



بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی اصفهان  
دانشکده مهندسی برق و کامپیوتر

## یادگیری تقویتی عمیق

گزارش پروژه  
پرنیان ملک‌زاده

استاد  
دکتر محمد حسین منشی

کلیه حقوق مادی مترتب بر نتایج مطالعات،  
ابتکارات و نوآوری های ناشی از تحقیق  
موضوع این گزارش پروژه متعلق به دانشگاه  
صنعتی اصفهان است.

## فهرست مطالب

| عنوان  | صفحه |
|--|------|
| فهرست مطالب .....  | پنج  |
| فهرست تصاویر .....   | شش   |
| فهرست جداول .....  | هفت  |
| چکیده .....  | ۱    |
| <b>فصل اول: مقدمه</b>  |      |
| ۱-۱ پیشینه تحقیق .....   | ۳    |
| ۱-۲ اهداف و دستاوردهای تحقیق .....   | ۳    |
| ۱-۳ ساختار گزارش .....   | ۳    |
| <b>فصل دوم: مفاهیم پایه در یادگیری تقویتی</b>                                  |      |
| ۱-۲ مسئله یادگیری تقویتی .....   | ۴    |
| ۲-۲ طبقه بندی الگوریتم های یادگیری تقویتی .....                                | ۵    |
| ۱-۲-۲ الگوریتم های مدل-مبنا .....  | ۵    |
| ۲-۲-۲ الگوریتم های بدون مدل .....  | ۶    |
| <b>فصل سوم: عوامل تأثیرگذار بر خاصیت بازتولیدی الگوریتم های یادگیری تقویتی</b> |      |
| ۱-۳ پیش زمینه فنی .....  | ۸    |
| ۲-۳ تحلیل آزمایشی .....  | ۹    |
| ۱-۲-۳ ابر پارامترها .....  | ۱۰   |
| ۲-۲-۳ دانه های تصادفی و تعداد آزمایش ها .....                                  | ۱۰   |
| ۳-۲-۳ ویژگی های محیط .....   | ۱۱   |
| ۴-۲-۳ معماری شبکه .....  | ۱۲   |
| ۵-۲-۳ کدهای پایه .....   | ۱۴   |
| ۳-۳ گزارش معیارهای ارزیابی .....   | ۱۴   |
| ۱-۳-۳ تفاوت بین مشاهده برخط و بهینه سازی سیاست .....                           | ۱۵   |
| ۲-۳-۳ فواصل اطمینان .....  | ۱۵   |
| ۳-۳-۳ تحلیل توان .....   | ۱۶   |
| ۴-۳-۳ اهمیت .....  | ۱۶   |
| ۴-۳ نتیجه کلی .....  | ۱۷   |
| پیوست .....  | ۱۸   |
| مراجع .....  | ۲۰   |

## فهرست تصاویر

|    |       |     |
|----|-------|-----|
| ۳  | ..... | ۱-۱ |
| ۵  | ..... | ۱-۲ |
| ۷  | ..... | ۲-۲ |
| ۱۱ | ..... | ۱-۳ |
| ۱۲ | ..... | ۲-۳ |
| ۱۳ | ..... | ۳-۳ |
| ۱۴ | ..... | ۳-۴ |
| ۱۹ | ..... | ۱-آ |

## فهرست جداول

|    |       |     |
|----|-------|-----|
| ۱۳ | ..... | ۱-۳ |
| ۱۳ | ..... | ۲-۳ |
| ۱۷ | ..... | ۳-۳ |

## چکیده

در سال‌های اخیر، پیشرفت‌های قابل توجهی در حل مسائل چالش برانگیز در حوزه‌های مختلف با استفاده از یادگیری تقویتی عمیق<sup>۱</sup> حاصل شده است. بازتولید<sup>۲</sup> کارهای موجود و ارزیابی دقیق بهبودهای ارائه شده توسط روش‌های جدید برای حفظ این پیشرفت‌ها بسیار مهم است. متأسفانه، بازتولید نتایج برای روش‌های پیشرفته یادگیری عمیق تقویتی اغلب به سادگی امکان پذیر نیست. به ویژه، غیرقابل پیش بینی بودن محیط‌های استاندارد معیار، همراه با واریانس<sup>۳</sup> ذاتی در روش‌ها، می‌تواند تفسیر نتایج گزارش شده را دشوار کند. بدون معیار<sup>۴</sup>‌های معنادار و استانداردسازی دقیق گزارش‌های آزمایشی، دشوار است که مشخص کنیم آیا بهبودها نسبت به حالت پیشرفته قبلی معنادار هستند یا خیر. در این مقاله، چالش‌های مربوط به قابلیت بازتولید، تکنیک‌های آزمایشی صحیح و روش‌های گزارش دهی را بررسی می‌کنیم. ما به وضوح نشان می‌دهیم که چقدر متغیر بودن متریک‌ها و نتایج گزارش شده هنگام مقایسه با مبناهای رایج می‌تواند متفاوت باشد و پیشنهاداتی برای راهنمایی جهت قابل بازتولیدتر شدن نتایج در یادگیری تقویتی عمیق ارائه می‌دهیم. هدف ما ایجاد بحث در مورد چگونگی تضمین پیشرفت مستمر در این حوزه با به حداقل رساندن تلاش‌های هدررفته ناشی از نتایج غیرقابل بازتولید و به راحتی قابل تفسیر است.

**واژه‌های کلیدی:** ۱- یادگیری تقویتی عمیق، ۲- قابلیت بازتولید، ۳- تکنیک‌های آزمایشی، ۴- محیط‌های معیار

<sup>1</sup> Deep Reinforcement Learning

<sup>2</sup> Reproducing

<sup>3</sup> Variance

<sup>4</sup> Benchmark

<sup>5</sup> Metric



## فصل اول

### مقدمه

یادگیری تقویتی (RL) شاخه‌ای از مطالعه است که به بررسی چگونگی تعامل<sup>۱</sup> یک عامل<sup>۲</sup> با محیط<sup>۳</sup> خود برای یادگیری یک سیاست<sup>۴</sup> که هدف آن بیشینه کردن پاداش‌های تجمعی<sup>۵</sup> مورد انتظار برای یک وظیفه<sup>۶</sup> است، می‌پردازد. در سال‌های اخیر، یادگیری تقویتی به دلیل نتایج امیدبخش در زمینه‌هایی مانند کنترل سیستم‌های پیوسته در رباتیک [۱]، بازی شطرنج [۲]، بازی‌های آتاری [۳] و بازی‌های ویدئویی رقابتی [۴] و [۵] توجه و علاقه زیادی را به خود جلب کرده است. شکل ۱-۱ رشد این زمینه را از طریق تعداد مقالات منتشر شده در هر سال نشان می‌دهد. برای حفظ پیشرفت سریع در تحقیقات یادگیری تقویتی، مهم است که کارهای موجود به راحتی قابل بازتولید و مقایسه باشند تا بتوان به دقت بهبودهای ارائه شده توسط روش‌های نوین را ارزیابی کرد. با این حال، بازتولید نتایج یادگیری عمیق تقویتی اغلب ساده نیست و ادبیات این حوزه نتایج متنوعی برای همان الگوریتم‌های پایه گزارش می‌دهد [۶]. بازتولیدپذیری می‌تواند تحت تأثیر عوامل خارجی (مانند ابرپارامترها یا کدهای پایه<sup>۷</sup> استفاده‌شده) و عوامل داخلی (مانند اثرات دانه‌های تصادفی<sup>۸</sup> یا خواص محیط) قرار گیرد. ما این منابع تفاوت در نتایج گزارش شده را از طریق یک مجموعه نماینده از آزمایش‌ها بررسی می‌کنیم. برای شفافیت، تمرکز خود را بر روی روش‌های گرادیان سیاست<sup>۱۱</sup> (PG) در کنترل پیوسته<sup>۱۲</sup> قرار می‌دهیم. روش‌های گرادیان سیاست با استفاده از تقریب‌زن<sup>۱۳</sup>های تابع شبکه عصبی<sup>۱۴</sup> به ویژه در کنترل پیوسته موفق بوده‌اند [۱]، [۷] و [۸] در محیط‌های گسسته با روش‌های مبتنی بر ارزش<sup>۱۵</sup> رقابت می‌کنند. توجه داریم که تنوع معیارها و عدم وجود تست‌های معناداری در ادبیات یادگیری تقویتی پتانسیل گزارش نتایج گمراه‌کننده<sup>۱۶</sup> را ایجاد می‌کند. ما فواید استفاده از تست‌های معنادار را با استفاده از تکنیک‌های رایج در یادگیری ماشین<sup>۱۷</sup> و روش‌های آماری نشان می‌دهیم.

<sup>1</sup> Interaction

<sup>2</sup> Agent

<sup>3</sup> Environment

<sup>4</sup> Policy

<sup>5</sup> Reward

<sup>6</sup> Commulative

<sup>7</sup> Task

<sup>8</sup> Hyperparameter

<sup>9</sup> Codebases

<sup>10</sup> Random Seeds

<sup>11</sup> Policy Gradient

<sup>12</sup> Continuous Control

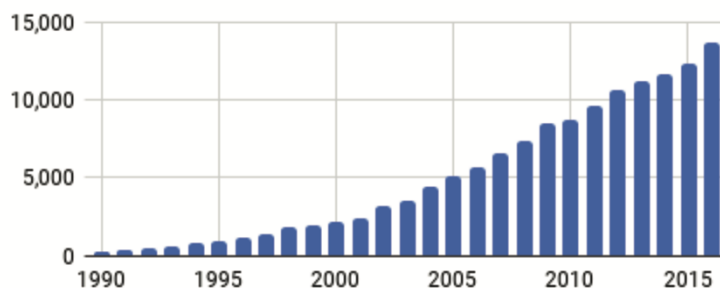
<sup>13</sup> Approximator

<sup>14</sup> Neural Network

<sup>15</sup> Value

<sup>16</sup> Misleading

<sup>17</sup> Machine Learning



شکل ۱-۱: رشد مقالات منتشر شده در زمینه یادگیری تقویتی. در این شکل، تعداد مقالات مرتبط با یادگیری تقویتی (محور y) در هر سال (محور x) که از جستجوهای Google Scholar استخراج شده‌اند، نشان داده شده است.

