

به نام خدا



دانشگاه صنعتی شریف
دانشکده مهندسی کامپیوتر

هوش مصنوعی

تمرین ۶: Markov Decision Process & RL

مهراد میلانلو

۹۹۱۰۵۷۷۵

سوال ۱

(آ) درست

اگر ضریب تخفیف به مقدار کافی کوچک باشد، ایجنت ترجیح می‌دهد پاداش ابتدایی را جمع کند و حریصانه و کوتاه‌نظر عمل می‌کند. زیرا پاداش‌های بعدی علی‌رغم بزرگ‌تر بودن، وقتی در γ ضرب شود، می‌تواند کوچک و کوچک‌تر و نزدیک به صفر باشد.

(ب) درست

هنگامی که پاداش منفی زندگی به‌اندازه‌ی کافی بزرگ باشد، ایجنت حریصانه عمل می‌کند و ترجیح می‌دهد یا زودتر به نقطه‌ی پایان برسد. زیرا با هر حرکتی که انجام می‌دهد، مقدار زیادی reward منفی دریافت می‌کند.

(ج) نادرست

ضریب منفی هنگامی که به‌توان‌های زوج برسد، مثبت و در غیر این صورت منفی است که نشان می‌دهد نمی‌توانیم پاداش منفی زندگی را با آن مدل کنیم. زیرا یک‌درمیان مثبت و منفی می‌شود.

(د) نادرست

دقیقا مشابه استدلال قبل، فرض کنید ضریب تخفیف منفی باشد. بنابراین در توان‌های زوج مثبت و در توان‌های فرد، منفی است. چون یکی در میان منفی و مثبت می‌شود، به‌وضوح نمی‌توان آن را با پاداش منفی زندگی مدل کنیم.

سوال ۲

(آ) تابع ارزش حالت‌ها را این‌گونه در نظر می‌گیریم:

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

در این رابطه، $T(s, \pi(s), s')$ احتمال آن است که در صورت انتخاب اعمال سیاست $\pi(s)$ از حالت s به حالت s' برسیم. (به دلیل stochastic بودن فرآیند). $R(s, \pi(s), s')$ مقدار reward حاصل از اعمال $\pi(s)$ برای رسیدن از s به s' است. همین‌طور γ مقدار ضریب تخفیف می‌باشد. بنابراین واضح است که در صورت انتخاب سیاست بهینه‌ی π ، $V^\pi(s)$ ماکسیمم مقدار utility با شروع از حالت s را نشان می‌دهد.

(ب) رابطه‌ای که در قسمت قبل نوشتیم، همان رابطه‌ی بلمن برای تابع ارزش حالت‌هاست که روی سیاست‌های مختلف عمل می‌کند. برای آن که بتوانیم ارزش حالت‌ها را آپدیت کنیم، از روش Time-Limited استفاده می‌کنیم. در واقع داریم:

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

(ج) برای حل این قسمت، ابتدا از روی نمودار فرآیند به‌دست می‌آوریم:

$$T(A, ab, B) = 1, T(B, ba, A) = 1, T(B, bc, C) = 1, T(C, ca, A) = 0.25, T(C, ca, C) = 0.75, T(C, cb, B) = 1$$

$$R(A, ab, B) = -8, R(B, bc, C) = -2, R(C, ca, C) = 4, R(B, ba, A) = 2, R(C, ca, A) = 4, R(C, cb, B) = 8$$

$$\gamma = 0.5$$

بنابراین طبق رابطه‌ی قسمت قبل خواهیم داشت:

$$V_1^{\pi_1}(A) = \sum_{s'} T(A, \pi_1(A), s') [R(A, \pi_1(A), s') + \gamma V_1^{\pi_1}(s')] = T(A, ab, B) [R(A, ab, B) + \gamma V_1^{\pi_1}(B)]$$

