.