### ۱-۱ پیشینه تحقیق

تاکنون چندین پژوهش به ارزیابی الگوریتم‌های یادگیری تقویتی پرداخته‌اند. در [۹] چندین الگوریتم یادگیری تقویتی مینا قرار گرفته و پیاده‌سازی‌های پایه را برای جامعه علمی فراهم می‌کنند. معیارهای ارزیابی قابل تعمیم یادگیری تقویتی در [۱۰] پیشنهاد شده است. در [۱۱] محیط یادگیری Arcade را بازبینی کرده و روش‌های ارزیابی بهتری را در این میناها پیشنهاد می‌دهند. با این حال، در حالی که موضوع بازتولیدپذیری و تمرین‌های تجربی در زمینه‌های مرتبط مورد بررسی قرار گرفته است؛ تا آنجا که می‌دانیم [۱۲]، [۱۳]، [۱۴]، [۱۵]، [۱۶] و [۱۷] اولین کاری است که این سوال مهم را در زمینه یادگیری تقویتی عمیق بررسی می‌کند.

### ۲-۱ اهداف و دستاوردهای تحقیق

در هر بخش از تحلیل، سوالاتی را در مورد عوامل کلیدی تأثیرگذار بر خاصیت بازتولیدی مطرح می‌کنیم و می‌بینیم که منابع متعدد عدم قطعیت<sup>۱</sup> در بازتولید و مقایسه الگوریتم‌های یادگیری تقویتی وجود دارد. به این منظور، نشان می‌دهیم که جزئیات دقیق روش تجربی می‌تواند بسیار مهم باشد. بر اساس آزمایش‌های انجام شده توسط ما به این نتیجه می‌رسیم که با توصیه‌ها، خطوط تحقیقات و نقاط بحثی برای کارهای آینده، یادگیری تقویتی عمیق قابل بازتولید بوده و همچنان اهمیت دارد.

### ۳-۱ ساختار گزارش

در فصل دوم، به طور کلی به بررسی مطالعه یادگیری تقویتی و الگوریتم‌های آن می‌پردازیم. در این بخش، به بررسی اصول و مبانی یادگیری و نیز الگوریتم‌های مهم در این حوزه می‌پردازیم. در فصل سوم، به طور کلی به بررسی عوامل تأثیرگذار بر خاصیت بازتولیدی الگوریتم‌های یادگیری تقویتی می‌پردازیم. در این بخش، عوامل خارجی و داخلی که می‌توانند بر فرایند بازتولید نتایج تأثیرگذار باشند را مورد بررسی قرار می‌دهیم و روش‌های مختلف برای کاهش عدم قطعیت در بازتولید نتایج را بررسی می‌کنیم.

<sup>1</sup> Nondeterminism

## فصل دوم

### مفاهیم پایه در یادگیری تقویتی

#### ۱-۲ مسئله یادگیری تقویتی

مسئله یادگیری تقویتی به عنوان یک چارچوب ساده برای یادگیری از تعامل به منظور دستیابی به هدف مطرح شده است. در این چارچوب، یادگیرنده و تصمیم گیرنده به عنوان عامل شناخته می شوند. هر چیزی که عامل با آن تعامل دارد و شامل تمام چیزهایی است که خارج از عامل قرار دارند، محیط نامیده می شود. این دو به طور مداوم با یکدیگر تعامل دارند، به این صورت که عامل اقداماتی<sup>۱</sup> را انتخاب می کند و محیط به این اقدامات پاسخ داده و وضعیت<sup>۲</sup>های جدیدی را به عامل ارائه می دهد. محیط همچنین پاداش هایی به شکل مقادیر عددی ارائه می دهد که عامل سعی در بیشینه سازی آن ها در طول زمان دارد. یک تعریف کامل از محیط، یک وظیفه را تعریف می کند که یک نمونه از مسئله یادگیری تقویتی است.

به طور خاص، عامل و محیط در هر دنباله از گام های زمانی<sup>۳</sup> مجزا، با هم تعامل دارند ( $t = 0, 1, 2, 3, \dots$ ):

در هر گام زمانی  $t$ ، عامل یک نمایش از وضعیت محیط دریافت می کند،  $S_t \in S$  که در آن  $S$  مجموعه ای از وضعیت های ممکن است. سپس بر اساس آن، یک اقدام انتخاب می کند،  $A_t \in A(S_t)$  که در آن  $A(S_t)$  مجموعه ای از اقدامات موجود در وضعیت  $S_t$  است. در گام زمانی بعد، عامل یک پاداش عددی، به عنوان نتیجه ای از اقدام خود دریافت می کند،  $R_{t+1} \in R \subset \mathbb{R}$  و خود را در وضعیت جدیدی می یابد،  $S_{t+1}$ .

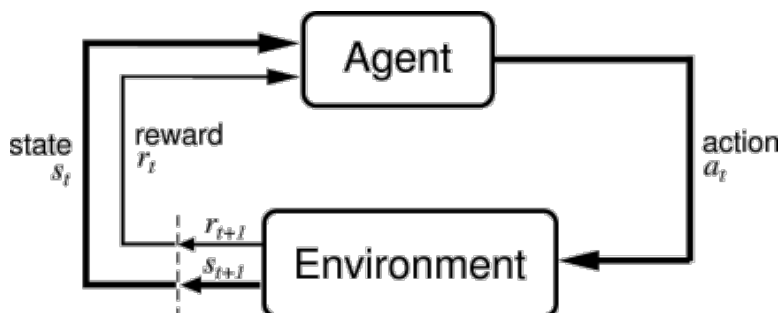
شکل ۱-۲ تعامل عامل و محیط را نشان می دهد. در هر گام زمانی، عامل یک نگاشت<sup>۴</sup> از وضعیت ها به احتمال انتخاب هر اقدام ممکن، پیاده سازی می کند. این نگاشت سیاست عامل نامیده می شود و به صورت  $(\pi_t)$  نمایش داده می شود، که در آن  $\pi_t(a|s)$  احتمال این است که  $A_t = a$  اگر  $S_t = s$  باشد.

<sup>1</sup> Actions

<sup>2</sup> State

<sup>3</sup> Timesteps

<sup>4</sup> Mapping



شکل ۲-۱: تعامل عامل-محیط در یادگیری تقویتی

روش‌های یادگیری تقویتی مشخص می‌کنند که چگونه عامل سیاست خود را به عنوان نتیجه‌ای از تجربه‌اش تغییر می‌دهد. هدف عامل، به طور خلاصه، بیشینه‌سازی مجموع کل پاداشی است که در طولانی مدت دریافت می‌کند.

## ۲-۲ طبقه‌بندی الگوریتم‌های یادگیری تقویتی

الگوریتم‌های یادگیری تقویتی به دو دسته کلی مدل-مبنا<sup>۲</sup> و بدون مدل<sup>۳</sup> تقسیم می‌شوند. منظور از مدل محیط، هر چیزی است که یک عامل می‌تواند از آن استفاده کند تا پیش‌بینی<sup>۴</sup> کند محیط چگونه به اقداماتش پاسخ می‌دهد. مدل، با داشتن یک حالت و یک عمل، پیش‌بینی می‌کند که حالت بعدی و پاداش بعدی چه خواهد بود. برخی مدل‌ها توصیفی از تمامی حالت‌های ممکن و احتمالات آنها ارائه می‌دهند که به آنها مدل‌های توزیعی<sup>۵</sup> می‌گوییم. مدل‌های دیگر فقط یکی از حالت‌ها را با توجه به احتمالات انتخاب می‌کنند که به آنها مدل‌های نمونه‌ای<sup>۶</sup> می‌گوییم. مدل‌ها می‌توانند برای تقلید<sup>۷</sup> یا شبیه‌سازی<sup>۸</sup> تجربه استفاده شوند. با داشتن یک حالت شروع و یک عمل، یک مدل نمونه‌ای یک انتقال<sup>۹</sup> ممکن را تولید می‌کند و یک مدل توزیعی تمامی انتقال‌های ممکن را با توجه به احتمالات آنها تولید می‌کند. با داشتن یک حالت شروع و یک سیاست، یک مدل نمونه‌ای می‌تواند یک روند<sup>۱۰</sup> کامل تولید کند و یک مدل توزیعی تمامی روندهای ممکن و احتمالات آنها را تولید می‌کند. در هر دو حالت، می‌گوییم مدل برای شبیه‌سازی محیط و تولید تجربه شبیه‌سازی شده استفاده می‌شود [۱۸]. در ادامه به توضیح این دو دسته و الگوریتم‌های مرتبط با هر کدام می‌پردازیم.

### ۲-۲-۱ الگوریتم‌های مدل-مبنا

در روش‌های مدل-مبنا، عامل به یک مدل از محیط دسترسی دارد یا آن را یاد می‌گیرد. عامل‌ها می‌توانند نتایج پیش‌بینی شده را در سیاست یادگرفته‌شده اعمال کنند. در این روش به دلیل پیش‌بینی‌های انجام شده، عواقب<sup>۱۱</sup> احتمالی سنجیده می‌شود و عامل می‌تواند تصمیمات<sup>۱۲</sup> بهتری بگیرد، اما یادگیری مدل محیط چندان آسان نیست و پیچیدگی‌هایی را ایجاد می‌کند.

<sup>1</sup> Experience

<sup>2</sup> Model-Based

<sup>3</sup> Model-Free

<sup>4</sup> Predict

<sup>5</sup> Distribution Models

<sup>6</sup> Sample Model

<sup>7</sup> Imitation

<sup>8</sup> Simulation

<sup>9</sup> Transition

<sup>10</sup> Episode

<sup>11</sup> Consequences

<sup>12</sup> Decisions

## ۲-۲-۲ الگوریتم های بدون مدل

در روش های بدون مدل، عامل به مدل محیط دسترسی ندارد و به جای آن، مستقیماً از تعامل با محیط یادگیری می کند. در این روش ها، دو رویکرد اصلی برای آموزش<sup>۱</sup> عامل وجود دارد:

بهینه سازی<sup>۲</sup> سیاست: روش های بهینه سازی سیاست، سیاست را به صورت واضحی به شکل  $\pi_\theta(a|s)$  نمایش می دهند. آن ها پارامترهای  $\theta$  را به دو روش بهینه سازی می کنند: یا به طور مستقیم با استفاده از گرادینت افزایشی<sup>۳</sup> بر روی هدف عملکرد<sup>۴</sup>  $J(\pi_\theta)$ ، یا به طور غیر مستقیم با بیشینه سازی<sup>۵</sup> تقریب های محلی  $J(\pi_\theta)$ . در اینجا،  $J(\pi_\theta)$  نشان دهنده تابع هدف عملکرد است که میزان عملکرد سیاست  $\pi_\theta$  را ارزیابی می کند و نشان می دهد که این سیاست چقدر خوب عمل می کند. این بهینه سازی تقریباً همیشه به صورت بر-سیاست<sup>۶</sup> انجام می شود، به این معنا که هر به روزرسانی<sup>۷</sup> فقط از داده هایی استفاده می کند که در حین اجرای نسخه جدید سیاست جمع آوری شده اند. بهینه سازی سیاست معمولاً شامل یادگیری یک تقریب زن برای تابع ارزش<sup>۸</sup> روی سیاست  $V_\phi(s)$  نیز می شود، که در تعیین نحوه به روزرسانی سیاست مورد استفاده قرار می گیرد.