$$\begin{aligned}
&= ۱ \times (-\lambda + \bullet/\delta \times ۲) \Rightarrow \boxed{V_{\gamma}^{\pi_1}(A) = -\gamma} \\
V_{\gamma}^{\pi_1}(B) &= \sum_{s'} T(B, \pi_1(B), s') [R(B, \pi_1(B), s') + \gamma V_{\gamma}^{\pi_1}(s')] \\
&= \frac{1}{\gamma} T(B, ba, A) [R(B, ba, A) + \gamma V_{\gamma}^{\pi_1}(A)] + \frac{1}{\gamma} T(B, bc, C) [R(B, bc, C) + \gamma V_{\gamma}^{\pi_1}(C)] \\
&= \frac{1}{\gamma} \times ۱ \times (۲ + \bullet/\delta \times ۲) + \frac{1}{\gamma} \times ۱ \times (-۲ + \bullet/\delta \times ۲) \Rightarrow \boxed{V_{\gamma}^{\pi_1}(B) = ۱} \\
V_{\gamma}^{\pi_1}(C) &= \sum_{s'} T(C, \pi_1(C), s') [R(C, \pi_1(C), s') + \gamma V_{\gamma}^{\pi_1}(s')] \\
&= \frac{1}{\gamma} \sum_{s'} T(C, ca, s') [R(C, ca, s') + \gamma V_{\gamma}^{\pi_1}(s')] + \frac{1}{\gamma} T(C, cb, B) [R(C, cb, B) + \gamma V_{\gamma}^{\pi_1}(B)] \\
&= \frac{1}{\gamma} (\frac{1}{\gamma} \times (۴ + \bullet/\delta \times ۲) + \frac{۳}{\gamma} \times (۴ + \bullet/\delta \times ۲)) + \frac{1}{\gamma} \times ۱ \times (\lambda + \bullet/\delta \times ۲) \Rightarrow \boxed{V_{\gamma}^{\pi_1}(C) = \gamma}
\end{aligned}$$

د) اکنون با توجه به تابع ارزش‌گذاری جدید، با استفاده از مرحله‌ی extraction policy از رابطه‌ی زیر استفاده می‌کنیم:

$$\begin{aligned}
\pi_{k+1}(s) &\leftarrow \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_k}(s')] \\
\pi_{\gamma}(A) &= \arg \max_a \sum_{s'} T(A, \pi_1(A), s') [R(A, \pi_1(A), s') + \gamma V_{\gamma}(s')] \\
&= \arg \max_a \{ ۱ \times (-\lambda + \bullet/\delta \times ۱) \} \Rightarrow \boxed{\pi_{\gamma}(A) = ab} \\
\pi_{\gamma}(B) &= \arg \max_a \sum_{s'} T(B, \pi_1(B), s') [R(B, \pi_1(B), s') + \gamma V_{\gamma}(s')] \\
&= \arg \max_a \{ ۱ \times (۲ + \bullet/\delta \times \gamma), ۱ \times (-۲ + \bullet/\delta \times -\gamma) \} \Rightarrow \boxed{\pi_{\gamma}(B) = bc} \\
\pi_{\gamma}(C) &= \arg \max_a \sum_{s'} T(C, \pi_1(C), s') [R(C, \pi_1(C), s') + \gamma V_{\gamma}(s')] \\
&= \arg \max_a \{ \frac{1}{\gamma} \times (۴ + \bullet/\delta \times (-\gamma)) + \frac{۳}{\gamma} \times (۴ + \bullet/\delta \times \gamma), ۱ \times (\lambda + \bullet/\delta \times ۱) \} \Rightarrow \boxed{\pi_{\gamma}(C) = cb}
\end{aligned}$$

۶

$$\begin{aligned}
V^{\pi}(s) &= \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V^{\pi}(s')] \\
\pi'(s) &= \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi}(s')] \\
\Rightarrow V^{\pi'}(s) &= \sum_{s'} T(s, \pi'(s), s') [R(s, \pi'(s), s') + \gamma V^{\pi'}(s')] \\
&= \sum_{s'} T(s, \pi'(s), s') [\gamma V^{\pi'}(s') - \gamma V^{\pi}(s')] + \sum_{s'} T(s, \pi'(s), s') [R(s, \pi'(s), s') + \gamma V^{\pi}(s')]
\end{aligned}$$

تهش به نتیجه می‌رسه احتمالا: به‌طور شهودی، سیاست π' به‌صورت حریصانه از π به‌دست می‌آید. یعنی از بین اکشن‌ها ماکسیمم می‌گیریم. واضح است که مقدار جدید تابع ارزش بیشتر یا مساوی است. در حالت تساوی نیز نشان‌دهنده‌ی این است که π' همان π بهینه است.

