

روشی یادگیر برای ترکیب وظایف در یادگیری تقویتی پیمانهای

سیدمحمدحسین میرهاشمی^۱، ناصر مزینی^۲ و محمدرضا جاهد مطلق^۳

^۱دانشگاه علم و صنعت ایران، mirhashemi@comp.iust.ac.ir

^۲دانشگاه علم و صنعت ایران، mozayani@iust.ac.ir

^۳دانشگاه علم و صنعت ایران، jahedmr@iust.ac.ir

چکیده - دسته ای از روش های یادگیری تقویتی سعی می کنند مسائل پیچیده را با تجزیه به مسائل کوچکتر حل کنند. به این صورت که هدف اصلی را به تعدادی زیرهدف یا وظیفه می شکنند و هریک را توسط یک یادگیر فرا می گیرند، سپس به ترکیب این وظایف یادگرفته شده می پردازند. بیشتر این روش ها به علت عدم استفاده مناسب از دانش موجود در این پیمانهای فراگرفته شده، در ترکیب آن ها با مشکل روبرو می شوند، و در نتیجه نمی توانند به خوبی به هدف اصلی دست بیابند. در این مقاله روشی یادگیر برای ترکیب وظایف ارائه شده است که به علت راهکار مناسبی که برای ترکیب پیمانها ارائه می کند، به جواب مناسبی برای مسئله اصلی می رسد. پیاده سازی این روش در یک مسئله که از منظر یادگیری تقویتی بسیار بزرگ و پیچیده است، عملکرد بسیار مناسبی از خود نشان داده و به کارایی بسیار بالایی می رسد. کلید واژه - ترکیب پیمانها، یادگیری تقویتی پیمانهای، یادگیری تقویتی.

ساده تر و مناسب یادگیری هستند با یادگیرهای تقویتی مرسوم حل کرده، و آنگاه سعی می کنند تا با ترکیب مناسب این پیمانهای یادگرفته شده به جواب قابل قبولی برای مسئله اصلی برسند. این سازوکار پیمانهای باعث انتساب نام یادگیری تقویتی پیمانهای^۱ به این روش ها شده است.

مشکل عمده ی روش های پیمانهای در ترکیب پیمانهای یادگرفته شده است. روش هایی که تا کنون ارائه شده اند یا در مسائلی خاص بوده تا بتوانند با استفاده از ویژگی های مسئله به حل این مشکل بپردازند، و یا سعی کرده اند تا با روش هایی ریاضی ترکیب پیمانها را عملی کنند. عده ی قلیلی هم که استفاده از یادگیری را برای ترکیب پیمانها در پیش گرفته اند، به خوبی از دانش موجود در پیمانها استفاده نکرده اند، و به همین

۱- مقدمه

یادگیری تقویتی یکی از شاخه های یادگیری ماشین است که به علت مزایای متعدد، بسیار مورد توجه قرار گرفته است. با این وجود یادگیری تقویتی از معایبی نیز رنج می برد که شاید مهمترین آن را بتوان مشکل نفرین ابعاد دانست. روش های یادگیری تقویتی به ابعاد محیط بسیار حساس هستند، و از آن جایی که محیط بیشتر مسائل واقعی بسیار پیچیده و بزرگ است، به کارگیری یادگیری تقویتی در آن ها با چالش جدی روبرو می شود.

دسته ای از روش ها برای غلبه بر این مشکل مسئله را به زیرمسائلی آسان تر تجزیه می کنند، سپس این زیرمسائل را که

¹ Modular Reinforcement Learning

در نظر گرفته می شود که فضای حالت آن، زیر فضای مربوط به آن زیرهدف، و پاداش آن همان پاداش مربوط به آن زیرهدف است: (S_i, A, P_i, R_i) . فضای اعمال همه ی این MDP ها یکسان و همان فضای اعمال MDP کلی است. تابع انتقال هر کدام (P_i) نیز به راحتی از روی تابع انتقال متعلق به MDP کلی قابل استخراج است، چراکه فضای حالت هر کدام، زیر فضایی از فضای حالت MDP کلی، و فضای اعمالش، یکسان با فضای اعمال MDP کلی است. هر کدام از این MDP ها که مربوط به یک زیرهدف عامل است را یک وظیفه، و یادگیری که آن را یاد می گیرد پیمانیه می نامیم، بنابراین عامل n وظیفه (پیمانیه) M_1, M_2, \dots, M_n دارد.