- الگوریتم PPO<sup>۹</sup> از طریق به روزرسانی هایی که به طور غیر مستقیم عملکرد را به حداکثر می رساند عمل می کند. این الگوریتم به جای اینکه مستقیماً هدف عملکرد  $J(\pi_\theta)$  را بهینه سازی کند، یک تابع هدف جایگزین را به حداکثر می رساند که یک تخمین محافظه کارانه<sup>۱۰</sup> از میزان تغییر  $J(\pi_\theta)$  در نتیجه به روزرسانی ها ارائه می دهد.
- الگوریتم TRPO<sup>۱۱</sup> سیاست را به گونه ای به روزرسانی می کند که تغییرات آن محدود به یک ناحیه اطمینان<sup>۱۲</sup> باشد. این کار باعث می شود تغییرات پایدارتر و کنترل شده تر باشند و از نوسانات شدید جلوگیری شود.
- الگوریتم ACKTR<sup>۱۳</sup> به جای اینکه مستقیماً هدف عملکرد  $J(\pi_\theta)$  را بهینه سازی کند، از نواحی اطمینان برای به روزرسانی های پایدار و کارآمد سیاست بهره می برد. نواحی اطمینان تضمین می کنند که تغییرات سیاست به شکلی پایدار و کنترل شده انجام شود.

یادگیری-Q<sup>۱۴</sup>: در این رویکرد، یک تقریب زن  $Q_\theta(s, a)$  برای تابع ارزش-عمل بهینه  $Q^*(s, a)$  یاد گرفته می شود. این بهینه سازی معمولاً بر اساس معادله ی بلمن<sup>۱۵</sup> و اکثراً به صورت بدون-سیاست<sup>۱۶</sup> انجام می شود، به این معنا که بدون توجه به اینکه عامل چگونه محیط را جستجو<sup>۱۷</sup> کرده است، هر به روزرسانی می تواند از داده های جمع آوری شده در هر زمان از

<sup>1</sup> Training

<sup>2</sup> Optimization

<sup>3</sup> Gradient Ascent

<sup>4</sup> Performance

<sup>5</sup> Maximizing

<sup>6</sup> On-Policy

<sup>7</sup> Update

<sup>8</sup> Value Function

<sup>9</sup> Proximal Policy Optimization

<sup>10</sup> Conservative

<sup>11</sup> Trust Region Policy Optimization

<sup>12</sup> Trust Region

<sup>13</sup> Actor-Critic using Kronecker-Factored Trust Region

<sup>14</sup> Q-Learning

<sup>15</sup> Bellman Equation

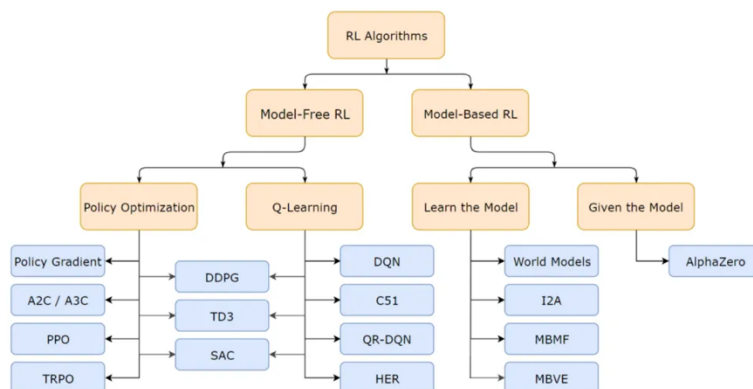
<sup>16</sup> Off-Policy

<sup>17</sup> Search

آموزش استفاده کند. سیاست مربوطه از طریق ارتباط بین  $Q^*$  و  $\pi^*$  به دست می‌آید و اقداماتی که توسط عامل یادگیری  $Q$  انجام می‌شود به صورت زیر است:

$$(s) = \arg \max_a Q_\theta(s, a) \quad (1-2)$$

میان‌یابی<sup>۱</sup> بین بهینه‌سازی سیاست و یادگیری-Q: به طرز عجیبی، بهینه‌سازی سیاست و یادگیری-Q با یکدیگر ناسازگار نیستند و تحت برخی شرایط حتی می‌توانند معادل باشند. الگوریتم‌هایی وجود دارند که در این طیف قرار می‌گیرند و می‌توانند به دقت بین نقاط قوت و ضعف هر دو روش تعادل برقرار کنند. این الگوریتم‌ها به بهره‌گیری از ویژگی‌های خوب هر دو روش کمک می‌کنند و می‌توانند در شرایط مختلف عملکرد بهتری داشته باشند. مثالی از این الگوریتم‌ها DDPG<sup>۲</sup> است. یک الگوریتم است که به طور همزمان یک سیاست تعیین‌کننده<sup>۳</sup> و یک تابع  $Q$ <sup>۴</sup> را می‌آموزد. این الگوریتم از هر کدام برای بهبود دیگری استفاده می‌کند. به عبارت دیگر، سیاست تعیین‌کننده برای بهبود تابع  $Q$  و تابع  $Q$  برای بهبود سیاست تعیین‌کننده به کار می‌رود. این تعامل به الگوریتم اجازه می‌دهد تا به طور موثرتری بهینه‌سازی را انجام دهد و عملکرد بهتری در محیط‌های پیچیده داشته باشد. شکل ۲-۲ تصویر مناسبی از این دسته بندی‌ها را نشان می‌دهد [۱۹].



شکل ۲-۲: دسته بندی در حوزه یادگیری تقویتی

<sup>1</sup> Interpolating

<sup>2</sup> Deep Deterministic Policy Gradient

<sup>3</sup> Deterministic Policy

<sup>4</sup> Q-Function

## فصل سوم

### عوامل تأثیرگذار بر خاصیت بازتولیدی الگوریتم‌های یادگیری تقویتی

#### ۳-۱ پیش‌زمینه فنی

این تحقیق به چندین الگوریتم یادگیری تقویتی بدون مدل مبتنی بر گرادیان سیاست پرداخته که پیاده‌سازی‌های آن‌ها به صورت عمومی در دسترس است و به طور گسترده در مقالات علمی به عنوان معیار مقایسه با روش‌های جدید به کار می‌روند. ما الگوریتم‌های بهینه‌سازی سیاست ناحیه اطمینان<sup>۱</sup> [۸]، گرادیان‌های سیاست قطعی عمیق<sup>۲</sup> [۱]، بهینه‌سازی سیاست پروگزیمال<sup>۳</sup> [۷]، و بازیگر-منتقد با استفاده از ناحیه اطمینان فاکتوری کروئکر<sup>۴</sup> [۲۰] را مورد بررسی قرار دادیم. این الگوریتم‌ها نتایج خوبی در وظایف کنترل پیوسته در محیط [۲۱] MuJoCo از [۲۲] OpenAI Gym نشان داده‌اند.

این الگوریتم‌ها عموماً به دنبال بهینه‌سازی پاداش تخفیفی مورد انتظار،

$$\rho(\theta, s_0) = \mathbb{E}_{\pi_\theta} [\sum_{t=0}^{\infty} \gamma^t r(s_t) | s_0]$$

(۱-۳)

با استفاده از تئوری گرادیان سیاست هستند که بیان می‌کند:

$$\frac{\delta \rho(\theta, s_0)}{\delta \theta} = \sum_s \mu_{\pi_\theta}(s | s_0) \sum_a \frac{\delta \pi_\theta(a | s)}{\delta \theta} Q_{\pi_\theta}(s | a)$$

(۲-۳)

اینجا

$$\mu_{\pi_\theta}(s | s_0) = \sum_{t=0}^{\infty} \gamma^t P(s_t = s | s_0)$$

(۳-۳)

[۲۳]

[۸] TRPO و [۷] PPO از محدودیت‌ها<sup>۵</sup> و تخمین مزیت<sup>۶</sup> برای انجام این به‌روزرسانی استفاده می‌کنند و مسئله بهینه‌سازی را به صورت زیر بازنویسی می‌کنند:

<sup>۱</sup> TRPO

<sup>۲</sup> DDPG

<sup>۳</sup> PPO

<sup>۴</sup> ACKTR

<sup>۵</sup> Constraints

<sup>۶</sup> Advantage

$$\max_{\theta} \mathbb{E}_t \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} A_t(s_t, a_t) \right]$$

(۴-۳)

اینجا ( $A_t$ ) تابع مزیت تعمیم یافته [۲۴] است. TRPO از روش نزول گرادینان مزدوج با یک محدودیت KL استفاده می کند:

$$\mathbb{E}_t \left[ KL[\pi_{\theta_{old}}(\cdot | s_t), \pi_{\theta}(\cdot, s_t)] \right] \leq \delta$$

(۵-۳)

PPO محدودیت را به عنوان یک جریمه (یا هدف برش) بازنویسی می کند.

الگوریتم های DDPG و ACKTR از روش بازیگر-منتقد<sup>۱</sup> استفاده می کنند. در این روش، ابتدا تابع ارزش<sup>۲</sup> که نشان دهنده مطلوبیت یک حرکت در یک وضعیت خاص است، تخمین زده می شود. سپس یک سیاست بهینه می شود که این تابع ارزش را بر اساس عرضه مونت کارلو<sup>۳</sup> رول اوت پیشینه کند.

تفاوت اصلی DDPG و ACKTR در نوع سیاست هایی است که استفاده می کنند. DDPG از سیاست های قطعی<sup>۴</sup> استفاده می کند، یعنی برای هر وضعیت، یک حرکت خاص را انتخاب می کند. اما ACKTR از سیاست های تصادفی<sup>۵</sup> بهره می برد که در آن ها برای هر وضعیت، یک توزیع احتمال<sup>۶</sup> بر روی حرکات وجود دارد. برای اطمینان از پایداری<sup>۷</sup> در سیاست های تصادفی، ACKTR از روش مناطق اعتماد فاکتور کرونگر<sup>۸</sup> استفاده می کند. این روش باعث می شود تغییرات در سیاست، محدود و کنترل شده باشند تا از ناپایداری ها جلوگیری شود.

