



## یادگیری تقویتی

نیم سال بهار ۱۴۰۱-۴۰۲

اساتید: دکتر رهبان، آقای حسنی

## الگوریتم‌های گرادیان سیاست

تمرین سری دوم

زمان تحویل: ۱۸ فروردین

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<sup>۱</sup> وابسته به حالت  $(b(s_t))$  بدون بایاس است.

پاسخ:

$$\begin{aligned}\nabla_{\theta} J(\theta) &= E_{\tau \sim p_{\theta}(\tau)} [(\nabla_{\theta} \log p_{\theta}(\tau)) (r(\tau) - b(s_t))] \\ &= E_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) r(\tau) \right] - E_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) b(s_t) \right]\end{aligned}$$

با استفاده از خاصیت خطی بودن امید ریاضی به راحتی می‌توان نشان داد:

$$E_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t=1}^T f(s_t, \mathbf{a}_t) \right] = \sum_{t=1}^T E_{(s_t, \mathbf{a}_t) \sim p_{\theta}(s_t, \mathbf{a}_t)} [f(s_t, \mathbf{a}_t)]$$

که در رابطه فوق منظور از  $p_{\theta}(s_t, \mathbf{a}_t)$  یک تابع توزیع حاشیه‌ای از متغیر تصادفی  $\tau$  است. حال با دانستن این موضوع داریم:

$$E_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) b(s_t) \right] = \sum_{t=1}^T E_{(s_t, \mathbf{a}_t) \sim p_{\theta}(s_t, \mathbf{a}_t)} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) b(s_t)]$$

همچنین می‌دانیم:

$$E_Y [E_{X|Y} [f(X, Y)]] = E_{X,Y} [f(X, Y)]$$

با توجه به نکته بالا می‌توان عبارت مربوط به بایاس را به صورت زیر نوشت:

$$\begin{aligned}\sum_{t=1}^T E_{s_t} [E_{a_t|s_t} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) b(s_t)]] &= \sum_{t=1}^T E_{s_t} \left[ \int \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) b(s_t) d\mathbf{a}_t \right] \\ &= \sum_{t=1}^T E_{s_t} \left[ \int \nabla_{\theta} \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) b(s_t) d\mathbf{a}_t \right] \\ &= \sum_{t=1}^T E_{s_t} \left[ b(s_t) \nabla_{\theta} \int \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) d\mathbf{a}_t \right] \\ &= \sum_{t=1}^T E_{s_t} [b(s_t) \nabla_{\theta}(1)] = 0\end{aligned}$$

<sup>1</sup>Baseline

در نهایت داریم:

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [(\nabla_{\theta} \log p_{\theta}(\tau)) (r(\tau) - b(s_t))] = E_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) r(\tau) \right]$$

بنابراین با اضافه کردن خط مبنا وابسته به حالت همچنان تخمین گرادیان بدون بایاس است. روند فوق را در حالتی که  $T \rightarrow \infty$  و یا ضریب تخفیف در عبارت  $r(\tau)$  وجود دارد، می توان به طور مشابه تکرار کرد و به همین نتیجه رسید.

(ب) هدف از اضافه کردن خط مبنا به تخمین گرادیان، کاهش واریانس تخمین است. خط مبنا بهینه (که کمترین واریانس تخمین گرادیان را ایجاد می کند) را به دست آورید.

**پاسخ:**

$$\text{Var}[x] = E[x^2] - E[x]^2$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b)]$$

$$\text{Var} = E_{\tau \sim p_{\theta}(\tau)} [(\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b))^2] - E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b)]^2$$

همانطور که در قسمت (آ) دیدیم می توان گفت امید ریاضی تخمین گرادیان نسبت به خط مبنا بدون بایاس است. (در عبارت مربوط به معذور امید ریاضی ترم خط مبنا حذف می شود و مستقل از  $b$  است) حال برای یافتن بهترین خط مبنا می که واریانس گرادیان را کاهش می دهد می توان از واریانس نسبت به بایاس مشتق گرفت.

