



هوش مصنوعی

زمستان ۱۴۰۱

استاد: محمدحسین رهبان

گردآورندگان: حسین گلی، محمد مشتاقی و علی ثالثی

مهلت ارسال: ۲۹ فروردین

فرآیندهای مارکوف و یادگیری تقویتی

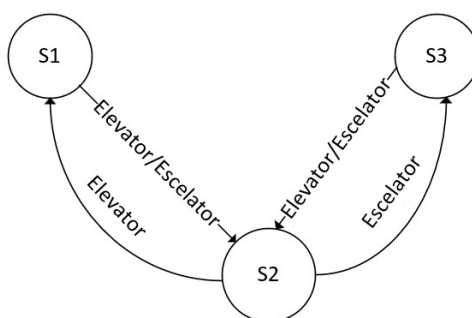
تمرین سوم

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- در طول ترم امکان ارسال با تاخیر پاسخ همه‌ی تمارین تا سقف ۷ روز و در مجموع ۱۵ روز، وجود دارد. پس از گذشت این مدت، پاسخ‌های ارسال شده پذیرفته نخواهند بود. همچنین، به ازای هر روز تأخیر غیر مجاز ۱۲ درصد از نمره تمرین به صورت ساعتی کسر خواهد شد.
- همکاری و هم‌فکری شما در انجام تمرین مانعی ندارد اما پاسخ‌های ارسال شده هر کس حتماً باید توسط خود او نوشته شده باشد.
- در صورت هم‌فکری و یا استفاده از هر منابع خارج درسی، نام هم‌فکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.
- لطفاً تصویری واضح از پاسخ سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.

سوالات نظری (۱۵۰ نمره)

۱. (۲۰ نمره) موارد صحیح و غلط را مشخص کنید و در صورت غلط بودن استدلال کنید.
 - الف) در مسائل یادگیری تقویتی عامل لزوماً باید با محیط به صورت برخط تعامل کند و با استفاده از پاداش‌هایی که به صورت برخط می‌گیرد سیاست خود را بهبود دهد.
 - ب) سیاست بهینه در MDP یکتا است.
 - ج) الگوریتم REINFORCE یک الگوریتم Model Based است.
 - د) در یک Infinite Horizon MDP به ازای تمامی سیاست و حالات داریم: $V^*(s) \geq V^\pi(s)$
۲. (۳۰ نمره) سو که در یک فروشگاه زنجیره‌ای گم شده است و با ۲ تصمیم مواجه است. به سمت بلوک شرقی یا بلوک غربی فروشگاه برود. برای فهم بهتر مسئله به MDP زیر توجه کنید که حالت‌ها و اکشن‌های مسئله را بهتر توصیف می‌کند. S_2 مکان فعلی او در فروشگاه است. S_1, S_3 به ترتیب حالات رفتن به بلوک شرقی و غربی هستند. همچنین او اکشن‌های پله برقی و آسانسور را در اختیار دارد.

شکل ۱: MDP



در واقع ما فرض کرده‌ایم سو مدت طولانی در فروشگاه سرگردان خواهد بود و مسئله را عملاً می‌توان با یک Infinite Horizon MDP مدل‌سازی کرد.

(آ) فرض کنید جدول Q-values به شما داده شده است. مطابق جدول زیر سیاست بهینه را بدست آورید.

Q-values	Elevator	Escalator
S_1	$+1/5$	$-0/5$
S_2	$+2/0$	$-0/3$
S_3	$+0/8$	$+0/9$

(ب) یک مدل سازی ساده می تواند MDP دو حالتی که سو در حالت S_2 است یا خیر باشد. برای این MDP یک state-diagram رسم کنید. فرض کنید با اجرا الگوریتم Q-learning به مقادیر زیر برای Q-values رسیده ایم.

Q^*	Elevator	Escalator
S_2	$+0/08$	$+2/08$
S_4	$+1/18$	$+1/38$

(ج) سیاست بهینه را بدست آورید.

(د) به چه دلیل سیاست بهینه در این مدل سازی با مدل سازی قبلی تفاوت کرده است؟ راهنمایی: از فرض مارکو کمک بگیرید.

۳. (۵۰ نمره) در بسیاری از کاربردهای یادگیری تقویتی ما نمی توانیم داده با استفاده از سیاست های مختلف جمع آوری کنیم (مانند سیستم های درمانی).

به همین دلیل می خواهیم در این سوال روش های off policy evaluation را باهم بررسی کنیم. فرض کنید در مسئله ما سیاست احتمالاتی است.

$$\pi(s, a) = P(a|s, \pi)$$

تابع هدف ما $R(s, a)$ با استفاده از یک سیاست π_1 به جای π است.

$$\begin{aligned} \mathbb{E}_{s \sim p(s), a \sim \pi_1(s, a)} R(s, a) &= \sum R(s, a) p(s) p(a|s) \\ &= \sum R(s, a) p(s) \pi_1(a|s) \end{aligned}$$

مشکل ما اینجا است که با روش هایی که در کلاس بررسی کردیم (روش های نمونه گیری) می توانیم R را برای سیاست π به دست آوریم. می خواهیم روش های که به ما این قابلیت را می دهد با استفاده از داده های جمع آوری شده از سیاست π تابع R را برای π_1 تخمین بزنند استفاده کنیم.

