

مقاله پژوهشی

توسعه منطق خودتطبیقی سیستم‌های خودتطبیق به کمک یادگیری تقویتی عمیق

Doi: 10.30508/kdip.2023.383007.1060

کاظم نیک فرجام^۱

۱- مربی و عضو هیات علمی دانشگاه آزاد اسلامی واحد بیرجند، بیرجند، ایران

تاریخ دریافت: ۱۴۰۱/۱۱/۰۴

تاریخ پذیرش: ۱۴۰۲/۰۵/۱۷

صفحه: ۰۰ - ۰۰

چکیده

یک سیستم خودتطبیق می‌تواند ساختار و رفتار خود را در زمان اجرا، بر اساس درکش از محیط و از خودش و نیازمندی‌هایش، اصلاح کند. یکی از عناصر کلیدی در توسعه این سیستم‌ها، منطق خودتطبیقی آن است که زمان و نحوه تطبیق سیستم را رمزگذاری می‌کند. هنگام توسعه منطق تطبیق، مهندسان با چالش عدم قطعیت زمان طراحی مواجه هستند. برای تعریف زمان تطبیق سیستم، باید تمام حالات محیطی بالقوه را پیش بینی کنند. پیش بینی تمام تغییرات محیطی بالقوه اغلب به دلیل اطلاعات ناقص در زمان طراحی، غیرممکن است. یادگیری تقویتی برخط، با یادگیری اثربخشی عملیات تطبیق، از طریق تعامل سیستم با محیط در زمان اجرا، مشکل عدم قطعیت زمان طراحی را برطرف، و توسعه منطق خودتطبیقی را به طور خودکار درمی‌آورد. عناصر یادگیری تقویتی، در حلقه MAPE-K سیستم‌های خودتطبیق ادغام می‌شود. روش‌های یادگیری تقویتی برخط موجود در سیستم‌های خودتطبیق، دانش آموخته شده را در قالب تابع ارزش نمایش می‌دهند و دو نقص دارند که درجه خودکارسازی و توسعه را محدود می‌کند: نیازمند تنظیم دقیق نرخ اکتشاف به صورت دستی هستند. از سوی دیگر برای تقویت توسعه‌پذیری، ممکن است نیاز به کمی‌سازی حالت‌های محیط به صورت دستی باشد. در این مقاله برای خودکارسازی فعالیت‌های فوق از یادگیری تقویتی عمیق، استفاده شد. در این یادگیری، دانش در قالب یک شبکه عصبی، در وزن‌های شبکه عصبی پنهان است. نتایج آزمایش‌ها، از سرعت همگرایی بالای یادگیری حکایت دارد.

واژگان کلیدی: یادگیری تقویتی عمیق، عدم قطعیت، منطق خود تطبیق، سیستم خودتطبیق.

۱- مقدمه

خودتطبیقی، باعث تسهیل در توسعه سیستم شده و سیستم را قادر می‌سازد، نیازمندی‌های کیفی خود را در محیط‌های پویا حفظ کند و در زمان اجرا به شکلی انعطاف‌پذیری عمل نماید (اشوف و زیسمن، ۲۰۱۱؛ صالحی و تحویل‌داری، ۲۰۰۸). سیستم خودتطبیق می‌تواند ساختار، پارامترها و رفتار خود را در زمان اجرا بر اساس درک خود از محیط، و از خود و نیازمندی‌هایش، تغییر دهد.

به عنوان مثال یک فروشگاه وب برخط خودتطبیق را در نظر بگیرید که مهندسان حفظ کارایی سیستم را به عنوان هدف یادگیری بیان کرده‌اند. بنابراین باید کارایی سیستم تحت بارهای کاری^۱ در حال تغییر در زمان اجرا، حفظ شود. سیستم در مواجهه با افزایش ناگهانی بارکاری، به صورت برخط، با غیرفعال کردن ویژگی‌های اختیاری، خود را با شرایط جدید بوجود آمده، تطبیق می‌دهد. یادگیری تقویتی برخط، عملیات توسعه منطق خودتطبیقی (که قبلاً دستی توسط مهندس انجام می‌شود)، را به طور خودکار انجام می‌دهد. به عنوان مثال؛ موتور توصیه سیستم را که یک ویژگی اختیاری است، در هنگام بارکاری زیاد (درخواست‌های زیاد کاربران)، غیرفعال می‌کند، تا منابع کمتری استفاده شده و بتواند به همه درخواست‌های کاربران با حداقل تاخیر پاسخ دهد (کلین و ماگو، ۲۰۱۴). عدم قطعیت زمان طراحی یک چالش مهم در توسعه یک سیستم خودتطبیق است. تعریف اینکه چگونه سیستم باید هنگام مواجهه با یک وضعیت محیطی جدید سازگار شود، نیازمند درک تأثیر دقیق عملیات تطبیقی است و احتمال دارد، در زمان طراحی شناخته شده، نباشد. یکی از روش‌های نوظهور برای رسیدگی به عدم قطعیت زمان طراحی در سیستم‌های

خودتطبیق، استفاده از یادگیری تقویتی برخط است (آمویی، صالحی و تحویل‌داری، ۲۰۰۸؛ مصطفی و ژانگ، ۲۰۱۳؛ عرب نژاد، ۲۰۱۷). با استفاده از این نوع یادگیری، سیستم خودتطبیق می‌تواند از داده‌های عملیاتی واقعی و نتایج بازخوردی که فقط در زمان اجرا در دسترس است، بیاموزد.

برای توسعه سیستم خودتطبیق، مهندسان باید منطق خودتطبیقی را طوری توسعه دهند که زمان و نحوه تطبیق سیستم را به شکل مناسبی نمایش دهد. برای مثال، مهندسان، ممکن است قوانینی به شکل (رویداد-شرایط-عمل) تعریف نموده، همچنین تعیین کند در پاسخ به یک تغییر محیطی، معین کدام عمل تطبیقی اجرا شود. توسعه منطق خودتطبیق مستلزم درک درستی از شرایط سیستم و محیط آن دارد و اینکه عملیات تطبیق بر کیفیت سیستم چه تأثیری می‌گذارند (پلیتو و همکاران، ۲۰۱۴؛ چن و باشمن، ۲۰۱۷). نگرانی مهم دیگر، پیش‌بینی تغییرات بالقوه محیطی بوده که سیستم ممکن است در زمان اجرا با آن مواجه شود. باید تعریف شود که چگونه سیستم می‌تواند خود را در پاسخ به این تغییرات محیطی تطبیق دهد. با این حال، به دلیل عدم قطعیت‌های موجود در محیط و سیستم در حین اجرا، پیش‌بینی تمام تغییرات محیطی بالقوه در زمان طراحی در بیشتر موارد غیرممکن است (رامیرز و جانسون، ۲۰۱۲). علاوه بر این، در حالی که اثر عملیات تطبیق بر روی سیستم ممکن است شناخته شده باشد، اما پیش‌بینی دقیق اثر عملیات تطبیق به دلیل ساده کردن فرضیات ساخته شده در زمان طراحی دشوار است (جمشیدی و کامارا، ۲۰۱۹).

یادگیری ماشین علم طراحی ماشین‌هایی است که با استفاده از داده‌هایی که به آن‌ها داده می‌شود (نمونه‌ها) و

1- Work load

بر تصميم‌گيري سيستم خودتطبيق تأثير مي‌گذارند. به عنوان مثال، يادگيرنده‌اي كه در حين اجرا مصرف انرژي باتري‌ها، مصرف CPU .. را پيش بيني مي‌كند.

● به روز نگه داشتن مدل‌هاي زمان اجرا^{۱۱}: پشتيباني از سيستم با به روز نگه داشتن مدل‌هاي زمان اجرا. به عنوان مثال به روز نگه داشتن مدل كارايي و مدل قابليت اطمينان سيستم خودتطبيق.

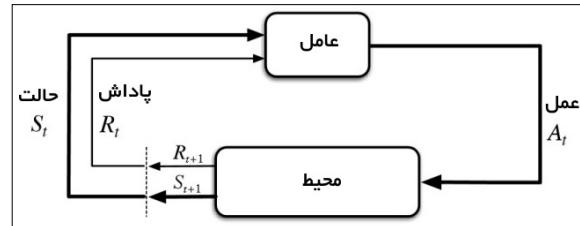
● کاهش فضاهاي تطبيق بزرگ^{۱۲}: پشتيباني از سيستم با کاهش تعداد زيادي از گزینه‌هاي تطبيق (فضاي تطبيق بزرگ) به طوري كه سيستم بتواند سريعتر تصميمات تطبيقی کارآمدتری اتخاذ کند. مثلاً استفاده از يك يادگيرنده براي پيش بيني ويژگي‌هاي كيفي گزینه‌هاي تطبيق براي انتخاب گزینه‌هاي مناسب و تسريع در تجزيه و تحليل سيستم.

● تشخيص و پيش بيني ناهنجاري‌ها^{۱۳}: شناسايي يا پيش بيني ناهنجاري‌ها در رفتار سيستم يا محيط آن كه مربوط به تطبيق است. مثلاً، يك يادگيرنده جريان غيرعادي ترافيك را در يك سيستم مديریت ترافیک خودتطبيق تشخيص مي‌دهد و تهديدات سايري را در يك شبکه ارتباطی شناسايي مي‌كند.

