



هوش مصنوعی

پاییز ۱۴۰۱

استاد: محمدحسین رهبان

گردآورندگان: محمدجواد هزاره-امیرحسین جوادی

مهلت ارسال: ۵ بهمن

فرآیندهای مارکف و یادگیری تقویتی

تمرین ششم

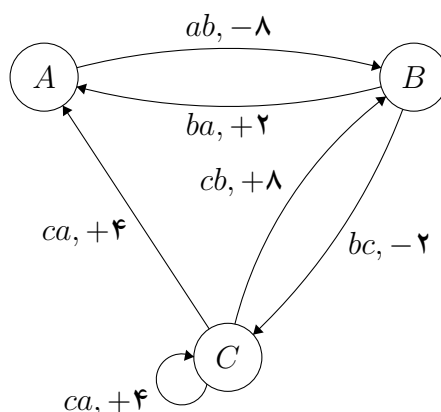
- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- امکان ارسال با تاخیر برای این تمرین وجود ندارد.
- همکاری و همفکری شما در انجام تمرین مانعی ندارد اما پاسخ ارسالی هر کس حتما باید توسط خود او نوشته شده باشد.
- در صورت همفکری و یا استفاده از هر منابع خارج درسی، نام همفکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.

سوالات نظری (۱۰۰ نمره)

۱. (۱۰ نمره) درستی یا نادرستی گزاره‌های زیر را در رابطه با یک فرآیند تصمیم‌گیری مارکف^۱ مشخص کنید و توضیحی کوتاه در رابطه با آن ارائه دهید.

- (آ) ضریب تخفیف^۲ کوچک و نزدیک به صفر به رفتار حریصانه و کوتاه‌نظر^۳ منجر می‌شود.
- (ب) پاداش منفی زندگی^۴ با اندازه‌ی زیاد (بسیار منفی) به رفتار حریصانه و کوتاه‌نظر منجر می‌شود.
- (ج) همواره می‌توان پاداش منفی زندگی را با استفاده از ضریب تخفیف منفی مدل کرد.
- (د) همواره می‌توان ضریب تخفیف منفی را با پاداش منفی زندگی مدل کرد.

۲. (۲۵ نمره) فرآیند تصمیم‌گیری مارکف که در شکل ۱ آمده است را با ضریب تخفیف $\gamma = 0.5$ در نظر بگیرید که در آن حالت‌ها با حروف A، B و C نشان داده شده‌اند. روی هر یال حروف کوچک نوشته شده که یکی از کنش‌های موجود است و یال مربوطه گذار متناظر با انجام آن کنش را نشان می‌دهد. عدد صحیح روی هر یال نیز پاداش کسب شده از آن کنش است. تمام گذارها با احتمال ۱ به وقوع می‌پیوندند و تنها گذار از حالت C به A تصادفی است که احتمال رفتن به حالت A برابر $\frac{1}{3}$ و احتمال رفتن به حالت C برابر $\frac{2}{3}$ است.



شکل ۱: گراف فرآیند تصمیم‌گیری مارکف.

^۱ Markov Decision Process

^۲ discount

^۳ shortsighted

^۴ negative living reward

با در نظر گرفتن این فرآیند به سوال‌های زیر پاسخ دهید.

- (آ) برای یک فرآیند تصمیم‌گیری مارکف به همراه ضریب تخفیف، تابع ارزش حالت‌ها^۵ یا همان $V^\pi(s)$ را توصیف کنید.
- (ب) رابطه‌ی بلمن را برای تابع ارزش حالت‌ها بنویسید.
- (ج) سیاست اولیه‌ی π_1 را در نظر بگیرید که به صورت تصادفی و با احتمال برابر در هر حالت یکی از کنش‌های موجود در آن حالت را انتخاب می‌کند. حال فرض کنید تابع ارزش‌گذاری اولیه را به صورت $V_1(A) = V_1(B) = V_1(C) = 2$ در نظر بگیریم. یک مرحله از الگوریتم ارزیابی سیاست^۶ را اجرا کنید تا به تابع ارزش $V_2(s)$ برای حالت‌های مختلف برسید.
- (د) براساس تابع ارزش‌گذاری جدید و به صورت حریصانه سیاست قطعی جدید π_2 را بدست آورید.
- (ه) سیاست قطعی π را در نظر بگیرید. اثبات کنید اگر سیاست جدید π' به صورت حریصانه از V^π بدست آمده باشد، آنگاه π' بهتر یا مساوی π است، یا به عبارتی برای تمام حالت‌ها داریم $V^{\pi'}(s) \geq V^\pi(s)$. همچنین اثبات کنید اگر تساوی برای تمام حالت‌ها رخ دهد آنگاه π' حتما سیاست بهینه است.
۳. (۲۰ نمره) صفحه‌ی 2×3 زیر را در نظر بگیرید. فرض کنید حرکت خود را از خانه‌ی شماره‌ی ۱ شروع می‌کنیم و با رسیدن به خانه‌ی شماره‌ی ۶ بازی تمام می‌شود و با رسیدن به این خانه ۱۰ امتیاز مثبت دریافت می‌کنیم. همچنین در تمام حرکت‌هایی که منجر به رسیدن به خانه‌ی شماره‌ی ۶ نمی‌شوند پاداش ۱- دریافت می‌کنیم.

۴	۵	۶
۱	۲	۳

شکل ۲: جدول بازی.

