



هوش مصنوعی

بهار ۱۴۰۴

استاد: احسان تن قطاری

دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

طراحان: امیدرضا معصومی، مبین باقری، محمد شفیق زاده، دنیا جعفری، سپهر ذوالفقاری، علیرضا ملک حسینی

مهلت ارسال: ۳۱ مرداد

یادگیری تقویتی

تمرین پنجم

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- در طول ترم امکان ارسال با تاخیر پاسخ همه‌ی تمارین سقف ۴ روز و در مجموع ۱۰ روز، وجود دارد. پس از گذشت این مدت، پاسخ‌های ارسال شده پذیرفته نخواهند بود. همچنین، به ازای هر ساعت تأخیر غیر مجاز نیم درصد از نمره‌ی تمرین کم 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.
- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- در طول ترم امکان ارسال با تاخیر پاسخ همه‌ی تمارین سقف ۴ روز و در مجموع ۱۰ روز، وجود دارد. پس از گذشت این مدت، پاسخ‌های ارسال شده پذیرفته نخواهند بود. همچنین، به ازای هر ساعت تأخیر غیر مجاز نیم درصد از نمره‌ی تمرین کم 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 (۱۰۰ نمره)

۱. (نمره) درستی یا نادرستی عبارت‌های زیر را با ذکر دلیل مشخص کنید.
 - (آ) اگر ضریب تخفیف γ شرط $0 < \gamma < 1$ را ارضا کند، می‌توان تضمین کرد که روش value iteration همگرا می‌شود
 - (ب) اگر برتر بودن یک سیاست را به معنی این در نظر بگیریم که موجب به دست آوردن reward بیشتری می‌شود، سیاست‌هایی که توسط روش value iteration به دست می‌آیند، از سیاست‌هایی که توسط روش policy iteration به دست می‌آیند، برتر هستند.
 - (ج) Q-learning می‌تواند تابع بهینه Q^* را بدون اینکه حتی یک بار سیاست بهینه را اجرا کند بیاموزد.
 - (د) اگر یک MDP مدل انتقال T ای داشته باشد که برای همه‌ی سه‌تایی‌های $T(s, a, s')$ احتمال غیرصفر اختصاص دهد، آنگاه Q-learning شکست خواهد خورد.
 - (ه) فرض کنید عامل ۱ تابع مطلوبیت U_1 و عامل ۲ تابع مطلوبیت U_2 را دارد. اگر

$$U_1 = k_1 U_2 + k_2$$

که در آن $k_1 > 0$ و $k_2 > 0$ باشد، آنگاه عامل ۱ و عامل ۲ ترجیحات یکسانی دارند.

(و) به این سوال پاسخ کوتاه دهید. یک تغییر متداول در یادگیری Q ، گنجاندن پاداش‌ها از گام‌های زمانی بیشتر در عبارت X است. بنابراین، عبارت معمول ما

$$r_t + \gamma \cdot \max_{a_{t+1}} q(s_{t+1}, a_{t+1}; w)$$

به این شکل تبدیل می‌شود:

$$r_t + \gamma \cdot r_{t+1} + \gamma^2 \cdot \max_{a_{t+2}} q(s_{t+2}, a_{t+2}; w)$$

به نظر شما مزایای استفاده از پاداش‌های بیشتر در این تخمین چیست؟

حل.

(۱) درست - برای $0 < \gamma < 1$ ، عملگر بلمن بهینگی یک انقباض با ضریب γ در نورم ℓ_∞ است؛ لذا طبق قضیه Value Iteration، به V همگرا می‌شود.

(۲) نادرست - اجرای کامل هر دو روش به π منجر می‌شود؛ بنابراین سیاست‌های خروجی یکی ذاتاً «برتر» از دیگری نیستند. Policy Iteration به صورت یکنواخت سیاست را بهبود می‌دهد، ولی Value Iteration چنین تضمینی برای سیاست‌های میانی ندارد.

(۳) درست - Q-learning یک روش off-policy است و بدون اجرای سیاست بهینه، با قاعده

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

و با کاوش کافی و نرخ‌های یادگیری مناسب، به Q همگرا می‌شود.

(۴) نادرست - نامنفی و غیرصفر بودن همه $T(s' | s, a)$ باعث شکست نمی‌شود؛ Q-learning به T نیاز ندارد و در MDPهای کاملاً تصادفی نیز (تحت شرایط استاندارد) همگرا است.

