



یادگیری تقویتی

پاییز ۱۴۰۲

مدرس: محمدحسین رهبان

دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

زمان: ۱۰۰ دقیقه

تاریخ: ۱۳ اردیبهشت ۱۴۰۳

امتحان میان ترم

سوالات (۱۰۰ نمره)

۱. (۲۰ نمره) به سوالات زیر پاسخ دهید.

(آ) فرض کنید در الگوریتم Policy Gradient ساده مانند REINFORCE به جای دنبال کردن سیاست $\pi_\theta(\cdot|s)$ سیاست دیگری را دنبال می‌کنیم که هدفش تقویت اکتشاف (exploration) است. به این منظور با احتمال ϵ یک عمل تصادفی و با احتمال $1 - \epsilon$ یک عمل از توزیع سیاست π_θ نمونه‌برداری می‌کنیم. گرادیان مورد استفاده برای به روزرسانی سیاست π_θ به چه صورت خواهد بود؟ راه حل را به صورت کامل شرح دهید. توجه داشته باشید که در گرادیان سیاست داریم:

$$\nabla_\theta J(\theta) \approx \sum_i \left(\sum_t \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) \right) \left(\sum_t r(s_t^i, a_t^i) \right) \quad (1)$$

پاسخ: کافی است قرار دهیم: $\pi'_\theta(a|s) = \epsilon U(a) + (1 - \epsilon) \pi_\theta(a|s)$. لذا $\nabla_\theta \log \pi'_\theta = \frac{(1-\epsilon)\nabla_\theta \pi_\theta}{\epsilon U(a) + (1-\epsilon)\pi_\theta}$. اگر کسی پاسخی با این مضمون نوشته که الگوریتم را به صورت off-policy در نظر گرفته نیز نمره کامل در صورت صحیح بودن پاسخ داده شود. توجه کنید در این حالت بایستی از مفهوم importance sampling استفاده کرده باشد.

(ب) فرض کنید به یک مجموعه از trajectory های نمونه‌برداری شده از سیاست‌های قبلی دسترسی داریم و هدف این است که تغییری در الگوریتم گرادیان سیاست ساده مانند REINFORCE ایجاد کنیم که بتواند از این داده‌ها نیز استفاده نماید. می‌خواهیم در گرادیان تابع هدف به گونه‌ای تغییر ایجاد کنیم که این مسئله محقق شود. چگونه این کار را انجام دهیم؟ راه حل را شرح دهید. پاسخ: در این حالت باید از روش وزن‌دهی نمونه استفاده کنیم تا وزن گرادیان‌هایی که از trajectory های قبلی بدست آمده تنظیم شود:

$$\nabla_\theta J = \mathbb{E} \left(\prod_{t=1}^T \frac{\pi_\theta(a_t | s_t)}{\pi_\theta^{\text{old}}(a_t | s_t)} \right) \left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right)$$

(ج) در اجرای یک روش گرادیان سیاست مشاهده شده است که میزان بازده (return) نوسان‌های زیادی دارد. علت احتمالی این مسئله چیست؟ دو راه حل برای این مشکل به اختصار مطرح کنید.

پاسخ: روش گرادیان سیاست در حالت ایده‌آل باید در هر گام میزان بازده را بهبود دهد. با این حال، در صورتی که تغییرات سیاست در یک گام به روزرسانی بالا باشد، ممکن است به دلیل تفاوت $p_\theta(\tau)$ و $p_{\theta'}(\tau)$ نتوان بیشتر شدن $J(\theta) - J(\theta')$ را تضمین نمود. لذا برای حل این مشکل، از روش‌های trust region می‌توان استفاده نمود. همچنین، یک دلیل محتمل دیگر می‌تواند بالا بودن واریانس تخمین گرادیان سیاست دانست. برای حل این مشکل می‌توان از راهکارهایی مانند اضافه کردن baseline و همچنین استفاده از تخمین علی پاداش باقی‌مانده تا انتهای اپیزود استفاده کرد.

۲. (۲۰ نمره) فرض کنید در یک MDP پاداش، برخلاف مورد متداول بررسی شده در کلاس، یک متغیر تصادفی وابسته به حالت جاری است. (منظور از وابستگی این متغیر تصادفی به حالت جاری این است که در مدل گرافی توصیف کننده MDP، یک پال از متغیر حالت در هر زمان به متغیر پاداش در آن زمان متصل است.) لذا معادله بهینگی Bellman به صورت زیر در آمده است:

$$V^*(s) = \max_a \mathbb{E}_{s',R} (R + \gamma V^*(s')). \quad (۲)$$

ثابت کنید الگوریتم Value Iteration در این حالت، به نقطه ثابت این معادله همگرا می‌شود.

