

سوال اول

الف) این عبارت درست است. حل معادلات Bellman در روش value iteration، در واقع پیدا کردن مقدار fixed-point برای این معادلات است. با توجه به مطالب تدریس شده، مقدار بدست آمده توسط value iteration همان جواب Bellman equations است که V^* می باشد. در نتیجه اگر مقدار تابع policy به ازای یک T بزرگ به اندازه کافی تغییر نکند، می توان مطمئن بود که مقادیر V به V^* همگرا شده است و در نتیجه مقادیر تابع policy نیز همگرا شده است.

ب) این عبارت نادرست است. طبق مطالب گفته شده درون اسلایدها، به ازای انتخاب درست مقدار α ، حتی به ازای sub-optimal policy، مقادیر Q نهایتاً به مقادیر Q برای optimal policy همگرا خواهد شد. در نتیجه، می توان با استفاده policy نادرست در ابتدا، مقادیر Q را به ازای π^* یاد گرفت.

ج) اگر $\text{learning rate} = 1$ باشد، این بدین معنی است که در هر مرحله از Q-learning، مقدار بدست آمده برای sample به جای Q قبلی قرار می گیرد. با توجه به اینکه MDP قطعی است، میانگین گیری از حرکات صورت گرفته خیلی معنا ندارد زیرا به ازای هر دو s, a ، می دانیم که یک حرکت رخ می دهد. در نتیجه، در صورت استفاده از $\alpha = 1$ برای یک deterministic MDP، مقادیر Q به مقادیر بهینه همگرا خواهند شد.

د) پیچیدگی زمانی اجرای هر iteration از بخش iterative الگوریتم value iteration از $O(|S|^2|A|)$ می باشد. همچنین این پیچیدگی برای اجرای هر iteration از بخش iterative الگوریتم policy iteration، در بخش policy evaluation از $O(|S|^2)$ است. حال دقت کنید که جواب این سوال، بستگی به تعریف iteration دارد. اگر منظور طراح سوال از iteration در هر دو الگوریتم، مرحله آپدیت کردن V باشد، این گزاره صحیح خواهد بود. اما اگر منظور طراح سوال از iteration برای value iteration، مرحله آپدیت کردن V و برای policy iteration، مرحله آپدیت کردن π باشد، این گزاره نادرست خواهد شد.

ه) نادرست است. فرض کنید در یک grid exploration هستیم که مقدار living reward به طور پیش فرض برابر 1- است. به طور شهودی، محرک تلاش می کند که کمترین زمان را در محیط چرخه بزند و زودتر به پاداش بزرگ برسد. اما اگر همه rewardها را با مقدار 1 جمع کنیم، چرخ زدن در محیط دیگر ضرری به محرک نخواهد زد و در نتیجه محرک از پذیرفتن ریسک گذر از خطر دوری می کند. دقت کنید که تحلیل صورت گرفته برای gridهایی می باشد که دارای terminal state می باشند.

در صورتی که بازی تا ابد ادامه داشته باشد، مقدار V برای تمامی stateها با مقدار ثابت $\sum_{i=0}^{\infty} c\gamma^i = \frac{c}{1-\gamma}$ جمع می شود. در این صورت، حرکتی که $Q(s,a)$ را قبل از افزایش تمامی rewardها با مقدار c بیشینه می کرد، هنوز هم مقدار $Q'(s,a)$ را بیشینه خواهد کرد. در نتیجه تفاوتی ایجاد نخواهد شد.

Value Iteration $\rightarrow V_{K+1}(s) = \max_a \sum_{s'} T(s, a, s') \{ R(s, a, s') + \gamma V_K(s') \}$ (Cell 2)

$$V_0 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

مقدار صفر است. چون هیچ کاری نداریم و هیچ کاری نداریم. پس mark صفر است.

first iter:

$$V_1 \{1,1\} = \max \{0, 0, 0, 0\} = 0$$

$$V_1 \{2,1\} = \max \{0, 0, 0, 0\} = 0$$

$$V_1 \{1,2\} = \max \{0, 1(-\gamma + 0) + 0, 1(0 + 0) + 0, 1(0 + 0) + 0, 1(0 + 0) + 0\} = 0 \rightarrow$$

$$V_1 \{2,2\} = \max \{0, 1(0 + 0) + 0, 1(0 + 0) + 0, 1(0 + 0) + 0, 1(0 + 0) + 0\} = 0 \rightarrow$$

$$\rightarrow V_1 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Second iter:

$$V_2 \{1,1\} = \max \{0, 0, 0, 0\} = 0$$

$$V_2 \{2,1\} = 0, 1(0 + \gamma) + 0, 1(0 + 0) + 0, 1(0 + 0) + 0 = 0, 1 \rightarrow$$

