

## یادگیری تقویتی

\* یو می

گروه مهندسی برق و کامپیوتر

دانشگاه ایالتی اوهایو

کلمبوس، اوهایو، 42310، ایالات متحده آمریکا

\*mei.103@osu.edu

مفهوم یادگیری ماشینی با هم، و بررسی نحوه کارکرد آنها در شرایط مختلف است.

که به عنوان تقطیر سیاست نیز شناخته می‌شود، DQN علاوه بر اصلاح تکنیک‌های بهینه‌سازی مختلفی برای یادگیری تقویتی وجود دارد. شکل دادن به RS پاداش یکی از آنها است که برای تسریع روند آموزش استفاده می‌شود. هدف تسریع یادگیری استراتژی‌های مؤثر با تغییر سیگنال‌های پاداش است، بنابراین به عامل کمک می‌کند ساختار و ویژگی‌های محیط را بهتر درک کند و در عین حال مقایسه RS راهنمایی صریح‌تری ارائه دهد. این مقاله کارایی تقطیر دانش را با می‌کند و این دو روش را برای دستیابی به نتایج بهتر ترکیب می‌کند.

## II. کار مرتبط

این بخش ادبیات موجود در مورد تقطیر دانش و شکل‌دهی پاداش را در حوزه یادگیری تقویتی عمیق مرور می‌کند. تفاوت‌ها و شباهت‌های بین این تکنیک‌ها نیز مورد بررسی قرار می‌گیرد و ببینی در مورد نقاط قوت و محدودیت‌های آنها ارائه می‌کند.

[7] (OD3) مطالعه قبلی روشی به نام یادگیری تقویتی عمیق روی دستگاه را برای انتقال مؤثر دانش مقطر برای کنترل دستگاه‌های لبه در سیستم‌های انتقال دانش و فشرده OD3 محاسباتی لبه محدود با منابع معرفی کرد. هدف سازی مدل‌های خط‌مشی به طور همزمان در طول آموزش بر روی دستگاه‌های را با پیاده‌سازی آن بر روی OD3 لبه با تقطیر دانش است. این مقاله عملکرد یک سیستم تعبیه شده تجاری با منابع سخت افزاری محدود ارزیابی کرد. نتایج علیرغم استفاده از یک شبکه خط‌مشی بسیار OD3 تجربی نشان داد که کوچکتر، به عملکرد قابل مقایسه با راه حل‌های مبتنی بر ابر دست یافت. علاوه بر این، زمان تمرین برای خط‌مشی لبه به طور قابل توجهی در مقایسه با تمرین را در سیستم‌های محاسباتی لبه‌ای OD3 از ابتدا کاهش یافت. این نتایج اثربخشی که منابع محدود هستند ثابت کرد و بر اهمیت نه تنها تقطیر سیاست کارآمد بلکه بر حفظ منابع نیز تأکید کرد.

یادگیری تقویتی در حوزه‌های مختلف به کار گرفته شده است، اما تحقیقات کنونی اغلب فاقد رویه‌های لازم برای اجرای مؤثر یادگیری تقویتی، مانند شکل دادن به پاداش، یادگیری برنامه درسی، و سایر تکنیک‌ها برای تجزیه و وظایف پیچیده به اجزای کوچکتر هستند. برای سهولت تلاش‌های مهندسی، این مطالعه استاندارد Q اثربخشی این تکنیک‌ها را در تقویت یک عامل یادگیری عمیق بررسی کرد [8]. برجسته‌ترین تکنیک‌ها سپس در سه محیط مختلف، در کنار و اولیه اعمال شدند DQN یک عامل