۳- مروری بر کارهای گذشته

اولین کار در یادگیری تقویتی پیمانیه ای را می توان [۱] دانست. در این کار ابتدا پیمانیه ها که از نوع یادگیر Q - [۲] هستند یاد گرفته می شوند؛ سپس یک یادگیر دیگر، که آن هم از نوع یادگیر Q - است، به کار گرفته می شود که می بایست در هر حالت یکی از پیمانیه ها را انتخاب کند. انتخاب یک پیمانیه به معنی انتخاب عمل پیشنهادی آن پیمانیه و اجرا شدن آن توسط عامل است. به بیان دیگر حالت این یادگیر همان حالت مسئله، ولی عمل آن انتخاب یک پیمانیه است. این روش یادگیری Q - سلسله مراتبی^۴ نامیده شده است.

مشکل آشکار روش فوق همان مشکل همیشگی بزرگی فضای حالت است. هر چند اگر تعداد وظایف کمتر از تعداد اعمال باشد، این روش تا حدی فضای حالت-عمل را کاهش می دهد، اما معمولاً خود فضای حالت به تنهایی چنان بزرگ است که کاهش اعمال تاثیر عملی در نتیجه کار نخواهد داشت.

راه کار دیگری که برای غلبه بر مشکل بزرگی فضای حالت در ادبیات موضوع مطرح شده است، استفاده از سازوکاری ساده و غیر یادگیر برای ترکیب اعمال پیشنهادی پیمانیه ها است. مفهوم

دلیل موفق به غلبه بر مشکل بزرگی فضای حالت نشده اند. در این مقاله روشی یادگیر برای ترکیب پیمانیه های یاد گرفته شده ارائه می شود که از یادگیری انجام شده توسط پیمانیه ها حداکثر استفاده را می برد. به علاوه سعی می کند وابسته به خواص محیط نباشد تا جامعیت خود را برای حل مسائل مشابه از دست ندهد. روش پیشنهادی در یک محیط که برای یادگیرهای تقویتی ساده بسیار بزرگ و پیچیده است، پیاده سازی شده و نتایج بسیار ارزنده ای کسب کرده است.

۲- فرمول بندی مسئله

مسئله به صورت یک فرآیند تصمیم مارکوف^۲ (MDP) در نظر گرفته شده است. یک MDP توسط چهار تایی (S, A, P, R) تعریف می شود که در آن S فضای حالت و A فضای اعمال است. $P(s'|s, a)$ تابع توزیع احتمال رسیدن به حالت s' با انجام عمل a در حالت s است و تابع انتقال نامیده می شود. $R(s, a)$ امید ریاضی پاداش دریافتی با انجام عمل a در حالت s است و تابع پاداش نامیده می شود. یک سیاست^۳ $\pi(s)$ تابعی است که برای هر حالت یک عمل را پیشنهاد می کند. هدف نهایی یافتن سیاست بهینه (π^*) است که با انجام اعمال پیشنهادی آن، عامل به حداکثر مجموع پاداش در طی فعالیتش برسد.

برای عامل n زیرهدف تعریف می شود، برای هر زیرهدف یک سیگنال پاداش جداگانه در نظر گرفته می شود، بنابراین n سیگنال پاداش R_1, R_2, \dots, R_n داده شده است. اگر تنها اجزائی از محیط که به یک زیرهدف مربوط است را در نظر بگیریم، زیرفضایی از محیط به وجود می آید که مخصوص آن زیرهدف است، به عبارت دیگر پاداش مربوط به آن زیرهدف تنها به حالت عامل در آن زیرفضا مربوط است. بنابراین n زیرفضای S_1, S_2, \dots, S_n که هر کدام مربوط به یکی از اهداف است نیز داده شده است. با کمک آنچه بیان شد برای هر زیرهدف یک MDP جداگانه