نه، وضعی حرکت می‌کند. به سمت چپ می‌رود.

$$V_2 \{1,2\} = \max \{0, 1(0 + \gamma) + 0, 1(-\gamma + 0) + 0, 1(0 + 0) + 0, 1(0 + 0) + 0\} = 0, 1 \uparrow$$

$$V_2 \{2,2\} = \max \{0, 1(0 + \gamma) + 0, 1(0 + \gamma) + 0, 1(0 + 0) + 0, 1(0 + 0) + 0\} = 0, 1 \rightarrow$$

$$\rightarrow V_2 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

third iter :

$$V_{\pi}[1,1] = 0,1(0+2,11\delta) + 0,1(0+2,41\delta) + 0,1(0+0) = 2,2811 \uparrow$$

وضع حرکت به بالا مطلوب است.

$$V_{\pi}[2,1] = 0,1(0+2,44\delta) + 0,1(0+2,11\delta) + 0,1(0+0) = 2,2911 \rightarrow$$

وضع حرکت به راست مطلوب است.

$$V_{\pi}[1,2] = \max \left\{ \begin{aligned} &0,1(0+2,44\delta) + 0,1(-2+0) + 0,1(0+0), \\ &0,1(-2+0) + 0,1(0+2,44\delta) + 0,1(0+2,11\delta), \\ &0,1(0+2,11\delta) + 0,1(0+0) + 0,1(-2+0), \\ &0,1(0+0) + 0,1(0+2,44\delta) + 0,1(0+2,11\delta) \end{aligned} \right\} = 2,4392 \uparrow$$

$$V_{\pi}[2,2] = 0,1(2+0) + 0,1(0+2,44\delta) + 0,1(0+2,11\delta) = 2,4044 \rightarrow$$

وضع حرکت به راست مطلوب است.

$$\rightarrow V_{\pi} = \begin{matrix} 2,2911 & \rightarrow & 2,4044 & \rightarrow & 0 \\ \uparrow & & \uparrow & & \\ 2,2811 & & 2,4392 & & 0 \end{matrix}$$

ب) جتنی که کار به کار دیگری π^* در این حالت است و π^* از π بهتر است Q-learning می باشد. π به π^* ساده ترین این می باشد، Passive Q-learning و Off policy learning.

$$\text{sample} = R(s, a, s') + \gamma V^{\pi}(s') \quad (ج)$$

$$V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha (\text{sample} - V^{\pi}(s))$$

$$\text{مسیر اول} \rightarrow (1,1) \rightarrow (1,2) \rightarrow (2,2) \rightarrow (2,1) \quad V^{\pi}(1,1) \rightarrow \text{sample} = 0 + \gamma \cdot 0 = 0$$

$$V(1,1) \leftarrow 0 + \alpha(0 + \gamma \cdot 0 - 0) = 0$$

$$V^{\pi}(1,2) \rightarrow V(1,2) \leftarrow 0 + \alpha(0 + \gamma \cdot 0 - 0) = 0$$

$$V(2,2) \leftarrow 0 + \alpha(2 + \gamma \cdot 0 - 0) = 0,2$$

$$\text{مسیر دوم} \rightarrow (1,1) \rightarrow (1,2) \rightarrow (1,1)$$

$$V(1,1) \leftarrow 0 + \alpha(0 + \gamma \cdot 0 - 0) = 0$$

$$V(1,2) \leftarrow 0 + \alpha(-2 + \gamma \cdot 0 - 0) = -0,2$$

اینها فقط مثال این ترتیب می باشد نه به ترتیب است :

$$\text{مسیر ۱} \rightsquigarrow (1,1) \rightarrow (1,2) \rightarrow (1,3)$$

حالت ۱ :

$$V(1,1) \leftarrow 0 + \alpha(0 + \gamma x_0 - 0) = 0$$

$$V(1,2) \leftarrow 0 + \alpha(-\Delta + \gamma x_0 - 0) = -0,1\Delta$$

$$\text{مسیر ۲} \rightsquigarrow (1,1) \rightarrow (1,2) \rightarrow (2,2) \rightarrow (2,3)$$

$$V(1,1) \leftarrow 0 + \alpha(0 + \gamma x_0 - 0,1\Delta - 0) = -0,1\Delta$$

$$V(1,2) \leftarrow -0,1\Delta + \alpha(0 + \gamma x_0 + 0,1\Delta) = -0,1\Delta$$

$$V(2,2) \leftarrow 0 + \alpha(\Delta + \gamma x_0 - 0) = 0,1\Delta$$

$$V^{\pi} : \begin{array}{ccc} 0 & 0,1\Delta & 0 \\ -0,1\Delta & -0,1\Delta & 0 \end{array}$$

