



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. به اسکرین‌شات از منابع یا پاسخ افراد دیگر نمره‌ای تعلق نمی‌گیرد.
- تمام پاسخ‌های خود را در یک فایل با فرمت RL_HW#[SID]_[Fullname].zip روی کوئرا قرار دهید.
- برای ارسال هر تمرین تا ساعت ۲۳:۵۹ روز ددلاین فرصت دارید. علاوه بر آن، در هر تمرین می‌توانید تا سقف پنج روز از تأخیر مجاز باقیمانده‌ی خود استفاده کنید.

سوال ۱: گرادیان سیاست (۲۵ نمره)

- (آ) نشان دهید تخمین گرادیان با خط مبنا^۱ وابسته به حالت $(b(s_t))$ بدون بایاس است.
- (ب) هدف از اضافه کردن خط مبنا به تخمین گرادیان، کاهش واریانس تخمین است. خط مبنا بهینه (که کمترین واریانس تخمین گرادیان را ایجاد می‌کند) را به دست آورید.
- (ج) فرض کنید احتمال یک مسیر در MDP تحت سیاست π و پاداش تخفیف داده شده به صورت زیر تعریف شده‌اند.

$$\Pr_{\mu}^{\pi}(\tau) = \mu(s_0) \pi(a_0 | s_0) P(s_1 | s_0, a_0) \pi(a_1 | s_1) \dots \quad (۱)$$

$$R(\tau) := \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \quad (۲)$$

هدف از روش گرادیان سیاست بیشینه‌سازی رابطه زیر است.

$$\max_{\theta \in \Theta} V^{\pi_{\theta}}(\rho) \quad (۳)$$

که در آن

$$V^{\pi}(\rho) := \mathbb{E}_{s_0 \sim \rho} [V^{\pi}(s_0)] \quad (۴)$$

ثابت کنید:

$$\nabla V^{\pi_{\theta}}(\mu) = \mathbb{E}_{\tau \sim \Pr_{\mu}^{\pi_{\theta}}} \left[\sum_{t=0}^{\infty} \gamma^t Q^{\pi_{\theta}}(s_t, a_t) \nabla \log \pi_{\theta}(a_t | s_t) \right] \quad (۵)$$

¹Baseline

(د) برای سیاست π می‌توان یک اندازه‌گیری احتمال^۲ برای بازدید حالت‌ها به صورت زیر تعریف کرد.

$$d_{s_0}^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \Pr^\pi(s_t = s \mid s_0) \quad (۶)$$

به جای s_0 می‌توان توزیع اولیه μ را در نظر گرفت که داریم:

$$d_\mu^\pi(s) = \mathbb{E}_{s_0 \sim \mu} [d_{s_0}^\pi(s)] \quad (۷)$$

با استفاده از تعاریف فوق می‌توان نشان داد:

$$\mathbb{E}_{\tau \sim \Pr^\pi} \left[\sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^\pi} \mathbb{E}_{a \sim \pi_\theta(\cdot \mid s)} [f(s, a)] \quad (۸)$$

با توجه به نتایج بخش‌های قبل روابط (۹) و (۱۰) را ثابت کنید.

$$\nabla V^{\pi_\theta}(\mu) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta(\cdot \mid s)} [Q^{\pi_\theta}(s, a) \nabla \log \pi_\theta(a \mid s)] \quad (۹)$$

$$\nabla V^{\pi_\theta}(\mu) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta(\cdot \mid s)} [A^{\pi_\theta}(s, a) \nabla \log \pi_\theta(a \mid s)] \quad (۱۰)$$

سوال ۲: الگوریتم‌های مبتنی بر ارزش برای مسائل با فعالیت‌های پیوسته (۲۰ نمره)

بسیاری از الگوریتم‌های یادگیری تقویتی از چهارچوب policy iteration برای یافتن سیاست مناسب استفاده می‌کنند. همانطور که در جلسات کلاس دیده‌ایم این الگوریتم شامل دو مرحله policy improvement و policy evaluation می‌باشد. در الگوریتم Q-learning نیز به صورت تکرار شونده اقدام به بهبود تخمین ارزش و یافتن سیاست با روابط ذیل می‌کنیم:

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \left[r + \gamma \max_{a'} Q(s', a') \right] \quad (۱۱)$$

$$\pi(s) = \operatorname{argmax}_a Q(s, a) \quad (۱۲)$$

(آ) عنوان کنید به چه دلیل این روابط در مسائلی با فضای فعالیت‌های پیوسته قابل اجرا نیستند؟

(ب) برای حل این مشکل و حل مسائل بهینه‌سازی $\max_{a'} Q(s, a')$ و $\operatorname{argmax}_a Q(s, a)$ روش‌های متنوعی ارائه شده‌اند. با افزایش ابعاد فضای فعالیت‌ها و پیچیده شدن مسئله بهینه‌سازی روش‌های مبتنی بر نمونه برداری کارایی خود را از دست می‌دهند. با این حال در صورت آگاهی از تابع هدف (روش‌های white-box) می‌توانیم از اطلاعاتی که گرادینان تابع هدف در اختیار ما قرار می‌دهد برای تسریع بهینه‌سازی استفاده کنیم. در این قسمت به بررسی برخی روش‌های در این دسته خواهیم پرداخت.

۱. یکی از روش‌های برای یافتن بهینه سراسری تابع Q فرض یک تابع محدب با بهینه به صورت فرم بسته برای این تابع بر حسب پارامتر \mathbf{a} است. برای مثال تابع ارزش به فرم زیر را در نظر بگیرید.

$$Q_\phi(\mathbf{s}, \mathbf{a}) = -\frac{1}{2} (\mathbf{a} - \mu_\phi(\mathbf{s}))^T P_\phi(\mathbf{s}) (\mathbf{a} - \mu_\phi(\mathbf{s})) + V_\phi(\mathbf{s})$$

در این حالت مقادیر $\operatorname{argmax}_a Q(s, a)$ و $\max_a Q(s, a)$ را به دست بیاورید و بیان کنید در نظر گرفتن چنین فرم‌های ساده‌ای برای تابع ارزش چه مزایا و معایبی همراه خواهد داشت.

²Probability Measure

۲. رویکرد دیگر برای حل مسئله بهینه‌سازی یادگیری بهینه‌سازی است. به این ترتیب که مدل پارامتری در طول آموزش اقدام به یادگیری حل مسئله بهینه‌سازی می‌کند. این رویکرد در روش deep deterministic policy gradient مورد استفاده قرار گرفته شده است و مدل پارامتری برای حل مسئله بهینه‌سازی یک شبکه عصبی در نظر گرفته شده است.

آ. در ابتدا الگوریتم، تابع زیان و معماری الگوریتم **DDPG** را بیان کنید.

ب. نحوه استفاده از شبکه‌های actor و critic در این الگوریتم را با الگوریتم REINFORCE مقایسه کنید. عبور گرادین از شبکه critic برای آموزش شبکه actor چه مزیت‌هایی به همراه دارد؟

ج. تابع هدف در بیشتر الگوریتم‌های یادگیری تقویتی بیشینه کردن امید مجموع پاداش دریافتی است که در ادامه نمایش داده شده است. در این رابطه شبکه actor با μ نمایش داده شده که با مجموعه پارامترهای θ مدل شده است.

$$\begin{aligned} J(\mu_\theta) &= \int_S \rho^\mu(s) r(s, \mu_\theta(s)) ds \\ &= \mathbb{E}_{s \sim \rho^\mu} [r(s, \mu_\theta(s))] \end{aligned}$$

در جلسات پیشین درس دیدیم که بیشینه کردن این تابع هدف معادل با بیشینه کردن تابع ارزش $Q(s, a)$ یا $V(s)$ است. با توجه به این اطلاعات و پاسخ قسمت قبلی ثابت کنید که گرادین تابع هدف برای پارامترهای شبکه actor به شکل زیر به دست خواهد آمد.

$$\begin{aligned} \nabla_\theta J(\mu_\theta) &= \int_S \rho^\mu(s) \nabla_\theta \mu_\theta(s) \nabla_a Q^\mu(s, a) \Big|_{a=\mu_\theta(s)} ds \\ &= \mathbb{E}_{s \sim \rho^\mu} \left[\nabla_\theta \mu_\theta(s) \nabla_a Q^\mu(s, a) \Big|_{a=\mu_\theta(s)} \right] \end{aligned}$$

برای این منظور می‌توانید گام‌های زیر را دنبال کنید.