### ۳-۲ تحلیل آزمایشی

در این بخش، به بررسی عواملی می پردازیم که بر قابلیت بازتولید روش های پیشرفته یادگیری تقویتی تأثیر گذار هستند. برای این منظور، مجموعه ای از آزمایش ها طراحی و اجرا شد تا اثرات این عوامل بررسی شوند. موارد زیر مورد توجه قرار گرفتند:

۱. تنظیم ابر پارامترها: بررسی تأثیر تنظیمات نامناسب ابر پارامترها بر عملکرد الگوریتم ها.
۲. دانه های تصادفی و تعداد تکرارها: ارزیابی اثر دانه های تصادفی و تعداد دفعات تکرار آزمایش بر نتایج.
۳. ویژگی های محیط: بررسی چگونگی تأثیر گذاری ویژگی های خاص محیط بر عملکرد الگوریتم ها.
۴. معماری شبکه: ارزیابی تأثیر انتخاب معماری شبکه برای تقریب تابع ارزش و سیاست بر عملکرد.
۵. کدهای پایه: بررسی تفاوت های عملکرد الگوریتم ها به دلیل تفاوت در کدهای پایه، با ثابت نگه داشتن سایر عوامل.

برای اکثر آزمایش ها، از پیاده سازی های موجود در OpenAI Baselines برای الگوریتم های ACKTR [۷]، [۲۰] PPO، DDPG، MuJoCo [۲۱] و Hopper-v1 و HalfCheetah-v1 از مجموعه [۲۱] استفاده شد. TRPO [۷]، [۲۵] استفاده شد. محیط های مورد استفاده شامل Hopper-v1 و HalfCheetah-v1 از مجموعه [۲۱] MuJoCo

<sup>1</sup> Actor-Critic

<sup>2</sup> Q Function

<sup>3</sup> Monte Carlo Rollout

<sup>4</sup> Deterministic

<sup>5</sup> Stochastic

<sup>6</sup> Probability Distribution

<sup>7</sup> Stability

<sup>8</sup> Kronecker Factored Trust Regions



در [۲۲] OpenAI Gym بودند که دینامیک‌های متفاوتی دارند و به بررسی اثرات مختلف کمک می‌کنند. برای اطمینان از نتایج منصفانه، هر آزمایش پنج بار با دانه‌های تصادفی متفاوت اجرا شده است (همه آزمایش‌ها از مجموعه یکسانی از دانه‌های تصادفی استفاده کردند). در این گزارش، جز در مواردی که به طور مشخص ذکر شده، از تنظیمات پیش فرض استفاده شده و تنها ابرپارامترهای مورد بررسی تغییر یافته‌اند. در همه آزمایش‌ها از تقریب‌زن‌های پرسپترون چندلایه<sup>۲</sup> استفاده شد. اندازه لایه‌های مخفی<sup>۳</sup> و توابع فعال‌سازی<sup>۴</sup> به صورت (M, N)، فعال‌سازی) نمایش داده شده‌اند. برای تنظیمات پیش فرض، ابرپارامترهای تحت بررسی به نوبت تغییر یافتند. برای DDPG از ساختار شبکه (۶۴، ۶۴) (ReLU) برای هر دو بخش عامل و منتقد، برای TRPO و PPO از ساختار (۶۴، ۶۴) (tanh) برای سیاست، و برای ACKTR از ساختار (۶۴، ۶۴) (tanh) برای عامل و (۶۴، ۶۴) (ELU) برای منتقد استفاده شده است.

### ۳-۲-۱ ابرپارامترها

انتخاب دقیق ابرپارامترها می‌تواند نقش بسیار مهمی در بهینه‌سازی نتایج الگوریتم‌ها داشته باشد. با این حال، در بسیاری از مقالات مرتبط، انتخاب تنظیمات بهینه ابرپارامترها یکسان نیست و محدوده مقادیر مورد بررسی به ندرت گزارش می‌شود. علاوه بر این، انتخاب نادرست ابرپارامترها می‌تواند به مقایسه ناعادلانه‌ای با الگوریتم‌های پایه منجر شود. در اینجا، به بررسی جنبه‌های مختلف انتخاب ابرپارامتر و تأثیر آن‌ها بر عملکرد الگوریتم‌ها می‌پردازیم.

### ۳-۲-۲ دانه‌های تصادفی و تعداد آزمایش‌ها

یکی از دغدغه‌های اصلی در یادگیری تقویتی عمیق، تفاوت نتایج به خاطر تصادفی بودن محیط یا فرآیند یادگیری مثل مقداردی اولیه تصادفی وزن‌هاست [۲۰]؛ [۲۶] بنابراین، حتی اگر نتایج چند آزمایش با دانه‌های کاملاً متفاوت را میانگین بگیریم، ممکن است به نتایج گمراه‌کننده‌ای برسیم. ما این موضوع را از طریق یک آزمایش بررسی می‌کنیم. نتایج: ما ۱۰ آزمایش را با همان تنظیمات ابرپارامتر انجام دادیم، فقط دانه‌های تصادفی را در هر ۱۰ آزمایش تغییر دادیم. سپس آزمایش‌ها را به دو گروه ۵ تایی تقسیم کرده و میانگین هر گروه را محاسبه کردیم. همانطور که در شکل ۳-۱ می‌بینید، عملکرد الگوریتم‌ها می‌تواند بسیار متفاوت باشد. نشان می‌دهیم که تفاوت بین اجراها به قدری زیاد است که فقط با تغییر دانه‌های تصادفی، توزیع‌های آماری کاملاً متفاوتی به دست می‌آید. متأسفانه در نتایج گزارش شده اخیر، گاهی فقط بهترین N آزمایش از میان چندین آزمایش انتخاب شده [۲۰]؛ [۲۶] یا نتایج روی تعداد کمی آزمایش ( $N < 5$ ) میانگین گرفته شده است [۲۷]؛ [۲۰] آزمایش ما نشان می‌دهد که این کار می‌تواند گمراه‌کننده باشد. به خصوص برای HalfCheetah، امکان دارد منحنی‌های یادگیری به دست آمده اصلاً در یک توزیع قرار نگیرند، فقط به این خاطر که میانگین اجراهای مختلف با همان ابرپارامترها اما دانه‌های متفاوت گرفته شده است. اگرچه نمی‌توان تعداد دقیقی از آزمایش‌ها را توصیه کرد، شاید بتوان از روش‌های تحلیل توان<sup>۵</sup> برای پیشنهاد یک تعداد تقریبی استفاده کرد، اما نیاز به بررسی بیشتری برای حل این مشکل است.

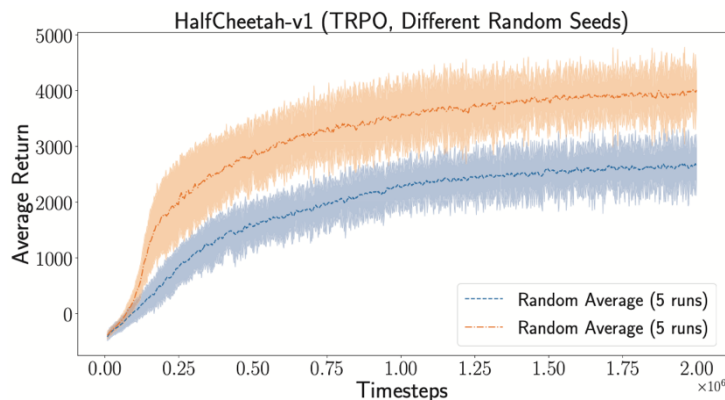
<sup>1</sup> Dynamic

<sup>2</sup> Multi Layer Perceptron

<sup>3</sup> Hidden Layers

<sup>4</sup> Activation Functions

<sup>5</sup> Power Analysis



شکل ۳-۱: اجرای TRPO بر روی HalfCheetah-v1 با استفاده از همان تنظیمات ابر پارامترها به طور میانگین بر روی دو مجموعه از ۵ بذر تصادفی مختلف. میانگین آزمون  $t$  دو نمونه‌ای در کل توزیع آموزشی به نتیجه  $t = -9.0916$  و  $p = 0.0016$  رسید.

### ۳-۲-۳ ویژگی‌های محیط

برای ارزیابی اینکه انتخاب محیط چگونه می‌تواند بر نتایج تاثیر بگذارد، از مجموعه‌ای از ابر پارامترهای پیش فرض در الگوریتم‌های مختلف استفاده کردیم و بررسی کردیم که هر الگوریتم در مجموعه‌ای از وظایف کنترل پیوسته چقدر خوب عمل می‌کند. برای این آزمایش‌ها، از محیط‌های روبرو از OpenAI Gym استفاده کردیم: Hopper-v1، HalfCheetah-v1، Walker2d-v1 و Swimmer-v1. انتخاب محیط نقش مهمی در نشان دادن عملکرد یک الگوریتم جدید نسبت به الگوریتم‌های موجود دارد. در وظایف کنترل پیوسته، محیط‌ها معمولاً دارای رفتارهای تصادفی، مسیرهای کوتاه‌تر یا ویژگی‌های دینامیکی متفاوت هستند. ما نشان می‌دهیم که به دلیل این تفاوت‌ها، عملکرد الگوریتم‌ها می‌تواند در محیط‌های مختلف متفاوت باشد و بهترین الگوریتم برای همه محیط‌ها همیشه مشخص نیست. بنابراین، ارائه نتایج برای طیف گسترده‌ای از محیط‌ها اهمیت دارد و نباید فقط محیط‌هایی را انتخاب کنیم که کار جدید ما نسبت به محیط‌های دیگر بهتر عمل کند.