خلاصه. تقطیر دانش تکنیکی است که دانش را از یک مدل پیچیده تر، که به عنوان مدل معلم نامیده می‌شود، به یک مدل کمتر پیچیده، به نام مدل دانش آموز، تقطیر می‌کند تا نیازها را کاهش دهد و کارایی مدل دوم را بهبود بخشد. علیرغم دستاوردهای اجرای سیاست‌هایی که از طریق این روش‌های (DRL)، متعدد یادگیری تقویتی عمیق پیشرفته در مقیاس بزرگ به دست می‌آیند به دلیل نبود منابع کافی مانع می‌شود. تقطیر دانش مدل‌های یادگیری عمیق تقویتی هنوز به اندازه تقطیر دانش سایر مدل‌های یادگیری عمیق شناخته شده مورد توجه قرار نگرفته است. این مقاله مروری بر تقطیر دانش ارائه می‌کند و اجزای مختلف، روش‌شناسی و کاربردهای بالقوه آن را مورد بحث DRL قرار می‌دهد. این یک تجزیه و تحلیل جامع از کارهای مرتبط را ارائه می‌دهد و پیشرفت‌های ایجاد شده در این زمینه را برجسته می‌کند. یک روش پایه برای تقطیر دانش یادگیری تقویتی عمیق با ارزیابی اثربخشی آن از طریق آزمایشات پیشنهاد شده است. پیاده‌سازی، Reward Shaping، علاوه بر این، یکی دیگر از تکنیک‌های بهینه‌سازی و با تقطیر دانش مقایسه می‌شود و بینش‌هایی را در مورد کاربرد آن‌ها در حوزه‌های مختلف ارائه می‌دهد. نتایج پتانسیل تقطیر دانش و شکل‌دهی پاداش را به عنوان تکنیک‌های قدرتمند برای افزایش عملکرد عوامل یادگیری تقویتی نشان می‌دهد.

کلمات کلیدی-یادگیری عمیق؛ یادگیری تقویتی؛ تقطیر دانش؛ عملکرد از دست دادن

## مقدمه

روشهای (KD) و تقطیر دانش (DRL) هر دو یادگیری تقویتی عمیق محبوب در یادگیری ماشین هستند. یادگیری تقویتی عمیق با سایر روش‌های یادگیری عمیق متفاوت است زیرا خط‌مشی به روز شده شامل یک شبکه عصبی ها DRL یکی از ساده‌ترین (Deep Q-Network (DQN). عمیق است [1,2] است، ایده آن شبیه به مینی بچ است. به جای به روزرسانی هر بار یک اقدام دنباله‌ای از حالت، اقدام و پاداش را به عنوان داده ورودی ذخیره می‌کند و DQN سپس از روش ضرر هوبر برای محاسبه پاداش مورد انتظار استفاده می‌کند. ضرر با استفاده از پاداش انتظار محاسبه می‌شود که در نهایت برای به روز رسانی شبکه استفاده می‌شود.

به عنوان یک تکنیک قوی برای حل مسائل پیچیده DRL، در سال‌های اخیر تصمیم‌گیری ظاهر شده است [3]. با این وجود، موفقیت یادگیری تقویتی عمیق به شدت به توانایی آموزش مدل‌های بزرگ و محاسباتی گران‌قیمت بستگی دارد، که اغلب به منابع و زمان قابل توجهی نیاز دارد. این محدودیت در سناریوهایی که منابع محاسباتی محدود هستند یا تصمیم‌گیری در زمان واقعی مورد نیاز است، چالش‌هایی ایجاد می‌کند که می‌تواند در بسیاری از کاربردهای عملی بازدارنده باشد [4]. تقطیر دانش یک راه حل امیدوارکننده برای رسیدگی ارائه می‌دهد

سه بخش عمده در تقطیر دانش وجود دارد: مدل معلم، مدل دانش آموز و مدل مقایسه.

با پیکربندی 128 نورون که در سه لایه Deep Q-Network مدل معلم از توزیع شده اند استفاده می کند.

ساخته شده است که شامل 96 نرون است که DQN مدل دانش آموز بایک در سه لایه پخش شده اند که تقریباً سه چهارم اندازه مدل معلم است. شبکه از پیش آموزش دیده الگوی معلم به مدل دانش آموز به ارث می رسد.

علاوه بر این، داده های مدل اضافی برای مقایسه ثبت می شود. این مدل مقایسه به طور یکسان با مدل دانش آموز ساخته شده است و فرآیند آموزشی مشابه مدل معلم را طی می کند.

2) رویکرد 2

که معادله T و دانش را نیز می توان با تغییر دمای DRL راندمان تقطیر داده شده را برآورده می کند تحت تاثیر قرار داد. معلم احتمالات هدف نرم  $\diamond\diamond$  و احتمالات هدف نرم دانش آموز  $\diamond\diamond\diamond$  را می توان از لاجیت معلم  $\diamond\diamond\diamond$  و دانش آموز محاسبه کرد که با نشان داده می شوند  $\diamond\diamond\diamond$  و متقابلا.

