

به نام خدا



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دانشکده مهندسی کامپیوتر

مبانی و کاربردهای هوش مصنوعی ترم پاییز ۱۴۰۱

پاسخنامه تمرین سوم

سوال ۱)

۱- نادرست است. مطابق با اسلایدها، فقط به حالت کنونی وابسته است:

- “Markov” generally means that given the present state, the future and the past are independent

۲- نادرست است. مطابق با اسلایدها، در direct evaluation پس از انجام چند sample با میانگین گرفتن ارزش را آپدیت می کنیم.

۳- نادرست است. یک مثال نقض در نظر می گیریم. یک MDP با دو حالت ترمینال A و B را در نظر بگیرید. انتقال به حالت A پاداش 1 و انتقال به حالت B پاداش 10 دارد. همه انتقال های دیگر پاداش صفر دارند. فرض کنید A یک قدم به سمت شمال از حالت شروع باشد. همچنین B دو قدم به سمت جنوب از حالت شروع باشد. (اقدامات همیشه موفق هستند). حالا اگر $\gamma < 0.1$ باشد سیاست بهینه عامل را یک قدم به سمت شمال از حالت شروع به A می برد، اما اگر $\gamma > 0.1$ سیاست بهینه عامل را دو گام به سمت جنوب از حالت شروع به B می رساند.

۴- درست است. Q-learning نیازی به مدلی از محیط ندارد و model-free است. همچنین اگر غیر بهینه (suboptimal) عمل شود، باز هم Q-learning به سیاست بهینه همگرا می‌شود و off-policy است.

سوال ۲) سرعت هلیکوپتر، موقعیت و زاویه با زمین هلیکوپتر، فاصله تا نزدیکترین مانع، سرعت باد و وضعیت آب و هوا و ...

حتی اگر از ویژگی‌های گفته شده استفاده کنیم و همه ضرایب را پیدا کنیم هم باز ممکن است به مشکل برخوردیم. مثلاً فرض کنید که هلیکوپتر به شکل برعکس در حال پرواز است. ممکن است با موانع فاصله داشته باشیم و همچنین همه چیز خوب باشد اما قطعاً پرواز به شکل برعکس مطلوب نیست.

* در این سوال پاسخ‌های قابل قبول دیگر نیز درست هستند.

سوال ۳)

الف) بله، همچنان یک سیاست بهینه دریافت می‌شود. این روش عملاً مشابه value iteration است، زیرا value iteration شامل یک مرحله ارزیابی (evaluation) نیز می‌شود. در این مورد، ما مقادیر را بر اساس بهترین سیاست فعلی خود به روز می‌کنیم و این کار را تا زمان همگرایی یک سیاست ادامه می‌دهیم.

فرمول value iteration به صورت زیر است:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

و فرمول policy evaluation به صورت زیر است:

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$

با داشتن یک سیاست ثابت π که در مرحله policy improvement به روز رسانی می‌کنیم:

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

که بسیار شبیه به policy extraction است که در value iteration انجام می‌شود.

ب) بله. با تکرارهای کافی، عامل می‌تواند همچنان Q-value ها را از اکشن‌های خود محاسبه کند. با این حال، یادگیری به طور قابل توجهی بیشتر از حالت $\epsilon < 1$ طول می‌کشد، زیرا عامل نمی‌تواند مقادیر Q ای که یاد گرفته را exploit کند، بنابراین اشتباهات مشابه را مکرراً انجام می‌دهد تا زمانی که راه حل صحیح به طور تصادفی پیدا شود.

سوال ۴) طبق روابط زیر محاسبات را انجام می‌دهیم.

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

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

$$V_1(A) = \max \{1 [1 + 0]\} = 1$$

$$V_1(B) = \max \{1 [1 + 0], 1[-2 + 0]\} = 1$$

$$V_1(C) = \max \{(0.5)[1 + 0] + (0.5)[4 + 0], 1[-2 + 0]\} = 2.5$$

$$V_1(D) = \max \{1[1 + 0], 1[-2 + 0]\} = 1$$

$$V_1(E) = \max \{1[1 + 0], 1[-2 + 0]\} = 1$$

$$V_1(F) = \max \{1[10 + 0], 1[-2 + 0]\} = 10$$

$$V_1(G) = 0$$

$$V_2(A) = \max \{1 [1 + 1]\} = 2$$

$$V_2(B) = \max \{1 [1 + 2.5], 1[-2 + 1]\} = 3.5$$

$$V_2(C) = \max \{(0.5)[1 + 1] + (0.5)[4 + 1], 1[-2 + 2.5]\} = 3.5$$

$$V_2(D) = \max \{1 [1 + 1], 1[-2 + 2.5]\} = 2$$

$$V_2(E) = \max \{1[1 + 10], 1[-2 + 1]\} = 11$$

$$V_2(F) = \max \{1[10 + 0], 1[-2 + 1]\} = 10$$

$$V_2(G) = 0$$

پاسخ‌های نهایی به صورت زیر است:

$$V_2(B) = 3.5$$

$$Q_2(B, \text{right}) = 1 [1 + 2.5] = 3.5$$

$$Q_2(B, \text{left}) = 1[-2 + 1] = -1$$

	A	B	C	D	E	F	G
V_1	1	1	2.5	1	1	10	0
V_2	2	3.5	3.5	2	11	10	0

ب) طبق به روز رسانی در Q-learning داریم:

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \left[r + \gamma \max_{a'} Q(s', a') \right]$$

$$Q(C, \text{jump}) = (0.5)(0) + (0.5)(4 + 0) = 2$$

$$Q(E, \text{right}) = (0.5)(0) + (0.5)(1 + 0) = 0.5$$

$$Q(F, \text{left}) = (0.5)(0) + (0.5)(-2 + \max \{Q(E, \text{right}), Q(E, \text{left})\}) = 0.5(-2 + \max \{0.5, 0\}) = -0.75$$

$$Q(E, \text{right}) = (0.5)(0.5) + (0.5)(1 + \max \{Q(F, \text{right}), Q(F, \text{left})\}) = 0.5(1 + \max \{0, -0.75\}) = 0.75$$

	Q(C, left)	Q(C, jump)	Q(E, left)	Q(E, right)	Q(F, left)	Q(F, right)
Initial	0	0	0	0	0	0
Transition 1		2				
Transition 2				0.5		
Transition 3					-0.75	
Transition 4				0.75		

یک state-action تنها زمانی به روز می شود که یک انتقال از آن انجام گیرد. $Q(C, \text{left}), Q(E, \text{left}), Q(F, \text{right})$. هرگز تجربه نشده اند و بنابراین این مقادیر به روز نمی شوند.

سوال ۵)

اولین مرحله از policy iteration, انجام policy evaluation است که مطابق با فرمول زیر انجام می شود:

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$

در ابتدا از $k=0$ و $i=0$ شروع می کنیم که مقدار $V_0^{\pi_0}$ برای هر وضعیت در صورت سوال داده شده است. ابتدا مقدار V_1 را برای سیاست ابتدایی محاسبه می کنیم:

Iteration 1:

Policy evaluation (calculate utilities):

$$\begin{aligned} V_1^{\pi_0}(\text{High}) &= T(\text{High}, \text{study}, \text{High}) \cdot (R(\text{High}, \text{study}, \text{High}) + \gamma V_0(\text{High})) \\ &\quad T(\text{High}, \text{study}, \text{Low}) \cdot (R(\text{High}, \text{study}, \text{Low}) + \gamma V_0(\text{Low})) \\ &= 1.0 \times (0 + \gamma 0) + 0 = 0 \end{aligned}$$

$$\begin{aligned} V_1^{\pi_0}(\text{Low}) &= T(\text{Low}, \text{study}, \text{High}) \cdot (R(\text{Low}, \text{study}, \text{High}) + \gamma V_0(\text{High})) \\ &\quad T(\text{Low}, \text{study}, \text{Low}) \cdot (R(\text{Low}, \text{study}, \text{Low}) + \gamma V_0(\text{Low})) \\ &= 0.3 \times (0 + \gamma 0) + 0.7 \times (0 + \gamma 0) = 0 \end{aligned}$$

یک بار دیگر ارزیابی را انجام می دهیم تا ببینیم که آیا value همگرا می شود یا خیر:

$$\begin{aligned} V_2^{\pi_0}(\text{High}) &= T(\text{High}, \text{study}, \text{High}) \cdot (R(\text{High}, \text{study}, \text{High}) + \gamma V_1(\text{High})) \\ &\quad T(\text{High}, \text{study}, \text{Low}) \cdot (R(\text{High}, \text{study}, \text{Low}) + \gamma V_1(\text{Low})) \\ &= 1.0 \times (0 + \gamma 0) + 0 = 0 \end{aligned}$$

$$\begin{aligned} V_2^{\pi_0}(\text{Low}) &= T(\text{Low}, \text{study}, \text{High}) \cdot (R(\text{Low}, \text{study}, \text{High}) + \gamma V_1(\text{High})) \\ &\quad T(\text{Low}, \text{study}, \text{Low}) \cdot (R(\text{Low}, \text{study}, \text{Low}) + \gamma V_1(\text{Low})) \\ &= 0.3 \times (0 + \gamma 0) + 0.7 \times (0 + \gamma 0) = 0 \end{aligned}$$

می بینیم که مقدار V تغییر نکرد، پس همگرا شده. در مرحله بعد، policy improvement را مطابق زیر انجام می دهیم:

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

توجه داشته باشید که V در فرمول بالا و محاسبات زیر، برابر با V ای است که در قسمت policy evaluation بدست آوردیم. مطابق با روند زیر، سیاست را بهبود می بخشیم.