(۵) درست - دو تابع مطلوبیت زمانی ترجیحات یکسانی را نمایش می‌دهند که یکی از آن‌ها تبدیل مثبتی از دیگری باشد، یعنی:

$$U_1 = k_1 U_2 + k_2 \quad k_1 > 0 \text{ که در آن}$$

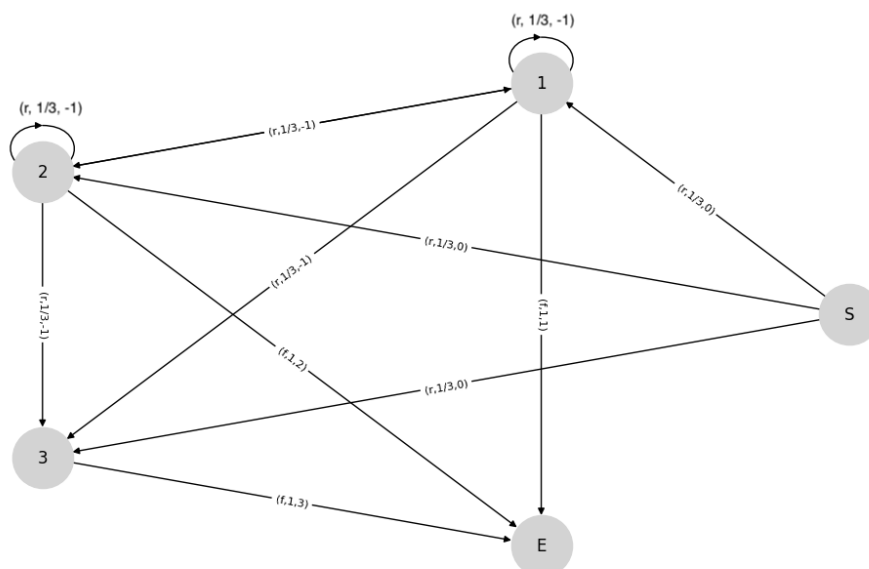
از آنجا که در صورت سوال $k_1 > 0$ و $k_2 \in \mathbb{R}$ داده شده‌اند، تابع U_1 تبدیل آفین مثبتی از U_2 است. بنابراین، عامل ۱ و عامل ۲ ترجیحات یکسانی دارند.

(۶) گنجاندن پاداش‌ها از چندین گام زمانی، تخمین «واقع‌بینانه‌تری» از پاداش کل واقعی ارائه می‌دهد، زیرا درصد بیشتری از آن از تجربه واقعی ناشی می‌شود. این روش می‌تواند به پایدارسازی فرآیند آموزش کمک کند، در حالی که همچنان امکان آموزش در هر گام زمانی (بوت‌استرپینگ) را فراهم می‌کند. این نوع روش، «یادگیری تفاضل زمانی» گامی-N نامیده می‌شود.

۲. (نمره) یک بازی جدید در نظر بگیرید که هربار یک عدد رندوم با توزیع احتمال یکنواخت بین ۱ تا ۳ تولید می‌شود. شما هربار می‌توانید به مقدار عدد تولید شده امتیاز بگیرید و بازی تمام شود، یا امتیاز ۱- بگیرید و ادامه دهید. $\gamma = 0.9$ در نظر بگیرید.

(آ) برای این بازی یک MDP ارائه دهید و اجزای مختلف آن را مشخص کنید.

(ب) سیاست بهینه‌ای که پس از ۳ مرحله اجرای value iteration بدست می‌آید را حساب کنید.



Iteration	State	Action		Updated Value(max)
		r	f	
۱	۱	- ۱	۱	۱
	۲	- ۱	۲	۲
	۳	- ۱	۳	۳
۲	۱	۰.۸	۱	۱
	۲	۰.۸	۲	۲
	۳	۰.۸	۳	۳
۳	۱	۰.۸	۱	۱
	۲	۰.۸	۲	۲
	۳	۰.۸	۳	۳

(ج) با داشتن سیاست اولیه پایان بازی برای اعداد ۱ و ۲ و ادامه برای عدد ۳، با اجرای ۳ مرحله از policy iteration سیاست بهینه را پیدا کنید.

حل.

(آ) محیط از S شروع می شود و با احتمال یکسان به هر سه حالت عدد می رود. در هر عدد، می توان اکشن f را انجام داد که بازی با همان مقدار پاداش تمام می شود. یا می توان اکشن r را انجام داد که با احتمال یکسان می تواند به هر حالت عدد برود.

(ب) با توجه به رابطه، ارزش هر اکشن را در جدول یادداشت می کنیم.

$$v_{k+1}(s) \leftarrow \max_{a \in \mathcal{A}(s)} \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma v_k(s')]$$

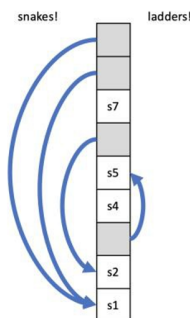
(ج) ابتدا ارزش هر حالت را با سیاست فعلی حساب می کنیم.

$$v_{k+1}(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma v_k(s')]$$

Iteration	State	Value	Policy Action
۱	۱	۱	f
	۲	۲	f
	۳	-۰.۱۴۳	f
۲	۱	۱	f
	۲	۲	f
	۳	۳	f
۳	۱	۱	f
	۲	۲	f
	۳	۳	f

سپس سیاست را با ارزش های بدست آمده بروزرسانی می کنیم.

$$\pi'(s) \leftarrow \arg \max_{a \in \mathcal{A}(s)} \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma v_{\pi}(s')]$$



شما در حال بازی نسخه‌ای ساده‌شده از یک بازی کلاسیک هستید که «مار و پله» نامیده می‌شود.