[illegible]

و DRL پیش‌فرض در عملکرد استاندارد نرم‌افزار 1 است. رفتار تقطیر نیز هنگامی که دما به مقادیر دیگر تغییر می‌کند کنترل می‌شود، Knowledge در حالی که سایر شرایط ثابت می‌مانند.

رویکرد 3) 3)

در شرایط مختلف بررسی شده است. DRL در آزمایشات قبلی، رفتار تقطیر با این حال، آن را با سایر تکنیک های بهینه سازی مقایسه نشده است. شکل دادن به پاداش خطی یکی از آنهاست که با تنظیم یک موقعیت هدف، پاداش خروجی را به صورت خطی تغییر می دهد. هنگامی که موقعیت فعلی به موقعیت هدف، CartPole-v1 نزدیکتر می شود، پاداش به تدریج افزایش می یابد. در محیط پاداش را می توان به معادله زیر نشان داد:

$$\diamondsuit_2 = (\diamondsuit_1 \diamondsuit_2 \diamondsuit_3 - |\diamondsuit_1 \diamondsuit_2|) / (\diamondsuit_1 \diamondsuit_2 \diamondsuit_3 - \diamondsuit_1 \diamondsuit_2) \quad (5)$$

نشان دهنده موقعیت مشاهده شده است، در حالی که  $\phi_0$  موقعیت شروع است. نقطه وسط بین نقطه شروع و نقطه پایان مثبت است. معادله  $r = 1$  فوق فقط زمانی کار می کند که  $\phi > \phi_0$  ، و در غیر این صورت این معادله باید برای محیط  $r = 0$  وقتی  $\phi = \phi_0$  ، و در غیر این صورت های مختلف اصلاح شود، اما رابطه خطی همچنان ادامه دارد.

[illegible]

کمک می‌کند و مقاومت آن را در برابر  $\delta$  اتلاف هویز به به حداقل رساندن MSE نویزدار است افزایش می‌دهد، زیرا تلفات هویز مشابه  $Q$  پرت زمانی که رفتار می‌کند. در MAE بین  $-1$  و  $1$  است یا در غیر این صورت  $\delta$  زمانی که از محاسبه می‌شود  $B$  یک دسته از انتقال حافظه تکرار نمونه

هدف اصلی از شامل شدن دما، نرم کردن نتیجه از دست دادن است. دمای

محدودیت اعمال می شود. IEEE Xplore از UTC دایلود شده در 27 فوریه 2024 ساعت 09:16:54 JAWAHARLAL NEHRU TECHNOLOGICAL UNIV. استفاده مجاز مجاز محدود به نشان‌دهنده مدت  $\gamma$  از توزیع داده‌ها و خط بهترین تناسب، که در آن محور آزمایش و نتیجه IV.