(الف)

(۳) معادلات بحین ضرورت زیر حتمی :

$$1) \forall s : V^*(s) = \max_a Q^*(s,a)$$

$$2) \forall s,a : Q^*(s,a) = \sum_{s'} T(s,a,s') \{ R(s,a,s') + \gamma V^*(s') \}$$

می توان این دو معادله را بصورت زیر بنویسید :

$$\forall s,a : Q^*(s,a) = \sum_{s'} T(s,a,s') \{ R(s,a,s') + \gamma \max_{a'} Q^*(s',a') \}$$

۱. با حل این معادلات ، ماتریس Q تشکیل می دهیم (با ابعاد $|S| \times |A|$) ، به نحایی $Q(s,a)$ ها را برای هر حالت s و a می یابیم.برای حل این معادلات ، می توان از روش iterative استفاده کرد . در ابتدا به ازای هر حالت s ، به صورت ماتریس Q ، ماتریس Q^* ، نقطه ثابت این معادله خواهیم داشت .

$$\text{Iteration : } Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \{ R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \}$$

در روش های RL ، به توان T و R دسترسی نداریم ، در نتیجه این مرحله iteration ، به صورت سری به انجام می رسد .

subject:

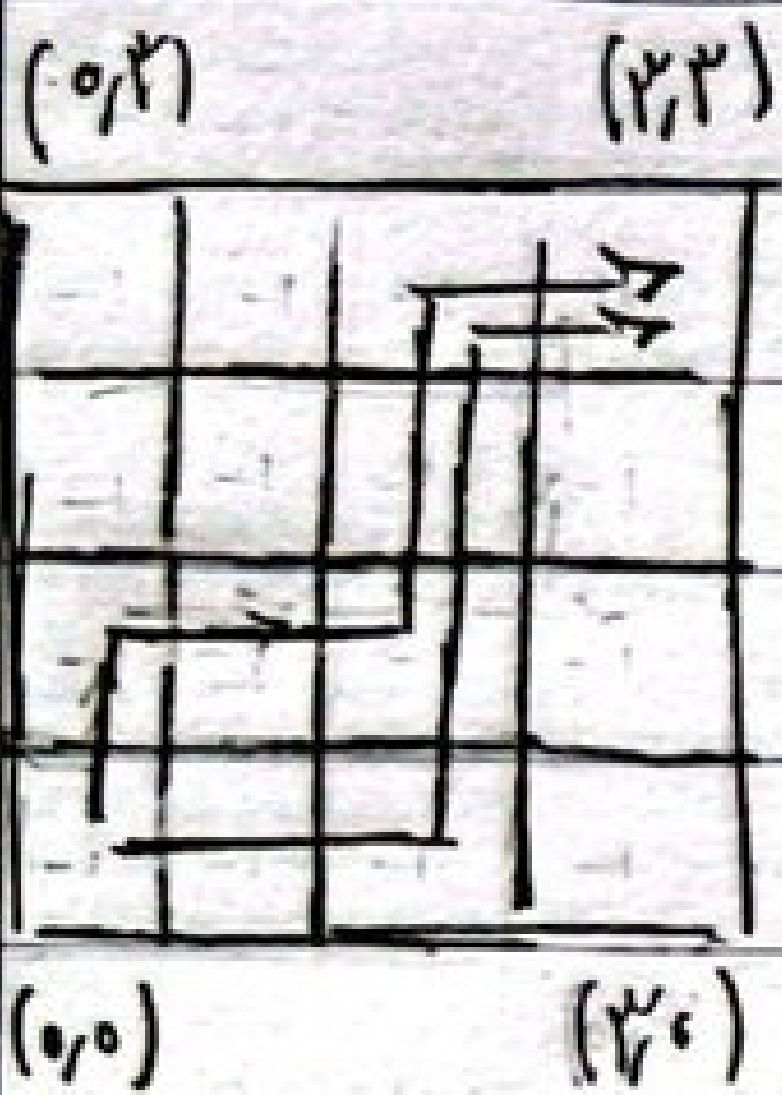
date:

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

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

✓ به مقدار این کار، از exp moving avg مقدار می گیریم:

در یادگیری بدون نظارت، می توان از روش های Active و Passive برای یادگیری مقدار Q استفاده کرد.
 در روش Passive یک fixed policy داریم که تنها optimal نیست، هدف ما یافتن آن است.
 در روش Active به جای آن، تغییر می دهیم و در واقع exploration صورت می گیرد.