Policy improvement:

For state "high":

- Utility of choosing "study": $1 \cdot (0 + 0.5 \cdot 0) + 0 \cdot (0 + 0.5 \cdot 0) = 0$
- Utility of choosing "exam": $0.9 \cdot (10 + 0.5 \cdot 0) + 0.1 \cdot (-10 + 0.5 \cdot 0) = 8$
- Utility of choosing "Netflix": $1.0 \cdot (1 + 0.5 \cdot 0) = 1$

So set $\hat{\pi}_1(\text{high}) \leftarrow \text{exam}$.

For state "low":

- Utility of choosing "study": $0.3 \cdot (0 + 0.5 \cdot 0) + 0.7 \cdot (0 + 0.5 \cdot 0) = 0$
- Utility of choosing "exam": $0.05 \cdot (10 + 0.5 \cdot 0) + 0.95 \cdot (-10 + 0.5 \cdot 0) = -9$
- Utility of choosing "Netflix": $1.0 \cdot (1 + 0.5 \cdot 0) = 1$

So set $\hat{\pi}_1(\text{low}) \leftarrow \text{Netflix}$.

در iteration دوم، باز هم مانند قبل عمل می کنیم:

مرحله اول: policy evaluation

ابتدا توجه کنید چون بر روی value همگرا شدیم، در ابتدای این پیمایش خواهیم داشت:

$$V^{\pi_1}_0(\text{High}) = 0$$

$$V^{\pi_1}_0(\text{Low}) = 0$$

حال تا سه مرحله همگرایی را بررسی می کنیم. توجه کنید که سیاست استفاده شده در این مراحل از سیاست بدست آمده در policy improvement است.

Iteration 2:

Policy evaluation:

$$\begin{aligned}V^{\pi^1}_1(\text{High}) &= T(\text{High}, \text{exam}, \text{High}) \cdot (R(\text{High}, \text{exam}, \text{High}) + \gamma V_0(\text{High})) \\&\quad T(\text{High}, \text{exam}, \text{Low}) \cdot (R(\text{High}, \text{exam}, \text{Low}) + \gamma V_0(\text{Low})) \\&= 0.9 \times (10 + \gamma 0) + 0.1 \times (-10 + \gamma 0) = 8\end{aligned}$$

$$\begin{aligned}V^{\pi^1}_1(\text{Low}) &= T(\text{Low}, \text{Netflix}, \text{High}) \cdot (R(\text{Low}, \text{Netflix}, \text{High}) + \gamma V_0(\text{High})) \\&\quad T(\text{Low}, \text{Netflix}, \text{Low}) \cdot (R(\text{Low}, \text{Netflix}, \text{Low}) + \gamma V_0(\text{Low})) \\&= 0.0 \times (1 + \gamma 0) + 1.0 \times (1 + \gamma 0) = 1\end{aligned}$$

$$\begin{aligned}V^{\pi^1}_2(\text{High}) &= T(\text{High}, \text{exam}, \text{High}) \cdot (R(\text{High}, \text{exam}, \text{High}) + \gamma V_1(\text{High})) \\&\quad T(\text{High}, \text{exam}, \text{Low}) \cdot (R(\text{High}, \text{exam}, \text{Low}) + \gamma V_1(\text{Low})) \\&= 0.9 \times (10 + \gamma 8) + 0.1 \times (-10 + \gamma 1) = 12.6 - 0.95 = 11.65\end{aligned}$$

$$\begin{aligned}V^{\pi^1}_2(\text{Low}) &= T(\text{Low}, \text{Netflix}, \text{High}) \cdot (R(\text{Low}, \text{Netflix}, \text{High}) + \gamma V_1(\text{High})) \\&\quad T(\text{Low}, \text{Netflix}, \text{Low}) \cdot (R(\text{Low}, \text{Netflix}, \text{Low}) + \gamma V_1(\text{Low})) \\&= 0.0 \times (1 + \gamma 8) + 1.0 \times (1 + \gamma 1) = 1.5\end{aligned}$$