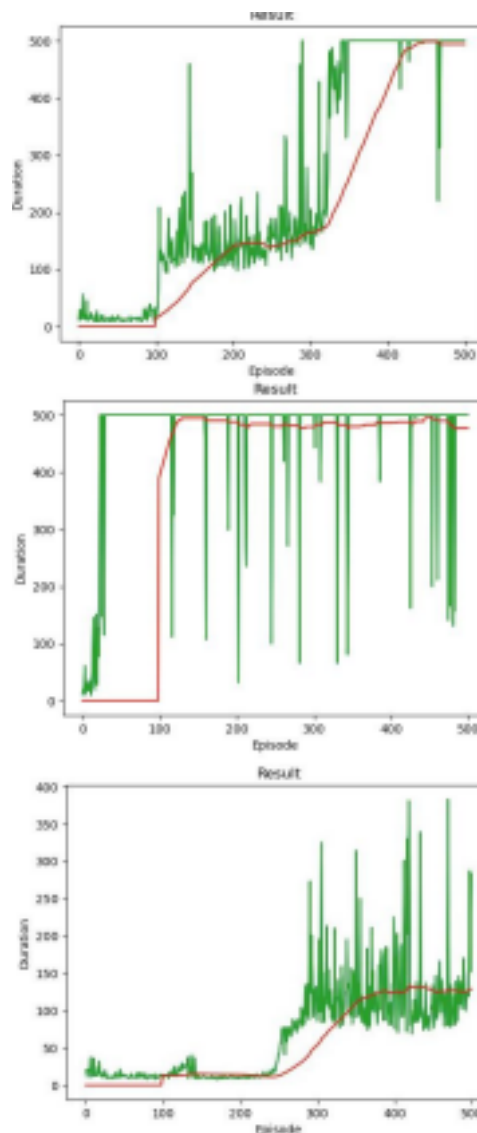
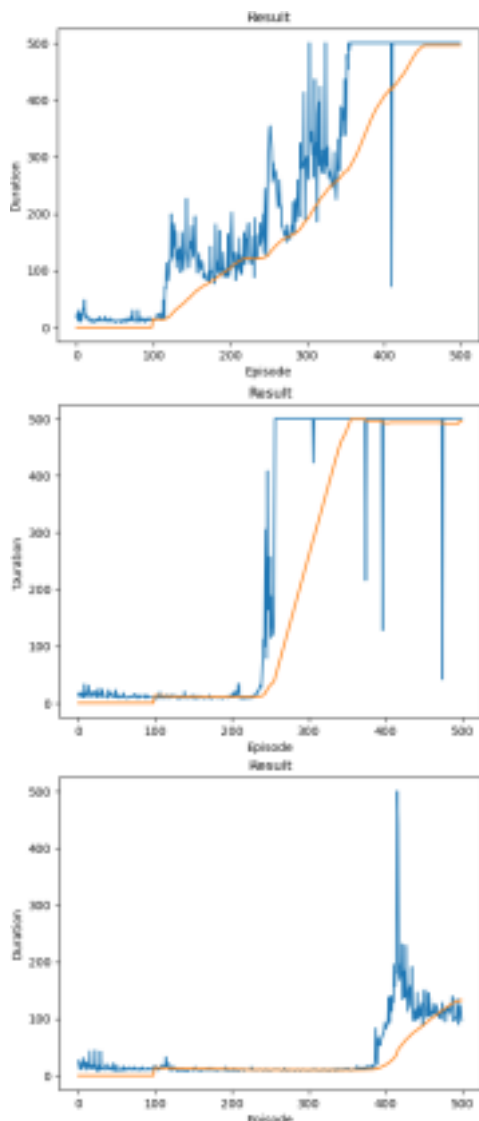
الف. آزمایشات

چهار آزمایش زیر از یادگیری تقوینی عمیق بر چهار جنبه مختلف تمرکز دارد: تقطیر دانش، تغییر تابع از دست دادن، تغییر دما، و شکل‌دهی پاداش. نتایج تجربی از طریق نمودار ارائه شده است

1) آزمایش 1

شکل 1 نتایج را پس از ارزیابی و مقایسه مدل معلم، مدل دانش آموز و مدل

مقایسه مربوطه از نظر عملکرد آنها نشان می دهد، در حالی که سایر شرایط را بدون تغییر نگه می دارد.



(شکل اعتبارات: اصلی) شکل 2. نتایج مدل های مقایسه معلم، دانش آموز و در آزمایش 2 با

(شکل 1. نتایج مدل های مقایسه معلم، دانش آموز و مقایسه در آزمایش 1 (شکل اعتبارات: اصلی)

مدت زمان بعد از 100 قسمت شروع به رشد می کند، در حدود 350 قسمت به بالاترین امتیاز خود یعنی 500 می رسد و تثبیت می شود بعد از آن. تقریباً 250 قسمت طول می کشد تا مازول دانش آموز به بالاترین امتیاز خود دست یابد و در محیط به ثبات برسد. با این حال، مازول مقایسه نمی تواند بالاترین امتیاز را در 500 قسمت به دست آورد.