برای اثبات همگرایی این الگوریتم کافی است Lipschitz بودن اپراتور Bellman جدید τ را بررسی نماییم.

$$\tau V_1(s) - \tau V_2(s) = \max_a \sum_{s',r} p(s'|s, a) p(r|s) (r + \gamma V_1(s')) - \max_a \sum_{s',r} p(s'|s, a) p(r|s) (r + \gamma V_2(s'))$$

. فرض کنید پاسخ بهینه‌سازی اول بالا، عملی مانند a_1 باشد. داریم:

$$\tau V_1(s) - \tau V_2(s) \leq \sum_{s',r} p(s'|s, a_1) p(r|s) (r + \gamma V_1(s')) - \sum_{s',r} p(s'|s, a_1) p(r|s) (r + \gamma V_2(s')).$$

لذا بدست می‌آید:

$$\tau V_1(s) - \tau V_2(s) \leq \sum_{s',r} \gamma p(s'|s, a_1) p(r|s) (V_1(s') - V_2(s')) \leq \gamma \max_{s'} V_1(s') - V_2(s') \sum_{s',r} p(s'|s, a) p(r|s).$$

لذا:

$$\tau V_1(s) - \tau V_2(s) \leq \gamma \max_{s'} V_1(s') - V_2(s').$$

به این ترتیب اثبات می‌شود این اپراتور حالت contraction mapping دارد و همگرایی تضمین می‌شود.

۳. (۱۵ نمره) رضا هر روز سه حالت دارد، یا آسوده است، یا نگران است و یا خسته است. او هر روز دو کار می‌تواند انجام دهد، یا درس می‌خواند یا درس نمی‌خواند. اطلاعات مربوط به پاداش‌ها و احتمالات انتقال هر اکشن در شکل آمده است. $\gamma = 0.8$.

Reward	Probability	Next state	Action	First state
۴	۰/۸	Tired	Study	Worried
۶	۰/۲	Relaxed	Study	Worried
-۲	۱	Worried	study Don't	Worried
۴	۰/۷	Tired	Study	Tired
۳	۰/۳	Relaxed	Study	Tired
۰	۰/۲	Tired	Don't study	Tired
-۲	۰/۳	Worried	Don't study	Tired
۵	۰/۵	Relaxed	Don't study	Tired
۱۰	۰/۴	Relaxed	Study	Relaxed
۸	۰/۶	Tired	Study	Relaxed
۲	۰/۵	Relaxed	Don't study	Relaxed
-۸	۰/۵	Worried	Don't study	Relaxed

(آ) ارزش هر حالت را با استفاده از Value Iteration تا سه مرحله حساب کنید.

(ب) ارزش هر اکشن در هر حالت یا همان Q-Value را حساب کنید.

(ج) رضا در کدام حالت‌های خود باید درس بخواند؟ در کدام حالت‌ها بهتر است درس نخواند؟ (نیازی به محاسبه تا همگرایی نیست و سه مرحله کافی است)

۴. (۱۵ نمره) به سوالات زیر در حوزه Deep Q-Learning پاسخ دهید.

(آ) دو دلیل استفاده از Replay Buffer در این الگوریتم را به اختصار توضیح دهید.
 پاسخ: این بافر برای ذخیره‌سازی انتقال بین حالت‌های متوالی قبلی و استفاده از آنها در کمیته‌سازی تابع زیان در گام‌های بعدی استفاده می‌شود. در صورتی که این کار انجام نشود، مدل دچار فراموشی می‌شود. همچنین، این کار باعث data efficient شدن الگوریتم نیز می‌شود، چرا که جمع‌آوری داده‌های قبلی، هزینه محاسباتی قابل توجهی برای عامل داشته است و با ذخیره آنها این هزینه محاسباتی صرفه‌جویی می‌شود.

(ب) چرا در Replay Buffer مقدار target Q-value را نگهداری نمی‌کنیم؟
 پاسخ: به این دلیل که target در طول زمان ثابت نیست و در هر لحظه به مقدار Q function که به صورت لگ‌دار در طول زمان Q function اصلی را دنبال می‌کند، وابسته است.