- شما در امتداد یک مسیر یک‌بعدی حرکت می‌کنید که خانه‌های s_1, s_2, s_4, s_5, s_7 را شامل می‌شود.
- شما دو عمل دارید: climb و quit.
- اگر از حالت s_i عمل climb را انتخاب کنید، با احتمال 0.5 یک خانه بالا می‌روید و با احتمال 0.5 دو خانه بالا می‌روید.
- اما! اگر روی خانه‌ای فرود بیایید که پله‌ای از آن بالا می‌رود، بلافاصله و با احتمال 1 به خانه مقصد پله منتقل می‌شوید.
- و اگر روی خانه‌ای فرود بیایید که ماری از آن پایین می‌آید، بلافاصله و با احتمال 1 به خانه مقصد مار سقوط می‌کنید.
- بنابراین، برای مثال، اگر در s_5 باشید و climb کنید:
 - با احتمال 0.5 یک خانه بالا می‌روید و به خانه خاکستری s_6 می‌رسید؛ چون s_6 مار دارد، بلافاصله به s_2 سقوط می‌کنید.
 - با احتمال 0.5 دو خانه بالا می‌روید و به s_7 می‌رسید.
 - اگر در s_7 باشید و climb کنید، ابتدا به خانه خاکستری s_8 می‌روید و چون s_8 مار دارد، به s_1 منتقل می‌شوید.
- اگر عمل quit را انتخاب کنید، بازی تمام می‌شود و دیگر نمی‌توانید حرکتی انجام دهید.
- هر اپیزود جدید از s_1 شروع می‌شود.
- پاداش انتخاب climb در هر حالت برابر با 0 است.
- پاداش انتخاب quit در حالت s_i برابر با i است.

- (آ) سیاست بهینه در افق زمانی یک چیست؟ برای هر یک از حالات s_1, s_2, s_4, s_5, s_7 بنویسید که عمل بهینه چیست.
- (ب) اگر مقادیر Q تمام حالات را با صفر مقداردهی اولیه کنید و یک دور value iteration بدون تنزیل ($\gamma = 1$) انجام دهید، تابع Q حاصل چه خواهد بود؟ مقدار $Q(s, \text{quit})$ و $Q(s, \text{climb})$ را برای هر یک از حالات s_1, s_2, s_4, s_5, s_7 مشخص کنید.
- (ج) اگر $\gamma = 1$ (یعنی هیچ تنزیلی نداریم)، سیاست بهینه در افق نامتناهی چیست؟ برای هر یک از حالات s_1, s_2, s_4, s_5, s_7 بنویسید که عمل بهینه چیست.
- (د) اکنون تنزیل ($\gamma < 1$) را در نظر بگیرید. نابرابری‌ای بنویسید که شامل مقادیر عددی، γ ، $Q(s_2, \text{climb})$ و $Q(s_7, \text{climb})$ باشد و شرط را مشخص کند که تحت آن عمل بهینه در حالت s_5 ، انتخاب quit باشد.

حل.

(آ) سیاست بهینه در افق ۱.

	s_1	s_2	s_4	s_5	s_7
عمل بهینه	quit	quit	quit	quit	quit

(ب) یک دور Value Iteration با $V^{(0)} \equiv 0$ و $\gamma = 1$.

	s_1	s_2	s_4	s_5	s_7
$Q(s, \text{quit})$	۱	۲	۴	۵	۷
$Q(s, \text{climb})$	۰	۰	۰	۰	۰

(ج) افق نامتناهی بدون تنزیل ($\gamma = 1$).

	s_1	s_2	s_4	s_5	s_7
عمل بهینه	climb	climb	climb	climb	quit

(د) شرط بهینگی quit در s_5 با تنزیل. می‌خواهیم $Q(s_5, \text{quit}) > Q(s_5, \text{climb})$.

$$\begin{aligned}
 Q(s_5, \text{quit}) &= 5, \\
 Q(s_5, \text{climb}) &= R(s_5, \text{climb}) + \frac{1}{4}\gamma \max_{a'} Q(s_2, a') + \frac{1}{4}\gamma \max_{a'} Q(s_7, a') \\
 &= 0 + \frac{1}{4}\gamma \max(Q(s_2, \text{quit}), Q(s_2, \text{climb})) + \frac{1}{4}\gamma Q(s_7, \text{quit}) \\
 &= \frac{\gamma}{4} (\max(Q(s_2, \text{quit}), Q(s_2, \text{climb})) + 7).
 \end{aligned}$$