در هر خانه چهار کنش ممکن وجود دارد: بالا، پایین، چپ و راست. فرض کنید کنش‌هایی که باعث خارج شدن از صفحه می‌شوند مجاز نیستند. هر کنش نیز به صورت قطعی انجام شده و به خانه‌ی مربوطه می‌رویم. حال فرض کنید جدول زیر را برای $Q(s, a)$ داریم:

$Q(1, \text{راست}) = 3$			$Q(1, \text{بالا}) = 4$
$Q(2, \text{راست}) = 8$	$Q(2, \text{چپ}) = 3$		$Q(2, \text{بالا}) = 6$
	$Q(3, \text{چپ}) = 7$		$Q(3, \text{بالا}) = 9$
$Q(4, \text{راست}) = 5$		$Q(4, \text{پایین}) = 2$	
$Q(5, \text{راست}) = 8$	$Q(5, \text{چپ}) = 5$	$Q(5, \text{پایین}) = 6$	

شکل ۳: جدول Q-value ها

با در نظر گرفتن این جدول و توضیح مسئله به سوال‌های زیر پاسخ دهید.

- (آ) باتوجه به داشتن دانش کامل در رابطه با محیط، می‌توان از رابطه‌ی بلمن برای بروزرسانی Q-value ها استفاده کرد. فرض کنید از سیاست حریصانه استفاده می‌کنیم و با در نظر گرفتن این سیاست، ابتدا رابطه‌ی بلمن برای بروزرسانی Q-value ها را نوشته و سپس مقدار بروز شده‌ی (چپ، ۳) $Q(3, \text{چپ})$ را حساب کنید.
- (ب) حال فرض کنید مدل محیط را نداریم و جدول Q-value های داده شده از روش یادگیری تفاوت زمانی^۷ بدست آمده است. توضیح دهید چرا در این صورت استفاده از سیاست حریصانه هوشمندانه نیست و با برقراری تعادل بین چه مواردی می‌توان سیاست بهتری داشت؟

^۵state-value function

^۶policy evaluation

^۷Temporal Difference Learning

(ج) توضیح دهید چرا به جای استفاده از ارزش حالت‌ها یا همان V-values^۱ از ارزش کنش‌ها یا همان Q-values استفاده شده است.

(د) یکی از روش‌های حل مشکل قسمت (ب) استفاده از سیاست تصادفی softmax است. در این روش احتمال انجام دادن کنش a از حالت s که آن را با $\pi(s, a)$ نشان می‌دهیم به صورت زیر محاسبه می‌شود:

$$\pi(s, a) = \frac{e^{Q(s, a)}}{\sum_b e^{Q(s, b)}}$$

با در نظر گرفتن این سیاست و جدول داده شده برای Q-value ها، احتمال انجام هر کنش در حالت‌های مختلف را بدست آورید. همین‌طور توضیح دهید چرا استفاده از این روش معقولانه است و مشکل قسمت (ب) را برطرف می‌کند.

(ه) حال می‌خواهیم با استفاده از الگوریتم SARSA^۲ مقدار Q-value ها را بروزرسانی کنیم. فرض کنید از خانه ۲ مسیر زیر را نمونه‌برداری کرده‌ایم.

$$2 \rightarrow 5 \rightarrow 6$$

با در نظر گرفتن رابطه‌ی زیر برای بروزرسانی به روش SARSA مقدار $Q(1, \text{بالا})$ و $Q(5, \text{راست})$ را بروزرسانی کنید. ($\alpha = 0.2, \gamma = 0.8$)

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

۴. (۲۰ نمره) یک MDP با دو استیت A و B ، با دو اکشن (۱) و (۲)، و استیت ترمینال (T) با $V(T) = 0$. reward function و transition function ناشناخته است اما نمونه‌های زیر را دیده‌ایم.

- (a) $A \rightarrow B : a_1 = 1, r_1 = -3$
- (b) $B \rightarrow A : a_2 = 1, r_2 = 4$
- (c) $A \rightarrow A : a_3 = 2, r_3 = -4$
- (d) $A \rightarrow B : a_4 = 1, r_4 = -3$
- (e) $A \rightarrow T : a_5 = 2, r_5 = 1$

که هر \rightarrow یک تغییر از حالت مبدا به مقصد با انجام action و reward مشخص شده است.

(آ) مقدار $Q(s, a)$ را بعد از مشاهده این نمونه‌ها تعیین کنید.

(ب) یک سیاست deterministic با توجه به سمپل‌ها معرفی کنید که از سیاست رندوم بهتر است. توضیح دهید.

(ج) سیاست رندوم را با π_{random} و سیاست طراحی شده را با π^* نام‌گذاری کنید. چه انتظار در مورد مقدار نهایی value estimation در زمانی که الگوریتم Q-Learning با سیاست π^* شروع شود نسبت به وقتی با π_{random} شروع شود دارید؟ هر کدام از این سیاست‌ها به چه مشکلاتی ممکن است بینجامد؟

۵. (۲۵ نمره) فرض کنید ما با نرخ اکتشاف ϵ شروع می‌کنیم. به این معنی که هرگاه مدل یک action را انتخاب کند، با احتمال ϵ به صورت تصادفی و با احتمال $1 - \epsilon$ action انتخاب شده انجام می‌شود. اگر فرض کنیم که محیط به اندازه کافی کاوش شده‌است، ممکن است بخواهیم پس از مدتی میزان اکتشاف را کاهش دهیم. یک الگوریتم برای کاهش این نرخ اکتشاف ارائه دهید. اگر حریف استراتژی‌اش را تغییر دهد، آیا روش شما کار می‌کند؟ چرا؟ اگر نه، یک heuristic ارائه دهید که بتواند با تغییرات در استراتژی حریف سازگار شود.

^۱state-action-reward-state-action