● جمع‌آوری دانش قبلی كه در دسترس نيست^{۱۴}: جمع‌آوری دانش اوليه ناشناخته زمان اجرا براي پشتيباني از تطبيق. مثلاً يادگيرنده يك مدل براي محاسبه سودمندی پيكربندي‌هاي مختلف سيستم (بعد از انجام عمليات تطبيق) مي‌سازد يا يادگيرنده‌اي كه سياست‌هاي مديریتی را بدون هيچ دانش قبلی شناسايي مي‌كند.

● همچنين يادگيري ماشين در سيستم‌هاي خودتطبيق اهداف زير را دنبال مي‌كند (غيبی امید و همكاران، ۲۰۲۱)

تجربيات خودشان عمل كنند، بدون آنكه همه اقدامات با بهره‌گيري از برنامه‌نويسي به آن‌ها ديكته شود. الگوريتم‌هاي يادگيري ماشين سه دسته: يادگيري با نظارت^۱، يادگيري نظارت نشده^۲ و يادگيري تقويتي^۳ است.



شكل (۱): عناصر يك مساله يادگيري تقويتي (ساتن، ۲۰۰۰)

در شكل شماره (۱)، عناصر مساله يادگيري ارايه شده است، كه شامل: محيط^۴؛ جهان فزيكي عامل؛ حالت^۵؛ موقعيت كنوني عامل؛ پاداش^۶؛ بازخورد از محيط؛ سياست^۷؛ روشي براي نگاهت حالت عامل به عمل و ارزش^۸؛ پاداش آينده كه عامل با اقدام به يك عمل در يك حالت خاص به آن دست مي‌يابد.

تاكنون يادگيري ماشين در سيستم‌هاي خودتطبيق در موارد زير استفاده شده است (غيبی امید، دنی، ونيز، و كوين، ۲۰۲۱).

● به روزرسانی و تغيير قوانين و سياست‌هاي تطبيق^۹: به روزرسانی يا تغيير قوانين تطبيق براي پشتيباني از سيستم در هنگام مواجهه با شرايط عملياتی متغير. به طور مثال، سيستم توسط يك يادگيرنده تقويتي به طور مداوم مورد بررسی قرار می‌گیرد، تا سياست‌هاي تطبيقی را براي مقابله با حجم کاری در حال تغيير در سيستم به روز کند.

● پيش بيني و تجزيه و تحليل استفاده از منابع^{۱۰}: بررسی وضعيت منابعی كه توسط سيستم استفاده مي‌شوند و

- 1- Supervised learning.
- 2- Unsupervised learning.
- 3- Reinforcement Learning.
- 4- Environment
- 5- State
- 6- Reward
- 7- Policy
- 8- Value
- 9- Update/Change adaptation rules/policies
- 10- Predict/Analyze resource usage
- 11- Keep runtime models up-to-date
- 12- Reduce large adaptation space
- 13- Detect/Predict anomalies
- 14- Collect unavailable prior knowledge

آن هنوز محدود مانده است. یکی از موانع پیاده‌سازی این نوع یادگیری، لزوم و تکیه این روش برای جستجوی اطلاعات (اکتشاف^۲) در محیط است. برای حل مساله یادگیری تقویتی، باید یک مساله مهم‌تر یعنی؛ توازن بین جست‌وجو (اکتشاف) و استخراج (بهره‌برداری) را مورد بررسی قرار داد. جست‌وجو (اکتشاف)، به معنای یافتن اطلاعات بیشتر درباره محیط و استخراج، به معنای بهره‌برداری از اطلاعات شناخته شده قبلی برای پیشینه‌سازی پاداش است.

هدف عامل یادگیرنده تقویتی پیشینه‌سازی پاداش انباره‌ای (تجمعی) است. بنابراین به نظر می‌رسد، بهتر باشد از اطلاعات شناخته شده قبلی بهره‌برداری کنیم، تا مجموع پاداش‌ها افزایش یابد. با این حال می‌توان در یک تله افتاد، ممکن است اطلاعات ناشناخته‌ای در محیط وجود داشته باشند، که پاداش بسیار بالایی دارند. اما تاکنون استخراج نشده‌اند. مهندسان باید قوانینی را تعیین کنند تا به برقراری توازن و مدیریت این دو مساله کمک کند. از طرفی یادگیری تقویتی تمرکز زیادی روی کارایی زنده سیستم دارد، و زمان برای اکتشاف محدود است. به این دلیل هم نیازمند پیدا کردن یک تعادل مناسب بین اکتشاف چیزهای جدید و بهره‌برداری از دانش اندوخته شده داریم. به این کار موازنه اکتشاف و استخراج^۳ گفته می‌شود. مشکل اصلی اکتشاف این است که تعداد حالات، ممکن است بسیار زیاد یا حتی نامتناهی باشد. از طرفی ممکن است تعداد خروجی‌ها بسیار زیاد و متنوع باشند که در این حالت نیازمند نمونه برداری بسیار گسترده‌ای برای تخمین خروجی نهایی سیستم هستیم. بنابراین در مسئله یادگیری تقویتی، نیازمند یک راهکار هوشمندانه برای اکتشاف هستیم. تصمیم‌گیری‌های تصادفی بدون استفاده از یک توزیع احتمال برآورد شده، معمولاً کارایی بسیار ضعیفی دارد. در دنیای واقعی، کار آسانی نیست که به طور پیوسته، بهترین عملیات

۱- بهبود ویژگی‌های کیفی سیستم؛ بهینه‌سازی و حفظ ویژگی‌های کیفی سیستم خودتطبیق در حین اجرا. مثلاً پایین نگه داشتن زمان پاسخ و بالا نگه داشتن قابلیت اطمینان سیستم خودتطبیق تحت تغییر بارکاری^۲ و وقوع رویدادهای غیرمنتظره.

۲- حفظ تعادل بین ویژگی‌های کیفی سیستم و منابع^۳: به طور مثال، نگهداری تأخیر در یک محدوده مشخص و در عین حال حداقل استفاده از منبع محاسباتی مورد نیاز (مانند CPU) برای دستیابی به آن تأخیر.

۳- تعادل بین ویژگی‌های کیفی سیستم و هزینه^۴: به طور مثال، پایین نگه داشتن نرخ شکست یک گردش کار مبتنی بر سرویس زیر یک آستانه مشخص در حالی که هزینه (عملیاتی، مالی و ...) استفاده از خدمات به حداقل برسد.

۴- بهبود در تخصیص منابع سیستم^۵: بهبود (بهینه‌سازی، مدیریت و ...) در تخصیص منابع سیستم. به عنوان مثال، مدیریت استفاده از CPU و حافظه سیستم تحت بارهای کاری نامشخص زمان اجرا در سیستم خودتطبیق.

۵- دفاع در برابر تهدیدات سایبری^۶: خودکارسازی دفاع سایبری در سیستم خودتطبیق با شناسایی و مدیریت تهدیدات (از قبیل؛ نفوذهای، ناهنجاری‌ها و غیره).

یادگیری تقویتی برخط می‌تواند اثربخشی عملیات تطبیق را از طریق تعامل با محیط بیاموزد، سیستم به طور خودکار منطق خود تطبیقی را به کمک یادگیری تقویتی در زمان اجرا یاد می‌گیرد و نیازی به تنظیم دستی توسط مهندس ندارد. به جای اینکه مهندسان مجبور باشند منطق خودتطبیقی، و مسئله یادگیری را به شکلی اعلانی، برحسب اهداف یادگیری که سیستم باید به آن دست یابد؛ بیان کنند. سیستم به طور خودکار با کمک یادگیری تقویتی برخط این عملیات را می‌آموزد.

با وجود پتانسیل‌های یادگیری تقویتی، پیاده‌سازی آن می‌تواند دشوار باشد و به همین علت کاربردهای

- 1- Improve quality
- 2- Work Load
- 3- Balance qualities with resources
- 4- Balance qualities with cost
- 5- Improve resource allocation
- 6- Protect against cyber threats
- 7- Explore
- 8- Exploration vs Exploitation trade-off

خودتطبيق (برنامه وب خودتطبيق) نشان می‌دهد.

۲- مبانی نظری

همانطور که ذکر شد، یادگیری تقویتی به صورت موثر کمک می‌کند تا مهندسی منطق خودتطبيقی سیستم، به طور خودکار انجام شود. به طور کلی، یادگیری تقویتی از طریق تعاملات عامل با محیط خود، اثربخشی اقدامات عامل را می‌آموزد (ساتن و بارتو، ۲۰۱۸).

در این نوع یادگیری عامل در مرحله زمانی t ، عملی را در حالت محیطی S_t انجام می‌دهد. در نتیجه، حالت سیستم در مرحله زمانی $t+1$ به S_{t+1} تبدیل می‌شود و عامل برای اجرای عمل، پاداش r_{t+1} را دریافت می‌کند. هدف یادگیری تقویتی، بهینه‌سازی پاداش‌های تجمعی است. هنگام بکارگیری یادگیری تقویتی برای سیستم‌های خودتطبيق، «عمل» به معنای عمل تطبیقی مشخص (مانند اصلاح ساختار، پارامترها یا رفتار سیستم)، و «عامل» نقش منطق خودتطبيقی را بر عهده می‌گیرد. محیط، سیستمی که در زمان اجرا باید تطبیق داده شود را شامل می‌شود.