نتایج: همانطور که در شکل ۴ نشان داده شده است، در محیط‌هایی با دینامیک پایدار مثل HalfCheetah-v1، الگوریتم DDPG عملکرد بهتری نسبت به بقیه دارد. اما با ناپایداری شدن دینامیک مثل در Hopper-v1 عملکرد DDPG به شدت کاهش می‌یابد. چون DDPG یک روش بدون سیاست است، نویز ناشی از اکتشاف<sup>۲</sup> می‌تواند باعث شکست<sup>۳</sup>های ناگهانی در محیط‌های ناپایدار شود. بنابراین، یادگیری یک تخمین درست از مقدار  $Q$  برای بازده<sup>۴</sup>های مورد انتظار دشوار است، خصوصاً وقتی که بسیاری از مسیرهای اکتشافی منجر به شکست می‌شوند. چون شکست‌ها در چنین وظایفی با مسیرهای کوتاه‌تر مشخص می‌شوند، یک بهینه محلی<sup>۵</sup> در این حالت فقط زنده ماندن تا حداکثر طول مسیر (که معادل هزار گام زمانی و پاداش مشابه به خاطر پاداش بقا در حالت Hopper-v1 است) می‌باشد. همانطور که در شکل ۲-۳ مشاهده می‌شود، DDPG در Hopper دقیقاً همین کار را انجام می‌دهد. این یک مثال واضح است که نشان می‌دهد تنها نشان دادن HalfCheetah پایدار و مناسب هنگام گزارش آزمایش‌های مبتنی بر DDPG ناعادلانه است. علاوه بر این، محیط Swimmer-v1 را در نظر بگیریم که در شکل ۳-۲ نشان داده شده است. در اینجا، TRPO به طور قابل توجهی عملکرد بهتری نسبت به سایر الگوریتم‌ها دارد. به دلیل دینامیک

<sup>1</sup> Trajectory

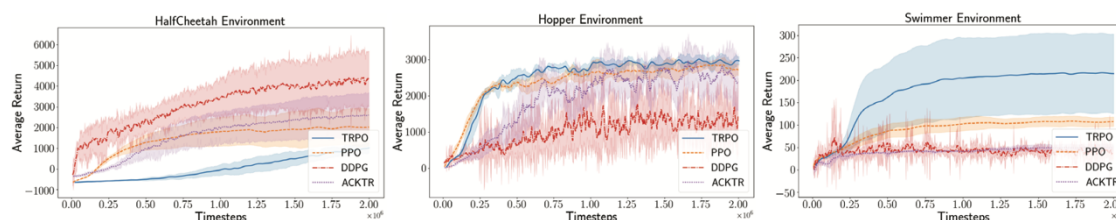
<sup>2</sup> Exploration

<sup>3</sup> Failure

<sup>4</sup> Return

<sup>5</sup> Local Optimal

محیط شبیه به آب، یک بهینه محلی برای سیستم این است که جمع شود و بدون شنا کردن صحیح تکان بخورد. اما این وضعیت منجر به بازده حدود ۱۳۰ می‌شود. با رسیدن به یک بهینه محلی، نمودارهای یادگیری می‌توانند بهینه‌سازی موفقیت‌آمیز سیاست را در طول زمان نشان دهند، در حالی که در واقع بازده‌های به دست آمده نمایانگر رفتار مطلوب نیستند، همانطور که در ویدیوی<sup>۱</sup> از سیاست یادگرفته شده نشان داده شده است. بنابراین، مهم است که نه تنها بازده‌ها، بلکه نمایش‌هایی از سیاست یادگرفته شده در عمل نشان داده شود. بدون درک اینکه ارزیابی بازده‌ها چه چیزی را نشان می‌دهد، ممکن است نتایج گمراه‌کننده‌ای گزارش شود که در واقع تنها بهینه‌های محلی را بیشینه می‌کند، نه رسیدن به رفتار مطلوب.



شکل ۳-۲: عملکرد چندین الگوریتم گرادین سیاست در مجموعه محیط‌های معیار MuJoCo

### ۳-۲-۴ معماری شبکه

مطالعات قبلی نشان داده‌اند که معماری شبکه<sup>۲</sup> سیاست می‌تواند نتایج را در الگوریتم‌های TRPO و DDPG به طور قابل توجهی تحت تأثیر قرار دهد [۶]. علاوه بر این، برخی تابع‌های فعال‌ساز مانند واحد خطی تصحیح شده<sup>۳</sup> ممکن است به دلیل مشکل "رلو مرده"<sup>۴</sup> باعث کاهش عملکرد یادگیری شوند. [۲۸] بنابراین، ما معماری شبکه و توابع فعال‌ساز را برای هر دو تقریب‌زن سیاست و ارزش بررسی کردیم. در ادبیات موضوع، بررسی‌های مشابهی نشان داده‌اند که عملکرد تقریب‌زن‌های خطی، RBF و شبکه‌های عصبی متفاوت است. [۲۹] جداول ۱ و ۲ خلاصه‌ای از عملکرد نهایی ارزیابی تمام تغییرات معماری پس از آموزش روی ۲ میلیون نمونه (یعنی ۲ میلیون گام زمانی در محیط) را نشان می‌دهند. ما سه معماری متداول پرسپترون چندلایه (MLP) را بررسی کردیم: (۶۴، ۶۴)، (۱۰۰، ۵۰، ۲۵) و (۴۰۰، ۳۰۰). علاوه بر این، توابع فعال‌ساز شبکه‌های ارزش و سیاست را در میان تانژانت هیپربولیک<sup>۵</sup>، ReLU و Leaky ReLU تغییر دادیم. نتایج نشان می‌دهد که عملکرد می‌تواند با تغییرات ساده در توابع فعال‌ساز شبکه‌های سیاست یا ارزش به شدت تحت تأثیر قرار گیرد. معمولاً ReLU یا Leaky ReLU در میان محیط‌ها و الگوریتم‌ها بهترین عملکرد را دارند. این اثرات در الگوریتم‌ها یا محیط‌های مختلف یکسان نیستند. این ناهمگونی نشان می‌دهد که معماری شبکه چقدر به روش الگوریتم وابسته است. به عنوان مثال، استفاده از یک شبکه بزرگ با PPO ممکن است نیاز به تنظیم ابرپارامترهای دیگری مانند میزان محدودیت ناحیه اطمینان یا نرخ یادگیری<sup>۶</sup> داشته باشد تا با تغییر معماری جبران شود. این بازی پیچیده ابرپارامترها یکی از دلایل دشواری تکرارپذیری روش‌های گرادین سیاست فعلی است. انتخاب معماری مناسب برای نتایج پایه‌ای صحیح بسیار مهم است. این همچنین نیاز احتمالی به الگوریتم‌های بی‌توجه به ابرپارامتر را پیشنهاد

<sup>1</sup> Video [Link](#)

<sup>2</sup> Network Architecture

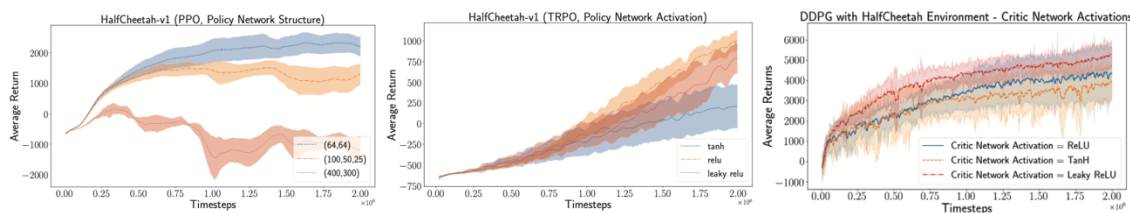
<sup>3</sup> ReLU

<sup>4</sup> Dying Relu

<sup>5</sup> Hyperbolic Tangent (tanh)

<sup>6</sup> Learning Rate

می‌کند - یعنی الگوریتم‌هایی که انطباق ابرپارامتر را بخشی از طراحی در نظر می‌گیرند - تا بتوان مقایسه‌های منصفانه‌ای بدون تگرانی از تنظیمات نامناسب برای کار در دست انجام داد.



شکل ۳-۳: اهمیت معماری شبکه سیاست و توابع فعال‌سازی: PPO (چپ)، TRPO (وسط) و DDPG (راست).

| Algorithm     | Environment    | 400,300     | 64,64      | 100,50,25  | tanh       | ReLU       | LeakyReLU  |
|---------------|----------------|-------------|------------|------------|------------|------------|------------|
| TRPO<br>[8]   | Hopper-v1      | 2980 ± 35   | 2674 ± 227 | 3110 ± 78  | 2674 ± 227 | 2772 ± 211 | -          |
|               | HalfCheetah-v1 | 1791 ± 224  | 1939 ± 140 | 2151 ± 27  | 1939 ± 140 | 3041 ± 161 | -          |
| TRPO<br>[9]   | Hopper-v1      | 1243 ± 55   | 1303 ± 89  | 1243 ± 55  | 1303 ± 89  | 1131 ± 65  | 1341 ± 127 |
|               | HalfCheetah-v1 | 738 ± 240   | 834 ± 317  | 850 ± 378  | 834 ± 317  | 784 ± 352  | 1139 ± 364 |
| TRPO<br>[7]   | Hopper-v1      | 2909 ± 87   | 2828 ± 70  | 2812 ± 88  | 2828 ± 70  | 2941 ± 91  | 2865 ± 189 |
|               | HalfCheetah-v1 | -155 ± 188  | 205 ± 256  | 306 ± 261  | 205 ± 256  | 1045 ± 114 | 778 ± 177  |
| PPO<br>[7]    | Hopper-v1      | 61 ± 33     | 2790 ± 62  | 2592 ± 196 | 2790 ± 62  | 2695 ± 86  | 2587 ± 53  |
|               | HalfCheetah-v1 | -1180 ± 444 | 2201 ± 323 | 1314 ± 340 | 2201 ± 323 | 2971 ± 364 | 2895 ± 365 |
| DDPG<br>[25]  | Hopper-v1      | 1419 ± 313  | 1632 ± 459 | 2142 ± 436 | 1491 ± 205 | 1632 ± 459 | 1384 ± 285 |
|               | HalfCheetah-v1 | 5579 ± 354  | 4198 ± 606 | 5600 ± 601 | 5325 ± 281 | 4198 ± 606 | 4094 ± 233 |
| DDPG<br>[27]  | Hopper-v1      | 600 ± 126   | 593 ± 155  | 501 ± 129  | 436 ± 48   | 593 ± 155  | 319 ± 127  |
|               | HalfCheetah-v1 | 2845 ± 589  | 2771 ± 535 | 1638 ± 624 | 1638 ± 624 | 2771 ± 535 | 1405 ± 511 |
| DDPG<br>[9]   | Hopper-v1      | 506 ± 208   | 749 ± 271  | 629 ± 138  | 354 ± 91   | 749 ± 271  | -          |
|               | HalfCheetah-v1 | 850 ± 41    | 1573 ± 385 | 1224 ± 553 | 1311 ± 271 | 1573 ± 385 | -          |
| ACKTR<br>[20] | Hopper-v1      | 2577 ± 529  | 1608 ± 66  | 2287 ± 946 | 1608 ± 66  | 2835 ± 503 | 2718 ± 434 |
|               | HalfCheetah-v1 | 2653 ± 408  | 2691 ± 231 | 2498 ± 112 | 2621 ± 381 | 2160 ± 151 | 2690 ± 231 |

جدول ۳-۱: نتایج تغییرات معماری سیاست ما در میان پیاده‌سازی‌ها و الگوریتم‌های مختلف. میانگین نهایی ± خطای استاندارد در عرض ۵ آزمایش از بازده‌ها در طول ۱۰۰ مسیر آخر پس از ۲ میلیون نمونه آموزشی. برای ACKTR، از تابع فعال‌ساز ELU استفاده می‌کنیم.