سوال ۳

(آ)

رابطه‌ی بلمن برای به‌روزرسانی Q-Value این‌گونه است:

$$Q(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q(s', a')]$$

اکنون با توجه به این‌که ضریب تخفیف $\gamma = 0.9$ است، مقدار به‌روزشده‌ی $Q(3, Left)$ را محاسبه می‌کنیم:

$$\begin{aligned} Q(3, Left) &= T(3, Left, 2) [R(3, Left, 2) + 0.9 \max_{a'} Q(2, a')] = 1 \times [-1 + 0.9 \times \max\{3, 6, 8\}] = -1 + 7/2 \\ &\Rightarrow Q(3, Left) = -6/2 \end{aligned}$$

(ب) در روش حریصانه در هر استیتی که باشیم، تنها حالت‌های همسایه که بیشترین مقدار reward را داشته باشند در نظر می‌گیریم در صورتی که ممکن است راهی وجود داشته باشد که در ابتدا پاداش کم و در ادامه پاداش بسیار بیشتری بدهد. از سوی دیگر عدم توجه به پاداش‌های همسایه‌ها و حرکات تصادفی بار محاسباتی را به‌شدت انجام می‌دهد. بنابراین ایجاد تعادل بین exploitation و exploration باعث سیاست‌های بهتری می‌شود.

(ج) می‌دانیم برای آپدیت کردن π با استفاده از V^π نیاز به دانستن T و R داریم که هنگام یادگیری مسئله ممکن است به مقادیر آن‌ها دسترسی نداشته باشیم. اما به‌روزرسانی با استفاده از $Q(s, a)$ تنها نیاز به استیت و اکشن نیاز دارد و به همین دلیل از Q-Value ها استفاده می‌کنیم.

(د)

در سیاست تصادفی softmax با استفاده از احتمالات محاسبه‌شده که در جدول زیر آورده شده‌است، اکشن مورد نظر را انتخاب می‌کنیم. همان‌طور که گفتیم، ایراد اصلی حریصانه عمل کردن، تصمیمات کوتاه‌نظرانه و عدم توجه به reward های مراحل جلوتر است. که در این روش این اتفاق نمی‌افتد و در محاسبه‌ی احتمالات، همسایه‌های دیگر نیز در نظر گرفته می‌شوند. بنابراین هم سنگینی محاسباتی روش رندوم را ندارد و سریع‌تر به هدف می‌رسد و هم مشکل حریصانه عمل کردن را ندارد.

$\pi(s, a) = \frac{e^{Q(s, a)}}{\sum_b e^{Q(s, b)}}$			
$\pi(1, U) = \frac{e^3}{e^3 + e^6} \approx 0.33$	$\pi(1, R) = \frac{e^7}{e^7 + e^8} \approx 0.47$	-	-
$\pi(2, U) = \frac{e^6}{e^6 + e^7 + e^8} \approx 0.118$	$\pi(2, R) = \frac{e^8}{e^6 + e^7 + e^8} \approx 0.876$	-	$\pi(2, L) = \frac{e^7}{e^6 + e^7 + e^8} \approx 0.006$
$\pi(3, U) = \frac{e^8}{e^7 + e^8} \approx 0.88$	-	-	$\pi(3, L) = \frac{e^7}{e^7 + e^8} \approx 0.12$
-	$\pi(4, R) = \frac{e^5}{e^5 + e^7} \approx 0.95$	$\pi(4, D) = \frac{e^7}{e^5 + e^7} \approx 0.05$	-
-	$\pi(5, R) = \frac{e^8}{e^8 + e^5 + e^6} \approx 0.84$	$\pi(5, D) = \frac{e^6}{e^8 + e^5 + e^6} \approx 0.12$	$\pi(5, L) = \frac{e^5}{e^8 + e^5 + e^6} \approx 0.04$

(ه)

با استفاده از رابطه‌ی داده‌شده، برای به‌روز رسانی مقادیر مورد نظر داریم: ($\gamma = 0.8$, $\alpha = 0.2$)

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R_{ss'}^a + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

$$Q(2, U) = Q(2, U) + 0.2 [R(2, U, 5) + 0.8 \max_{a'} Q(5, a') - Q(2, U)] = 6 + 0.2 [-1 + 0.8 \times 8 - 6]$$

$$\Rightarrow Q(2, U) = 5.88$$

$$Q(5, R) = Q(5, R) + 0.2 [R(5, R, 6) + 0.8 \max_{a'} Q(6, a') - Q(5, R)] = 8 + 0.2 [10 + 0.8 \times 0 - 8]$$

$$\Rightarrow Q(5, R) = 8.4$$

سوال ۴

(آ)

با استفاده از همان رابطه‌ای که در سوال قبل نوشتیم، مقدار $Q(s, a)$ را بعد از مشاهده‌ی هر نمونه به‌روزرسانی می‌کنیم: $\alpha = 0.1, \gamma = 0$