$$\begin{aligned}V^{\pi^1}_3(\text{High}) &= T(\text{High}, \text{exam}, \text{High}) \cdot (R(\text{High}, \text{exam}, \text{High}) + \gamma V_2(\text{High})) \\&\quad T(\text{High}, \text{exam}, \text{Low}) \cdot (R(\text{High}, \text{exam}, \text{Low}) + \gamma V_2(\text{Low})) \\&= 0.9 \times (10 + \gamma 11.65) + 0.1 \times (-10 + \gamma 1.5) = 13.3175\end{aligned}$$

$$\begin{aligned}V^{\pi^1}_3(\text{Low}) &= T(\text{Low}, \text{Netflix}, \text{High}) \cdot (R(\text{Low}, \text{Netflix}, \text{High}) + \gamma V_2(\text{High})) \\&\quad T(\text{Low}, \text{Netflix}, \text{Low}) \cdot (R(\text{Low}, \text{Netflix}, \text{Low}) + \gamma V_2(\text{Low})) \\&= 0.0 \times (1 + \gamma 11.65) + 1.0 \times (1 + \gamma 1.5) = 1.75\end{aligned}$$

Policy improvement:

For state "high":

- Utility of choosing "study": $1 \cdot (0 + 0.5 \cdot 13.3175) + 0 \cdot (0 + 0.5 \cdot 1.75) = 6.66$
- Utility of choosing "exam": $0.9 \cdot (10 + 0.5 \cdot 13.3175) + 0.1 \cdot (-10 + 0.5 \cdot 1.75) = 14.99 - 0.91 = 14.08$
- Utility of choosing "Netflix": $1 \cdot (1 + 0.5 \cdot 13.3175) + 0 \cdot (0 + 0.5 \cdot 1.75) = 7.66$

So set $\hat{\pi}_2(\text{high}) \leftarrow \text{exam}$

For state "low":

- Utility of choosing "study": $0.3 \cdot (0 + 0.5 \cdot 13.3175) + 0.7 \cdot (0 + 0.5 \cdot 1.75) = 1.99 + 0.6125 = 2.6$
- Utility of choosing "exam": $0.05 \cdot (10 + 0.5 \cdot 13.3175) + 0.95 \cdot (-10 + 0.5 \cdot 1.75) = 0.83 - 8.66 = -7.83$
- Utility of choosing "Netflix" $1 \cdot (1 + 0.5 \cdot 1.75) = 1.875$

So set $\hat{\pi}_2(\text{low}) \leftarrow \text{study}$

می بینیم که سیاست تغییر کرد فلذا سیاست ها هنوز همگرا نشده اند.

سوال ۶)

الف) با استفاده از Q-learning و محاسبه سه نمونه زیر داریم:

$$Q(A, \text{Go}) \leftarrow (1 - \alpha)Q(A, \text{Go}) + \alpha(r + \gamma \max Q(B, a')) = 0.5(0) + 0.5(2) = 1$$

$$Q(C, \text{Stop}) \leftarrow (1 - \alpha)Q(C, \text{Stop}) + \alpha(r + \gamma \max Q(A, a')) = 0.5(0) + 0.5(0 + 1) = 0.5$$

$$Q(C, \text{Go}) \leftarrow (1 - \alpha)Q(C, \text{Go}) + \alpha(r + \gamma \max Q(A, a')) = 0.5(0) + 0.5(2 + 1) = 1.5$$

ب) محاسبات را طبق روابط زیر انجام می دهیم.

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$

$$\text{difference} = \left[r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$$

$$w_i \leftarrow w_i + \alpha [\text{difference}] f_i(s, a)$$

پس از اولین به روزرسانی:

$$Q(A, \text{Go}) = w_1 f_1(A, \text{Go}) + w_2 f_2(A, \text{Go}) = 0$$

$$\text{difference} = [r + \max Q(B, a')] - Q(A, \text{Go}) = [4 + 0] - 0 = 4$$

$$w_1 = w_1 + \alpha(\text{difference})f_1 = 0 + (0.5)(4)(1) = 2$$

$$w_2 = w_2 + \alpha(\text{difference})f_2 = 0 + (0.5)(4)(1) = 2$$

پس از دومین به روز رسانی:

$$Q(B, \text{Stop}) = w_1f_1(B, \text{Stop}) + w_2f_2(B, \text{Stop}) = 2(1) + 2(-1) = 0$$

$$Q(A, \text{Go}) = w_1f_1(A, \text{Go}) + w_2f_2(A, \text{Go}) = 2(1) + 2(1) = 4$$

$$\text{difference} = [r + \max Q(A, a)] - Q(B, \text{Stop}) = [0 + 4] - 0 = 4$$

$$w_1 = w_1 + \alpha(\text{difference})f_1 = 2 + (0.5)(4)(1) = 4$$

$$w_2 = w_2 + \alpha(\text{difference})f_2 = 2 + (0.5)(4)(-1) = 0$$