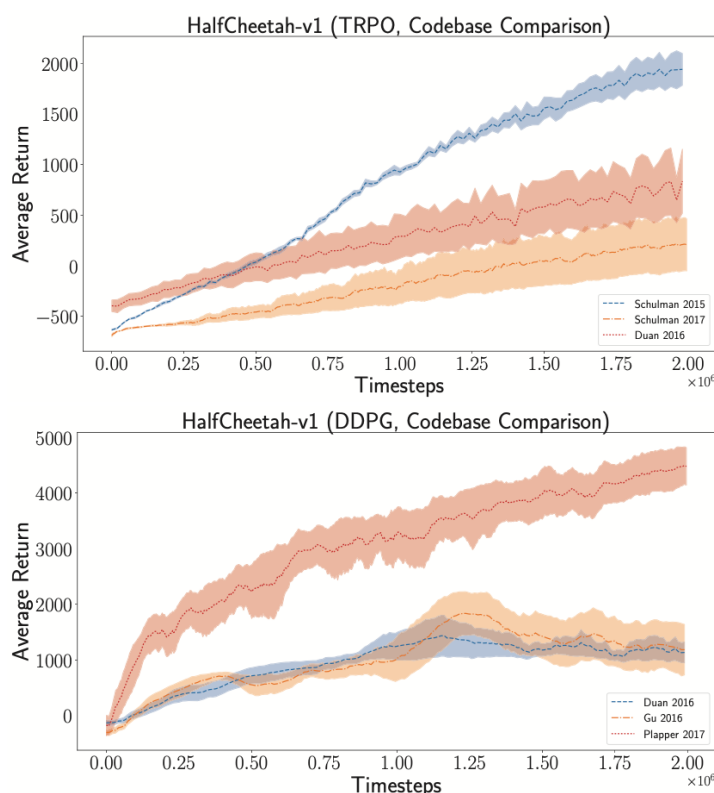
| Algorithm     | Environment    | 400,300    | 64,64      | 100,50,25  | tanh       | ReLU       | LeakyReLU  |
|---------------|----------------|------------|------------|------------|------------|------------|------------|
| TRPO<br>[8]   | Hopper-v1      | 3011 ± 171 | 2674 ± 227 | 2782 ± 120 | 2674 ± 227 | 3104 ± 84  | -          |
|               | HalfCheetah-v1 | 2355 ± 48  | 1939 ± 140 | 1673 ± 148 | 1939 ± 140 | 2281 ± 91  | -          |
| TRPO<br>[7]   | Hopper-v1      | 2909 ± 87  | 2828 ± 70  | 2812 ± 88  | 2828 ± 70  | 2829 ± 76  | 3047 ± 68  |
|               | HalfCheetah-v1 | 178 ± 242  | 205 ± 256  | 172 ± 257  | 205 ± 256  | 235 ± 260  | 325 ± 208  |
| PPO<br>[7]    | Hopper-v1      | 2704 ± 37  | 2790 ± 62  | 2969 ± 111 | 2790 ± 62  | 2687 ± 144 | 2748 ± 77  |
|               | HalfCheetah-v1 | 1523 ± 297 | 2201 ± 323 | 1807 ± 309 | 2201 ± 323 | 1288 ± 12  | 1227 ± 462 |
| DDPG<br>[25]  | Hopper-v1      | 1419 ± 312 | 1632 ± 458 | 1569 ± 453 | 971 ± 137  | 852 ± 143  | 843 ± 160  |
|               | HalfCheetah-v1 | 5600 ± 601 | 4197 ± 606 | 47134374   | 3908 ± 293 | 4197 ± 606 | 5324 ± 280 |
| DDPG<br>[27]  | Hopper-v1      | 523 ± 248  | 343 ± 34   | 345 ± 44   | 436 ± 48   | 343 ± 34   | -          |
|               | HalfCheetah-v1 | 1373 ± 678 | 1717 ± 508 | 1868 ± 620 | 1128 ± 511 | 1717 ± 508 | -          |
| DDPG<br>[9]   | Hopper-v1      | 1208 ± 423 | 394 ± 144  | 380 ± 65   | 354 ± 91   | 394 ± 144  | -          |
|               | HalfCheetah-v1 | 789 ± 91   | 1095 ± 139 | 988 ± 52   | 1311 ± 271 | 1095 ± 139 | -          |
| ACKTR<br>[20] | Hopper-v1      | 152 ± 47   | 1930 ± 185 | 1589 ± 225 | 691 ± 55   | 500 ± 379  | 1930 ± 185 |
|               | HalfCheetah-v1 | 2653 ± 408 | 2691 ± 231 | 2498 ± 112 | 2621 ± 381 | 2160 ± 151 | 2691 ± 231 |

جدول ۳-۲: نتایج تغییرات معماری تابع ارزش (Q یا V) در میان پیاده‌سازی‌ها و الگوریتم‌های مختلف. میانگین نهایی ± خطای استاندارد در عرض ۵ آزمایش از بازده‌ها در طول ۱۰۰ مسیر آخر پس از ۲ میلیون نمونه آموزشی. برای ACKTR، به جای Leaky ReLU، از توابع فعال‌ساز ELU استفاده می‌کنیم.

### ۳-۲-۵ کدهای پایه

در بسیاری از موارد، نویسندگان نسخه‌های خود را از الگوریتم‌های مبنا برای مقایسه پیاده‌سازی می‌کنند. ما پیاده‌سازی مبناهای OpenAI از الگوریتم TRPO که در (Schulman et al. ۲۰۱۷) استفاده شده، کد اصلی [۸] TRPO، و پیاده‌سازی TRPO با Tensorflow در [۹] را بررسی می‌کنیم. همچنین پیاده‌سازی‌های [۹] rllab Theano، [۲۷] rllabplusplus و [۲۵] OpenAI baselines از DDPG را مقایسه می‌کنیم. هدف ما توجه به تفاوت‌های عملکرد به دلیل جزئیات پیاده‌سازی در الگوریتم‌ها است. ما زیرمجموعه‌ای از آزمایش‌های معماری خود را با استفاده از پیاده‌سازی‌های OpenAI baselines و همان ابرپارامترها انجام می‌دهیم.

نتایج: ما متوجه شدیم که تفاوت‌های پیاده‌سازی که اغلب در مقالات منعکس نمی‌شوند، می‌توانند تأثیرات چشمگیری بر عملکرد داشته باشند. این را می‌توان در عملکرد نهایی ارزیابی ما پس از آموزش روی ۲ میلیون نمونه در جداول ۱ و ۲، و همچنین یک مقایسه نمونه در شکل ۳-۴ مشاهده کرد. این نشان می‌دهد که جزئیات پیاده‌سازی باید شرح داده شوند، کدها همراه با مقالات منتشر شوند، و عملکرد آزمایش‌های مبنا در کارهای جدید با کدهای اصلی مقالات مبنا مطابقت داشته باشد.



شکل ۳-۴: مقایسه کد پایه TRPO با مجموعه پیش‌فرض ابرپارامترهای ما (همانطور که در آزمایش‌های دیگر استفاده شده است).

### ۳-۳ گزارش معیارهای ارزیابی

در این بخش، برخی از معیارهای ارزیابی که معمولاً در مطالعات یادگیری تقویتی استفاده می‌شوند را بررسی می‌کنیم. در عمل، الگوریتم‌های یادگیری تقویتی معمولاً با نمایش نمودارها یا جداولی از پاداش تجمعی متوسط (بازده متوسط) و اخیراً، پاداش حداکثر به دست آمده در تعداد معینی از گام‌های زمانی ارزیابی می‌شوند. به دلیل ناپایدار بودن بسیاری از این الگوریتم‌ها، صرفاً گزارش پاداش‌های حداکثری برای مقایسه منصفانه کافی نیست؛ حتی گزارش بازده متوسط نیز می‌تواند

گمراه کننده باشد زیرا دامنه عملکرد در بین آزمایش ها و دانه های تصادفی متفاوت است. این گزارش ها ممکن است به تنهایی، تصویر واضحی از عملکرد یک الگوریتم ارائه ندهند. با این حال، هنگامی که با فواصل اطمینان<sup>۱</sup> ترکیب شوند، ممکن است برای گرفتن تصمیمی آگاهانه<sup>۲</sup> کافی باشند به شرطی که تعداد آزمایش ها کافی باشد. از این رو، ما استفاده از روش های بوت استرپ<sup>۳</sup> و آزمون معناداری مشابه یادگیری ماشین (مانند مطالعات [۳۰] و [۳۱]) را برای ارزیابی عملکرد الگوریتم ها بررسی می کنیم.

### ۳-۳-۱ تفاوت بین مشاهده برخط<sup>۴</sup> و بهینه سازی سیاست

یک تفاوت مهم در گزارش نتایج یادگیری تقویتی مشاهده برخط و بهینه سازی سیاست است. در مشاهده برخط، عامل در طول فرآیند یادگیری بهینه سازی بازده ها را هدف قرار می دهد و مسیر عامل پایانی ندارد. در این حالت، ارزیابی ها می توانند شامل پاداش های تجمعی متوسط در طول کل فرآیند یادگیری باشند. این روش به تعادل بین اکتشاف و بهره برداری<sup>۵</sup> توجه دارد و معمولاً از پاداش های تجمعی متوسط در طول زمان برای ارزیابی استفاده می شود، همانطور که در مطالعات [۳۲] آمده است. همچنین، ممکن است از روش های ارزیابی برون خط<sup>۶</sup> نیز استفاده شود، همانطور که در مطالعه [۳۳] دیده می شود. در مقابل، بهینه سازی سیاست رویکرد دیگری است که در آن ارزیابی با استفاده از یک سیاست هدف در حالت برون خط انجام می شود. در این دیدگاه، ارزیابی ها باید در طول کل مسیر وظیفه با یک سیاست هدف واحد انجام شوند تا بازده های متوسطی که سیاست هدف می تواند به دست آورد، مشخص شود. این روش به ما اجازه می دهد تا عملکرد یک سیاست مشخص را در شرایطی ثابت و کنترل شده ارزیابی کنیم.

### ۳-۳-۲ فواصل اطمینان

روش بوت استرپ به عنوان یک روش محبوب برای به دست آوردن بینشی از توزیع یک جمعیت<sup>۷</sup> با استفاده از یک نمونه<sup>۸</sup> کوچکتر شناخته می شود. این تکنیک به خصوص در آزمایش های A/B بسیار پر کاربرد است و می توانیم از ایده های این حوزه بهره بگیریم. معمولاً یک تخمین گر<sup>۹</sup> بوت استرپ با باز نمونه گیری<sup>۱۰</sup> مکرر با جایگزینی ایجاد می شود تا میانگین آماری معتبر و حدود اطمینان به دست آید. با استفاده از این روش، می توانیم به بازه اطمینان<sup>۱۱</sup> ۹۵٪ نتایج بخش محیط های خود دست پیدا کنیم.