(ج) چرا در محاسبه مقدار target Q-value از میانگین متحرک نمایی (exponential moving average) وزن‌های شبکه Q به صورت زیر استفاده می‌کنیم؟ در صورتی که بخواهیم مقدار τ را بر اساس یک زمان‌بندی در طول آموزش تغییر دهیم، چه زمان‌بندی برای این کار مناسب‌تر است؟
 پاسخ: دلیل این مسئله آن است که نمی‌خواهیم اهداف یادگیری بین گام‌های متوالی تغییرات زیادی داشته باشد. در غیر اینصورت بهینه‌سازی تابع زیان همگرا نخواهد شد. برای زمان‌بندی τ بهتر است این مقدار از یک شروع شده و به تدریج به صفر میل داده شود.

$$\theta' = \tau\theta + (1 - \tau)\theta'. \quad (۳)$$

۵. (۱۵ نمره) رابطه گرادیان سیاست را در نظر داشته و به سوالات زیر پاسخ دهید:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \quad (۴)$$

(آ) با حذف جملات غیر علی از رابطه بالا به عبارت ساده‌تری برسید.

$$\frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \sum_{t'=t}^T r(s_{i,t'}, a_{i,t'}) \right).$$

(ب) نشان دهید رابطه بدست آمده از قسمت قبل، چگونه باعث کاهش واریانس تخمین گرادیان سیاست می‌شود.

(ج) نشان دهید عبارات به دست آمده در قسمت اول، به تخمین گرادیان سیاست بایاس اضافه نمی‌کند.

۶. (۱۵ نمره) در رابطه با روش‌های Temporal Difference و Monte Carlo به سوالات زیر پاسخ دهید.

(آ) در چه شرایطی نمی‌توان از روش‌های Monte Carlo استفاده نمود؟

در شرایطی که محیط episodic نباشد.

(ب) دو روش را از حیث Bias/Variance Trade-off مقایسه نمایید.

MC تخمین بدون سوگیری از $v_{\pi}(S_t)$ به دست می‌دهد در حالیکه تخمین TD دارای سوگیریست مگر آنکه $E[v_t(S_{t+1}) | S_{t+1}] = v_{\pi}(S_{t+1})$ ؛ اما تخمین TD واریانس کمتری دارد زیرا بر مبنای مشاهدات بیشتریست.

(ج) یک محیط با دو حالت A و B را در نظر بگیرید. دو دنباله‌ی نمونه از کنش‌ها و پاداش‌ها به صورت زیر داده شده است.

$$A \xrightarrow{3} A \xrightarrow{2} B \xrightarrow{1} \text{Terminate}$$

$$A \xrightarrow{3} A \xrightarrow{3} A \xrightarrow{2} B \xrightarrow{1} \text{Terminate}$$

مقدار تابع Value برای حالات A و B را با روش TD دو قدمه و Monte Carlo به طور جداگانه به دست آورید. مقدار اولیه‌ی تابع را ۰ و γ را برابر ۱ در نظر بگیرید.

• MC:

$$G_{1,3} = 1, G_{1,2} = 3, G_{1,1} = 6$$

$$G_{2,4} = 3, G_{2,3} = 3, G_{2,2} = 6, G_{2,1} = 9$$

$$G(A) = 3 + 6 + 3 + 6 + 9 = 27, N(A) = 5 \rightarrow V(A) = 5.4$$

$$G(B) = 1 + 1, N(B) = 2 \rightarrow V(B) = 1$$

• TD-2: مقدار α به دانشجویان گفته نشده و بسته به مقداری که خودشان در نظر گرفته‌اند پاسخ بررسی می‌شود.

$$V(A) = \alpha(3 + 2)$$

$$V(A) = V(A) + \alpha(2 + 1 - V(A)) = \alpha(3 + 5(1 - \alpha))$$

$$V(B) = \alpha$$

$$V(A) = \alpha(6 + 3(1 - \alpha) + 5(1 - \alpha)^2)$$

$$V(A) = \alpha(5 + 6(1 - \alpha) + 3(1 - \alpha)^2 + 5(1 - \alpha)^3)$$

$$V(A) = \alpha(3 + 5(1 - \alpha) + 6(1 - \alpha)^2 + 3(1 - \alpha)^3 + 5(1 - \alpha)^4)$$

$$V(B) = \alpha + \alpha(1 - \alpha)$$

برای $\alpha = \frac{1}{2}$ ، $V(A) = 3.84375 = \frac{123}{32}$ و $V(B) = 0.75 = \frac{3}{4}$ ، $\alpha = 1$ برای $V(B) = 1$ و $V(A) = 3$.