$$\begin{aligned} \frac{d \text{Var}}{db} &= \frac{d}{db} E[g(\tau)^2 (r(\tau) - b)^2] \\ &= \frac{d}{db} (E[g(\tau)^2 r(\tau)^2] - 2E[g(\tau)^2 r(\tau) b] + b^2 E[g(\tau)^2]) \\ &= -2E[g(\tau)^2 r(\tau)] + 2bE[g(\tau)^2] \\ &= 0 \end{aligned}$$

با حل معادله نسبت به  $b$  داریم:

$$b = \frac{E[g(\tau)^2 r(\tau)]}{E[g(\tau)^2]}$$

(ج) فرض کنید احتمال یک مسیر در MDP تحت سیاست  $\pi$  و پاداش تخفیف داده شده به صورت زیر تعریف شده اند.

$$\text{Pr}_{\mu}^{\pi}(\tau) = \mu(s_0) \pi(a_0 | s_0) P(s_1 | s_0, a_0) \pi(a_1 | s_1) \dots \quad (1)$$

$$R(\tau) := \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \quad (2)$$

هدف از روش گرادیان سیاست بیشینه سازی رابطه زیر است.

$$\max_{\theta \in \Theta} V^{\pi_{\theta}}(\rho) \quad (3)$$

که در آن

$$V^{\pi}(\rho) := \mathbb{E}_{s_0 \sim \rho} [V^{\pi}(s_0)] \quad (4)$$

ثابت کنید:

$$\nabla V^{\pi_{\theta}}(\mu) = \mathbb{E}_{\tau \sim \text{Pr}_{\mu}^{\pi_{\theta}}} \left[ \sum_{t=0}^{\infty} \gamma^t Q^{\pi_{\theta}}(s_t, a_t) \nabla \log \pi_{\theta}(a_t | s_t) \right] \quad (5)$$

برای هر حالت شروع  $s_0$  داریم:

$$\begin{aligned}
 \nabla V^{\pi_\theta}(s_0) &= \\
 &= \nabla \sum_{a_0} \pi_\theta(a_0 | s_0) Q^{\pi_\theta}(s_0, a_0) \\
 &= \sum_{a_0} (\nabla \pi_\theta(a_0 | s_0)) Q^{\pi_\theta}(s_0, a_0) + \sum_{a_0} \pi_\theta(a_0 | s_0) \nabla Q^{\pi_\theta}(s_0, a_0) \\
 &= \sum_{a_0} \pi_\theta(a_0 | s_0) (\nabla \log \pi_\theta(a_0 | s_0)) Q^{\pi_\theta}(s_0, a_0) \\
 &\quad + \sum_{a_0} \pi_\theta(a_0 | s_0) \nabla \left( r(s_0, a_0) + \gamma \sum_{s_1} P(s_1 | s_0, a_0) V^{\pi_\theta}(s_1) \right) \\
 &= \sum_{a_0} \pi_\theta(a_0 | s_0) (\nabla \log \pi_\theta(a_0 | s_0)) Q^{\pi_\theta}(s_0, a_0) + \gamma \sum_{a_0, s_1} \pi_\theta(a_0 | s_0) P(s_1 | s_0, a_0) \nabla V^{\pi_\theta}(s_1) \\
 &= \mathbb{E}_{\tau \sim \text{Pr}_{s_0}^{\pi_\theta}} [Q^{\pi_\theta}(s_0, a_0) \nabla \log \pi_\theta(a_0 | s_0)] + \gamma \mathbb{E}_{\tau \sim \text{Pr}_{s_0}^{\pi_\theta}} [\nabla V^{\pi_\theta}(s_1)]
 \end{aligned}$$

با استفاده از خاصیت خطی بودن امید ریاضی

$$\begin{aligned}
 \nabla V^{\pi_\theta}(\mu) &= \\
 &= \mathbb{E}_{\tau \sim \text{Pr}_{\mu}^{\pi_\theta}} [Q^{\pi_\theta}(s_0, a_0) \nabla \log \pi_\theta(a_0 | s_0)] + \gamma \mathbb{E}_{\tau \sim \text{Pr}_{\mu}^{\pi_\theta}} [\nabla V^{\pi_\theta}(s_1)] \\
 &= \mathbb{E}_{\tau \sim \text{Pr}_{\mu}^{\pi_\theta}} [Q^{\pi_\theta}(s_0, a_0) \nabla \log \pi_\theta(a_0 | s_0)] + \gamma \mathbb{E}_{\tau \sim \text{Pr}_{\mu}^{\pi_\theta}} [Q^{\pi_\theta}(s_1, a_1) \nabla \log \pi_\theta(a_1 | s_1)] + \dots
 \end{aligned}$$

که عبارت آخر به صورت بازگشتی به دست آمده است و برابر با عبارتی است که به دنبال آن بودیم.

