

باسمه تعالی

تکلیف ۳ درس هوش مصنوعی - دکتر جوانمردی

سید امیرمهدی میرشریفی - ۹۸۳۱۱۰۵

سوال ۱ (۲۰ نمره)

صحیح یا غلط بودن موارد زیر را با ذکر دلیل بیان کنید.

- ۱- در فرایند تصمیم گیری مارکوف نتیجه هر عمل به حالت کنونی و قبلی وابسته است.
- ۲- در direct evaluation بعد از انجام هر انتقال، ارزش آن حالت را دوباره محاسبه می کنیم.
- ۳- اگر تنها تفاوت بین دو MDP، مقدار discount factor باشد، این دو قطعاً سیاست بهینه (optimal policy) یکسان دارند.
- ۴- Q-learning یک شیوه یادگیری تقویتی model-free و off-policy است.

پاسخ:

- ۱- غلط - در تصمیم گیری مارکوف نتیجه هر عمل ما صرفاً به حالت کنونی وابسته است.
- ۲- غلط - اگر مقصود جمله این است که بعد از هر نمونه نتیجه را به روز رسانی می کنیم، جواب غلط است در این روش ما در تمام بخش های نمونه مان، دنباله هایی که شروعشان state مدنظرمان هست را به عنوان یک نمونه در نظر می گیریم و ارزش آن را با میانگیری آن ها حساب می کنیم.
- ۳- غلط - در مسئله خوردن دات ها این قضیه نقض خواهد شد.
- ۴- درست

سوال ۲ (۱۰ نمره)

فرض کنید می خواهیم از روش یادگیری Q تخمینی در یک هلیکوپتر کوچک برای خودداری از برخورد با درختان و ساختمان ها استفاده کنیم. ابتدا مشخص کنید از کدام یک از ویژگی های محیط را برای تابع ارزش خطی استفاده می کنید و سپس دو حالتی که بر اساس ویژگی های گفته شده مشابه هستند اما ارزش بسیار متفاوتی دارند را بیان کنید.

پاسخ:

ویژگی ها همچون ، ارتفاع ، فاصله تا درخت ، سرعت ، جهت ، تراکم درختان اطراف و برای ساخت تابع ارزش خطی میتوان از پارامتر های سرعت (عامل خنثی یا مثبت) ، پراکندگی (عامل منفی) و فاصله تا درخت (عامل منفی) استفاده کرد. اگر هلی کوپتر در حالتی باشد که پراکندگی محیط زیاد است اما فاصله اش تا نزدیک ترین درخت متوسط باشد ارزش آن برابر خواهد شد با حالتی که پراکندگی کم است اما فاصله اش تا نزدیک ترین درخت بسیار بسیار کم باشد. به این صورت شرایط متفاوت اما ارزش برآورد شده شان یکی خواهد بود.

سوال ۳ (۱۰ نمره)

الف) اگر در طول اجرای policy iteration، فقط یک iteration از policy evaluation را به جای اجرای آن تا زمان همگرایی انجام دهیم، آیا همچنان به سیاست بهینه می‌رسیم؟ توضیح دهید.

ب) فرض کنید در Q-learning، مقدار $\epsilon=1$ باشد. در این صورت آیا تضمینی وجود دارد که به سیاست بهینه همگرا شود؟ توضیح دهید.

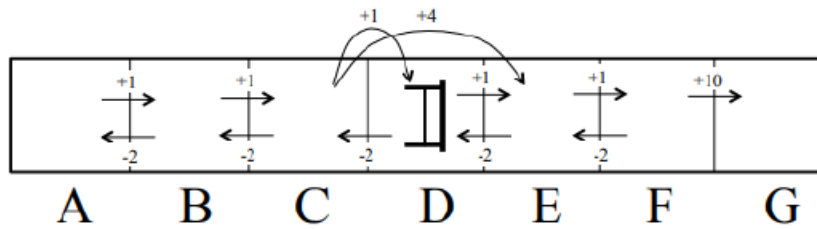
پاسخ:

الف) وقتی یک دور را بخواهیم مدام اجرا کنیم یعنی با یک سیاست یکسان هر دفعه مقادیر حالات را حساب کنیم ، بنابراین بعد از دور اول که مقادیر با سیاست ها همگرا شد سراغ دور دوم میرویم با همان سیاست ها که طبق همگرایی مرحله آخر مجددا همگرا است بنابراین لزوما بهینه نیست.