## آزمایش 2

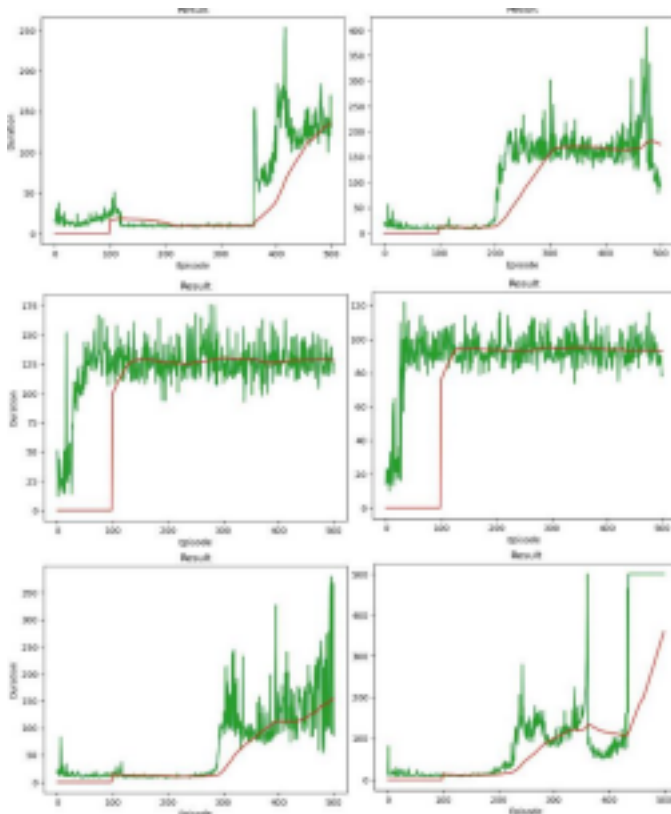
هنگامی که تابع حداکثر دما درگیر باشد، عملکرد تلفات اصلاح می شود. همانطور که در شکل 2 نشان داده شده است تحت تأثیر DRL بنابراین عملکرد قرار می گیرد.

854

محدودیت اعمال می شود. IEEE Xplore از UTC دانلود شده در 27 فوریه 2024 ساعت 09:16:54: استفاده مجاز مجاز محدود به

## آزمایش 3

سپس دما به 5 و 10 تغییر می کند، در حالی که نمودارهای داده ها نتیجه مدل معلم، دانش آموز و مقایسه را در شکل 3 نشان می دهند.



و 10 (شکل اعتبارات: T=5 شکل 3. نتایج مدل های مقایسه معلم، دانش آموز و در آزمایش 3 با اصلی)

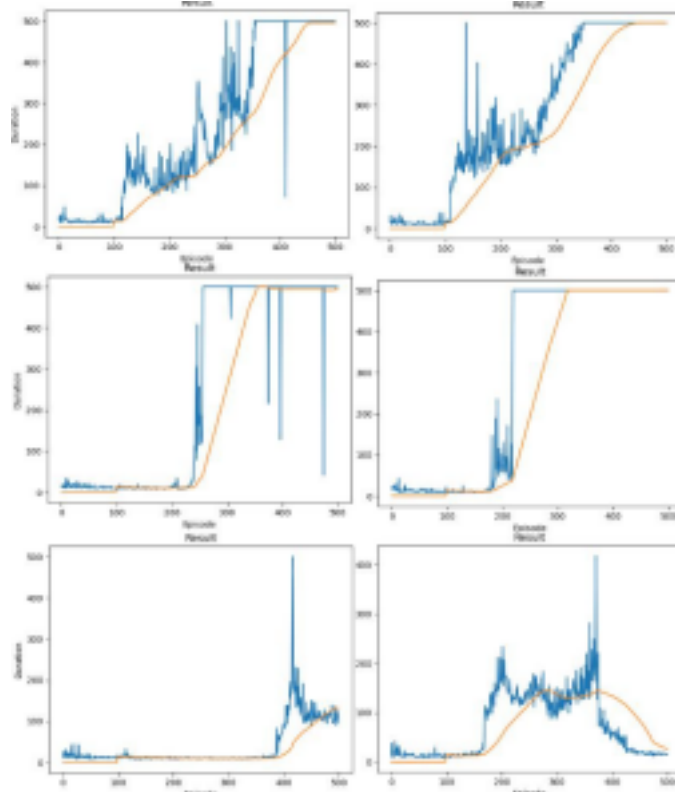
مدل دانش آموز به دلیل تأثیر دمای بالا در مقایسه با مدل معلم شروع به کاهش جزئی در اوج امتیاز می کند، اما همچنان نسبت به مدل مقایسه بسیار بهتر عمل می کند و بالاترین سرعت رشد را حفظ می کند. به طور کلی، تقطیر دانش هنوز در دماهای مختلف قابل اعتماد است.

#### آزمایش 4

آخرین آزمایش شکل دهی پاداش خطی را به سه مدل اصلی اضافه می کند.