(د) برای سیاست  $\pi$  می‌توان یک اندازه‌گیری احتمال<sup>۲</sup> برای بازدید حالت‌ها به صورت زیر تعریف کرد.

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

به جای  $s_0$  می‌توان توزیع اولیه  $\mu$  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:

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

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

$$\nabla V^{\pi_\theta}(\mu) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^\pi} \mathbb{E}_{a \sim \pi_\theta(\cdot | s)} [Q^{\pi_\theta}(s, a) \nabla \log \pi_\theta(a | s)] \quad (۹)$$

$$\nabla V^{\pi_\theta}(\mu) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^\pi} \mathbb{E}_{a \sim \pi_\theta(\cdot | s)} [A^{\pi_\theta}(s, a) \nabla \log \pi_\theta(a | s)] \quad (۱۰)$$

<sup>۲</sup>Probability Measure

با توجه به توضیحات داده شده در صورت سوال برای اثبات رابطه ۹ کافی است به جای  $f(s, a)$  عبارت  $Q^{\pi_\theta}(s, a) \nabla \log \pi_\theta(a | s)$  را قرار دهیم تا نتیجه به دست آید.

روند اثبات رابطه ۱۰ بسیار مشابه با اثبات بخش (آ) است.

## 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:

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \left[ r + \gamma \max_{a'} Q(s', a') \right] \quad (11)$$

$$\pi(s) = \operatorname{argmax}_a Q(s, a) \quad (12)$$

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

**پاسخ:**

پیوسته بودن طیف فعالیت‌های ممکن، امکان بررسی تمام فعالیت‌ها برای یافتن فعالیت مناسب را سلب می‌کند و به این ترتیب یافتن max و argmax تابع ارزش دیگر به صورت بدیهی امکان پذیر نیست. برای حل این مسئله بهینه‌سازی و یافتن یک بهینه محلی خوب در زمان مناسب روش‌های مختلفی ارائه شده است که در ادامه سوال مشاهده کرده‌اید.

(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 گرادیان 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  $a$  است. برای مثال تابع ارزش به فرم زیر را در نظر بگیرید.

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

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.

**پاسخ:**

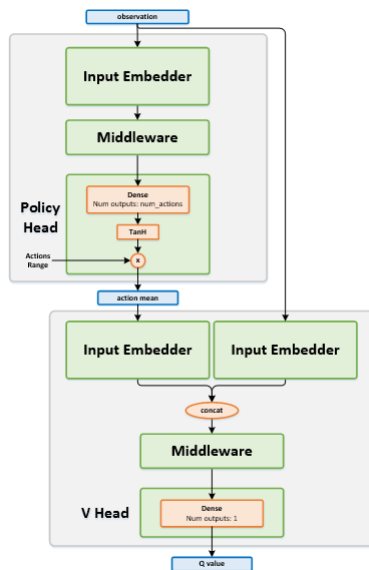
به دلیل فرم محدب تابع ارزش بر حسب پارامتر  $a$  می‌توانیم بهینه سراسری را به فرم بسته به شکل زیر محاسبه کنیم. با مشتق‌گیری از تابع ارزش بر حسب  $a$  و برابر قرار دادن با صفر خواهیم داشت (در رابطه اولیه بیان شده  $x$  یک بردار و  $A$  یک ماتریس مربعی است):

$$\begin{aligned} \frac{\partial x^T A x}{\partial x} &= 2Ax \\ \nabla_a Q(s, a) &= -P_\phi(s)(a - \mu_\phi(s)) = 0 \\ P_\phi(s)^{-1} (a - \mu_\phi(s)) &= 0 \rightarrow a^* = \mu_\phi(s) \\ Q_\phi(s, a^*) &= -\frac{1}{2} (\mu_\phi(s) - \mu_\phi(s))^T P_\phi(s) (\mu_\phi(s) - \mu_\phi(s)) + V_\phi(s) = V_\phi(s) \end{aligned}$$