ب) وقتی اپسیلون برابر یک باشد یعنی همیشه رندوم عمل کند و در نهایت که همگرا میشود ، بهینه هم هست

سوال (۴) (۲۰ نمره)

مدل MDP شکل زیر را در نظر بگیرید. در خانه D یک مانع وجود دارد و عامل ما می تواند به صورت deterministic به چپ (left) یا راست (right) حرکت کند. خانه G نیز حالت ترمینال است. اگر عامل در خانه C باشد، نمی تواند به راست حرکت کند و باید بپرد (jump). پرش ممکن است موفقیت آمیز باشد و به خانه E برسد و $T(C, \text{jump}, E) = 0.5$ و همچنین ممکن است شکست خورده و به مانع خانه D برخورد کند و $T(C, \text{jump}, D) = 0.5$. به سوالات زیر با فرض $\gamma = 1$ پاسخ دهید.



الف) دو مرحله value iteration را انجام دهید و مقادیر زیر را محاسبه کرده و جدول را کامل کنید. Iteration صفر برای تخصیص مقدار اولیه صفر برای همه مقادیر است.

- 1) $V_2(B)$
- 2) $Q_2(B, \text{right})$
- 3) $Q_2(B, \text{left})$

پاسخ:

In v1:

$$V_a = v_b + 1 = 0 + 1 = 1$$

$$V_b = \max((v_c + 1), v_a - 2) = \max(1, -1) = 1$$

$$V_c = \max((v_d + 1)/2 + (v_e + 4)/2, -2 + 1) = \max(5/2, -1) = 5/2$$

$$V_d = \max(v_c - 2, v_e + 1) = \max(0/2, 1) = 1$$

$$V_e = \max((v_f + 1), v_d - 2) = \max(1, -2) = 1$$

$$V_f = \max((v_g + 10), v_e - 2) = \max(10, -1) = 10$$

In v2:

$$V_a = v_b + 1 = 1 + 1 = 2$$

$$V_b = \max((v_c + 1), v_a - 2) = \max(3.5, 0) = 3.5$$

$$V_c = \max((v_d + 1)/2 + (v_e + 4)/2, -2 + 3.5) = \max(7/2, 1.5) = 7/2$$

$$V_d = \max(v_c - 2, v_e + 1) = \max(1.5, 2) = 2$$

$$V_e = \max((v_f + 1), v_d - 2) = \max(11, 0) = 11$$

$$V_f = \max((v_g + 10), v_e - 2) = \max(10, -9) = 10 \Rightarrow V(B) = 3.5 / Q(B, \text{left}) = 0 / Q(B, \text{right}) = 3.5$$

	A	B	C	D	E	F	G
V_1	1	1	5/2	1	1	10	0
V_2	2	3.5	7/2	2	11	10	0

ب) با توجه به اپیزود زیر و چهار انتقال انجام شده، به روز رسانی Q-learning را اعمال کنید و Q-value های به دست آمده را در خانه های جدول پر کنید. خانه هایی که تحت تاثیر قرار نمی گیرند را خالی بگذارید.

Episode											
s	a	r	s	a	r	s	a	r	s	a	r
C	jump	+4	E	right	+1	F	left	-2	E	right	+1

بر فرض آن که آلفا برابر نیم باشد:

	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	۰	۲	۰	۰	۰	۰
Transition 2	۰	۲	۰	/۵	۰	۰
Transition 3	۰	۲	۰	/۵	-۰,۷۵	۰
Transition 4	۰	۲	۰	/۷۵	-۱	۰

سوال ۵) (۲۰ نمره)

فرض کنید یک عامل هوشمند داریم که میتواند درس بخواند! به جهت سادگی، این عامل فقط دارای حالت وضعیت درسی {High, Low} میباشد. در هر کدام از این حالت ها، عامل میتواند یکی از عمل های زیر را انجام دهد:

- درس بخواند
- امتحان بدهد
- به تماشای Netflix بنشیند.