ب) فرض کنیم در 2 iteration، از مسیرهای زیر حرکت کردیم:

این مسیر را با یک tempdiff مقایسه می کنیم. در بخش اول با هم مقایسه می کنیم. Q به صورت زیر می شود:

$$Q((0,0), \uparrow) = -0.1, Q((1,1)) = Q((1,2)) = Q((2,2)) = Q((2,1)) = -0.1$$

$$Q((0,1), \rightarrow) = -0.1$$

$$Q((1,1), \rightarrow) = -0.1$$

$$Q((2,1), \uparrow) = -0.1$$

$$Q((2,2), \uparrow) = -0.1$$

$$Q((2,3), \rightarrow) = 1$$

$$Q((0,0), \rightarrow) = -0.1$$

$$Q((1,0), \rightarrow) = -0.1$$

$$Q((2,0), \uparrow) = 0 + 0.1(-1 + 0.9 \max Q(2,1)) - 0.1 = -0.1$$

$$Q((2,1), \uparrow) = -0.1 + 0.1(-1 + 0.9 \max Q(2,2) + 0.1) = -0.19$$

$$Q((2,2), \uparrow) = -0.1 + 0.1(-1 + 0.9 \max Q(2,3) + 0.1) = -0.1$$

$$Q((2,3), \rightarrow) = 1 + 0.1(10 + 0.9 \max Q(2,2) - 1) = 1.9$$

ح) در روش ϵ -greedy، هر گاه با احتمال ϵ به صورت random عمل می کنیم. این روش می تواند explore و exploit را به هم ترکیب کند.

Simulated annealing: در این روش، به صورت random عمل می کنیم.

with prob ϵ act randomly

with prob $1-\epsilon$ act on current policy

$$\text{argmax}_a \hat{Q}(s, a)$$

1) explore
2) exploit

$$a' = \arg \max_a Q^*(s, a) \rightarrow v^*(s) = Q^*(s, a')$$

(الف) (۴)

$$\pi(s) = \arg \max_a \tilde{Q}(s, a) \rightarrow \tilde{Q}(s, \pi(s)) \geq \tilde{Q}(s, a') \quad (۳)$$

$$\max_{s, a} |Q^*(s, a) - \tilde{Q}(s, a)| \leq \epsilon \rightarrow Q^*(s, a') - \tilde{Q}(s, a') \leq \epsilon \quad (۱)$$

$$\rightarrow \tilde{Q}(s, \pi(s)) - Q^*(s, \pi(s)) \leq \epsilon \quad (۲)$$

$$(۱) \leadsto v^*(s) - \tilde{Q}(s, a') \leq \epsilon$$

$$(۲) \leadsto v^*(s) - \tilde{Q}(s, \pi(s)) \leq \epsilon$$

$$(۲) \leadsto v^*(s) - Q^*(s, \pi(s)) \leq 2\epsilon$$

$$v^*(s) - v_{\pi}(s) = f(s)$$

(ب)

$$\rightarrow v^*(s) - Q^*(s, \pi(s)) + Q^*(s, \pi(s)) - v_{\pi}(s) \leq 2\epsilon + Q^*(s, \pi(s)) - v_{\pi}(s)$$

$$Q^*(s, \pi(s)) = \sum_{s'} T(s, \pi(s), s') \{ R(s, \pi(s), s') + \gamma v^*(s') \} = R(s, a) + \gamma \sum T(s, \pi(s), s') v^*(s')$$

$$v_{\pi}(s) = \sum_{s'} T(s, \pi(s), s') \{ R(s, \pi(s), s') + \gamma v_{\pi}(s') \} = R(s, a) + \gamma \sum T(s, \pi(s), s') v_{\pi}(s')$$

$$\rightarrow Q^*(s, \pi(s)) - v_{\pi}(s) = \gamma \sum T(s, \pi(s), s') \{ v^*(s') - v_{\pi}(s') \} = \gamma E \{ v^*(s') - v_{\pi}(s') \}$$

$$E \{ v^*(s') - v_{\pi}(s') \} \leq \max_{s'} \{ v^*(s') - v_{\pi}(s') \} = v^*(s'') - v_{\pi}(s'') = f(s'')$$

$$\leadsto \forall s \exists s'' \text{ such that : } f(s) \leq 2\epsilon + \gamma f(s'')$$

$$\leadsto f(s) \leq 2\epsilon \sum_{i=0}^n \gamma^i + \gamma^{n+1} f(s'')$$

این به صند می‌رسد چنان که لا، حدسین
 $v^*(s)$ و $v_{\pi}(s)$ کردن در حدسین.