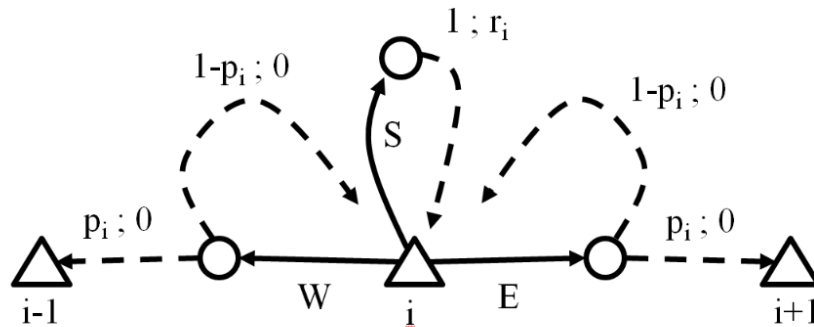
پس شرط بهینگی quit در s_5 برابر است با

$$5 > \frac{\gamma}{4} (\max\{2, Q(s_2, \text{climb})\} + 7).$$

۴. (نمره)

در امتداد یک بزرگراه اصلی، N شهر وجود دارد که با شماره‌های ۱ تا N شماره‌گذاری شده‌اند. شما تاجری از شهر ۱ هستید (شروع شما از آنجاست). هر روز، می‌توانید یکی از کارهای زیر را انجام دهید: به شهر مجاور بروید (حرکت به سمت شرق یا غرب)، یا در شهر فعلی بمانید و به تجارت بپردازید (عمل Stay). اگر تصمیم بگیرید از شهر i سفر کنید، با احتمال p_i با موفقیت به شهر بعدی می‌رسید، اما با احتمال $1 - p_i$ گرفتار طوفان می‌شوید و در این صورت روزتان هدر می‌رود و هیچ‌جا نمی‌روید. اگر تصمیم بگیرید در شهر i بمانید و تجارت کنید، پاداشی برابر $r_i > 0$ دریافت می‌کنید؛ روزهای سفر پاداشی برابر با صفر دارند، چه موفق به رفتن به شهر دیگر شوید و چه نشوید.

نمودار زیر، اقدامات و انتقال‌ها از شهر i را نشان می‌دهد. پیکان‌های پررنگ نشان‌دهنده اقدامات هستند؛ پیکان‌های خط‌چین انتقال‌های حاصل را با برچسبی شامل احتمال و پاداش (به همین ترتیب) نمایش می‌دهند.



- (آ) اگر برای همه i ، داشته باشیم $r_i = 1$ ، $p_i = 1$ و $\gamma = 0.5$ باشد، مقدار $V_{\text{stay}}(1)$ در حالتی که همیشه stay انتخاب می‌شود چقدر است؟
- (ب) اگر برای همه i ، داشته باشیم $r_i = 1$ ، $p_i = 1$ و $\gamma = 0.5$ باشد، مقدار بهینه $V^*(1)$ بودن در شهر ۱ چقدر است؟
- (ج) اگر مقادیر r_i و p_i اعداد مثبت معلوم باشند و discount factor تقریباً برابر با یک باشد ($\gamma \approx 1$)، سیاست بهینه را توصیف کنید.
- (د) فرض کنید از الگوریتم Value Iteration استفاده می‌کنیم. مقدار حالت s پس از k مرحله مقدارگذاری است و همه مقادیر اولیه صفر هستند.
- اگر مقدار بهینه برای بودن در شهر ۱ مثبت باشد یعنی $V^*(1) > 0$ ، بیشترین مقدار k که در آن ممکن است $V_k(1) = 0$ باشد چیست؟
- (ه) اگر همه r_i و p_i مثبت باشند، بیشترین مقدار k که در آن ممکن است $V_k(s) = 0$ برای برخی حالت‌ها باشد چیست؟
- (و) فرض کنید r_i ها و p_i ها را نمی‌دانیم، بنابراین تصمیم می‌گیریم از Q-Learning استفاده کنیم.
- فرض کنید دنباله زیر از حالت‌ها، کنش‌ها و پاداش‌ها تجربه شده است:
- $(s=1, a=\text{stay}, r=4)$
 - $(s=1, a=\text{east}, r=0)$
 - $(s=2, a=\text{stay}, r=6)$
 - $(s=2, a=\text{west}, r=0)$
 - $(s=1, a=\text{stay}, r=4)$

اگر نرخ یادگیری برابر 0.5 ، discount factor برابر ۱ و همه مقادیر اولیه $Q(s, a) = 0$ باشند، مقدارهای نهایی $Q(s, a)$ را در جدول زیر وارد کنید. هر سطر باید مقادیر q-value پس از انتقال مشخص شده در ستون اول را نشان دهد. مقادیر بدون تغییر را می‌توان خالی گذاشت.

(s, a, r, s')	$Q(1, S)$	$Q(1, E)$	$Q(2, W)$	$Q(2, S)$
initial	0	0	0	0
(1, S, 4, 1)				
(1, E, 0, 2)				
(2, S, 6, 2)				
(2, W, 0, 1)				
(1, S, 4, 1)				

حل.