این عامل را می توانیم با MDP زیر مدل کنیم:

S	A	S'	T(S, A, S')	R (S, A, S')
High	study	High	1.0	0
High	study	Low	0.0	0
High	exam	High	0.9	10
High	exam	Low	0.1	-10
High	Netflix	High	1.0	1
High	Netflix	Low	0.0	1
Low	study	High	0.3	0
Low	study	Low	0.7	0
Low	exam	High	0.05	10
Low	exam	Low	0.95	-10
Low	Netflix	High	0.0	1
Low	Netflix	Low	1.0	1

Policy iteration را برای این MDP تا دو iteration اعمال کنید. سیاست اولیه را عمل study برای هر حالت در نظر بگیرید. Utility اولیه هر دو حالت را برابر با صفر در نظر بگیرید. مقدار $\gamma = 0.5$ است. راه حل خود را شرح دهید. آیا در انتها سیاست ها همگرا می شوند؟ توضیح دهید.

پاسخ:

$\Pi_0 = \text{study}$ and $v_1 = 0$ and $v_2 = 0$

$$V_1 = 1 \cdot (0+0) + 0 \cdot (0+0) = 0$$

$$V_2 = .3 \cdot (0+0) + 0 \cdot (0+0) = 0$$

\Rightarrow Now check Q for Each V :

$$Q(V_1, \text{study}) = 0$$

$$Q(V_1, \text{exam}) = .9 \cdot (10+0) + .1 \cdot (-10+0) = 8 \Rightarrow \text{best action} = \text{exam}$$

$$Q(V_1, \text{Netflix}) = 1+0=1$$

$$Q(V_2, \text{study}) = 0$$

$$Q(V_2, \text{exam}) = .05 \cdot (10+0) + .95 \cdot (-10+0) = -9 \Rightarrow \text{best action} = \text{netflix}$$

$$Q(V_2, \text{Netflix}) = 1+0=1$$

$\Rightarrow \Pi_1 = \text{exam, netflix}$

$$V_1 = 8$$

$$V_2 = 1+0=1$$

\Rightarrow Now check Q for Each V :

$$Q(V_1, \text{study}) = 1 \cdot (0 + 8/2) = 4$$

$$Q(V_1, \text{exam}) = .9 \cdot (10 + 8/2) + .1 \cdot (-10 + 1/2) = 12.6 - .95 = 11.65 \Rightarrow \text{best action} = \text{exam}$$

$$Q(V_1, \text{Netflix}) = 1 + 8/2 = 5$$

$$Q(V2, \text{study}) = .3 * (8/2) + .7 * (1/2) = 1.2 + .35 = -1.55$$

$$Q(V2, \text{exam}) = .05 * (10 + 8/2) + .95 * (-10 + 1/2) = .7 - 9.025 = -8.325 \Rightarrow \text{best action} = \text{study}$$

$$Q(V1, \text{Netflix}) = 1 - 1/2 = -.5$$

$\Rightarrow \Pi_2 = \text{for } v_1 = \text{exam} \text{ \& for } v_2 = \text{study}$

برای حل این مسئله همانطور که در بالا نوشته شده است اول با سیاست این که هر دو حالت درس جلو میرویم. ارزش حالت ها را به دست می آوریم. سپس بهترین اکشن را به عنوان سیاست بعدی پیدا میکنیم و همین مراحل را دو بار انجام میدهیم تا در نهایت برای حالت یک سیاست آزمون و حالت دو سیاست درس انتخاب میشود. برای این که بفهمیم همگرا شده است یا نه نیاز به یک مرحله دیگر دارد زیرا تا آخرین مرحله همگرا نبود.