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.  
**پاسخ:**

این الگوریتم در واقع توسعه روش‌های مبتنی بر ارزش برای مسائل با طیف پیوسته از فعالیت‌ها است. ایده کلی این روش ارائه معماری actor-critic است که در آن شبکه actor وظیفه یادگیری عملگر  $\arg\max_a Q(s, a)$  را بر عهده دارد و مابقی اجزا همانند روش DQN قابل انجام خواهد بود. نحوه قرار گیری دو شبکه در شکل قابل مشاهده است.



به این ترتیب تابع هدف شبکه actor

$$J(\theta) = \mathbb{E} \left[ Q(s, a) |_{s=s_t, a_t=\mu(s_t)} \right]$$

و شبکه critic

$$J(\theta) = -\mathbb{E} \left[ \left( Q(s_t, a_t | \theta^Q) - y_t \right)^2 \right]$$

قابل بیان است.

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

**پاسخ:**

شبکه critic در الگوریتم REINFORCE به عنوان یک baseline قابل یادگیری برای کاهش واریانس گرادینان شبکه actor مورد استفاده قرار می‌گیرد در حالی که در الگوریتم DDPG شبکه actor وظیفه بیشینه کردن خروجی شبکه critic را بر عهده دارد. این موضوع به صورت مستقیم در توابع زیان و گرادینان توابع زیان دو الگوریتم برای شبکه actor مشهود است. در الگوریتم REINFORCE به شکل زیر:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{s, a \sim \rho^{\mu}} \left[ \nabla_{\theta} \log \pi_{\theta}(\mathbf{a} | \mathbf{s}) \hat{A}^{\pi}(\mathbf{s}, \mathbf{a}) \right]$$

که تابع مزیت به شکل زیر تعریف می‌شود:

$$\hat{A}^{\pi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \hat{V}_{\phi}^{\pi}(\mathbf{s}') - \hat{V}_{\phi}^{\pi}(\mathbf{s})$$

و در الگوریتم DDPG به این شکل:

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

خواهند بود.

ج. تابع هدف در بیشتر الگوریتم‌های یادگیری تقویتی بیشینه کردن امید مجموع پاداش دریافتی است که در ذیل نمایش داده شده است. در این رابطه شبکه actor با  $\mu$  نمایش داده شده که با مجموعه پارامترهای  $\theta$  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 گرادیان 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  $\theta$  که به دنبال آن امکان جابه‌جایی عملگرهای گرادیان و انتگرال را خواهیم داشت، به یک رابطه بازگشتی برای گرادیان 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 گرادیان تابع ارزش بر حسب پرامتر  $\theta$  می‌رسیم.

- در گام آخر از این گرادیان بر حسب  $S$  امیدریاضی می‌گیریم.

**پاسخ:**

در مرحله اول طبق روابط عنوان شده باید به یک رابطه بازگشتی برسیم:

$$\begin{aligned}\nabla_{\theta} V^{\mu_{\theta}}(s) &= \nabla_{\theta} Q^{\mu_{\theta}}(s, \mu_{\theta}(s)) \\ &= \nabla_{\theta} \left( r(s, \mu_{\theta}(s)) + \int_S \gamma p(s' | s, \mu_{\theta}(s)) V^{\mu_{\theta}}(s') ds' \right) \\ &= \nabla_{\theta} \mu_{\theta}(s) \nabla_a r(s, a) \Big|_{a=\mu_{\theta}(s)} + \nabla_{\theta} \int_S \gamma p(s' | s, \mu_{\theta}(s)) V^{\mu_{\theta}}(s') ds' \\ &= \nabla_{\theta} \mu_{\theta}(s) \nabla_a r(s, a) \Big|_{a=\mu_{\theta}(s)} \\ &\quad + \int_S \gamma \left( p(s' | s, \mu_{\theta}(s)) \nabla_{\theta} V^{\mu_{\theta}}(s') + \nabla_{\theta} \mu_{\theta}(s) \nabla_a p(s' | s, a) \Big|_{a=\mu_{\theta}(s)} V^{\mu_{\theta}}(s') \right) ds' \quad (1) \\ &= \nabla_{\theta} \mu_{\theta}(s) \nabla_a \left( r(s, a) + \int_S \gamma p(s' | s, a) V^{\mu_{\theta}}(s') ds' \right) \Big|_{a=\mu_{\theta}(s)} \\ &\quad + \int_S \gamma p(s' | s, \mu_{\theta}(s)) \nabla_{\theta} V^{\mu_{\theta}}(s') ds' \\ &= \nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\mu_{\theta}}(s, a) \Big|_{a=\mu_{\theta}(s)} + \int_S \gamma p(s \rightarrow s', 1, \mu_{\theta}) \nabla_{\theta} V^{\mu_{\theta}}(s') ds'.\end{aligned}$$

سپس به صورت تکرار شونده این رابطه بازگشتی را جایگذاری می‌کنیم:

$$\begin{aligned}
&= \nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\mu_{\theta}}(s, a)|_{a=\mu_{\theta}(s)} \\
&\quad + \int_S \gamma p(s \rightarrow s', 1, \mu_{\theta}) \nabla_{\theta} \mu_{\theta}(s') \nabla_a Q^{\mu_{\theta}}(s', a)|_{a=\mu_{\theta}(s')} ds' \\
&\quad + \int_S \gamma p(s \rightarrow s', 1, \mu_{\theta}) \int_S \gamma p(s' \rightarrow s'', 1, \mu_{\theta}) \nabla_{\theta} V^{\mu_{\theta}}(s'') ds'' ds' \\
&= \nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\mu_{\theta}}(s, a)|_{a=\mu_{\theta}(s)} \\
&\quad + \int_S \gamma p(s \rightarrow s', 1, \mu_{\theta}) \nabla_{\theta} \mu_{\theta}(s') \nabla_a Q^{\mu_{\theta}}(s', a)|_{a=\mu_{\theta}(s')} ds' \\
&\quad + \int_S \gamma^2 p(s \rightarrow s', 2, \mu_{\theta}) \nabla_{\theta} V^{\mu_{\theta}}(s') ds' \tag{2}
\end{aligned}$$

⋮

$$= \int_S \sum_{t=0}^{\infty} \gamma^t p(s \rightarrow s', t, \mu_{\theta}) \nabla_{\theta} \mu_{\theta}(s') \nabla_a Q^{\mu_{\theta}}(s', a)|_{a=\mu_{\theta}(s')} ds'.$$

و در آخر نسبت به S امید ریاضی می‌گیریم:

$$\begin{aligned}
\nabla_{\theta} J(\mu_{\theta}) &= \nabla_{\theta} \int_S p_1(s) V^{\mu_{\theta}}(s) ds \\
&= \int_S p_1(s) \nabla_{\theta} V^{\mu_{\theta}}(s) ds \\
&= \int_S \int_S \sum_{t=0}^{\infty} \gamma^t p_1(s) p(s \rightarrow s', t, \mu_{\theta}) \nabla_{\theta} \mu_{\theta}(s') \nabla_a Q^{\mu_{\theta}}(s', a)|_{a=\mu_{\theta}(s')} ds' ds \\
&= \int_S \rho^{\mu_{\theta}}(s) \nabla_{\theta} \mu_{\theta}(s) \nabla_a Q^{\mu_{\theta}}(s, a)|_{a=\mu_{\theta}(s)} ds, \tag{3}
\end{aligned}$$

### 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 مسأله پیش برویم. در این مسأله 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  $\pi$  روی یک MDP  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \rho, \gamma)$  ارزش عر 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  $\mu$  و  $\nu$  دارند را به صورت زیر تعریف می‌کنیم:

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

(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:

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

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

$$J(\pi) = \mathbb{E}_{\rho, \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \tag{۱۵}$$

که در آن  $\rho$  توزیع اولیه روی وضعیت شروع است.