بطور کلی سه رویکرد به یادگیری تقویتی وجود دارد:

الف) رویکرد ارزش محور: در این رویکرد، هدف بهینه‌سازی تابع ارزش است. تابع ارزش، تابعی است که پاداش بیشینه آینده را مشخص می‌کند که عامل در هر حالت دریافت می‌کند. ارزش هر حالت برابر است با ارزش کل پاداشی که عامل می‌تواند انتظار داشته باشد در آینده با آغاز از آن حالت جمع‌آوری کند. عامل از این تابع ارزش برای انتخاب آنکه کدام حالت در هر گام انتخاب شود، استفاده می‌کند. عامل حالتی با بیشترین ارزش را انتخاب می‌کند.

ب) رویکرد سیاست محور: در یادگیری تقویتی سیاست محور، هدف، بهینه‌سازی تابع سیاست بدون استفاده از تابع ارزش است. سیاست، چیزی است که رفتار عامل را در یک زمان داده شده، تعیین می‌کند. عامل یک تابع سیاست را می‌آموزد. این امر به او کمک می‌کند تا هر حالت را به بهترین عمل ممکن نگاشت کند. دو دسته سیاست وجود دارد. سیاست قطعی که برای یک

انتخاب و انجام شوند؛ چراکه محیط پیوسته در حال تغییر است. یکی دیگر از مشکلاتی که سر راه یادگیری تقویتی وجود دارد، زمان و منابع محاسباتی که لازم است تا اطمینان حاصل کنیم، یادگیری به درستی انجام شده است. هر چه محیط بزرگ‌تر باشد، به زمان و منابع بیشتری برای فرآیند آموزش الگوریتم نیاز است.

روش‌های یادگیری تقویتی برخط موجود برای توسعه سیستم‌های خودتطبيقی اغلب از جدول جستجو برای نشان دادن دانش آموخته شده، استفاده می‌کنند. لذا مهندسان سیستم را ملزم می‌کنند تا به صورت دستی حالات محیط را (برای توسعه پذیری در صورتی که محیط دارای تعداد حالت‌های زیادی است)، کمی‌سازی (گسسته‌سازی) کنند. این فعالیت دستی ممکن است گران و به طور بالقوه غیرقابل اعتماد باشد (جمشیدی و کامار، ۲۰۱۹). ایده اصلی در این مقاله این است که فعالیت‌های دستی فوق‌الذکر به طور خودکار با استفاده از یادگیری تقویتی مبتنی بر سیاست^۱ (یا خط‌مشی)، به عنوان یک نوع متفاوت از یادگیری تقویتی برای توسعه سیستم‌های خودتطبيق انجام شوند (نچمن و نوروزی، ۲۰۱۷؛ ساتن، ۲۰۰۰). همچنین موازنه بین اکتشاف و استخراج به طور خودکار در سیستم انجام شود. در یادگیری تقویتی مبتنی بر سیاست، به جای جدول، برای نگهداری دانش آموخته شده، از یک شبکه عصبی مصنوعی استفاده شده و به کمک یک یادگیرنده تقویتی مبتنی بر سیاست دانش در حین اجرا و گام به گام آموخته شده و در قالب وزن‌های شبکه عصبی نمایش می‌یابد (نوروزی و نچمن، ۲۰۱۷). رویکرد پیشنهادی از لحاظ مفهومی، رسمی و فنی، یادگیری تقویتی مبتنی بر سیاست را در مدل مرجع سیستم‌های خودتطبيق شناخته شده، ادغام و یکپارچه می‌کند. این کار، استفاده از یادگیری تقویتی برخط را برای توسعه سیستم‌های خودتطبيق ساده می‌کند، زیرا نیاز به کمی کردن حالت‌های محیطی به صورت دستی ندارد. همچنین نیاز به تنظیم دقیق نرخ اکتشاف به طور دستی نیست. در این مقاله امکان‌سنجی و کاربرد یادگیری تقویتی برخط مبتنی بر سیاست را، روی یک سیستم اطلاعاتی

1- Policy based Reinforcement Learning

2- Value based

3- Policy based

۲۰۱۴؛ مصطفی و ژانگ (۲۰۱۴).

دوم اینکه، اندازه جدول جستجو مستقیماً به تعداد حالت‌های محیطی که باید ذخیره شوند، بستگی دارد و اندازه جدول با افزایش تعداد متغیرهای حالت به طور تصاعدی افزایش می‌یابد. و جهت نگهداری، نیازمند حافظه بالایی است. در نتیجه، راه حل جدولی از مقیاس‌پذیری ضعیفی رنج می‌برد. فرآیند یادگیری برای اینکه بتواند به طور مؤثر یاد بگیرد، باید برای همه ورودی‌های جدول، تمام داده‌های لازم را جمع‌آوری کند (کلین ماگو، ۲۰۱۴؛ مصطفی و ژانگ، ۲۰۱۴).

یک روش متداول برای حل این محدودیت‌ها، کمی‌سازی حالت‌های محیطی پیوسته با تعریف تعداد کمی از حالت‌های محیطی گسسته است. این کمی‌سازی یک فعالیت دستی است که باید توسط مهندس سیستم انجام شود. بنابراین ممکن است گران و به طور بالقوه غیرقابل اعتماد (خطای انسانی) باشد. همان طور که در مطالب قبلی اشاره شد، یادگیری تقویتی در حین اجرا و به صورت تدریجی میزان اثربخشی اقدامات عامل را از طریق تعامل عامل با محیط یاد می‌گیرد و بدین ترتیب منطق خودتطبیقی سیستم به طور خودکار حین اجرا ساخته می‌شود (ساتن و بارتو، ۲۰۱۸).

حالت‌های محیطی می‌تواند به صورت درشت دانه^۳، یا خیلی ریزدانه^۴ کمی‌سازی شوند، و ممکن است به دلیل اندازه بسیار بزرگ جدول جستجو منجر به مقیاس‌پذیری ضعیف شود. دانستن وسعت فضای حالت به این معنی است که توسعه دهندگان در زمان طراحی می‌توانند تمام تغییرات محیطی بالقوه‌ای که سیستم ممکن است در زمان اجرا با آن مواجه شود، را پیش‌بینی کنند. از طرفی نیز احتمال دارد، وسعت فضای حالات (یعنی مجموعه همه حالات‌ها)، به دلیل عدم قطعیت در زمان طراحی، ناشناخته بماند. از این رو تعیین کران‌های پایین و بالای حالت‌های گسسته امکان‌پذیر نیست.

حالت داده شده همیشه عمل مشابهی را باز می‌گرداند. و سیاست تصادفی که برای هر یک از اعمال یک توزیع احتمالی در نظر می‌گیرد. همان طور که مشهود است، سیاست مستقیماً بهترین عمل برای هر حالت را بیان می‌کند.

ج) رویکرد مدل محور: در یادگیری تقویتی مدل محور، ابتدا محیط مدل می‌شود. این یعنی مدلی از رفتار محیط ساخته می‌شود. مساله مهم در این رویکرد نیاز به مدلی متفاوت برای ارائه هر محیط است.

رویکردهای موجود که از یادگیری تقویتی برای ساخت سیستم‌های خودتطبیق استفاده می‌کنند، اکثراً از یادگیری تقویتی مبتنی بر ارزش (یا مقدار) استفاده می‌کنند. یادگیری تقویتی مبتنی بر مقدار، از یک تابع مقدار برای نمایش دانش آموخته شده، استفاده می‌کند. همیشه عملی که بالاترین ارزش را در یک حالت معین دارد، انتخاب می‌شود. دو نوع مشهور یادگیری تقویتی مبتنی بر مقدار عبارتند از Q-Learning و SARSA که این دو الگوریتم در نحوه به روزرسانی تابع ارزش متفاوت هستند (ساتن و بارتو، ۲۰۱۸). اما رویکردهای یادگیری تقویتی موجود برای سیستم‌های خودتطبیق از دور روش مختلف برای نمایش تابع ارزش (مقدار) استفاده می‌کنند، که در ادامه توضیح داده شده است.