سوال (۶) (۲۰ نمره)

یک مسئله MDP را در نظر بگیرید که در آن سه حالت $[A, B, C]$ و دو اکشن حرکت (Go) و توقف (Stop) وجود دارد. نمونه های جدول زیر از انجام اکشن های مختلف در این MDP تولید شده اند. همچنین فرض کنید $\gamma = 1$ و $\alpha = 0.5$. به سوالات زیر پاسخ دهید.

الف) Q-learning را روی نمونه های زیر اجرا کنید. با فرض مقدار اولیه Q-value صفر، مقادیر Q-value زیر که از Q-learning به دست آمده اند را محاسبه کنید.

s	a	s'	r
A	Go	B	2
C	Stop	A	0
B	Stop	A	-2
B	Go	C	-6
C	Go	A	2
A	Go	A	-2

1) $Q(C, \text{Stop})$

2) $Q(C, \text{Go})$

پاسخ:

$$\text{First } (A, \text{Go}, B, 2) \Rightarrow Q(A, \text{Go}) += .5 * (2 + \max Q(B) - Q(A, \text{Go})) \Rightarrow Q(A, \text{Go}) = 1$$

$$\text{Then } (C, \text{stop}, A, 0) \Rightarrow Q(C, \text{stop}) += .5 * (0 + \max Q(A) - Q(C, \text{stop})) \Rightarrow Q(C, \text{stop}) = 0.5$$

$$\text{Then } (B, \text{stop}, A, -2) \Rightarrow Q(B, \text{stop}) += .5 * (-2 + \max Q(A) - Q(B, \text{stop})) \Rightarrow Q(B, \text{stop}) = -0.5$$

$$\text{Then } (B, Go, C, -6) \Rightarrow Q(B, Go) = 0.5 * (-6 + \max Q(C) - Q(B, Go)) \Rightarrow Q(B, Go) = -0.275$$

$$\text{Then } (C, GO, A, 2) \Rightarrow Q(C, Go) = 0.5 * (2 + \max Q(A) - Q(C, Go)) \Rightarrow Q(C, Go) = 1.5$$

$$\text{Finally } (A, GO, A, -2) \Rightarrow Q(A, Go) = 0.5 * (-2 + \max Q(A) - Q(A, Go)) \Rightarrow Q(A, Go) = 0$$

ب) در این بخش از دو ویژگی داده شده زیر استفاده کنید.

- $f_1(s, a) = 1$
- $f_2(s, a) = \begin{cases} 1 & a = Go \\ -1 & a = Stop \end{cases}$

وزن های w_1 و w_2 را با توجه به نمونه های مشاهده شده پس از اولین به روزرسانی (با استفاده از نمونه اول) و دومین

بروزرسانی (با استفاده از نمونه دوم) بنویسید. وزن های اولیه صفر هستند.

s	a	s'	r
A	Go	B	4
B	Stop	A	0

پاسخ:

با توجه به آن که آلفا بیان نشده است مقدار آن را نیم در نظر میگیریم:

$$Q(S, A) = w_1 * f_1(s, a) + w_2 * f_2(s, a)$$

$$W_1 = W_2 = 0$$

$$\text{First: } Q(A, Go) = 0 + 0 = 0$$

$$R + \max Q(B) - Q(A, Go) = 4 + 0 - 0 = 0 \rightarrow w_1 = w_1 + \alpha * \text{difference} * f_1 = 0 + .5 * 4 * 1 = 2$$

$$W_2 = w_2 + \alpha * \text{difference} * f_2 = 0 + .5 * 1 = 0.5$$

$$W_1 = W_2 = 2$$

Then:

$$Q(B, Stop) = 2 * 1 + 2 * -1 = 0$$

$$\text{Difference} = 0 + \max Q(A) - Q(B, stop) = 0 + 3 - 0 = 3 \Rightarrow w_1 = 2 + 0.5 * 3 * 1 = 3.5$$

$$W_2 = 2 + .5 * 3 * (-1) = 0.5$$

$$W_1 = 3.5 \quad W_2 = 0.5$$