جدول ۳ نشان دهنده میانگین بوت استرپ و حدود اطمینان ۹۵٪ در آزمایش های محیطی ما است. حدود اطمینان می تواند بسته به الگوریتم ها و محیط ها بسیار متغیر باشند. ما مشاهده کردیم که الگوریتم های TRPO و PPO از پایدارترین ها هستند و حدود اطمینان کوچکی از بوت استرپ دارند. در مواردی که حدود اطمینان بسیار بزرگ است، ممکن است نیاز باشد تعداد آزمایش ها را افزایش دهیم (یعنی اندازه نمونه را بیشتر کنیم).

<sup>1</sup> Confidence Intervals

<sup>2</sup> Informed

<sup>3</sup> Bootstrap

<sup>4</sup> Online

<sup>5</sup> Exploitation

<sup>6</sup> Offline

<sup>7</sup> Population

<sup>8</sup> Sample

<sup>9</sup> Estimator

<sup>10</sup> Resampling

<sup>11</sup> Confidence Interval

### ۳-۳-۳ تحلیل توان

یکی دیگر از روش‌های تعیین نیاز به بیشتر شدن حجم نمونه، تحلیل توان<sup>۱</sup> بوت‌استرپ است [۳۴]. اگر از نمونه موجود استفاده کنیم و به آن یک افزایش<sup>۲</sup> یکنواخت<sup>۳</sup> اعمال کنیم (مثلاً با ضریب ۱.۲۵ به طور یکنواخت افزایش دهیم)، می‌توانیم شبیه‌سازی‌های بوت‌استرپ متعددی را اجرا کنیم و مشخص کنیم که چه درصدی از این شبیه‌سازی‌ها منجر به مقادیر آماری معنادار با این افزایش می‌شوند. اگر درصد کمی از مقادیر معنادار باشند، به نمونه بزرگتری نیاز داریم (باید آزمایشات بیشتری انجام شود). ما این کار را در تمامی آزمایش‌های محیطی انجام دادیم و به این نتیجه رسیدیم که در تنظیمات ناپایداری، درصد توان بوت‌استرپ به نتایج غیرمعنادار در آزمایش افزایش متمایل است. برعکس، در آزمایش‌های پایداری (مثلاً TRPO در Hopper-v1) با نمونه‌های کوچک، آزمایش افزایش نشان داد که نیازی به آزمایش‌های بیشتر برای ایجاد مقایسه‌های معنادار نیست.

### ۴-۳-۳ اهمیت<sup>۴</sup>

یک عامل مهم در انتخاب الگوریتم یادگیری تقویتی، اهمیت دستاوردهای گزارش شده بر اساس یک معیار مشخص است. چندین پژوهش به استفاده از معیارهای اهمیت برای ارزیابی قابلیت اطمینان معیارهای ارزیابی در یادگیری ماشین پرداخته‌اند. با این حال، تحقیقات کمی در یادگیری تقویتی به ارزیابی اهمیت معیارهای گزارش شده می‌پردازند. بر اساس نتایج تجربی ما که نشان می‌دهد عملکرد الگوریتم‌ها می‌تواند به شدت با تغییرات تصادفی در دانه‌های تصادفی متفاوت باشد، واضح است که به یک معیار برای ارزیابی اهمیت دستاوردهای عملکرد الگوریتم و اطمینان از معیارهای گزارش شده نیاز است. در حالی که تحقیقات و بررسی‌های بیشتری برای تعیین بهترین معیارها برای ارزیابی الگوریتم‌های یادگیری تقویتی لازم است، ما یک مجموعه اولیه از معیارها را بر اساس نتایج یادگیری ماشین بررسی می‌کنیم.

در یادگیری نظارت شده<sup>۵</sup>، معیارهای اهمیتی مانند  $t$ -test،  $k$ -fold  $t$ -test، باز نمونه گیری تصحیح شده<sup>۶</sup> و دیگر معیارها برای مقایسه نتایج یادگیری ماشین مورد بحث قرار گرفته‌اند [۳۱] و [۳۵]؛ با این حال، فرضیاتی که مربوط به داده‌های اصلی با معیارهای اصلاح شده است، الزاماً در یادگیری تقویتی قابل اعمال نیستند. برای بررسی آزمون‌های اهمیت اصلاح شده مناسب برای یادگیری تقویتی به تحقیقات بیشتری نیاز است.

با این وجود، ما چندین معیار اهمیت را بررسی می‌کنیم که نشان می‌دهند آیا یک الگوریتم جدید واقعاً به اندازه الگوریتم‌های پیشرفته عمل می‌کند یا خیر. ما  $t$ -test<sup>۷</sup> ۲ نمونه‌ای ساده (مرتب‌سازی تمام بازده‌های نهایی ارزیابی در میان  $N$  آزمایش تصادفی با دانه‌های تصادفی مختلف)، آزمون کولموگروف-اسمیرنوف<sup>۸</sup> [۳۶] و تفاوت درصدی بوت‌استرپ با فواصل اطمینان ۹۵٪ را در نظر می‌گیریم. تمامی معیارهای محاسبه شده در مواد تکمیلی قابل مشاهده هستند. به طور کلی، ما دریافتیم که مقادیر اهمیت با انتظارات ما مطابقت دارند. به عنوان مثال، مقایسه عملکرد Walker۲d-v1 بین ACKTR و DDPG را در نظر بگیرید. عملکرد ACKTR کمی بهتر است، اما این عملکرد به دلیل همپوشانی فواصل اطمینان دو الگوریتم معنادار نیست:  $p = ۰.۳۳۴$ ،  $KS = ۰.۴۰$ ،  $p = ۰.۶۹۷$ ، تفاوت درصدی بوت‌استرپ  $(-۰.۴۴\% - ۰.۸۰\%)$ ،  $۱۱۱.۷۲\%$ .

<sup>1</sup> Power

<sup>2</sup> Lift

<sup>3</sup> Uniform

<sup>4</sup> Significance

<sup>5</sup> Gain

<sup>6</sup> Supervised Learning

<sup>7</sup> Corrected

<sup>8</sup> Kolmogorov-Smirnov Test

| Environment    | DDPG              | ACKTR             | TRPO               | PPO               |
|----------------|-------------------|-------------------|--------------------|-------------------|
| HalfCheetah-v1 | 5037 (3664, 6574) | 3888 (2288, 5131) | 1254.5 (999, 1464) | 3043 (1920, 4165) |
| Hopper-v1      | 1632 (607, 2370)  | 2546 (1875, 3217) | 2965 (2854, 3076)  | 2715 (2589, 2847) |
| Walker2d-v1    | 1582 (901, 2174)  | 2285 (1246, 3235) | 3072 (2957, 3183)  | 2926 (2514, 3361) |
| Swimmer-v1     | 31 (21, 46)       | 50 (42, 55)       | 214 (141, 287)     | 107 (101, 118)    |

جدول ۳-۳: میانگین بوت‌استرپ و حدود اطمینان ۹۵٪ برای زیرمجموعه‌ای از آزمایش‌های محیط. در این آزمایش‌ها از ۱۰ هزار تکرار بوت‌استرپ و روش محوری استفاده شده است.

### ۴-۳ نتیجه کلی:

از طریق آزمایش‌های متمرکز بر روش‌های گرادینان سیاست برای کنترل پیوسته، ما مشکلات قابلیت بازتولید در یادگیری تقویتی عمیق را بررسی کردیم. دریافتیم که هم منابع درونی (مانند دانه‌های تصادفی و ویژگی‌های محیط) و هم منابع بیرونی (نظیر ابرپارامترها و کدهای پایه) عدم قطعیت می‌توانند به دشواری‌های قابلیت بازتولید الگوریتم‌های پایه کمک کنند. بر اساس نتایج آزمایش‌ها، برخی توصیه‌های کلی را می‌توان ارائه داد: ابرپارامترها می‌توانند اثرات متفاوتی بر الگوریتم‌ها و محیط‌ها داشته باشند. بنابراین یافتن مجموعه‌ای که حداقل عملکرد گزارش شده اولیه الگوریتم‌های پایه را تکرار کند، مهم است. همچنین پیاده‌سازی‌های جدید برای مقایسه باید با نتایج کد پایه اصلی مطابقت داشته باشند. با توجه به واریانس بالا در میان آزمایش‌ها و دانه‌های تصادفی، باید آزمایش‌های متعددی با دانه‌های مختلف انجام شود. استفاده از آزمون‌های معنی‌دار مناسب برای تعیین اینکه آیا بازده‌های بالاتر نشان‌دهنده عملکرد بهتر هستند یا نه، ضروری است. گزارش همه جزئیات از جمله ابرپارامترها، پیاده‌سازی، تنظیمات آزمایشی و روش‌های ارزیابی برای هم روش‌های مقایسه پایه و هم کارهای جدید، مهمترین گام برای قابلیت بازتولید است. بدون انتشار پیاده‌سازی‌ها و جزئیات مربوطه، تلاش‌های بیهوده برای تکرار کارهای برتر، جامعه را آزار خواهد داد و پیشرفت را کند می‌کند.



## پیوست اول

### شبکه های عصبی و توابع فعال سازی

توابع فعال سازی، برای تعیین خروجی شبکه های عصبی، مانند بله یا خیر استفاده می شوند. این توابع، یک نگاشت بین مقادیر حاصل (خروجی) و بازه ۰ تا ۱، یا ۱- تا ۱ (بسته به تابع) ایجاد می کنند. توابع فعال سازی، به دو دسته خطی<sup>۱</sup> و غیر خطی<sup>۲</sup> تقسیم بندی می شوند. توابع فعال سازی غیر خطی، پر استفاده ترین توابع فعال سازی هستند، زیرا این ویژگی، تعمیم یا انطباق مدل با انواع داده ها و تمایز بین خروجی ها را برای مدل آسان می کند. در اینجا، به معرفی برخی از توابع فعال سازی غیر خطی پر کاربرد می پردازیم:

- **تابع Sigmoid:** منحنی این تابع، شبیه به شکل حرف S است. دلیل اصلی استفاده از این تابع این است که برد آن، بازه (۰،۱) است. از آنجایی که احتمال هر چیزی مقداری بین ۰ و ۱ است، بنابراین این تابع به ویژه برای مدل هایی که باید احتمال را به عنوان خروجی پیش بینی کنیم، انتخاب درستی است و استفاده می شود. این تابع همچنین می تواند باعث شود که شبکه عصبی در زمان آموزش، گیر کند.
- **تابع Tanh:** منحنی این تابع نیز، شبیه به شکل حرف S است و محدوده آن، بازه (۱،-۱) است. مزیت این تابع این است که ورودی های منفی، به مقادیر به شدت منفی و ورودی های صفر، به مقادیر نزدیک به صفر، نگاشت می شوند. تابع Tanh، عمدتاً برای طبقه بندی<sup>۳</sup> استفاده می شود.
- **تابع ReLU:** از آنجایی که این تابع در تقریباً تمام شبکه های عصبی پیچشی<sup>۴</sup> و یادگیری عمیق استفاده می شود، در حال حاضر پر استفاده ترین تابع فعال سازی است. همان طور که در شکل ۱-آ می بینید، این تابع برابر ۰ است وقتی ورودی آن کمتر از ۰ (منفی) باشد و برابر ورودی اش است وقتی ورودی آن بیشتر یا مساوی ۰ (نامنفی) باشد. مسئله اصلی این تابع این است که تمام مقادیر منفی، بلافاصله ۰ می شوند، بنابراین نگاشت خوبی از مقادیر منفی ایجاد نمی شود که این موضوع توانایی مدل را برای برازش<sup>۵</sup> یا آموزش صحیح از داده ها کاهش می دهد. این مشکل، در واقع همان مشکل رلو مرده است که به آن اشاره کردیم.
- **تابع Leaky Relu:** این تابع، تلاشی برای حل مشکل رلو مرده است. همان طور که در شکل ۱-آ می بینید، وجود نشت<sup>۶</sup> به افزایش محدوده تابع کمک می کند. معمولاً، مقدار a (شیب خط در قسمت منفی) برابر ۰.۰۱ در نظر گرفته می شود.