• در ابتدا با استفاده از تعریف تابع ارزش به شکل زیر

$$r(s, \mu_\theta(s)) + \int_S \gamma p(s' | s, \mu_\theta(s)) V^{\mu_\theta}(s') ds'$$

و فرض پیوستگی توابع $p(s' | s, a)$, $\mu_\theta(s)$, $V^{\mu_\theta}(s)$ و مشتق آنها نسبت به θ که به دنبال آن امکان جابه‌جایی عملگرهای گرادین و انتگرال را خواهیم داشت، به یک رابطه بازگشتی برای گرادین تابع $V^{\mu_\theta}(s)$ بر حسب $V^{\mu_\theta}(s')$ برسید.

• در گام بعدی با جایگذاری متوالی رابطه بازگشتی به دست آماده به صورت حدی و با فرض محدود بودن نرم تابع ارزش و finite horizon بودن مسئله، به یک فرم بسته برای گرادین تابع ارزش بر حسب پارامتر θ می‌رسیم.

• در گام آخر از این گرادین بر حسب S امیدریاضی می‌گیریم.

سوال ۳: روش بهینه‌سازی جدید با تغییر Trust Region (۳۰ نمره)

در این سوال سعی می‌کنیم تا با تغییر Trust Region الگوریتم بهینه‌سازی TRPO یک روش جدید و البته قابل اتکاتر برای به دست آوردن سیاست بهینه به دست آوریم. سعی می‌کنیم تا گام به گام به سمت حل مسئله پیش برویم. در این مسئله فضای \mathcal{X} را برای سادگی یک مجموعه بسته و کراندار اندازه‌گیری‌پذیر در نظر می‌گیریم به صورتی که $\mathcal{P}(\mathcal{X})$ مجموعه‌ی تمام توزیع‌های احتمال موجود روی \mathcal{X} است. همچنین با داشتن یک سیاست π روی یک MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \rho, \gamma)$ ارزش هر وضعیت یا حرکت-وضعیت را با $V^\pi(s)$ و $Q^\pi(s, a)$ نشان دهیم تابع مزیت را به صورت $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$ می‌توان تعریف کرد. در ابتدا با چند تعریف شروع می‌کنیم:

(آ) هزینه انتقال: تابع هزینه‌ی انتقال را به صورت یک تابع $c: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ نشان می‌دهیم.

(ب) همچنین برای هر دو توزیع $\mu, \nu \in \mathcal{P}(\mathcal{X})$ مجموعه توزیع‌های توأمی که حاشیه‌های μ و ν دارند را به صورت زیر تعریف می‌کنیم:

$$\Gamma(\mu, \nu) = \{\gamma \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) : \gamma(A \times \mathcal{X}) = \mu(A), \gamma(\mathcal{X} \times B) = \nu(B)\} \quad (۱۳)$$

(ج) تابع اختلاف انتقال بهینه: برای هر دو توزیع $\mu, \nu \in \mathcal{P}(\mathcal{X})$ تابع بهینه‌ی هزینه‌ی انتقال را به صورت زیر تعریف می‌کنیم:

$$C(\mu, \nu) = \min_{\gamma \in \Gamma(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{X}} c(x, x') d\gamma(x, x') \quad (۱۴)$$

(د) **تابع هدف:** تابع هدف که پاداش کاهش‌یابنده را امید ریاضی محاسبه می‌کند به صورت زیر تعریف می‌کنیم:

$$J(\pi) = \mathbb{E}_{\rho, \pi} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \quad (15)$$

که در آن ρ توزیع اولیه روی وضعیت شروع است.

حال می‌توانیم گام به گام به سمت حل مسئله حرکت کنیم:

(آ) ابتدا نشان دهید در صورتی که $\pi, \tilde{\pi} \in \Pi$ دو سیاست دلخواه باشند، خواهیم داشت:

$$J(\tilde{\pi}) = J(\pi) + \int_S \int_{\mathcal{A}} A^\pi(s, a) d\tilde{\pi}(a|s) d\rho_{\tilde{\pi}}(s) \quad (16)$$

که در آن $\rho_{\tilde{\pi}}(s)$ ، توزیع وضعیت آینده کاهشی، به صورت زیر است:

$$\rho_{\tilde{\pi}}(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}[s_t = s | \tilde{\pi}, \rho] \quad (17)$$