(آ)

$$\forall i \in \{1, \dots, N\}, \quad V^{\text{stay}}(i) = r_i + \gamma V^{\text{stay}}(i)$$

$$V^{\text{stay}}(i) = 1 + 0.5 V^{\text{stay}}(i)$$

$$V^{\text{stay}}(i) = 2$$

$$V^{\text{stay}}(1) = 2$$

(ب) برای تمام شهرها (حالت‌ها) $i = 1, \dots, N$ ، معادلات بلمن به صورت زیر نوشته می‌شوند:

$$V^*(i) = \max \left\{ \underbrace{r_i + \gamma V^*(i)}_{\text{stay}}, \underbrace{p_i \gamma V^*(i-1) + (1-p_i) \gamma V^*(i)}_{\text{left}}, \underbrace{p_i \gamma V^*(i+1) + (1-p_i) \gamma V^*(i)}_{\text{right}} \right\}$$

با توجه به اینکه $p_i = 1$ ، معادله به شکل زیر ساده می‌شود:

$$V^*(i) = \max \left\{ \underbrace{1 + \gamma V^*(i)}_{\text{stay}}, \underbrace{\gamma V^*(i-1)}_{\text{left}}, \underbrace{\gamma V^*(i+1)}_{\text{right}} \right\}$$

از آنجا که همه شهرها پاداش یکسان دارند ($r_i = 1$) و مقدار $V^*(i)$ در تمام حالت‌ها یکسان است، ماکزیمم مقدار همیشه با عمل stay به دست می‌آید. در نتیجه:

$$V^*(i) = 1 + \gamma V^*(i) \implies V^*(i) = \frac{1}{1-\gamma} = 2 \quad (\text{برای } \gamma = 0.5)$$

(ج) سیاست بهینه این است که همیشه به سمت شهری با بیشترین پاداش حرکت کنید. پس از رسیدن به آن شهر، برای همیشه در آنجا بمانید و به تجارت ادامه دهید.

(د) با فرض $r_i > 0$ ، بزرگترین مقدار k برابر با ۰ خواهد بود، زیرا:

$$V_1(s) = \max\{r_i + 0, \dots\} > 0$$

اگر فرض $r_i > 0$ را نداشته باشیم، آنگاه بزرگترین k ممکن برابر با $N-1$ خواهد بود.

۱. از آنجا که $V^*(1) > 0$ ، حداقل یکی از r_i ‌ها باید positive strictly باشد.
۲. پس از یک تکرار:

$$V_1(i) > 0$$

۳. پس از دو تکرار:

$$V_2(i-1) > 0$$

۴. نهایتاً پس از i تکرار:

$$V_i(1) > 0$$

۵. اگر برای همه $i < j$ داشته باشیم $r_j = 0$ ، آنگاه:

$$V_j(1) = 0 \quad \forall j < i$$

۶. در بدترین حالت وقتی $i = N$:

$$V_{N-1}(1) = 0 \quad \text{امکان پذیر است}$$

اما:

$$V_N(1) > 0$$

(ه) با فرض اینکه $r_i > 0$ باشد، بزرگترین مقدار ممکن برای k برابر با ۰ خواهد بود، زیرا:

$$V_1(s) = \max\{r_i + 0, \dots\} > 0$$

(و) پس از مشاهده انتقال $(1, S, 4, 1)$ ، مقدار $Q(1, S)$ را update می‌کنیم:

$$Q(1, S) \leftarrow 0.5[4 + 1 \cdot 0] + 0.5(0) = 2$$

پس از مشاهده انتقال $(1, E, 0, 2)$ ، مقدار $Q(1, E)$ را update می‌کنیم:

$$Q(1, E) \leftarrow 0.5[0 + 1 \cdot 0] + 0.5(0) = 0$$

پس از مشاهده انتقال $(2, S, 6, 2)$ ، مقدار $Q(2, S)$ را update می‌کنیم:

$$Q(2, S) \leftarrow 0.5[6 + 1 \cdot 0] + 0.5(0) = 3$$

پس از مشاهده انتقال $(2, W, 0, 1)$ ، مقدار $Q(2, W)$ را update می‌کنیم:

$$Q(2, W) \leftarrow 0.5[0 + 1 \cdot 2] + 0.5(0) = 1$$

پس از مشاهده انتقال $(1, S, 4, 1)$ ، مقدار $Q(1, S)$ را update می‌کنیم:

$$Q(1, S) \leftarrow 0.5[4 + 1 \cdot 2] + 0.5(2) = 4$$

۵. (نمره) یک مدل MDP با افق نامتناهی و ضریب تخفیف γ با مشخصات زیر در نظر بگیرید:

$$T(s, a, s') = P(s' | f(s, a)), \quad R(s, a, s') = R(s, a)$$

که در آن $f : S \times A \rightarrow Y$ یک تابع قطعی به مجموعه حالت‌های پس از تصمیم Y است. دنباله حالت‌ها s_t ، اقدامات a_t ، حالت‌های پس از تصمیم y_t و پاداش‌ها r_t به صورت زیر است:

$$\begin{array}{ccccccc} (s_0, a_0) & \xrightarrow{f} & y_0 & \xrightarrow{P} & (s_1, a_1) & \xrightarrow{f} & y_1 & \xrightarrow{P} & (s_2, a_2) \\ \downarrow & & \downarrow & & \downarrow & & \downarrow & & \downarrow \\ & & r_0 & & r_1 & & r_2 & & \end{array}$$