(آ) تخمین گر زیر را در نظر بگیرید:

$$\mathbb{E}_{s \sim p(s), a \sim \pi(s, a)} \frac{\pi_1(s, a)}{\hat{\pi}(s, a)} R(s, a)$$

نشان دهید اگر $\pi = \hat{\pi}$ باشد، این تخمین گر سازگار است.

(ب) مورد بالا را برای تخمین گر زیر نشان دهید.

$$\frac{\mathbb{E}_{s \sim p(s), a \sim \pi(s, a)} \frac{\pi_1(s, a)}{\hat{\pi}(s, a)} R(s, a)}{\mathbb{E}_{s \sim p(s), a \sim \pi(s, a)} \frac{\pi_1(s, a)}{\hat{\pi}(s, a)}}$$

(ج) مثالی نقضی برای این که نشان دهد تخمین گر بالا می تواند بایاس باشد ارائه کنید.

(د) اختیاری: درباره تخمین‌گر Doubly Robust تحقیق کنید و برای چه از این تخمین‌گر به جای تخمین‌گرهایی که بررسی کردیم استفاده می‌شود.

۴. (۵۰ نمره) حسین که به تازگی به ماشین زمان دست یافته است به چین قدیم سفر کرده و یک فرمانده محلی شده است. تابع پاداش به‌ازای هر حالت را در جدول زیر برای فرمانروایی حسین مشخص کرده‌ایم. (هر حالت مکان فرمانروایی اوست)

State	Reward
Mountain	$1/0$
Riverside	$+2/0$
Desert	$-1/0$

همچنین در جدول زیر transition function توسط مشاور کاردرستش به او داده شده است.

Start State	Action	Probability	State End
Mountain	Peace	$0/5$	Riverside
Mountain	Peace	$0/5$	Desert
Mountain	War	$0/1$	Mountain
Mountain	War	$0/2$	Desert
Mountain	War	$0/7$	Riverside
Riverside	Peace	$1/0$	Riverside
Riverside	War	$0/2$	Riverside
Riverside	War	$0/8$	Mountain
Desert	Peace	$1/0$	Desert
Desert	War	$1/0$	Mountain

(آ) حسین تصمیم گرفته است که صلح بهترین رویکرد برای موفقیت است. از این رو سیاست او به این صورت است که هیچ‌گاه به جنگ نمی‌رود. تابع ارزش این سیاست را به‌دست آورید. (بر حسب γ)

(ب) حسین پس از شکست‌های پی‌درپی متوجه شده است که باید سیاستش را عوض کند. با استفاده از الگوریتم Policy iteration یک گام سیاست حسین را بهبود دهید. (می‌توانید فرض کنید $\gamma = 0/9$ است)

(ج) مشاور حسین خائن بوده است. او دیگر از transition function هایی که او داده است نمی‌تواند استفاده کند و مجبور است از الگوریتم‌های RL استفاده کند. با توجه به عملکرد اخیر حسین Q Table زیر را پر کنید. (فرض کنید $\gamma = 1, \alpha = 0/5$)

Start State	Action	State End	Reward
Desert	War	Desert	$-2/0$
Desert	War	Riverside	$3/0$
Riverside	Peace	Mountain	$1/0$
Mountain	Peace	Riverside	$1/0$

Mountain-peace	Riverside-peace	Desert-war
.	.	.

۵. (اختیاری) پاسخ به این سوال اختیاری است. حل آن در این **لینک** قرار دارد. پیشنهاد می‌کنیم آن را مشاهده کنید.

در این تمرین ما با استفاده از cauchy sequence می‌خواهیم اثبات کنیم value iteration به یک جواب یکتا مستقل از نقطه شروع همگرا می‌شود.
Bellman Backup Operator مطابق زیر تعریف می‌شود.

$$V_{k+1} = BV_k = \max(R(s, a) + \gamma \sum p(s'|s, a) V_k^\pi(s'))$$

همچنین می‌دانیم که این اپراتور یک Contraction mapping است.

$$\|BV' - BV''\|_\infty \leq \gamma \|V' - V''\|_\infty$$

$$\|V_{n+1} - V_n\|_\infty \leq \gamma^n \|V_1 - V_0\|_\infty \quad (\text{آ})$$

$$\|V_{n+c} - V_n\|_\infty \leq \frac{\gamma^n}{1-\gamma} \|V_1 - V_0\|_\infty \quad c > 0 \quad (\text{ب})$$

(ج) cauchy sequence یک سری است که عناصر آن به با جلو رفتن سری به یک دیگر نزدیک می‌شوند. درواقع به ازای هر $\epsilon > 0$ یک k وجود دارد به طوری که اگر $m, n > k$ آنگاه $d(a_m, a_n) < \epsilon$ و در نتیجه سری همگرا است. با استفاده از توصیف گفته‌شده نشان دهید سری V_1, V_2, V_3, \dots همگرا است.
(د) نشان دهید که جواب Value iteration یکتا است.

سوالات عملی (۲۰۰ نمره)

۱. (۱۰۰ نمره) برای حل این سوال به نوت‌بوک Q_Tabular.ipynb مراجعه کنید.
۲. (۱۰۰ نمره) برای حل این سوال به نوت‌بوک RL_Chat.ipynb مراجعه کنید. پیشنهاد می‌کنیم از google colab استفاده کنید.