(ب) اکنون به جای استفاده از عبارت بالا به صورت مستقیم (که محاسبه‌ی آن هزینه‌ی زیادی در بر دارد) از شکل تغییر یافته‌ی زیر استفاده می‌کنیم:

$$L_\pi(\tilde{\pi}) = J(\pi) + \int_S \int_{\mathcal{A}} A^\pi(s, a) d\tilde{\pi}(a|s) d\rho_\pi(s) \quad (18)$$

حال استفاده از این تابع هدف برای بیشینه سازی معقول‌تر به نظر می‌رسد، در صورتی که تغییر سیاست به گونه‌ای نباشد که توزیع وضعیت آینده کاهشی سیاست آتی با سیاست فعلی تفاوت قابل توجهی داشته باشد. که این فرض بنا بر شرایط مسئله فرض معقولی است و برای سادگی می‌توان از آن در ادامه‌ی کار استفاده کرد.

(ج) حال در ادامه با توجه به معادله ۱۸ می‌توانیم مسئله‌ی بهینه‌سازی را با توجه به بیشینه کردن عبارت دوم در آن، به شکل زیر تعریف کرد:

$$\begin{aligned} & \sup_{\tilde{\pi} \in \Pi} \int_S \int_{\mathcal{A}} A^\pi(s, a) d\tilde{\pi}(a|s) d\rho_\pi(s) \\ & \text{s.t. } \tilde{\pi} \in \mathcal{T}_\epsilon := \left\{ \tilde{\pi} \in \Pi : \int_S C(\pi(\cdot|s), \tilde{\pi}(\cdot|s)) d\rho_\pi(s) \leq \epsilon \right\} \end{aligned} \quad (19)$$

همانطور که در مسئله‌ی بهینه‌سازی بالا مشاهده می‌شود \mathcal{T}_ϵ در واقع یک تعریف جدید از Trust Region مسئله است.

در ادامه باید گفت که در مسئله‌ی بهینه‌سازی جدید معرفی شده محاسبه‌ی تابع اختلاف انتقال بهینه‌ی مورد نظر باز هم از لحاظ محاسباتی دشواری‌های قابل توجهی دارد. اما برای حل آن از می‌توان از تکنیک تبدیل مسئله به دوگان آن استفاده کرد که در ادامه خواهیم دید تا چه حدی می‌تواند حل مسئله را بهبود بخشد. اما پیش از آغاز به صورتی گذرا تکنیک دوگان‌گیری از یک مسئله‌ی بهینه‌سازی را بیان می‌کنیم:

(آ) فرض کنید یک مسئله‌ی بهینه‌سازی به صورت زیر داده شده‌است:

$$\begin{aligned} & \min_{x \in X} f(x) \\ & \text{s.t. } g_i(x) \leq 0 \quad \forall i \in \{1, 2, \dots, k\} \end{aligned} \quad (20)$$

حال تابع لاگرانژی آن را به صورت زیر تعریف می‌کنیم:

$$L(x, \lambda_1, \lambda_2, \dots, \lambda_k) = f(x) + \lambda_1 g_1(x) + \lambda_2 g_2(x) + \dots + \lambda_k g_k(x) \quad (21)$$

همچنین تابع دوگان آن را به صورت زیر تعریف می‌کنیم:

$$g(\lambda_1, \dots, \lambda_k) = \begin{cases} \inf_{x \in X} L(x, \lambda_1, \dots, \lambda_k) & \text{if } \lambda_i \geq 0, \forall i \\ -\infty & \text{o.w.} \end{cases} \quad (22)$$

حال مسئله‌ی دوگان را به صورت زیر تعریف می‌کنیم:

$$\begin{aligned} & \max \quad g(\lambda_1, \dots, \lambda_k) \\ & \text{s.t. } \lambda_i \geq 0 \end{aligned} \quad (23)$$

(ب) به سادگی قابل مشاهده است که اگر پاسخ مسأله‌ی ۲۰ را با f^* و پاسخ مسأله‌ی ۲۳ را با g^* نمایش دهیم، آنگاه با توجه به تعریف بالا خواهیم داشت: $g^* \leq f^*$. که این نتیجه به دوگانی ضعیف معروف است.

