



لطفا نکات زیر را رعایت کنید:

- سوالات خود را از طریق پست مربوط به تمرین در Quera مطرح کنید.
- پاسخ ارسالی واضح و خوانا باشد.
- در هر کدام از سوالات، اگر از منابع خاصی استفاده کرده‌اید باید آن را ذکر کنید.
- اگر با افرادی همفکری کرده‌اید، نام ایشان را ذکر کنید.
- پاسخ ارسالی باید توسط خود شما نوشته شده باشد.
- تمام پاسخ‌های خود را در یک فایل با فرمت zip [Fullname]_[SID]_RL_HW# روی کوثر قرار دهید.
- برای ارسال هر تمرین تا ساعت ۲۳:۵۹ روز ددلاین فرصت دارید. علاوه بر آن، در هر تمرین می‌توانید تا سقف ۲ روز از تأخیر مجاز باقیمانده‌ی خود استفاده کنید و در مجموع ۵ روز تأخیر مجاز برای تمارین در اختیار دارید.

سوال ۱: هدف یادگیری تقویتی (نظری) (۱۰ نمره)

هدف یادگیری تقویتی پیشنهادی مجموعه پاداش در طول زمان است

$$J(\theta) = \mathbb{E}[\sum_t \gamma^t r(s_t, a_t)]. \quad (1)$$

الف) منشأ تصادفی بودن میزان جمع پاداش که منجر می‌شود بخواهیم به جای مقدار آن در یک اپیزود، به امید ریاضی آن توجه کنیم چه چیزهایی است؟

ب) فرض کنید الگوریتم مورد استفاده، یادگیری یک سیاست $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$ باشد. چرا نمی‌توان θ را مستقیم با استفاده از گرادین کاهشی بر روی تابع هزینه تعریف شده در ۱ بهینه کرد؟ پیشنهادی دهید که چگونه می‌توان تنظیمات یادگیری تقویتی (شیوه‌ی اعمال کنش و دریافت پاداش) را تغییر داد به طوری که امکان بهینه‌سازی با اعمال گرادین کاهشی به طور مستقیم بر روی رابطه‌ی ۱ وجود داشته باشد (راهنمایی: دلیلی که سبب می‌شود نتوان از رابطه‌ی ۱ به صورت مستقیم مشتق گرفت آن است که ما در حالت معمول در یادگیری تقویتی کنش را نمونه‌گیری کرده و انتخاب می‌کنیم و این عملیات قابلیت مشتق‌گیری ندارد).

سوال ۲: بایاس و واریانس گرادین سیاست (نظری) (۲۰ نمره)

الف) ثابت کنید تخمین رابطه‌ی ۲ از $\nabla_\theta J(\theta)$ نااریب است. چه فرض‌هایی را در اثباتتان در نظر گرفته‌اید؟

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_{i,t} | s_{i,t}) \right) \left(\sum_{t=0}^T \gamma^t r(s_{i,t}, a_{i,t}) \right) \quad (2)$$

ب) تعریف می‌کنیم

$$\rho_\pi(s) = \sum_{t=0}^{\infty} \gamma^t p(s_t = s) \quad (3)$$

. نشان دهید

$$\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] = \mathbb{E}_{\rho_\pi, \pi} \left[\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) \hat{Q}(s_t, a_t) \right]$$

که در آن $\hat{Q}(s_t, a_t) = \sum_{t'=t}^{\infty} \gamma^{t'-t} r(s_{t'}, a_{t'})$.

(ب) تابع $b(s) : \mathcal{S} \rightarrow \mathbb{R}$ را بیابید که واریانس تخمینگر نااریب

$$\nabla_{\theta} J(\theta) \approx \mathbb{E}_{\rho_{\pi, \pi}} \left[\sum_{t=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (\hat{Q}(s_t, a_t) - b(s_t)) \right] \quad (۴)$$

کمینه شود.

(ج) حال با توجه به رابطه‌ای که در قسمت قبل به دست آورده‌اید، نشان دهید که اگر به جای $\hat{Q}(s_t, a_t)$ از مقدار واقعی $Q(s_t, a_t)$ استفاده کنیم، خواهیم داشت $b(s_t) = V^{\pi}(s_t)$.

سوال ۳: TRPO (نظری) (۳۵ نمره)

(الف) دو سیاست π_{θ} و $\pi_{\theta'}$ را در نظر بگیرید. فرض کنید

$$\forall s_t : |\pi_{\theta'}(a_t | s_t) - \pi_{\theta}(a_t | s_t)| \leq \epsilon. \quad (۵)$$

ثابت کنید

$$|p_{\theta'}(s_t) - p_{\theta}(s_t)| \leq 2t\epsilon.$$

اگر از قضایای کمکی استفاده می‌کنید، آن‌ها را اثبات نمایید.