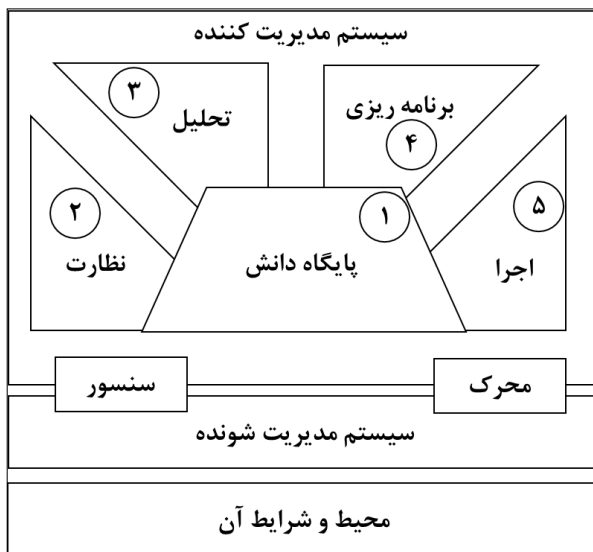
روش اول: تابع ارزش در قالب جدول جستجو^۱

اغلب رویکردهای موجود، مقادیر تابع ارزش را در یک جدول جستجو، ذخیره می‌کنند. اگر چه راه حل جدولی پیاده‌سازی ساده‌ای داشته و به خوبی قابل درک است (کلین و ماگو، ۲۰۱۴)، اما دو محدودیت کلیدی دارد. اول، به دلیل ماهیت گسسته جدول جستجو، تکنیک‌های حل جدولی به فضاهای حالت و عمل گسسته محدود می‌شوند و نمی‌توانند با فضاهای حالت و عمل پیوسته کنار بیایند. این بدان معنی است که حالات محیطی باید قابل شمارش باشند و نمی‌توانند توسط متغیرهای پیوسته (با ارزش واقعی) نمایش داده شوند (کلین و ماگو،

- 1- Model based
- 2- Value Function as Lookup Table
- 3- Coarse-grained
- 4- Fine-grained

حلقه بازخورد سیستم خودتطبيق استفاده می‌شود.

رویکرد یادگیری تقویتی برخط مبتنی بر سیاست رویکرد پیشنهادی، مطابق شکل شماره (۲)، یادگیری تقویتی را در مدل MAPE-K^۴ که یک مدل مرجع شناخته شده برای سیستم‌های خودتطبيق^۵ است، ادغام می‌کند. (فیلو و پورتر، ۲۰۱۷؛ لموس و همکاران، ۲۰۱۳).



شکل (۲): مدل MAPE-K برای سیستم خودتطبيق (آمویی و همکاران، ۲۰۰۸)

همان طور که در شکل شماره (۲) نشان داده شده، این مدل، ساختار مفهومی یک سیستم خودتطبيق را به دو عنصر اصلی تقسیم می‌کند: منطق سیستم (یا سیستم مدیریت شونده^۶) و منطق خودتطبيق (یا سیستم مدیریت کننده^۷). منطق خودتطبيقی به چهار فعالیت مفهومی اصلی تقسیم شده: (نظارت یا مانیتور)، تحلیل، برنامه‌ریزی و اجرا) که از یک پایگاه دانش مشترک استفاده می‌کنند. پایگاه دانش شامل: اطلاعات مربوط به سیستم مدیریت شده و محیط آن است (مانند: مدل‌های زمان اجرای کد شده). این چهار فعالیت شامل: نظارت بر

روش دوم: تقریب زدن تابع ارزش^۱

یک رویکرد جایگزین برای جلوگیری از کمی‌سازی فضای حالت، تقریب زدن تابع مقدار است. به عنوان مثال، با استفاده از تکنیک‌های خطی و غیرخطی (مانند شبکه‌های عصبی مصنوعی). این روش امکان مقابله با فضای حالت بزرگ را با تعمیم وضعیت‌های نادیده فراهم می‌کند. با وجود چنین تقریبی، یادگیری تقویتی مبتنی بر مقدار همچنان با معضل بهره‌برداری- اکتشاف مواجه است (ساتن و بارتو، ۲۰۱۸). برای بهینه‌سازی پاداش‌ها، عملیاتی انتخاب می‌شوند که اثربخشی خود را قبلاً نشان داده و درون پایگاه دانش ذخیره شده‌اند (بهره‌برداری) در عین حال، برای کشف چنین عملیاتی در وهله اول، باید عملیاتی را انتخاب کرد که قبلاً انتخاب نشده‌اند (اکتشاف). یک راه حل معمول برای معضل اکتشاف- بهره‌برداری، روش حریم‌بندی^۲ است. این روش در طول فرآیند یادگیری، حریم‌بندی به طور تصادفی یک عمل را با احتمال افسیلنی (خیلی کم) انتخاب و کاوش می‌کند. چالش مهندس سیستم، تنظیم دقیق تعادل بین نرخ بهره‌برداری و نرخ اکتشاف به منظور اطمینان از همگرایی فرآیند یادگیری است (عرب نژاد و جمشیدی، ۲۰۱۷). به عنوان مثال، مهندس ممکن است مکانیزمی را اجرا کند که در طول زمان اجرا مقدار افسیلن کاهش یافته و در نتیجه میزان اکتشاف کاهش یابد تا همگرایی را تسهیل کند. با این حال، در یادگیری برخط، چالش این است که چه زمانی و چگونه می‌توان دوباره این مقدار را در محیط‌های غیرثابت^۳، (یعنی محیط‌هایی که پویا هستند و در آن نتایج عملیات تطبیقی در طول زمان تغییر می‌کند)، افزایش داد؟ (ساتن و بارتو، ۲۰۱۸). بنابراین درجه خودکارسازی رویکردهای موجود محدود است. مهندس سیستم مجبور است این تنظیمات را دستی انجام دهد که به کار و تلاش زیادی نیاز داشته و به دلیل عدم قطعیت در زمان طراحی، انجام آنها دشوار است. برای غلبه بر این چالش‌ها از یادگیری تقویتی مبتنی بر سیاست در مدل

- 1- Value function Approximation.
- 2- e-Greedy
- 3- Non-stationary
- 4- Monitor-Analys-Planning-Execute-Knowledge
- 5- SelfAdaptive Systems
- 6- Maneged system
- 7- Maneging system

بهینه‌سازی را در بین نمونه‌های فرآیند توزیع می‌کنند. دوتریل و همکاران از Q-Learning برای مدیریت منابع ابری خودمختار (دوتریل و همکاران، ۲۰۱۱) استفاده کردند. آنها فرض می‌کنند که کران بالایی برای متغیرهای حالت وجود دارد. بو و همکاران از Q-Learning برای پیکربندی خودکار ماشین‌های مجازی ابری و برنامه‌های کاربردی استفاده کردند (بو، رایو و اپکسو، ۲۰۱۳) برای تسهیل مقیاس‌پذیری، سه حالت برای گسسته‌سازی تعریف می‌کنند، که نشان دهنده بازه‌های بالا، متوسط و پایین متغیر حالت مربوطه است. عرب نژاد و همکارانش از Q-Learning فازی و SARSA برای ایجاد ابری با خاصیت خود مقیاسی استفاده می‌کند (عرب نژاد و همکاران، ۲۰۱۷) حالات محیطی کمی‌سازی می‌شوند و در نتیجه به مجموعه‌های کوچکی از حالت‌های بیان شده در منطق فازی محدود می‌شوند. مزیت منطق فازی این است که بسیاری از حالت‌ها را می‌توان تنها با چند حالت فازی نشان داد. با این حال، رویکرد آنها هنوز نیازمند شناسایی عناصر فازی «گسسته» در مجموعه فازی است که یادگیری تقویتی بر اساس آن عمل می‌کند. کاپوروسیو و همکاران پیشنهاد استفاده از یادگیری تقویتی مبتنی بر ارزش را برای مونتاز خدمات چند عاملی می‌دهند (کاپوروسیو و آنجلو، ۲۰۱۶). عامل‌ها اطلاعات نظارت بر وضعیت را به اشتراک می‌گذارند و از نمایش جدولی تابع مقدار استفاده می‌کند. سیلواندر پیشنهاد استفاده از Q-Learning با تقریب تابعی از طریق یک شبکه عصبی عمیق^۴ برای بهینه‌سازی فرآیندهای تجاری دارد (سیلواندر، ۲۰۱۹). همه این رویکردها از الگوریتم حریصانه برای مکانیزم اکتشاف استفاده می‌کنند که نیاز به تنظیم دقیق نرخ اکتشاف به صورت دستی دارد. در مقابل، یادگیری ما نیازی به کنترل صریح نرخ اکتشاف ندارد، کاوش به طور خودکار از طریق انتخاب عملیات احتمالاتی انجام می‌شود.

مروری مفهومی بر معماری روش پیشنهادی

شکل شماره (۳)، معماری مفهومی رویکرد پیشنهادی را نشان داده و مشخص می‌کند، چگونه عناصر یادگیری

منطق سیستم و محیط از طریق حس‌گرها، تجزیه و تحلیل داده‌های نظارت شده در مرحله قبل برای تعیین نیاز به تطبیق سیستم، برنامه‌ریزی عملیات تطبیق، و در نهایت اجرای عملیات تطبیقی از طریق محرک‌ها^۱، که در نهایت باعث اصلاح منطق سیستم در زمان اجرا، می‌شود. ایده اساسی پشت یادگیری تقویتی مبتنی بر سیاست، بهینه‌سازی سیاست انتخاب عملیات تصادفی پارامتری شده^۳ است (ناچمن و نوروزی، ۲۰۱۲؛ ساتن ۲۰۰۰). سیاست انتخاب عمل، وضعیت‌ها را به یک توزیع احتمال در فضای عملیات (یعنی مجموعه عملیات تطبیقی ممکن) نگاشت می‌کند.