مقدار مورد انتظار با تخفیف پاداش‌ها تحت سیاست π به صورت زیر تعریف می‌شود:

$$V^\pi(s_0) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid a_t = \pi(s_t) \right]$$

و تابع ارزش بهینه به صورت $V^*(s) = \max_{\pi} V^\pi(s)$ برای حالت‌های پس از تصمیم:

$$W^\pi(y_0) = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} R(s_t, a_t) \mid y_0 \right], \quad W^*(y) = \max_{\pi} W^\pi(y)$$

(آ) W^* را بر حسب V^* به دست آورید.

(ب) V^* را بر حسب W^* بنویسید. راهنمایی: ابتدا V^* را بر حسب Q بنویسید و سپس آن را بر اساس W^* بنویسید.

(ج) می‌دانیم معادله بلمن فورد برای تابع ارزش به شکل زیر تعریف می‌شود.

$$V^*(s) = \max_a \left[R(s, a) + \gamma \sum_{s'} T(s, a, s') V^*(s') \right]$$

معادل این معادله برای W^* را ارائه دهید.

(د) جاهای خالی را در الگوریتم زیر طوری پر کنید که یک الگوریتم policy iteration به دست آید. این الگوریتم باید تضمین کند که π^* را پیدا می‌کند. همچنین هر مورد را توضیح دهید:

- سیاست $\pi^{(1)}$ را به صورت دلخواه مقداردهی اولیه کنید
- برای $i = 1, 2, 3, \dots$:
 - $W^{\pi^{(i)}}(y)$ را برای تمام $y \in Y$ محاسبه کن
 - سیاست جدید $\pi^{(i+1)}$ را طوری محاسبه کن که $\pi^{(i+1)}(s) = \arg \max_a$ _____
 - اگر _____ برای تمام $s \in S$ ، $\pi^{(i)}$ را برگردان

(ه) با توجه به دنباله مشاهده شده $s_t, a_t, y_t, s_{t+1}, a_{t+1}, y_{t+1}$ و نرخ یادگیری $\alpha \in (0, 1)$ ، قاعده بهروزرسانی برای تخمین W^* را مشابه قاعده بهروزرسانی الگوریتم Q-learning کامل کنید.

$$W(y_t) \leftarrow (1 - \alpha)W(y_t) + \alpha \underline{\hspace{2cm}}$$

حل.

(آ) W^* بر حسب V^*

$$W^*(y) = \sum_{s'} P(s' | y) V^*(s')$$

$$\begin{aligned} W^*(y_\cdot) &= \mathbb{E} [R(s_1, a_1) + \gamma R(s_2, a_2) + \gamma^2 R(s_3, a_3) + \dots | y_\cdot] \\ &= \sum_{s_1} P(s_1 | y_\cdot) \mathbb{E} [R(s_1, a_1) + \gamma R(s_2, a_2) + \dots | s_1] \\ &= \sum_{s_1} P(s_1 | y_\cdot) V^*(s_1) \end{aligned}$$

با استفاده از مستقل بودن V^* از زمان، جایگزین می‌کنیم $y_\cdot \rightarrow y$ و $s_1 \rightarrow s'$.

(ب) V^* بر حسب W^*

$$V^*(s) = \max_a [R(s, a) + \gamma W^*(f(s, a))]$$

$$\begin{aligned} V^*(s_\cdot) &= \max_{a_\cdot} Q(s_\cdot, a_\cdot) \\ &= \max_{a_\cdot} \mathbb{E} [R(s_\cdot, a_\cdot) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \dots | s_\cdot, a_\cdot] \\ &= \max_{a_\cdot} (R(s_\cdot, a_\cdot) + \mathbb{E} [\gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \dots | f(s_\cdot, a_\cdot)]) \\ &= \max_{a_\cdot} (R(s_\cdot, a_\cdot) + \gamma W^*(f(s_\cdot, a_\cdot))) \end{aligned}$$

با تعمیم و تغییر نام متغیرها $s \rightarrow s_\cdot, a \rightarrow a_\cdot$.