نتایج در شکل 4 نشان می دهد که شکل دهی پاداش قابلیت بهبود نرخ یادگیری را دارد، اما زمانی که تعداد کل قسمت ها نسبتاً کوچک باشد، به اندازه DRL تقطیر دانش مؤثر نیست.

با هم اعمال می شوند، مدل دانش آموز در اوایل 200 RS و KD وقتی قسمت به اوج امتیاز می رسد و ثابت می کند که همکاری این دو تکنیک عملکرد DRL را بیشتر افزایش می دهد.



شکل 4. نتیجه مدل های معلم، دانش آموز و مقایسه در آزمایش 4 با طرح اصلی و شکل دادن به پاداش (شکل اعتبارات: اصلی)

## V. بحث

اجرای تقطیر دانش با موفقیت کارایی شبکه دانشجویی را در دو جنبه زیر افزایش می دهد:

اول، مدل دانش آموز پایدارترین فرآیند آموزشی را نشان می دهد و به طور مداوم شبکه خود را در جهت دلخواه بهبود می بخشد. علاوه بر این، در مقایسه با مدل های دیگر، مدل دانش آموز نرخ یادگیری به طور قابل توجهی بالاتری را نشان می دهد.

دوم، تقطیر دانش نه تنها باعث صرفه جویی در زمان در طول آموزش می شود، بلکه دقت نتایج را نیز افزایش می دهد. ماژول دانشجویی بالاترین امتیاز را کسب می کند و از ماژول مقایسه ای که قادر به رسیدن به این سطح از عملکرد نبود عملکرد بهتری دارد.

از آنجایی که تقطیر هیچ تأثیری بر تابع تلفات ندارد، همچنان همانطور که انتظار می رود کار می کند و مدل دانشجویی عملکرد بهتری دارد. بدون شک، عملکرد پایدار تقطیر دانش، زمانی که تابع تلفات تغییر می کند، نسبت به سایر ترندهای تقطیر مفیدتر است.

نتیجه آزمایش همچنین ثابت می کند که وقتی الگوریتم مدل معلم اصلاح می شود، تقطیر دانش کارایی خود را از دست نمی دهد. این ویژگی باعث می شود که تقطیر دانش در برخی از سناریوها مفیدتر از بهینه سازهای سنتی باشد.

دمای بالاتر می تواند روند تمرین را تسریع کند، اما به قیمت از دست دادن برخی جزئیات در داده های ویژگی است. برعکس، دمای پایین تر تمام داده های ویژگی را حفظ می کند.

در مورد دید رایانه و تشخیص الگو، صفحات 5008-5017 IEEE/CVF مقالات کنفرانس 2021.

- [7] I. Jang, H. Kim, D. Lee, Y. S. Son, و S. Kim, تقویت یادگیری برای یادگیری تقویتی، IEEE، عمیق روی دستگاه در سیستم‌های محاسباتی لایه محدود منابع. دسترسی جلد. 8، ص 146588-146597، 2020.
- [8] A. Kanervisto, C. Scheller, Y. Schraner, و V. Hautamäki, تقویتی تقطیر برای بازی های ویدیویی. در کنفرانس در مورد بازی ها، IEEE 2021 تقویتی تقطیر برای بازی های ویدیویی. در کنفرانس صفحات 04-01، 2021.
- [9] Y. Hu, W. Wang, H. Jia, Y. Wang, و همکاران. آموزش استفاده از پاداش‌های، پیشرفت‌ها در سیستم‌های پردازش اطلاعات شکل‌دهی: رویکردی جدید در شکل‌دهی پاداش. پیشرفت‌ها در سیستم‌های پردازش اطلاعات، جلد. 33، صفحات 15931-15941، 2020.
- [10] CartPole-v1: نشانی اینترنتی [https://www.gymnasium.dev/environments/classic\\_control/cart\\_pole/](https://www.gymnasium.dev/environments/classic_control/cart_pole/)، آخرین دسترسی: 2023/09/24

نیز از طریق دماها و عملکردهای KD هم ترکیب شوند. سازگاری و اثربخشی تلفات مختلف بررسی می شود