بدین ترتیب عملیات با نمونه‌گیری از این توزیع احتمال انتخاب می‌شوند. همچنین از یک چرخه یادگیری با تعداد دفعات معین از قبل تعیین شده، برای به روزرسانی سیاست استفاده می‌شود. در پایان هر چرخه یادگیری، وضعیت و پاداش‌های دریافتی، بروز می‌شوند، طوری که توزیع احتمال حاصل به سمتی تغییر کند، که احتمال انتخاب عملیاتی را افزایش می‌دهد که منجر به پاداش تجمعی بالاتری می‌شوند. بیشتر مقالات موجود که از یادگیری تقویتی در سیستم‌های خودتطبیق استفاده می‌کنند، از یادگیری تقویتی مبتنی بر مقدار استفاده کرده‌اند. در این قسمت به چگونگی تعریف تابع مقدار یا همان تابع ارزش در کارهای گذشته می‌پردازیم. (آمویی و همکاران، ۲۰۰۸) از الگوریتم SARSA برای یادگیری عملیات تطبیقی مؤثر در یک برنامه وب خودتطبیق و از تابع مقدار در قالب جدول جستجو استفاده کردند. آنها با تعریف تنها دو مقدار برای هر یک از متغیرهای حالت، حالت‌های محیط را به طور بنیادی کمی می‌کنند. بدین منظور، از مقادیر آستانه‌ای تعریف شده توسط کارشناسان حوزه استفاده می‌شود. هوانگ و همکاران از Q-Learning برای بهینه‌سازی پویای تخصیص منابع در فرآیندهای تجاری استفاده کردند (کپارت و چیس، ۲۰۰۳) علاوه بر بهینه‌سازی تخصیص منابع برای یک نمونه فرآیند، آنها با در نظر گرفتن هزینه‌های منابع سراسری هنگام به روزرسانی جدول ارزش،

- 1- Sensors
- 2- Effectors
- 3- Parametrized stochastic action selection policy
- 4- Deep Network Learning

از دوشبکه عصبی استفاده می‌کند، یکی وظیفه عامل (بازیگر) و دیگری محاسبه پاداش‌ها (منتقد) را انجام می‌دهد.

هدف مورد استفاده برای بهینه‌سازی شبکه عصبی با رابطه ۵ از احتمال انجام عمل تحت سیاست فعلی تقسیم بر احتمال انجام عمل تحت سیاست قبلی استفاده می‌کند.

$$r(\theta) = \frac{\prod_{\theta(at|st)} / \prod_{\theta_{old}(at|st)}}{5}$$

این $r(\theta)$ زمانی بزرگ‌تر از ۱ خواهد بود که احتمال این عمل برای سیاست فعلی بیشتر از سیاست قدیمی باشد. زمانی که احتمال این عمل برای سیاست فعلی کمتر از سیاست قبلی باشد، بین ۰ و ۱ خواهد بود. برخلاف روش‌های گرادیان سیاست، PPO امکان می‌دهد توسط چند عامل یادگیرنده به طور موازی چندین دوره از شیب صعود را روی نمونه‌های مشاهده شده، اجرا کنیم؛ بدون اینکه به روزرسانی‌های سیاست بسیار مخربی ایجاد شود. این ویژگی این امکان را می‌دهد که اطلاعات بیشتری از محیط جمع‌آوری کرد و ناکارآمدی نمونه را کاهش داد. الگوریتم PPO با استفاده از یک تابع برش^۳ از به روزرسانی سیاست‌های خیلی بزرگ جلوگیری می‌کند. به روزرسانی سیاست‌های خیلی بزرگ ممکن است، به این معنی باشد که یادگیری تقویتی جواب بهینه سراسری را از دست داده و در یک بهینه محلی، گیر کرده است. برای نمایش مدل‌های بازیگر - منتقد، این الگوریتم از یک شبکه پرسپترون چند لایه با دو لایه پنهان ۶۴ نورونی استفاده می‌کند (تعداد نرون‌ها در لایه‌ی ورودی و خروجی شبکه به تعداد متغیرهای عملیات و حالت سیستم بستگی دارد).

قطعیت زمان طراحی (چگونگی تأثیر تطبیق بر کیفیت سیستم) نمی‌دانیم. به عنوان مثال، ممکن است درک دقیقی از نحوه عملکرد پیکربندی‌های مختلف سیستم تحت بارهای کاری مختلف نداشته باشیم.

برای انتخاب یک الگوریتم یادگیری تقویتی مبتنی بر سیاست، دونکته اصلی را در نظر می‌گیریم. ابتدا، همان‌طور که فرض می‌کنیم تابع انتقال T را نمی‌دانیم، باید از یک نوع یادگیری بدون مدل برای یادگیری تقویتی مبتنی بر سیاست استفاده کنیم. دوم، برای تسهیل یادگیری برخط، به الگوریتمی نیاز داریم که به طور مداوم سیاست را بدون انتظار برای نتیجه نهایی، (یعنی بدون انتظار برای رسیدن به حالت پایانی) به روز کند. الگوریتم‌های منتقد-بازیگر یک نوع الگوریتم بدون مدل از الگوریتم‌های یادگیری تقویتی مبتنی بر سیاست هستند که در آن دانش به طور مداوم بدون انتظار برای نتیجه نهایی به روزرسانی می‌شود. در این مقاله از بهینه‌سازی سیاست تقریبی^۲ PPO به عنوان الگوریتم منتقد-بازیگر استفاده می‌کنیم (اسچالمن و همکاران، ۲۰۱۷؛ وینز و ایگلسیا، ۲۰۱۵).

الگوریتم بهینه‌سازی سیاست پروگزیمال (PPO) از محبوب‌ترین الگوریتم‌های یادگیری تقویتی توسط تیم OpenAI در سال ۲۰۱۷ معرفی شد. این الگوریتم به جمع‌آوری دسته کوچکی از تجربیات در تعامل با محیط پرداخته و از آن دسته برای به روزرسانی سیاست تصمیم‌گیری استفاده می‌کند. الگوریتم PPO به جای تخصیص مقادیر به جفت‌های حالت-عمل، فضای سیاست‌ها را جستجو می‌کند. همچنین با محدود کردن تغییری که در سیاست خود در هر مرحله ایجاد می‌کند، به ثبات آموزش و بهینه‌سازی شبکه عصبی کمک می‌کند. رایج‌ترین پیاده‌سازی PPO از طریق مدل عامل منتقد است که

```

Algorithm 1 PPO, Actor-Critic Style
1 for iteration=1, 2, ... do
2   for actor=1, 2, ..., N do
3     Run policy  $\pi_{\theta_{old}}$  in environment for T timesteps
4     Compute advantage estimates  $\hat{A}_1, \dots, \hat{A}_T$ 
5   end for
6   Optimize surrogate L wrt  $\theta$ , with K epochs and minibatch size  $M \leq NT$ 
7    $\theta_{old} \leftarrow \theta$ 
8 end for
    
```

شکل شماره (۴): شبه کد الگوریتم PPO سبک عامل منتقد (اسچالمن و همکاران، ۲۰۱۷)

- 1- Actor-Critic.
- 2- Proximal Policy Optimization
- 3- Clipping

پياده سازي و ارزيابي تجربي

برای نشان دادن امکان سنجی و کاربرد یادگیری پیشنهادی، دامنه آزمایش‌ها این است که تحلیل کنیم، آیا سیستم قادر به یادگیری و بهبود منطق خودتطبیقی در زمان اجرا است یا خیر؟ در این مرحله یک تحلیل مقایسه‌ای با رویکردهای یادگیری تقویتی مبتنی بر مقدار انجام ندادیم، چنین مقایسه‌ای فراتر از محدوده مقاله فعلی می‌باشد. چنین مقایسه‌ای نیازمند تغییر و تحلیل دقیق طیف وسیعی از پارامترها برای رویکرد مبتنی بر مقدار است، از جمله تنظیم نرخ اکتشاف، و همچنین سطوح و اشکال مختلف گسسته‌سازی فضای حالت. به ویژه، باید مراقب بود که مقایسه ناعادلانه‌ای انجام نشود. مقایسه ممکن است به شدت تحت تأثیر گسسته‌سازی یا به اصطلاح کوانتیزاسیون انتخابی قرار گیرد. گسسته‌سازی بسیار ریزدانه، ممکن است به این معنی باشد که یادگیری مبتنی بر ارزش همگرایی بسیار کندی خواهد داشت. از طرفی گسسته‌سازی خیلی درشت ممکن است به این معنی باشد که یادگیری مبتنی بر ارزش قادر به تمایز و تفکیک بین حالت‌های مختلف سیستم نیست و در نتیجه نمی‌تواند پاداش‌های تجمعی را بهینه کند. برای آزمایش از یک سیستم وب خودتطبیق حراجی برخط RUBIS-Brownout استفاده می‌کنیم (کلین و ماگیو، ۲۰۱۴). هنگامی که کاربر، مورد خاصی را از بین کالاهای حراجی درخواست می‌کند، موتور توصیه برنامه، لیستی از موارد توصیه شده را بر اساس حراجی‌های گذشته به کاربر ارائه می‌دهد. با

توجه به نیاز موتور توصیه به منابع مختلف جهت راهنمایی و توصیه به کاربر، RUBIS Brownout باید بین دو نیازمندی کیفی موازنه برقرار کند: به حداکثر رساندن تجربه کاربر با ارائه توصیه‌های بیشتر، و در عین حال به حداقل رساندن تأخیر ملاحظه شده توسط کاربران.