(ج) معادله بلمن برای W^*

$$W^*(y) = \sum_{s'} P(s'|y) \max_a (R(s', a) + \gamma W^*(f(s', a)))$$

$$W^*(y) = \sum_{s'} P(s'|y) V^*(s') \quad (\text{از قسمت الف})$$

$$= \sum_{s'} P(s'|y) \max_a [R(s', a) + \gamma W^*(f(s', a))] \quad (\text{از قسمت ب})$$

Policy iteration (د)

$$\pi^{(i+1)}(s) = \arg \max_a [R(s, a) + \gamma W^{\pi^{(i)}}(f(s, a))] \quad ; \quad \pi^{(i)}(s) = \pi^{(i+1)}(s)$$

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E} [R(s, a) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \dots \mid s, a] \\ &= R(s, a) + \gamma \mathbb{E} [R(s_1, a_1) + \gamma R(s_2, a_2) + \dots \mid f(s, a)] \\ &= R(s, a) + \gamma W^\pi(f(s, a)) \end{aligned}$$

بهبود سیاست از $Q^{\pi^{(i)}}(s, a) = \arg \max_a Q^{\pi^{(i)}}(s, a)$ استفاده می‌کند. هنگامی که سیاست همگرا شد، الگوریتم پایان می‌یابد.

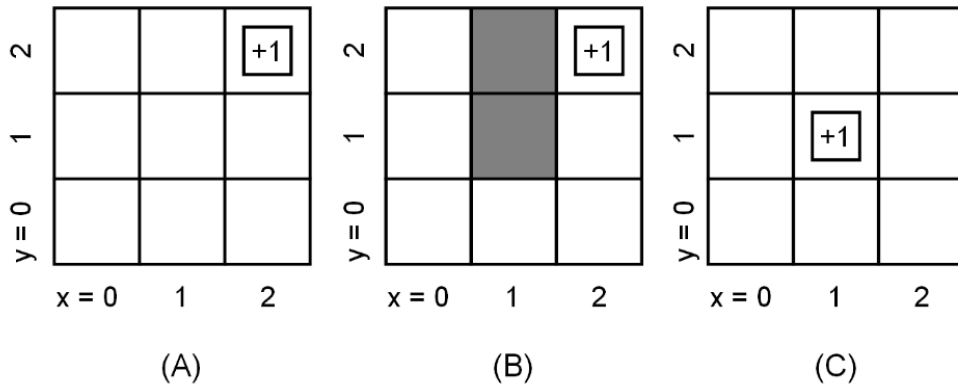
(ه) قاعده به‌روزرسانی برای W^*

$$W(y_t) \leftarrow (1 - \alpha)W(y_t) + \alpha \max_a (R(s_{t+1}, a) + \gamma W(f(s_{t+1}, a)))$$

$$\begin{aligned} W^*(y) &= \mathbb{E}_{s'} [\max_a (R(s', a) + \gamma W^*(f(s', a)))] \\ &\approx \max_a (R(s_{t+1}, a) + \gamma W(f(s_{t+1}, a))) \quad (\text{تخمین با یک نمونه}) \end{aligned}$$

تخمین فعلی را با نرخ یادگیری α به سمت مقدار هدف به‌روزرسانی کنید.

۶. (نمره) در مسئله زیر مربوط به agent gridworld، می‌تواند اعمال E، S، N، W را انجام دهد که عامل را به ترتیب به سمت شمال، جنوب، شرق و غرب یک خانه حرکت می‌دهند. هیچ نویزی وجود ندارد، بنابراین این اعمال همیشه agent را در جهت مورد نظر حرکت می‌دهند، مگر آنکه آن جهت به بیرون از شبکه یا به خانه‌ای مسدود (خاکستری) منتهی شود، که در آن صورت هیچ حرکتی انجام نمی‌شود. خانه‌هایی که دارای ۱+ هستند همچنین اجازه اجرای عمل X را می‌دهند که عامل را از شبکه خارج کرده و وارد حالت نهایی (terminal) می‌کند. پاداش برای تمام انتقال‌ها صفر است، به جز انتقال خروج که پاداش آن ۱+ است. discount factor برابر 0.5 فرض می‌شود.



- (آ) مقدارهای بهینه برای شبکه (A) را کامل کنید.
- (ب) سیاست بهینه برای شبکه (B) را مشخص کنید.
- (ج) فرض کنید برای هر حالت غیرنهایی $s = (x, y)$ مجموعه‌ای از ویژگی‌های حقیقی $f_i(s)$ داریم، و می‌خواهیم مقدار بهینه $V^*(s)$ را با رابطه خطی زیر تقریب بزنیم:

$$V(s) = \sum_i w_i \cdot f_i(s)$$

- اگر ویژگی‌ها $f_1(x, y) = x$ و $f_2(x, y) = y$ باشند، مقادیر w_1 و w_2 را طوری بدهید که سیاست به‌دست‌آمده از (one-step look-ahead) در شبکه (A) بهینه باشد.
- (د) آیا می‌توان مقادیر واقعی و بهینه V^* را برای شبکه (A) تنها با استفاده از این دو ویژگی نمایش داد؟ چرا یا چرا نه؟
- (ه) برای هر یک از مجموعه ویژگی‌های زیر مشخص کنید کدام یک (در صورت وجود) از های MDP شبکه‌ای بالا می‌توانند «حل شوند»؛ به این معنا که بتوان مقادیری (احتمالاً غیربهینه) برای $V(s)$ پیدا کرد که سیاست حاصل از نگاه یک‌مرحله‌ای به آینده، بهینه باشد.

i. $f_2(x, y) = y, f_1(x, y) = x$

ii. برای هر (i, j) ، ویژگی $f_{i,j}(x, y) = 1$ اگر $(x, y) = (i, j)$ ، و صفر در غیر این صورت.

iii. $f_1(x, y) = (x - 1)^2, f_2(x, y) = (y - 1)^2$ و $f_3(x, y) = 1$

حل.