اگرچه به نظر می‌رسد تقطیر دانش عملکرد بهتری نسبت به شکل‌دهی پاداش به عنوان یک تکنیک RS نمی‌تواند جایگزین KD، دارد CartPole در محیط توانایی KD. بهینه‌سازی شود، زیرا آنها حوزه‌های تخصصی متفاوتی دارند را می‌توان به RS. کاهش سخت افزار مورد نیاز در کنار بهبود کارایی را دارد طور مستقیم طراحی و پیاده سازی کرد و در زمان اضافی برای آموزش مدل دانش آموز صرفه جویی کرد

(MountainCar-v0) تکرار آزمایش های قبلی در محیط های مختلف نتایج مشابهی را به همراه دارد. تقطیر دانش مدل های خیره (Pendulum-v1) بهتری را ایجاد می کند و تحت تاثیر تغییر تابع تلفات و دما قرار نمی گیرد. اثر منجر به بهبودهای قابل توجهی در فرایند آموزش و عملکرد RS و KD ترکیبی کلی مدل می شود

در آزمایش‌های بالا، تقطیر در یادگیری تقویتی عمیق در توابع مختلف دما و افت اجرا شده است، در حالی که شکل‌دهی پاداش نیز اعمال شده است. با این وجود، هنوز تکنیک های بهینه سازی و تقطیر بسیار زیادی وجود دارد، مانند گرادینان خط مشی، یادگیری برنامه درسی و سلسله مراتب دستی. زنجیره پیچیده تر از آن تکنیک ها عملکرد مدل را بیشتر بهبود می بخشد. رویکرد دیگر استفاده یا Double Deep Q Learning از یک مدل معلم پیشرفته‌تر است، مانند Transfer Deep Recurrent Q Learning. تقطیر مانند Deep Recurrent Q Learning. نیز روش خوبی برای بررسی است

#### اره. نتیجه

این کار مفهوم تقطیر دانش و شکل‌دهی پاداش اعمال شده بر روی مدل‌های یادگیری تقویتی عمیق را بررسی کرده است. کارهای مرتبط، روش‌شناسی و کاربردهای این تکنیک مورد بحث قرار گرفت. نتایج تجربی ثابت کرد که تقطیر دانش می‌تواند عملکرد شبکه سیاست را افزایش دهد. این یافته‌ها همچنین پتانسیل شکل‌دهی پاداش را به عنوان ابزاری ارزشمند برای بهبود کارایی مدل‌های یادگیری تقویتی عمیق برجسته کردند. علاوه بر این، این مقاله بینش‌هایی را در مورد کاربرد و چالش‌های بالقوه مرتبط با تقطیر دانش و یادگیری تقویتی ارائه کرده است. تحقیقات و کاوش بیشتر در این زمینه‌ها احتمالاً به الگوریتم‌های یادگیری تقویتی کارآمدتر و مؤثرتر و روش‌های تقطیر قوی‌تر منجر می‌شود

#### منابع

- [1] V. François-Lavet, P. Henderson, R. Islam, M. G. Bellemare, و J. Pineau. مقدمه ای بر یادگیری تقویتی عمیق میانی و روندها در یادگیری ماشین، جلد. 11 (4-3)، صص 219-354، 2018.
- [2] International Journal of Computer Vision, vol. 129, 2021 1819-1789. جی. گو، بی. یو، اس جی میبانک و دی. تائو. تقطیر دانش: نظرسنجی
- [3] پی هندرسون، آر اسلام، پی باخمن، جی. پینو، دی. پرکاپ و دی. میگر. یادگیری تقویتی عمیق در مورد هوش مصنوعی، جلد 32 (1)، AAAI که مهم است. در مجموعه مقالات کنفرانس صفحات 8-1، 2018.
- [4] T. T. Nguyen, N. D. Nguyen, و S. Nahavandi, یادگیری تقویتی عمیق برای سیستم در سایبرنتیک، IEEE های چند عاملی: بررسی چالش‌ها، راه حل‌ها و برنامه‌ها. معاملات جلد. 50 (9)، صص 3826-3839، 2020.
- [5] Wang, L., & Yoon, K. J. تقطیر دانش و یادگیری دانش‌آموز-معلم برای هوش بصری: یک در تحلیل الگو و هوش ماشینی، جلد. 44 (6)، IEEE بررسی و دیدگاه‌های جدید. تراکنش‌های صص 3048-3068، 2021.
- [6] P. Chen, S. Liu, H. Zhao, و J. Jia, تقطیر دانش از طریق بررسی دانش. در مجموعه