

HW06 - Artificial Intelligence

Saeed Rostami

Student ID: 98106542

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.

(آ) هر چند γ کوچک‌تر شود، پاداش‌ها باید از reward های آینده کمتر شود. به منبر به رفتار حریصانه و کوتاه‌نظر می‌شود.

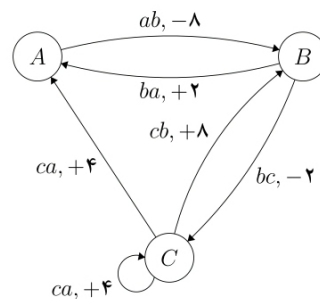
(ب) هر چند γ نزدیک به ۱ باشد، یا agent سعی کند زودتر به هدف برسد یا زودتر terminate شود.

که در هر دو صورت منبر به رفتار حریصانه و کوتاه‌نظر می‌شود.

(ج) خیر - زیرا "منفی شدن مدل کردن به دلیل اینست که living reward تا زیر صفر را به دست جمع می‌گذارد اما discount reward به دست می‌آید چند جملای درجه یک را در این می‌زنند باعث شود که شرایطی نداشته باشد living reward منفی را با آن مدل کنیم.

(د) خیر - پاسخ هاست تست قبل می‌باشد.

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 $\frac{1}{2}$ و احتمال رفتن به حالت C برابر $\frac{1}{2}$ است.



شکل ۱: گراف فرآیند تصمیم‌گیری مارکف.

(آ) برای یک فرآیند تصمیم‌گیری مارکوف به همراه ضریب تخفیف، تابع ارزش حالت‌ها^۵ یا همان $V^\pi(s)$ را توصیف کنید.

$$V_K^\pi(s) = E \left[r_{k+1} + \gamma r_{k+2} + \dots \mid s_k = s \right] \triangleq E \left[\sum_{i=0}^{\infty} \gamma^i r_{k+i} \mid s_k = s \right]$$

(ب) رابطه‌ی بین تابع ارزش حالت‌ها بنویسید.

$$\text{Bellman Eq} \triangleq V_{k+1}^\pi(s) = \sum_{\max a} T(s, a, s') [R(s, a, s') + \gamma V_k^\pi(s')]$$

(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.

$$V_2(A) = -8 + 0.5 V_1(B) = -7$$

$$V_2(B) = 0.5 [-2 + 0.5 V_1(C) + 2 + 0.5 V_1(A)] = 1$$

$$V_2(C) = 0.5 [4 + 0.5 (\frac{3}{4} V_1(C) + \frac{1}{4} V_1(A))] + 8 + 0.5 V_1(B) = 7$$

(د) براساس تابع ارزش‌گذاری جدید و به صورت حریصانه سیاست قطعی جدید π_2 را بدست آورید.

$$Q_2(B, ba) = 2 + 0.5(-7) = -1.5$$

$$Q_2(B, bc) = -2 + 0.5(7) = 1.5$$

$$\pi_2(B) = bc \quad (2)$$

$$\pi_2(A) = ab \quad (1)$$

$$Q_2(C, ca) = 4 + 0.5 \left[\frac{1}{4} V_2(A) + \frac{3}{4} V_2(B) \right] = 5.75$$

$$Q_2(C, cb) = 8 + 0.5(1) = 8.5 \Rightarrow \pi_2(C) = cb \quad (3)$$

(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.

$$\text{we know that for } Q^\pi(s, \pi'(s)) = \max_a Q^\pi(s, a) \geq Q^\pi(s, \pi(s)) = V^\pi(s)$$

$$\rightarrow Q^\pi(s, \pi'(s)) \geq V^\pi(s) \xrightarrow{V^{\pi'}(s) = Q(s, \pi'(s))} V^{\pi'}(s) \geq V^\pi(s)$$

$$\text{if } V^{\pi'}(s) = V^\pi(s) \Rightarrow Q^\pi(s, \pi'(s)) = Q^\pi(s, \pi(s)) \xrightarrow{Q^\pi(s, \pi'(s)) = \max_a Q^\pi(s, a)} \max_a Q^\pi(s, a) = Q^\pi(s, \pi(s))$$

$$\max_a Q^\pi(s, a) = Q^\pi(s, \pi(s))$$

$$\Rightarrow \pi(s) = \text{optimal policy}$$

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:

$Q(۱, \text{بالا}) = ۴$			$Q(۱, \text{راست}) = ۳$
$Q(۲, \text{بالا}) = ۶$	$Q(۲, \text{چپ}) = ۳$		$Q(۲, \text{راست}) = ۸$
$Q(۳, \text{بالا}) = ۹$	$Q(۳, \text{چپ}) = ۷$		
		$Q(۴, \text{پایین}) = ۲$	$Q(۴, \text{راست}) = ۵$
	$Q(۵, \text{چپ}) = ۵$	$Q(۵, \text{پایین}) = ۶$	$Q(۵, \text{راست}) = ۸$

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 $Q(۳, \text{چپ})$ را حساب کنید.
- (ب) حال فرض کنید مدل محیط را نداریم و جدول Q-value های داده شده از روش یادگیری تفاوت زمانی 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؟

$$Q_{k+1} = \sum_a \pi(s, a, s') [R(s, a, s') + \gamma \max_a Q_k(s', a)] \quad (۱)$$

$$\Rightarrow Q(3, \text{چپ}) = \frac{1}{2} [-1 + \gamma \max(6, 3, 8)] = 8\gamma - 1$$

(ب) این کار باعث می‌شود تا خیلی از خانه‌های مسئله explore نشود زیرا ز را هم با استفاده از یک سری randomness رت غیر optimal کنیم.