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

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

$$J(\tilde{\pi}) = J(\pi) + \int_S \int_{\mathcal{A}} A^\pi(s, a) d\tilde{\pi}(a|s) d\rho_\pi(s) \quad (۱۶)$$

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

$$\rho_\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}[s_t = s | \pi, \rho] \quad (۱۷)$$

**پاسخ:** به صورت زیر می‌نویسیم:

$$\begin{aligned} V^{\tilde{\pi}}(s) &= \mathbb{E}_{\tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s \right] = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\tilde{\pi}} [r(s_t, a_t) + V^\pi(s_t) - V^\pi(s_t)] \\ &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\tilde{\pi}} [r(s_t, a_t) + \gamma V^\pi(s_{t+1}) - V^\pi(s_t)] + V^\pi(s) \\ &= \mathbb{E}_{a \sim \tilde{\pi}(\cdot|s), s' \sim P(\cdot|s, a)} [Q^\pi(s', a) - V^\pi(s')] + V^\pi(s) \\ &= \mathbb{E}_{a \sim \tilde{\pi}(\cdot|s), s' \sim P(\cdot|s, a)} [A^\pi(s', a)] + V^\pi(s) \end{aligned}$$

با گرفتن امید ریاضی از دو سمت روی تمامی وضعیت‌ها به معادله‌ی مورد نظر می‌رسیم.  $\square$

(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:

$$L_\pi(\tilde{\pi}) = J(\pi) + \int_S \int_{\mathcal{A}} A^\pi(s, a) d\tilde{\pi}(a|s) d\rho_\pi(s) \quad (۱۸)$$

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 مسأله فرض معقولی است و برای سادگی می‌توان از آن در ادامه‌ی کار استفاده کرد.

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

$$\begin{aligned} \sup_{\tilde{\pi} \in \Pi} \quad & \int_S \int_{\mathcal{A}} A^\pi(s, a) d\tilde{\pi}(a|s) d\rho_\pi(s) \\ \text{s.t.} \quad & \tilde{\pi} \in \mathcal{T}_\epsilon := \left\{ \tilde{\pi} \in \Pi : \int_S C(\pi(\cdot|s), \tilde{\pi}(\cdot|s)) d\rho_\pi(s) \leq \epsilon \right\} \end{aligned} \quad (۱۹)$$

همانطور که در مسأله‌ی بهینه‌سازی بالا مشاهده می‌شود  $\mathcal{T}_\epsilon$  در واقع یک تعریف جدید از Trust Region مسأله است.

در ادامه باید گفت که در مسأله‌ی 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 مسأله به دوگان آن استفاده کرد که در ادامه خواهیم دید تا چه حدی می‌تواند حل مسأله را بهبود بخشد. اما پیش از آغاز به صورتی گذرا تکنیک دوگان‌گیری از یک مسأله‌ی بهینه‌سازی را بیان می‌کنیم:

(آ) فرض کنید یک مسأله‌ی بهینه‌سازی به صورت زیر داده شده است:

$$\begin{aligned} \min_{x \in X} \quad & f(x) \\ \text{s.t.} \quad & g_i(x) \leq 0 \quad \forall i \in \{1, 2, \dots, k\} \end{aligned} \quad (۲۰)$$

حال تابع لاگرانژیان آن را به صورت زیر تعریف می‌کنیم:

$$L(x, \lambda_1, \lambda_2, \dots, \lambda_k) = f(x) + \lambda_1 g_1(x) + \lambda_2 g_2(x) + \dots + \lambda_k g_k(x) \quad (۲۱)$$

همچنین تابع دوگان آن را به صورت زیر تعریف می‌کنیم:

$$g(\lambda_1, \dots, \lambda_k) = \begin{cases} \inf_{x \in X} L(x, \lambda_1, \dots, \lambda_k) & \text{if } \lambda_i \geq 0, \forall i \\ -\infty & \text{o.w.} \end{cases} \quad (۲۲)$$



حال مسأله‌ی دوگان را به صورت زیر تعریف می‌کنیم:

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

(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 دوگانگی ضعیف معروف است.