(۰/۹)

$$Q(s, a) \leftarrow Q(s, a) + \alpha [R_{ss'}^a + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

a:

$$Q(A, 1) \leftarrow Q(A, 1) + \alpha [R(A, 1, B) + \gamma \max_{a'} Q(B, a') - Q(A, 1)] = 0 + 0.1 [-3 + 0.9 \times 0 - 0] = -0.3$$

b:

$$Q(B, 1) \leftarrow Q(B, 1) + \alpha [R(B, 1, A) + \gamma \max_{a'} Q(A, a') - Q(B, 1)] = 0 + 0.1 [4 + 0.9 \times 0 - 0] = 0.4$$

c:

$$Q(A, 2) \leftarrow Q(A, 2) + \alpha [R(A, 2, A) + \gamma \max_{a'} Q(A, a') - Q(A, 2)] = 0 + 0.1 [-4 + 0.9 \times 0 - 0] = -0.4$$

d:

$$Q(A, 1) \leftarrow Q(A, 1) + \alpha [R(A, 1, B) + \gamma \max_{a'} Q(B, a') - Q(A, 1)] = -0.3 + 0.1 [-3 + 0.9 \times 0.4 - 0.3] = -0.534$$

e:

$$Q(A, 2) \leftarrow Q(A, 2) + \alpha [R(A, 2, T) + \gamma \max_{a'} Q(T, a') - Q(A, 2)] = -0.4 + 0.1 [1 + 0.9 \times 0 + 0.4] = -0.26$$

در نهایت مقادیر $Q(s, a)$ برابر است با:

$$Q(A, 1) = -0.534, Q(A, 2) = -0.26, Q(B, 1) = 0.4$$

(ب)

سیاست زیر را در نظر می‌گیریم:

$$\pi(s) \leftarrow \arg \max_a Q(s, a)$$

که نتیجه می‌دهد:

$$\pi(A) = 2, \pi(B) = 1$$

این سیاست با توجه به مقادیر Q-Value به‌دست آمده، از سیاست تصادفی بهتر است و اکشن‌های بهتری را انتخاب می‌کند.

(ج)

استفاده از سیاست π_{random} به وضوح امکان exploration بیشتری می‌دهد و احتمالا به جواب بهینه برای Q-Value ها خواهیم رسید. هرچند که از لحاظ زمانی ممکن است بیشتر طول بکشد. از طرف دیگر چون π^* از نمونه‌ها به‌دست آمده، آگه نمونه‌ها به مقدار میانگین نزدیک باشند و تعدادشان افزایش پیدا کند تغییر خاصی نداشته باشند، طبق محاسبات سیاست بهتری از رندوم است. اگرچه به دلیل حریصانه بودن، مانع exploration می‌شویم و ممکن است به جواب کاملاً بهینه نرسیم.

سوال ۵

برای آن‌که بتوانیم مقدار نرخ اکتشاف را در طول زمان کاهش بدهیم، می‌توانیم آن را تابعی از زمان تعریف کنیم. مانند توابعی که در روش Simulated Annealing تعریف می‌کردیم ($\epsilon = e^{\frac{\Delta E}{T}}$) که به‌وضوح با گذر زمان مقدار نرخ اکتشاف کاهش می‌یابد. همین‌طور می‌توان از الگوریتم‌هایی مانند softmax که در سوال‌های قبلی نیز داشتیم، استفاده کنیم. در گذر زمان، اختلاف Q-Value استیت‌هایی که پاداش بیشتری می‌دهند، با دیگر استیت‌ها بیشتر می‌شود و چون موقع به‌روز رسانی، مقدار جدید برابر است با نسبت $Q(s, a)$ به دیگر Q-Value ها، بنابراین احتمال حرکت‌هایی که پاداش بیشتری دارند بیشتر می‌شود.

برای حل مشکل دیگر، دقت کنید با تغییر استراتژی حریف، محیط یعنی استیت‌ها و پاداش‌ها نیز تغییر می‌کنند. در این صورت می‌توان از رویکرد featured-based استفاده کرد. به این صورت که توابع ارزش‌گذاری را بر اساس تعدادی feature تعریف می‌کنیم:

$$V(s) = \sum_i w_i f_i(s), Q(s, a) = \sum_i w_i f_i(s, a)$$

برای به دست آوردن وزن‌ها نیز، می‌توانیم از روش‌های لرنینگ مانند Gradient-Descent استفاده کنیم. به این صورت با تغییر وزن‌ها می‌توان عملکرد مناسبی در برابر تغییر استراتژی حریف داشت.