(ج) چون می‌خواهیم action بهینه در هر مرحله را انتخاب کنیم که در آن کن را به ما می‌دهد. اما طبق رابطه‌ی $\pi(s) = \arg \max_a Q(s, a)$ این policy را به ما می‌دهد. در نتیجه ما سببی Q-value ها را نتق می‌باشد.

$$s: 1 \Rightarrow \pi(1, \text{بالا}) = \frac{e^4}{e^4 + e^3}, \pi(1, \text{پایین}) = \frac{e^3}{e^3 + e^4}$$

$$s: 2 \Rightarrow \pi(2, \text{بالا}) = \frac{e^6}{e^6 + e^8}, \pi(2, \text{پایین}) = \frac{e^8}{e^6 + e^8 + e^3}$$

$$s: 3 \Rightarrow \pi(3, \text{بالا}) = \frac{e^9}{e^9 + e^7}, \pi(3, \text{چپ}) = \frac{e^7}{e^9 + e^7}$$

4	5	6
1	2	3

$$5:4 \rightarrow \pi(4, \underline{a}) = \frac{e^2}{e^2 + e^5}, \pi(4, \overline{a}) = \frac{e^5}{e^2 + e^5}$$

$$8:5 \rightarrow \pi(5, \underline{a}) = \frac{e^6}{e^6 + e^8 + e^5}, \pi(5, \overline{a}) = \frac{e^8}{e^6 + e^8 + e^5}, \pi(5, \underline{b}) = \frac{e^5}{e^6 + e^8 + e^5}$$

ه با استفاده از این مسئله درباره ی ران explanation را انجام دهیم.

$$Q(1, \underline{a}) = 4 + 0.2 \left[\cancel{R_{1,4}^{-1}} + 0.8 Q(1, \overline{a}) \right] - Q(1, \underline{a}) = 3.8 \quad (5)$$

$$Q(5, \overline{a}) = Q(5, \underline{a}) + 0.2 \left[R_{5,6} + 0.8 Q(6, \overline{a}) \right] - Q(5, \overline{a}) = 10$$

۴. (۲۰ نمره) یک MDP با دو استیت A و B ، با دو اکشن (۱) و (۲)، و استیت ترمینال (T) با $V(T) = 0$.
 transition function و reward 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 مبدأ 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؟

$$Q(s, a) = \sum_{s'} T(s, a, s') \underbrace{[R(s, a, s') + \gamma V_K(s')]}_{\text{sample}}$$

$$\rightarrow Q(A, 1) \Rightarrow \text{samples} = \{a, d\} = \frac{-3 - 3}{2} = -3$$

$$Q(A, 2) \Rightarrow \text{samples} = \{c, e\} = \frac{-4 + 1}{2} = -1.5$$

$$Q(B, 1) \Rightarrow \text{samples} = \{b\} = 4$$

$$Q(B, 2) = 0 \rightarrow \text{initid - value}$$

مناسب! Q -value های درست‌تر نسبت به Q های قبلی وارد می‌شود: $V^{\pi}(s) = \max_a Q(s, a)$

$$\boxed{\pi^*(A) = 2} \quad \boxed{\pi^*(B) = 1} \quad \leftarrow \text{سیاست انتخاب شده}$$

اگر با سیاست π_{random} این رسم می‌کشیم می‌بینیم که نیاز گذشتن مقدار خوبی Iteration، جواب ما به مقدار واقعی converge می‌شود. هر چند ممکن است زمان بسیاری طول بکشد تا converge شود.

اگر با π^* شروع کنیم امکان دارد چون یک سیاست را مناسب با نقطه ۵، Episode امتداد ندهیم، exploration به خوبی صورت نگیرد و باعث شود به value های واقعی converge نکرده. اما نسبت به π_{random} این مسئله درست است.

سرعت π_{random} \leftarrow اطمینان از convergence اما کمی time consuming می‌باشد.

سرعت π^* \leftarrow سریع‌تر converge می‌اند اما احتمال هست به مقدار واقعی converge نکند.

۵. (۲۵ نمره) فرض کنید ما با نرخ اکتشاف ϵ شروع می کنیم. 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 می شود. اگر فرض کنیم که محیط به اندازه کافی کاوش شده است، ممکن است بخواهیم پس از مدتی میزان اکتشاف را کاهش دهیم. یک الگوریتم برای کاهش این نرخ اکتشاف ارائه دهید. اگر حریف استراتژی اش را تغییر دهد، آیا روش شما کار می کند؟ چرا؟ اگر نه، یک heuristic ارائه دهید که بتواند با تغییرات در استراتژی حریف سازگار شود.

→ برای بخش اول، یک روش naive و ساده این است که یک $decay rate$ انتخاب کنیم و هر بار (بعد از k ϵ دارد روش ϵ -greedy با فرمول $\epsilon \leftarrow \epsilon - decay rate$ → آیدیت کنیم. حرکت منفی) به ظاهر روش بدی نیست و به خصوص اگر حریف استراتژی خود را تغییر ندهد روش خوبی است.

→ برای بخش دوم که حریف استراتژی خود را تغییر دهد بهتر است که دوباره آیدیت کردن میزان ϵ در معنی
$$C = \sum_k |Q^{\pi}(s) - sample|^2$$
 را حساب کنیم و اگر آیدیم این میزان $error$ ، از یک $Threshold$ پست تر است، هر بار $\epsilon \leftarrow \epsilon + decay rate$ → آیدیت کنیم.

→ در حقیقت در این روش، هر بار میزان اعمال $sample$ ، Q بین بینی شده از قبل را حساب می کنیم تا هیچگاه از حالت $convergent$ خارج نشویم. به زود تغییر حالت $convergent$ به صورت ناگهانی نشان از تغییر استراتژی حریف است.