(آ)

(با انجام عمل خروج) $V(2, 2) = 1$

$V(1, 2) = 0.5 \cdot V(2, 2) = 0.5$

$V(2, 1) = 0.5 \cdot V(2, 2) = 0.5$

$V(1, 1) = 0.5 \cdot V(2, 1) = 0.25$

$V(0, 2) = 0.5 \cdot V(1, 2) = 0.25$

$V(2, 0) = 0.5 \cdot V(2, 1) = 0.25$

$V(1, 0) = 0.5 \cdot V(2, 0) = 0.125$

$V(0, 1) = 0.5 \cdot V(1, 1) = 0.125$

$V(0, 0) = 0.5 \cdot V(1, 0) = 0.0625$

ماتریس نهایی مقادیر:

	$x = 0$	$x = 1$	$x = 2$
$y = 2$	0/25	0/5	1
$y = 1$	0/125	0/25	0/5
$y = 0$	0/0625	0/125	0/25

(ب)

	$x = 0$	$x = 1$	$x = 2$
$y = 2$	↓		✓
$y = 1$	↓		↑
$y = 0$	→	→	↑

(ج) • اگر $|w_1| > |w_2|$ ، عامل ترجیح می دهد در جهت افقی حرکت کند:

- $w_1 > 0$: به شرق

- $w_1 < 0$: به غرب

• اگر $|w_1| < |w_2|$ ، عامل ترجیح می دهد در جهت عمودی حرکت کند:

- $w_2 > 0$: به شمال

- $w_2 < 0$: به جنوب

• اگر $|w_1| = |w_2|$:

- هر دو مثبت: حرکت به صورت قطری به شمال شرق (رسیدن به (2, 2))

- هر دو منفی: حرکت به جنوب غرب (دور شدن از هدف)

در نتیجه، اگر $w_1 = w_2 = c > 0$ باشد و عامل از برخورد با دیوار اجتناب کند، مسیر بهینه را به سمت (2, 2) طی کرده و با انجام عمل X پاداش می گیرد.

(د) فرض می کنیم تابع ارزش را با فرم زیر تقریب می زنیم:

$$V(x, y) = w_1 x + w_2 y$$

و مقادیر بهینه برخی خانه ها در گرید A به صورت زیر است:

$$V(0, 2) = 0/25 \Rightarrow 0w_1 + 2w_2 = 0/25$$

$$V(1, 2) = 0/5 \Rightarrow 1w_1 + 2w_2 = 0/5$$

$$V(2, 2) = 1 \Rightarrow 2w_1 + 2w_2 = 1$$

از رابطه اول:

$$2w_2 = 0/25 \Rightarrow w_2 = 0/125$$

جای گذاری در رابطه دوم:

$$w_1 + 2(0/125) = 0/5 \Rightarrow w_1 = 0/25$$

بررسی در معادله سوم:

$$2(0/25) + 2(0/125) = 0/5 + 0/25 = 0/75 \neq 1$$

دستگاه معادلات ناسازگار است؛ یعنی هیچ انتخابی از w_1, w_2 وجود ندارد که همه مقادیر بهینه را دقیق بازسازی کند. تابع ارزش به طور غیرخطی به موقعیت و مسیر بهینه وابسته است. ویژگی های x, y اطلاعات کافی ندارند.

- (ه) i. $f_2(x, y) = y, f_1(x, y) = x$ فقط شبکه‌ی (A) قابل حل است. زیرا ویژگی‌ها agent را به یک جهت ثابت هدایت می‌کنند، اما در (B) وجود مانع و در (C) موقعیت مرکزی هدف باعث ناهماهنگی با این جهت می‌شود.
- ii. $f_{i,j}(x, y) = 1$ اگر $(x, y) = (i, j)$ ، وگرنه 0. هر سه شبکه (A)، (B)، (C) قابل حل هستند؛ زیرا برای هر خانه ویژگی جداگانه وجود دارد و می‌توان با انتخاب مناسب وزن‌ها جهت بهینه را تعیین کرد. (one-hot encoded)
- iii. $f_2(x, y) = (y - 1)^2, f_1(x, y) = (x - 1)^2$ فقط شبکه‌ی (C) قابل حل است؛ این تابع ارزش زمانی جهت درست ایجاد می‌کند که $w_1 = w_2 = c < 0$ و $w_3 = 0$ انتخاب شود. اما در شبکه‌های (A) و (B) هدف در گوشه قرار دارد.