(ج) همچنین می‌توان نشان داد در صورتی که بر روی مجموعه‌ی  $feasible$  مسأله‌ی ۲۰ شرایطی برقرار باشد آنگاه دوگانگی قوی است یعنی:  $f^* = g^*$ . از دوگانگی 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  $\pi$  تابع مزیت  $A^\pi$  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  $\pi$  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:

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

**پاسخ:** با نوشتن فرم دوگان به شکل زیر خواهیم داشت:

$$g(\lambda) = \int_S \int_A A^\pi(s, a) d\tilde{\pi}(a|s) d\rho_\pi(s) + \lambda \left( \epsilon - \int_S C(\pi(\cdot|s), \tilde{\pi}(\cdot|s)) d\rho(s) \right)$$

در واقع برای حل دوگان قصد داریم مسأله‌ی بهینه‌سازی زیر را حل کنیم:

$$\begin{aligned} \min \quad & g(\lambda) \\ \text{s.t.} \quad & \lambda \geq 0 \end{aligned}$$

اما برای حل مسأله‌ی بهینه‌سازی بالا به جای اینکه مستقیم آن را حل کنیم روی تابع هدف آن کران بالا می‌زنیم تا به یک تابع هدف بهینه‌سازی خوش دست برسیم که بتوان از الگوریتم‌های بهینه‌سازی معروف برای حل آن استفاده کرد. به این منظور لازم است تا تابع  $g$  را بازنویسی کنیم:

$$g(\lambda) = \lambda\epsilon + \left( \int_S \int_A (A^\pi(s, a) d\tilde{\pi}(a|s) - \lambda C(\pi(\cdot|s), \tilde{\pi}(\cdot|s))) d\rho_\pi(s) \right)$$

حال سعی می‌کنیم برای عبارت درون انتگرال یک کران بالا با استفاده از قضیه کانورویچ ارائه کنیم. به همین جهت می‌نویسیم:

$$\psi(\cdot) = \frac{A^\pi(s, \cdot)}{\lambda} \quad \phi(\cdot) = \inf_{a' \in A} \{c(\cdot, a') - \psi(a')\}$$

که نتیجه می‌دهد:

$$\phi(a_1) + \psi(a_2) \leq c(a_1, a_2) - \psi(a_2) + \psi(a_2) \leq c(a_1, a_2)$$

لذا برای توابع بالا شرط مربوط به قضیه کانتورویچ صادق است. پس می‌توان نوشت:

$$C(\pi(\cdot|s), \tilde{\pi}(\cdot|s)) \geq \int_{\mathcal{A}} \phi(a) d\pi(\cdot|s) + \int_{\mathcal{A}} \psi(a) d\tilde{\pi}(\cdot|s)$$

پس در ادامه می‌توانیم بنویسیم:

$$\begin{aligned} (A^\pi(s, a) d\tilde{\pi}(a|s) - \lambda C(\pi(\cdot|s), \tilde{\pi}(\cdot|s))) &\leq \left( A^\pi(s, a) d\tilde{\pi}(a|s) - \lambda \left( \int_{\mathcal{A}} \phi(a) d\pi(\cdot|s) + \int_{\mathcal{A}} \psi(a) d\tilde{\pi}(\cdot|s) \right) \right) \\ &\leq \sup_{a' \in \mathcal{A}} \{A^\pi(s, a) - c(\cdot, a')\} d\pi(a|s) \end{aligned}$$

همانطور که ملاحظه می‌شود وابستگی به  $\tilde{\pi}$  در این کران بالا از بین رفته‌است، لذا با جایگذاری این کران بالا نتیجه‌ی دلخواه مسئله به دست می‌آید.

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 مسئله‌ی ۲۴ یک مسئله‌ی 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  $\lambda^*$  بهینه برای این مسئله می‌توان با استفاده از آن سیاست (تقریباً) بهینه‌ی مربوط به مسئله‌ی اول ۱۹ را با تکنیک‌هایی به دست آورد که به جهت جلوگیری از اطناب از ذکر آن خودداری می‌کنیم. حال تنها مسئله‌ای 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<sup>۳</sup> میانگین بگیرید.)

<sup>3</sup>Sliding Window