(ج) همچنین می‌توان نشان داد در صورتی که بر روی مجموعه‌ی $feasible$ مسأله‌ی ۲۰ شرایطی برقرار باشد آنگاه دوگانی قوی است یعنی: $f^* = g^*$. از دوگانی قوی در حل مسأله‌ی بهینه‌سازی گذشته استفاده خواهد شد که البته با توجه به اینکه این موضوع از حوزه‌ی این تمرین خارج است برقراری آن را در این سوال فرض می‌کنیم.

اما اکنون به سراغ مسأله‌ی بهینه‌سازی ۱۹ می‌رویم:

(آ) پیش از ادامه برای به دست آوردن دوگان مناسب دو فرض که به صورت معمول برقرار است را در نظر می‌گیریم:

- فضای وضعیت S و کنش A مجموعه‌هایی بسته و کران‌دار هستند. به علاوه تابع پاداش r یک تابع پیوسته و امید ریاضی هر تابع پیوسته w روی S یک تابع پیوسته است.
- برای هر سیاست π تابع مزیت A^π یک تابع پیوسته است همچنین تابع هزینه انتقال c یک تابع پیوسته است که هزینه‌ی هر کنش با خودش برابر صفر است: $c(a, a) = 0$.

(ب) حال فرم دوگان مسأله‌ی ۱۹ را به دست آورید. سپس یک کران بالا برای آن ارائه کنید به صورتی که به π وابستگی نداشته باشد. (راهنمایی: یک کران بالا برای فرم دوگان را به صورت زیر می‌توان نوشت:

$$\begin{aligned} \min \quad & \lambda \epsilon + \int_S \int_A \max_{a' \in A} \{A^\pi(s, a') - \lambda c(a, a')\} d\pi(a|s) d\rho_\pi(s) \\ \text{s.t.} \quad & \lambda \geq 0 \end{aligned} \quad (24)$$

. همچنین برای به دست آوردن این کران می‌توانید از قضیه Kantorovich در مورد اختلاف انتقال بهینه استفاده کنید).

(ج) قضیه Kantorovich: می‌توان نشان داد که پاسخ اختلاف انتقال بهینه ۱۴ را می‌توان به صورت زیر نوشت:

$$C(\mu, \nu) = \sup_{\phi, \psi} \left\{ \int_{\mathcal{X}} \phi(x) d\mu(x) + \int_{\mathcal{X}} \psi(x') d\nu(x') \right\} \quad (25)$$

$$\text{subject to } \phi(x) + \psi(x') \leq c(x, x')$$

با عرض تبریک سوال تا اینجا به پایان رسید! (:

در ادامه به بیان روش به دست آوردن سیاست بهینه مبتنی بر این کران بالای به دست آمده می‌پردازیم.

می‌توان نشان داد که مسأله‌ی ۲۴ یک مسأله‌ی بهینه‌سازی محدب است که با توجه به تک پارامتره بودن می‌توان آن را با استفاده از حل‌کننده‌های مرسوم به سرانجام رساند. حال پس از پیدا کردن λ^* بهینه برای این مسأله می‌توان با استفاده از آن سیاست (تقریباً) بهینه‌ی مربوط به مسأله‌ی اول ۱۹ را با تکنیک‌هایی به دست آورد که به جهت جلوگیری از اطناب از ذکر آن خودداری می‌کنیم. حال تنها مسأله‌ای که باقی می‌ماند تخمین زدن تابع مزیت و محاسبه‌ی انتگرال‌هاست که می‌توان آن‌ها را با استفاده از یک شبکه‌ی عصبی یا روش های MC یا TD به دست آورد.

سوال ۴: پیاده‌سازی (۲۵ نمره)

هدف این بخش از تمرین پیاده‌سازی دو الگوریتم PPO و DDPG و مقایسه نتایج این دو الگوریتم در محیط Pendulum-v1 از کتابخانه gym است. با استفاده از نوت‌بوک داده شده این دو الگوریتم را پیاده‌سازی کنید. برای مقایسه، نمودارهای هزینه شبکه‌های actor و critic و نمودار پاداش در طول اپیزودها را بر هر دو الگوریتم رسم کنید. (برای کاهش نوسانات شدید نمودار پاداش در طول اپیزود و مقایسه بهتر نمودارها می‌توانید با روش پنجره لغزان^۳ میانگین بگیرید).

³Sliding Window