بنابراین، نیاز به استفاده از موتور توصیه‌رامی توان با تنظیم یک متغیر بنام متغیر تنظیم تطبیق (دیمر) در محدوده [۰، ۱] که نشان دهنده احتمال فعال شدن موتور توصیه به ازای هر درخواست است، تطبیق داد. بنابراین، مقدار کم این متغیر روی هر دو نیازمندی کیفی تأثیر می‌گذارد. نرخ بالای توصیه، تجربه کاربر را افزایش می‌دهد، اما در عین حال برای اجرای این توصیه و راهنمایی، نیاز به منابع، افزایش یافته و در نتیجه ممکن است تأخیر را افزایش دهد و زمان انتظار کاربر برای پاسخ گرفتن به درخواست‌هایش زیاد شود، باعث نارضایتی کاربر گردد. همان طور که در جدول شماره (۱) نشان داده شده، مساله یادگیری را در قالب فرآیند تصمیم مارکوف MDP تعریف می‌کنیم. تابع پاداش^۱ را به گونه‌ای تعریف می‌کنیم، که الگوریتم یادگیری تعادل و موازنه خوبی بین تأخیر کم و نرخ‌های توصیه بالا پیدا کند. تابع پاداش طوری تعریف شده که پاداش بیشتر، بهتر باشد. هدف سیستم به حداکثر رساندن پاداش تجمعی است. فرض می‌کنیم که رضایت کاربر برای تأخیرهای بالاتر از آستانه λ_{max} (که اینجا ۱۰ میلی ثانیه در نظر گرفته شده) کاهش می‌یابد، بنابراین تأخیرهای بالاتر از λ_{max} جریمه می‌شوند.

جدول (۱): یادگیری تقویتی برخط در برنامه وب خود تطبیق

State $St=(Ut, \alpha_t, \lambda_t)$	$Ut \in \mathbb{N}^+$ $\alpha_t \in [0, 1]$ $\lambda_t \in \mathbb{R}^+$	تعداد درخواست‌های کاربر در لحظه t نرخ توصیه در زمان t تأخیر در زمان t
Action $\alpha_t \in A$	$A = \delta \in [0, 1]$	متغیر تطبیقی (دیمر)
Reward $r(t) = \alpha_t * f(\lambda_t)$	$A_t \in [0, 1]$ $\lambda_t \in \mathbb{R}^+$ $f(\lambda_t) = 1$ if $\lambda_t < \lambda_{max}$ $f(\lambda_t) = 0$ if $\lambda_t > 2 * \lambda_{max}$ else $f(\lambda_t) = - \lambda_t / \lambda_{max} + 2$	$\lambda_{max} = 10 \text{ ms}$

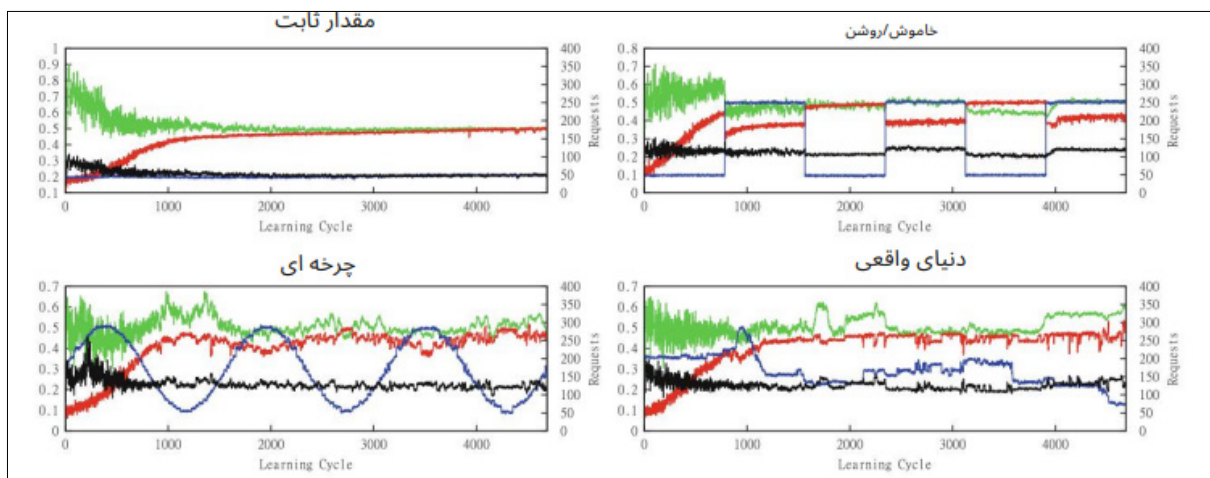
1- Dimmer value

2- Reward Function

مطابق جدول شماره (۱)، حالت سیستم به صورت سه‌گانه‌ای شامل تعداد درخواستهای کاربران، نرخ توصیه و میزان تاخیر در لحظه t تعریف می‌شود. نرخ توصیه همان مقدار متغیر تطبیق است که نسبت عکس با تاخیر دارد. اگر تاخیر در سیستم افزایش یابد، باید نرخ توصیه را کم کنیم و برعکس. بنابراین مقدار پاداش در هر لحظه بر اساس این دو ویژگی که اهداف تطبیقی سیستم است، تعیین می‌شود. مقدار پاداش در زمان t برابر با حاصل ضرب مقدار متغیر تطبیق که همان نرخ توصیه است (عددی بین صفر و یک)، در تابعی از تاخیر سیستم در زمان t است. اگر تاخیر سیستم زیر ده میلی ثانیه (آستانه تحمل تاخیر) باشد، پاداش دریافتی در زمان t برابر است با مقدار متغیر تطبیق یا همان نرخ توصیه و چنانچه تاخیر بالاتر از ۲۰ میلی ثانیه باشد، پاداش دریافتی از محیط صفر است (زیرا خروجی تابع تاخیر صفر است). چنانچه تاخیر از ده به سمت بیست میلی ثانیه حرکت کند میزان پاداش هم به صورت خطی کم می‌شود. یادگیرنده در چرخه سیستم خود تطبیقی یاد می‌گیرد، که در دراز مدت از این حالات (تاخیر بالای ده میلی ثانیه) پرهیز کند. هدف یادگیرنده‌ی تقویتی پیشینه‌کردن پاداش تجمعی است. جهت پیاده‌سازی

مطابق جدول شماره (۱)، حالت سیستم به صورت سه‌گانه‌ای شامل تعداد درخواستهای کاربران، نرخ توصیه و میزان تاخیر در لحظه t تعریف می‌شود. نرخ توصیه همان مقدار متغیر تطبیق است که نسبت عکس با تاخیر دارد. اگر تاخیر در سیستم افزایش یابد، باید نرخ توصیه را کم کنیم و برعکس. بنابراین مقدار پاداش در هر لحظه بر اساس این دو ویژگی که اهداف تطبیقی سیستم است، تعیین می‌شود. مقدار پاداش در زمان t برابر با حاصل ضرب مقدار متغیر تطبیق که همان نرخ توصیه است (عددی بین صفر و یک)، در تابعی از تاخیر سیستم در زمان t است. اگر تاخیر سیستم زیر ده میلی ثانیه (آستانه تحمل تاخیر) باشد، پاداش دریافتی در زمان t برابر است با مقدار متغیر تطبیق یا همان نرخ توصیه و چنانچه تاخیر بالاتر از ۲۰ میلی ثانیه باشد، پاداش دریافتی از محیط صفر است (زیرا خروجی تابع تاخیر صفر است). چنانچه تاخیر از ده به سمت بیست میلی ثانیه حرکت کند میزان پاداش هم به صورت خطی کم می‌شود. یادگیرنده در چرخه سیستم خود تطبیقی یاد می‌گیرد، که در دراز مدت از این حالات (تاخیر بالای ده میلی ثانیه) پرهیز کند. هدف یادگیرنده‌ی تقویتی پیشینه‌کردن پاداش تجمعی است. جهت پیاده‌سازی

شکل شماره (۵)، نتایج را برای انواع الگوهای بار کاری نشان می‌دهد. الگوی چرخه‌ای، الگوی (خاموش / روشن)، الگوی دنیای واقعی و الگوی ثابت.