$$\lim_{n \rightarrow \infty} : v^*(s) - v_{\pi}(s) \leq \frac{2\epsilon}{1-\gamma}$$

$$v^*(s_1) = \max \{ Q^*(s_1, \text{stay}), Q^*(s_1, q_0) \}$$

(c)

$$v^*(s_1) = r\varepsilon + \gamma v^*(s_1) \leadsto v^*(s_1) = \frac{r\varepsilon}{1-\gamma}$$

$$Q^*(s_1, \text{stay}) = 0 + \gamma v^*(s_1)$$

$$Q^*(s_1, q_0) = r\varepsilon + \gamma v^*(s_1) = \frac{r\varepsilon}{1-\gamma}$$

$$\textcircled{1} v^*(s_1) = Q^*(s_1, \text{stay}) \xrightarrow{\gamma=1.0} v^*(s_1) = Q^*(s_1, q_0) = \frac{r\varepsilon}{1-\gamma}$$

or $v^*(s_1) = 0$

$$\xrightarrow{Q^*(s_1, \text{stay}) = \frac{r\varepsilon\gamma}{1-\gamma}}$$

$$\|Q^* - \tilde{Q}\| < \varepsilon \xrightarrow{\text{up}} \tilde{Q}(s_1, \text{stay}) = Q^*(s_1, \text{stay}) + \varepsilon = \varepsilon \left(\frac{r\gamma}{1-\gamma} + 1 \right) = \varepsilon \times \frac{1+\gamma}{1-\gamma} \quad \textcircled{1}$$

$$\xrightarrow{\tilde{Q}(s_1, q_0) = Q^*(s_1, q_0) - \varepsilon = \varepsilon \left(\frac{r}{1-\gamma} - 1 \right) = \varepsilon \times \frac{1+\gamma}{1-\gamma}} \quad \textcircled{2}$$

$$\textcircled{2}, \textcircled{1} \leadsto \tilde{Q}(s_1, \text{stay}) = \tilde{Q}(s_1, q_0)$$

$$\hookrightarrow \pi(s_1) = \text{stay}$$

$$v_{\pi}(s_1) = 0 + \gamma v_{\pi}(s_1) \xrightarrow{v_{\pi}(s_1) = 0}$$

$$\xrightarrow{v_{\pi}(s_1) = v^*(s_1) - \frac{r\varepsilon}{1-\gamma}}$$

$$v^*(s_1) = 1 + \delta v^*(s_1) \rightarrow v^*(s_1) = \frac{1}{1-\delta}$$

$$\chi^{\mu}(\delta r) = 0$$

بجای γ به γ' و γ به γ ، $\gamma' \neq \gamma$ ، $Q^*(s_0, a_1) > Q^*(s_0, a_2)$ ← عمل a_1 است.

$$v_0(\delta_i) = 0, \forall i$$

$$v_i(s_0) = \max \left(0 + \delta v_i(s_1), \frac{\gamma}{1-\delta} + \delta v_i(s_r) \right) = \frac{\gamma}{1-\delta} \rightarrow \text{استجاب می کند} \quad a_r$$

$$v_1(\delta_1) = 1$$

$$\mathcal{V}_1(\delta_y) = 0$$

$$\rightarrow \left. \begin{aligned} v_n(s_1) &= 1 + \gamma + \gamma^r + \dots + \gamma^n = \frac{1 - \gamma^{n+1}}{1 - \gamma} \\ v_n(s_r) &= 0 \end{aligned} \right\} v_n(s_0) = \max \left\{ \gamma v_{n-1}(s_1), \frac{\gamma^r}{1 - \gamma} \right\}$$

$\rightarrow v_n(s_0)$ finds optimal action when $\rightarrow \gamma v_{n-1}(s_1) \} \frac{\gamma^r}{1-\gamma}$

$$\rightarrow 1 - 2^{-1} 2^{-2} 2^{-3} \dots$$

$$1 - r \geq r^n$$

$$\log(1-x) \neq n \log(x)$$

$$\rightarrow n^* \gg \frac{\log(1-\delta)}{\log \delta}$$