(ب) در صورتی که رابطه‌ی ۵ برقرار باشد، مرتبه‌ی اختلاف

$$\mathbb{E}_{p_{\theta'}(s_t)} [E_{a_t \sim \pi_{\theta}(a_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \right]] \quad (۶)$$

و

$$\mathbb{E}_{p_{\theta}(s_t)} [E_{a_t \sim \pi_{\theta}(a_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \right]] \quad (۷)$$

را در حالت افق محدود و نامتناهی به دست آورید.

(ج) در صورت برقراری رابطه‌ی ۵، کران بالایی بر روی $D_{KL}(\pi_{\theta'}(a_t | s_t) || \pi_{\theta}(a_t | s_t))$ بیان نموده و آن را ثابت کنید.

(د) بسط تیلور مرتبه‌ی اول

$$\sum_t \mathbb{E}_{s_t \sim p_{\theta}(s_t)} [E_{a_t \sim \pi_{\theta}(a_t | s_t)} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi_{\theta}(a_t | s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \right]] \quad (۸)$$

را حول $\theta' = \theta$ بنویسید.

سوال ۴: DPG (نظری) (۲۰ نمره)

در این سوال می‌خواهیم گرادینان سیاست را برای یک سیاست قطعی^۱ $\mu_{\theta} : \mathcal{S} \rightarrow \mathcal{A}$ به دست آوریم که $A = \mathbb{R}^m$ و S یک زیر فضای فشرده از \mathbb{R}^d است. محیط، مارکوف با تابع انتقال^۲ p و توزیع حالت شروع p_1 است. احتمال حضور در حالت s' پس از انجام t کنش مطابق سیاست μ با شروع از s را با $p(s \rightarrow s', t, \mu)$ نشان می‌دهیم. تعریف می‌کنیم

$$\rho_{\mu}(s') := \int_S \sum_{t=0}^{\infty} \gamma^{t-1} p_1(s) p(s \rightarrow s', t, \mu) ds. \quad (۹)$$

تابع هدف به صورت

$$J(\mu_{\theta}) = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r(s_t, a_t) | \mu \right] = \int_S \rho_{\mu}(s) r(s, \mu_{\theta}(s)) ds = \mathbb{E}_{s \sim \rho_{\mu}} [r(s, \mu_{\theta}(s))] \quad (۱۰)$$

^۱ Deterministic
^۲ Transition

است. فرض کنید $p_1(s)$ ، $\nabla_a r(s, a)$ ، $r(s, a)$ ، $\nabla_a p(s'|s, a)$ ، $p(s'|s, a)$ و $p_1(s)$ همگی نسبت به تمام پارامترها پیوسته‌اند. همچنین فرض کنید b وجود دارد به طوری که $\sup_{a,s} r(s, a) < b$ ، $\sup_{a,s,s'} p(s'|s, a) < b$ ، $\sup_s p_1(s) < b$ ، به علاوه، L نیز وجود دارد به طوری که $\sup_{a,s} \|\nabla_a r(s, a)\| < L$ و $\sup_{a,s,s'} \|\nabla_a p(s'|s, a)\| < L$ نشان دهید.

$$\nabla_\theta J(\mu_\theta) = \int_S \rho_\mu(s) \nabla_\theta \mu_\theta(S) \nabla_a Q^{\mu_\theta}(s, a)|_{a=\mu_\theta(s)} ds = \mathbb{E}_{s \sim \rho_\mu} [\nabla_\theta \mu_\theta(S) \nabla_a Q^{\mu_\theta}(s, a)|_{a=\mu_\theta(s)}]. \quad (11)$$

راهنمایی ۱: اثبات، شباهت‌هایی به قضیه‌ی استاندارد گرادین سیاست در کتاب درس دارد.
 راهنمایی ۲: برای جا به جایی ترتیب مشتق و انتگرال و همچنین ترتیب انتگرال‌ها در اثبات، توجه نمایید که با استفاده از شرایط گفته شده می‌توان نشان داد که $V^{\mu_\theta}(s)$ و $\nabla_\theta V^{\mu_\theta}(s)$ توابع پیوسته‌ای از θ و s هستند. به علاوه، می‌توان از فشرده بودن S نتیجه گرفت که برای هر θ ، $\|\nabla_\theta V^{\mu_\theta}(s)\|$ ، $\|\nabla_a Q^{\mu_\theta}(s, a)|_{a=\mu_\theta(s)}\|$ و $\|\nabla_\theta \mu_\theta(s)\|$ توابع کران‌داری از s هستند.