شکل (۵): رفتار یادگیری تقویتی عمیق مبتنی بر سیاست در سیستم وب خود تطبیق (رنگ بارکاری=سیاه تاخیر=سیاه متغیر تطبیق یا دیمر=سبز پاداش=قرمز)

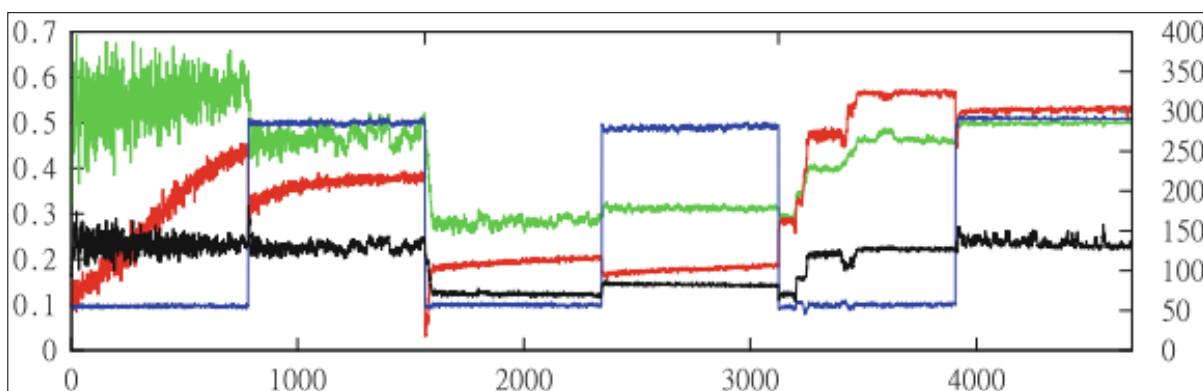
- 1- Constant workload pattern
- 2- On/Off workload pattern
- 3- Cyclic workload pattern

مشاهدات سراسری از سیستم، بسیار با الگوی بارکاری خاموش / روشن قابل مقایسه است، با این تفاوت که میانگین پاداش آهسته تر افزایش و کاهش می یابد، زیرا بارکاری به آرامی تغییر می کند. برای بارکاری دنیای واقعی، الگوریتم می تواند از حالت های تجربه شده قبلی بیاموزد و با تنظیم مقدار کم برای متغیر تطبیق، پاداش را تقریباً در همان سطح نگه دارد، حتی اگر بارکاری در طول زمان تغییر کند. مزیت مهم این یادگیری این است که اگر بار کاری مشابه ای، دوباره تکرار شود، این رویکرد قادر است به سرعت عملیات تطبیقی موثر را پیدا کند.

شکل شماره (۶) نشان می دهد چگونه این نوع یادگیری به طور خودکار در محیط های غیر ثابت^۱ عمل می کند. پس از چرخه یادگیری ۱۵۶۲ (رانش^۲ اول)، منابع محاسباتی ماشین را به نصف کاهش دادیم. این بدان معنی است که برای همان مقدار کم دیمر، سیستم تأخیر بیشتری را تجربه می کند، زیرا منابع محاسباتی کمتری در دسترس است. الگوریتم یاد می گیرد، مقدار متغیر تطبیق (دیمر) را به گونه ای کاهش دهد که آستانه تأخیر (تاخیر از ۱۰ میلی ثانیه بیشتر نشود) نقض نشود. پس از چرخه یادگیری ۳۱۲۵ (رانش دوم) منابع را یک و نیم برابر افزایش دادیم. مجدداً مقدار متغیر تطبیق (دیمر) بر این اساس تنظیم می شوند. این نوع یادگیری قادر است ویژگی غیرایستایی محیط را بدون انجام فرآیند نظارت صریح، تغییرات در منابع محاسباتی و بدون تغییر صریح نرخ اکتشاف، یاد بگیرد.

نتایج نشان می دهد که چگونه این نوع یادگیری تقویتی عامل منتقد مبتنی بر سیاست، سیستم را قادر می سازد تا به مرور زمان خود را تطبیق دهد. سیستم به طور خودکار متغیر تنظیمی برای تطبیق (دیمر) را، بسته به نوع الگوی بارکاری، تنظیم می کند، و در نتیجه تعادل بین تأخیر و افزایش تجربه کاربر با نرخ توصیه و پیشنهادات بیشتر ایجاد می کند (که در افزایش پاداش های تجمعی قابل مشاهده است). در ابتدای فرآیند یادگیری، مقدار کم متغیر تنظیم، واریانس بالایی را برای همه الگوهای بارکاری نشان می دهد، اما پس از مدتی واریانس عملیات تطبیق به وضوح کمتر می شود.

برای الگوی بارکاری ثابت، پاداش پس از حدود ۱۹۵۰ چرخه یادگیری به مقدار ۴۷ صدم همگرا شده و مقدار متغیر تطبیق در حدود نیم، تثبیت می شود و این حالت منجر به بالاترین نرخ توصیه بدون نقض آستانه تأخیر می شود. برای الگوی بارکاری خاموش / روشن، پاداش در طول زمان برای تنظیمات خاموش و همچنین روشن، افزایش می یابد. از تکرار دوم به بعد می توان همگرایی را مشاهده کرد. هنگام مقایسه یادگیری، برای دوره های خاموش و روشن به طور جداگانه، می توان مشاهده کرد که یادگیری تقویتی می تواند از دانش کسب شده در مورد بارهای کاری مشابه در طول زمان استفاده مجدد کند. الگوی بارکاری چرخه ای، مشابه بارکاری خاموش / روشن، نشان می دهد، چگونه یادگیری تقویتی می تواند به کمک پایگاه دانش از قبل به دست آمده تعمیم و گسترش یابد.



شکل (۶): رفتار الگوریتم یادگیری تقویتی عمیق مبتنی بر سیاست بر خط در محیط غیر ثابت (آبی = بارکاری، سیاه = تأخیر، سبز = متغیر تطبیق دیمر، قرمز = پاداش)

- 1- Non_Stationary
- 2- Drift 1

تهدیدات

برای مشاهده اینکه آیا یادگیری تقویتی مبتنی بر سیاست نتایج مورد انتظار را برآورده می‌کند یا خیر، از یک شبکه پرسپترون چند لایه به عنوان یک شبکه عصبی ساده برای نشان دادن سیاست استفاده کردیم. برای حل مشکل کمی‌سازی خوب فضای حالات و تنظیم مناسب نرخ اکتشاف، از تنظیمات پیش فرض برای فرآیندهای شبکه استفاده شد. همچنین، به دلیل ماهیت تصادفی، هر یک از آزمایش‌ها را چندین بار تکرار کردیم، تا اثرات تصادفی بودن را کاهش دهیم. یادگیری تقویتی مبتنی بر سیاست، می‌تواند فضای عمل بزرگی را در حین یادگیری برخط کنترل کند، مشروط بر اینکه فضای عملیات پیوسته باشد. با این حال، بر روی مجموعه‌ای از عملیات غیرپیوسته، یعنی گسسته، کار نمی‌کند. بنابراین نمی‌تواند به عملیاتی که قبلاً دیده نشده بودند، گسترش و تعمیم یابد. به طور معمول، سیستم‌های اطلاعاتی خودتطبیقی دارای فضاهای عملیاتی گسسته بزرگی هستند، مانند سیستم‌های خودتطبیق مبتنی بر ویژگی یا مبتنی بر معماری. به عنوان مثالی از سیستم خودتطبیق مبتنی بر ویژگی با فضای گسسته بزرگ، سیستمی را در نظر بگیرید که ۱۰ ویژگی اختیاری دارد. این ویژگی‌ها، می‌توانند به صورت پویا حین اجرا فعال و غیرفعال شوند. بنابراین فضای تطبیق سیستم یک فضای گسسته شامل ۱۰۲۴ عملیات تطبیقی است (یعنی دو به توان ده). این ۱۰۲۴ عملیات تطبیقی را نمی‌توان به عنوان یک متغیر پیوسته نشان داد. در این حالت، یک راه حل می‌تواند، جای‌گذاری عملیات گسسته در یک فضای پیوسته و استفاده از جستجوی نزدیک‌ترین همسایه برای یافتن نزدیک‌ترین عملیات گسسته باشد (دالاس، ایوان، سانگ و کوپینگ، ۲۰۱۵). عملکرد یادگیری ماشین تا حدود زیادی به مقدار داده موجود برای یادگیری بستگی دارد. هنگامی که از یادگیری تقویتی برای سیستم‌های خودتطبیقی استفاده می‌شود، ممکن است به چرخه‌های یادگیری زیادی نیاز باشد تا فرآیند یادگیری همگرا شود (مصطفی و ژانگ، ۲۰۱۴).

در آزمایش‌های انجام شده، یادگیری حدوداً به ۲۰۰۰ چرخه یادگیری (با داده‌های ۲۵۶۰۰۰ مرحله زمانی) برای همگرایی نیاز داشت. تا وقتی یادگیری تقویتی همگرا

نشود، سیستم به احتمال زیاد تطبیق‌های ناکارآمدی را اجرا می‌کند، زیرا هنوز مشاهدات کافی وجود ندارد. تطبیق‌های ناکارآمدی ممکن است منجر به اثرات منفی شود، زیرا در یک سیستم زنده اجرا می‌شوند (فیلپو و پورتر ۲۰۱۷). برای سرعت بخشیدن به همگرایی، یافتن برآوردهای اولیه خوب برای دانش آموخته شده (ساتن و بارنو، ۲۰۱۸؛ دوتریل ملخوف، ۲۰۱۱) یا انجام یادگیری برون خط از طریق شبیه‌سازی سیستم، می‌تواند استفاده شود (تزار، جانگ، داس و بنائی، ۲۰۰۷). در عین حال، یادگیری تقویتی برخط ممکن است برای سیستم‌هایی که در محیطی کار می‌کنند تا اثرات آزمایش و خطا یادگیری تقویتی، قابل تحمل نباشد (به عنوان مثال، اگر عملیات تطبیقی به محیط آسیب برساند)، قابل استفاده نباشد.