<sup>1</sup> Linear

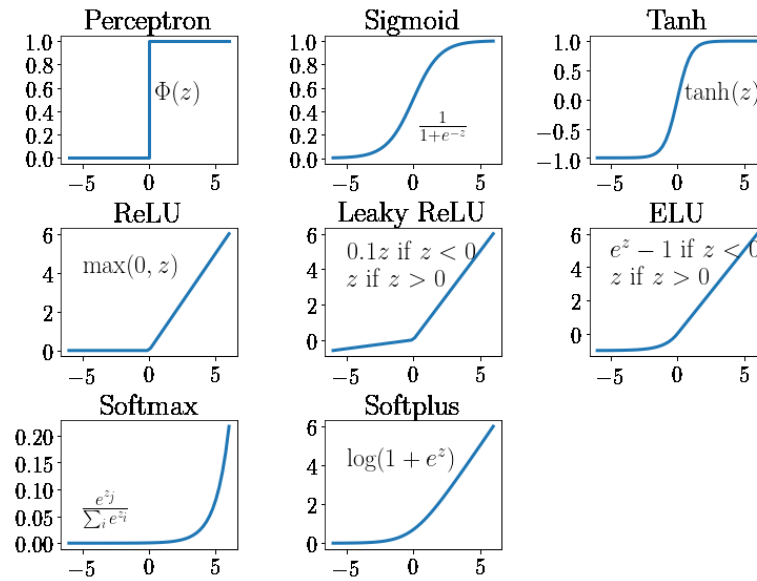
<sup>2</sup> Non-Linear

<sup>3</sup> Classification

<sup>4</sup> Convolutional

<sup>5</sup> Fitting

<sup>6</sup> Leak



شکل آ-۱: برخی توابع فعال‌سازی پرکاربرد، به همراه رابطه و نمودار آن‌ها

## مراجع

- [1] Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D., "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [2] Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al., "Mastering the game of go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484-489, 2016.
- [3] Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M., "Playing atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.
- [4] Vinyals, O., Ewalds, T., Bartunov, S., Georgiev, P., Vezhnevets, A. S., Yeo, M., Makhzani, A., Küttler, H., Agapiou, J., Schrittwieser, J., et al., "Starcraft ii: A new challenge for reinforcement learning," *arXiv preprint arXiv:1708.04782*, 2017.
- [5] Silva, V. d. N., and Chaimowicz, L., "Moba: a new arena for game ai," *arXiv preprint arXiv:1705.10443*, 2017.
- [6] Islam, R., Henderson, P., Gomrokchi, M., and Precup, D., "Reproducibility of benchmarked deep reinforcement learning tasks for continuous control," in *ICML Reproducibility in Machine Learning Workshop*, 2017.
- [7] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O., "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [8] Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P., "Trust region policy optimization," in *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, 2015.
- [9] Duan, Y., Chen, X., Houthoofd, R., Schulman, J., and Abbeel, P., "Benchmarking deep reinforcement learning for continuous control," in *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, 2016.
- [10] Whiteson, S., Tanner, B., Taylor, M. E., and Stone, P., "Protecting against evaluation overfitting in empirical reinforcement learning," in *2011 IEEE Symposium on Adaptive Dynamic Programming And Reinforcement Learning, ADPRL 2011*, Paris, France, Apr. 12-14, 2011, pp. 120-127.

- [11] Machado, M. C., Bellemare, M. G., Talvitie, E., Veness, J., Hausknecht, M., and Bowling, M., “Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents,” *arXiv preprint arXiv:1709.06009*, 2017.
- [12] Wagstaff, K., “Machine learning that matters,” *arXiv preprint arXiv:1206.4656*, 2012.
- [13] Boulesteix, A.-L., Lauer, S., and Eugster, M. J., “A plea for neutral comparison studies in computational sciences,” *PloS One*, vol. 8, no. 4, p. e61562, 2013.
- [14] Stodden, V., Leisch, F., and Peng, R. D., *Implementing reproducible research*, CRC Press, 2014.
- [15] Bouckaert, R. R., and Frank, E., “Evaluating the replicability of significance tests for comparing learning algorithms,” in *Proc. PAKDD*, 2004, pp. 3-12.
- [16] Bouckaert, R. R., “Estimating replicability of classifier learning experiments,” in *Proceedings of the 21st International Conference on Machine Learning (ICML)*, 2004.
- [17] Vaughan, R., and Wawerla, J., “Publishing identifiable experiment code and configuration is important, good and easy,” *arXiv preprint arXiv:1204.2235*, 2012.
- [18] Sutton, R. S., and Barto, A. G., “Reinforcement Learning: An Introduction”, 2nd ed., in progress, 2014, 2015.
- [19] Achiam, J., “Spinning Up Documentation”, Release, Feb. 07, 2020.
- [20] Wu, Y., Mansimov, E., Liao, S., Grosse, R., and Ba, J., “Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation,” *arXiv preprint arXiv:1708.05144*, 2017.
- [21] Todorov, E., Erez, T., and Tassa, Y., “Mujoco: A physics engine for model-based control,” in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2012)*, Vilamoura, Algarve, Portugal, Oct. 7-12, 2012, pp. 5026-5033.
- [22] Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W., “OpenAI gym,” *arXiv preprint arXiv:1606.01540*, 2016.
- [23] Sutton, R. S., McAllester, D. A., Singh, S. P., and Mansour, Y., “Policy gradient methods for reinforcement learning with function approximation,” in *Advances in Neural Information Processing Systems (NIPS)*, 2000.
- [24] Schulman, J., Moritz, P., Levine, S., Jordan, M., and Abbeel, P., “High-dimensional continuous control using generalized advantage estimation,” *arXiv preprint arXiv:1506.02438*, 2015.
- [25] Plappert, M., Houthoofd, R., Dhariwal, P., Sidor, S., Chen, R., Chen, X., Asfour, T., Abbeel, P., and Andrychowicz, M., “Parameter space noise for exploration,” *arXiv preprint arXiv:1706.01905*, 2017.
- [26] Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K., “Asynchronous methods for deep reinforcement learning,” in *International Conference on Machine Learning*, 2016, pp. 1928-1937.
- [27] Gu, S., Lillicrap, T., Ghahramani, Z., Turner, R. E., Schölkopf, B., and Levine, S., “Interpolated policy gradient: Merging on-policy and off-policy gradient estimation for deep reinforcement learning,” *arXiv preprint arXiv:1706.00387*, 2017.

- [28] Xu, B., Wang, N., Chen, T., and Li, M., “Empirical evaluation of rectified activations in convolutional network,” arXiv preprint arXiv:1505.00853, 2015.
- [29] Rajeswaran, A., Lowrey, K., Todorov, E., and Kakade, S., “Towards generalization and simplicity in continuous control,” arXiv preprint arXiv:1703.02660, 2017.
- [30] Kohavi, R., et al., “A study of cross-validation and bootstrap for accuracy estimation and model selection,” in IJCAI, vol. 14, 1995.
- [31] Nadeau, C., and Bengio, Y., “Inference for the generalization error,” in Advances in Neural Information Processing Systems (NIPS), 2000.
- [32] Hofer, L., and Gimbert, H., “Online reinforcement learning for real-time exploration in continuous state and action Markov decision processes,” *arXiv preprint arXiv:1612.03780*, 2016.
- [33] Mandel, T., Liu, Y.-E., Brunskill, E., and Popovic, Z., “Offline evaluation of online reinforcement learning algorithms,” in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2016.
- [34] Yuan, K.-H., and Hayashi, K., “Bootstrap approach to inference and power analysis based on three test statistics for covariance structure models,” *British Journal of Mathematical and Statistical Psychology*, vol. 56, no. 1, pp. 93-110, 2003.
- [35] Bouckaert, R. R., and Frank, E., “Evaluating the replicability of significance tests for comparing learning algorithms,” in *Proceedings of the 8th Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD)*, 2004, pp. 3-12, Springer.
- [36] Wilcox, R., “Kolmogorov–Smirnov test,” in *Encyclopedia of Biostatistics*, 2005.

# Deep Reinforcement Learning

**Parniyan Malekzadeh**

p.malekzadeh@ec.iut.ac.ir

June 15, 2024

Department of Electrical and Computer Engineering

Isfahan University of Technology, Isfahan 84156-83111, Iran

Degree: B.Sc.

Language: Farsi

**Supervisor: Prof. Mohammad Hossein Manshaei (manshaei@cc.iut.ac.ir)**

## **Abstract:**

In recent years, significant progress has been made in solving challenging problems across various domains using deep reinforcement learning (RL). Reproducing existing work and accurately judging the improvements offered by novel methods is vital to sustaining this progress. Unfortunately, reproducing results for state-of-the-art deep RL methods is seldom straightforward. In particular, non-determinism in standard benchmark environments, combined with variance intrinsic to the methods, can make reported results tough to interpret. Without significance metrics and tighter standardization of experimental reporting, it is difficult to determine whether improvements over the prior state-of-the-art are meaningful. In this paper, we investigate challenges posed by reproducibility, proper experimental techniques, and reporting procedures. We illustrate the variability in reported metrics and results when comparing against common baselines and suggest guidelines to make future results in deep RL more reproducible. We aim to spur discussion about how to ensure continued progress in the field by minimizing wasted effort stemming from results that are non-reproducible and easily misinterpreted.

**Key Words:** Deep Reinforcement Learning, Reproducibility, Experimental Techniques, Benchmark Environments.