۳- نتیجه‌گیری

یادگیری تقویتی مبتنی بر ارزش در سیستم‌های خودتطبیق بکار گرفته می‌شوند این مقاله به معرفی و ارزیابی یادگیری تقویتی برخط مبتنی بر سیاست برای تسهیل در مهندسی سیستم‌های خودتطبیقی پرداخت. برای تسهیل یادگیری برخط تقویتی، به الگوریتمی نیاز است که بتواند به طور مداوم، سیاست را بدون انتظار برای نتیجه نهایی، یعنی بدون انتظار برای رسیدن به حالت پایانی، به روز کند. الگوریتم‌های منتقد بازیگر یک نوع بدون مدل از الگوریتم‌های یادگیری تقویتی مبتنی بر سیاست هستند که در آن دانش به طور مداوم بدون انتظار برای نتیجه نهایی به روزرسانی می‌شود. این نوع یادگیرنده را در مدل چرخه بازخورد سیستم‌های خودتطبیق ترکیب شده که نقش تحلیلی را برای سیستم ایفا می‌کند. با توجه به تنوع داده‌ها در انواع الگوهای بار کاری (که لازمه هر الگوریتم تقویتی، حجم و تنوع داده‌های زیاد در حالات مختلف است)، می‌توان به دقت بالایی در تطبیق سیستم نسبت به رویکردهای مبتنی بر ارزش رسید. با طراحی یک تابع پاداش خوب که تمام ویژگی‌های هدف را دربردارد، توازن بالایی بین نرخ توصیه (بر اساس حجم درخواست‌ها) و میزان تاخیر برقرار کرد. رویکرد پیشنهادی این مقاله می‌تواند فضای عملیات بزرگ پیوسته را در حین یادگیری برخط کنترل کند (چون از شبکه عصبی به

آینده، این رویکرد می‌تواند برای مدیریت فضاهاى عملیاتی گسسته بزرگ گسترش خواهد یافت. تا انواع بیشتری از سیستم‌های خودتطبیق را به تصویر کشید.

جای راه‌جدولی بهره‌می‌برد). این رویکرد نه به حالت‌های محیطی گسسته‌شده به صورت دستی نیاز دارد، و نه به صورت دستی نیاز به تعیین نرخ اکتشاف مناسب همانند الگوریتم یادگیری تقویتی ارزش‌محور دارد. به عنوان کار

منابع

- Amoui, M., Salehie, M., Mirarab, S., & Tahvildari, L. (2008, March). Adaptive action selection in autonomic software using reinforcement learning. In *Fourth International Conference on Autonomic and Autonomous Systems (ICAS'08)* (pp. 175-181). IEEE.
- Arabnejad, H., Pahl, C., Jamshidi, P., & Estrada, G. (2017, May). A comparison of reinforcement learning techniques for fuzzy cloud auto-scaling. In *2017 17th IEEE/ACM international symposium on cluster, cloud and grid computing (CCGRID)* (pp. 64-73). IEEE.
- Aschoff, R., & Zisman, A. (2011). QoS-driven proactive adaptation of service composition. In *Service-Oriented Computing: 9th International Conference, ICSOC 2011, Paphos, Cyprus, December 5-8, 2011 Proceedings 9* (pp. 421-435). Springer Berlin Heidelberg.
- Barrett, E., Howley, E., & Duggan, J. (2013). Applying reinforcement learning towards automating resource allocation and application scalability in the cloud. *Concurrency and computation: practice and experience*, 25(12), 1656-1674.
- Bu, X., Rao, J., & Xu, C. Z. (2012). Coordinated self-configuration of virtual machines and appliances using a model-free learning approach. *IEEE transactions on parallel and distributed systems*, 24(4), 681-690.
- Caporuscio, M., D'Angelo, M., Grassi, V., & Mirandola, R. (2016). Reinforcement learning techniques for decentralized self-adaptive service assembly. In *Service-Oriented and Cloud Computing: 5th IFIP WG 2.14 European Conference, ESOCC 2016, Vienna, Austria, September 5-7, 2016, Proceedings 5* (pp. 53-68). Springer International Publishing.
- Chen, T., & Bahsoon, R. (2016). Self-adaptive and online qos modeling for cloud-based software services. *IEEE Transactions on Software Engineering*, 43(5), 453-475.
- D'Ippolito, N., Braberman, V., Kramer, J., Magee, J., Sykes, D., & Uchitel, S. (2014, May). Hope for the best, prepare for the worst: multi-tier control for adaptive systems. In *Proceedings of the 36th International Conference on Software Engineering* (pp. 688-699).
- Dulac-Arnold, G., Evans, R., van Hasselt, H., Sunehag, P., Lillicrap, T., Hunt, J., ... & Coppin, B. (2015). Deep reinforcement learning in large discrete action spaces. *arXiv preprint arXiv:1512.07679*.
- Dutreilh, X., Kirgizov, S., Melekhova, O., Malenfant, J., Rivierre, N., & Truck, I. (2011, May). Using

- reinforcement learning for autonomic resource allocation in clouds: towards a fully automated workflow. In *ICAS 2011, The Seventh International Conference on Autonomic and Autonomous Systems* (pp. 67-74).
- Filho, R. R., & Porter, B. (2017). Defining emergent software using continuous self-assembly, perception, and learning. *ACM Transactions on Autonomous and Adaptive Systems (TAAS)*, 12(3), 1-25.
- Jamshidi, P., Cámara, J., Schmerl, B., Käestner, C., & Garlan, D. (2019, May). Machine learning meets quantitative planning: Enabling self-adaptation in autonomous robots. In *2019 IEEE/ACM 14th International Symposium on Software Engineering for Adaptive and Self-Managing Systems (SEAMS)* (pp. 39-50). IEEE.
- Gheibi, O., Weyns, D., & Quin, F. (2021). Applying machine learning in self-adaptive systems: A systematic literature review. *ACM Transactions on Autonomous and Adaptive Systems (TAAS)*, 15(3), 1-37.
- Kephart, J. O., & Chess, D. M. (2003). The vision of autonomic computing. *Computer*, 36(1), 41-50.
- Klein, C., Maggio, M., Ārzén, K. E., & Hernández-Rodriguez, F. (2014, May). Brownout: Building more robust cloud applications. In *Proceedings of the 36th International Conference on Software Engineering* (pp. 700-711).
- De Lemos, R., Giese, H., Müller, H. A., Shaw, M., Andersson, J., Litoiu, M., ... & Wuttke, J. (2013). Software engineering for self-adaptive systems: A second research roadmap. In *Software Engineering for Self-Adaptive Systems II: International Seminar, Dagstuhl Castle, Germany, October 24-29, 2010 Revised Selected and Invited Papers* (pp. 1-32). Springer Berlin Heidelberg.
- Lorido-Botran, T., Miguel-Alonso, J., & Lozano, J. A. (2014). A review of auto-scaling techniques for elastic applications in cloud environments. *Journal of grid computing*, 12, 559-592.
- Mann, Z. Ā. (2017). Resource optimization across the cloud stack. *IEEE Transactions on Parallel and Distributed Systems*, 29(1), 169-182.
- Moustafa, A., & Zhang, M. (2014, June). Learning efficient compositions for QoS-aware service provisioning. In *2014 IEEE International Conference on Web Services* (pp. 185-192). IEEE.
- Nachum, O., Norouzi, M., Xu, K., & Schuurmans, D. (2017). Bridging the gap between value and policy based reinforcement learning. *Advances in neural information processing systems*, 30.
- Ramirez, A. J., Jensen, A. C., & Cheng, B. H. (2012, June). A taxonomy of uncertainty for dynamically adaptive systems. In *2012 7th International Symposium on Software Engineering for Adaptive and Self-Managing Systems (SEAMS)* (pp. 99-108). IEEE.
- Salehie, M., & Tahvildari, L. (2009). Self-adaptive software: Landscape and research challenges. *ACM transactions on autonomous and adaptive systems (TAAS)*, 4(2), 1-42.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Silvander, J. (2019). Business process optimization with reinforcement learning. In *Business Modeling and Software Design: 9th International Symposium, BMSD 2019, Lisbon, Portugal, July 1-3, 2019, Proceedings 9* (pp. 203-212). Springer International Publishing.
- Sutton, R. S., McAllester, D., Singh, S., & Mansour, Y. (1999). Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12.
- Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT press.
- Tesauro, G., Jong, N. K., Das, R., & Bennani, M. N. (2007). On the use of hybrid reinforcement learning for autonomic resource allocation. *Cluster Computing*, 10, 287-299.
- Iglesia, D. G. D. L., & Weyns, D. (2015). MAPE-K formal templates to rigorously design behaviors for self-adaptive systems. *ACM Transactions on Autonomous and Adaptive Systems (TAAS)*, 10(3), 1-31.

©Authors, Published by Journal of Intelligent Knowledge Exploration and Processing. This is an open-access paper distributed under the CC BY (license <http://creativecommons.org/licenses/by/4.0/>).