⁴ Hierarchical Q-Learning

² Markov Decision Process

³ Policy

برهمن اساس پیمانها را یادگیر سارسا^۷ که یک یادگیر داخل-سیاست است، تعیین کرده است. یک یادگیر داخل-سیاست^۸ سیاستی که در طول یادگیری به اجرا درمی‌آید را فرا می‌گیرد، بنابراین ارزش اعمال آن اشکال فوق را نخواهد داشت و انتظار می‌رود سامانه‌ی کلی نتیجه بهتری بدست آورد. با این روش کارایی عامل تا حدی افزایش می‌یابد، اما استفاده از ارزش اعمال که مشکل اصلی این روش‌هاست در این روش نیز وجود دارد. به علاوه همانطور که خودش بیان کرده همگرایی این روش اثبات نشده است.

سه روش بالا ([۳]، [۴] و [۶]) و روش‌های مانند آن‌ها که بر اساس ارزش اعمال به ترکیب پیمانها می‌پردازند، برای مسائل چند-هدفه ارائه شده‌اند. در مسائل چند-هدفه هر زیرهدف دارای اصالت است و هدف کلی عامل نیز چیزی جز ارضای حداکثری همه‌ی زیراهداف نیست. به همین دلیل ارزش اعمال از دید زیراهداف مختلف قابل مقایسه بوده و عامل می‌تواند بر اساس این مقایسه تصمیم‌گیری کرده و یک عمل را انتخاب کند. در مسائل چند-هدفه مجموع پاداش زیراهداف به عنوان پاداش به عامل تعلق می‌گیرد: $R = R_1 + R_2 + \dots + R_n$.

در دیگر مسائل که زیر اهداف دارای اصالت نیستند و تنها برای رسیدن به هدفی کلی تر طراحی شده‌اند، استفاده از این روش‌ها ناکارآمد خواهد بود، چراکه ارضای خود زیراهداف مدنظر نیست و بنابراین ارزش اعمال از منظر آن‌ها معیار مناسبی برای یافتن عمل منتخب نیست، بلکه تنها هدف کلی دارای اصالت بوده و زیراهداف تنها وسیله‌ای برای رسیدن به آن هستند. روش یادگیری-Q سلسله‌مراتبی [۱] و روشی که در این مقاله ارائه شده است برای این نوع مسائل نیز قابل به کارگیری هستند.

۴- روش پیشنهادی

برای حل مسئله‌ی کلی با استفاده از پیمانهای که هر یک

نهفته در این شیوه آنست که پیمانها عمده کار یادگیری را به انجام رسانده‌اند، حال کافی است که به نحوی ساده ترکیب شوند تا جواب کلی دلخواه حاصل شود.

کارلسون [۳] و هامفریس [۴] روش‌هایی از این نوع را ارائه می‌دهند. در هر دوی آن‌ها پیمانها از نوع یادگیر-Q هستند، و هر دو برای انتخاب عمل از ارزش اعمال استفاده می‌کنند. هامفریس بیان می‌کند که انتخاب عمل پیمانهای که ارزش عمل پیشنهادی‌اش در جدول-Q خودش، از ارزش اعمال پیشنهادی دیگر پیمانها در جدول-Q خودشان بیشتر است، کار درستی نیست؛ چراکه ممکن است با این کار پیمانهای دیگر کاهش ارزش بزرگتری را متحمل شود، و در نتیجه بیشینه پاداش مجموع حاصل نشود. بنابراین روشی را پیشنهاد می‌کند که در آن عمل پیشنهادی پیمانهای انتخاب می‌شود که در صورت انتخاب نشدن عملش، بیشترین کاهش ارزش را متحمل خواهد شد.

کارلسون از روشی به نام بزرگترین جرم^۵ برای انتخاب عمل استفاده می‌کند. در این روش در هر حالت عملی انتخاب می‌شود که مجموع ارزش‌های انتسابی به آن در جدول-Q پیمانها، بیشینه باشد. روش بزرگترین جرم اولین بار در [۵] پیشنهاد شده است.

در [۶] به روش‌هایی که عمل را بر اساس ارزششان در پیمانها انتخاب می‌کنند و پیمانها از نوع یادگیر-Q هستند اشکال گرفته شده است. یادگیر-Q یک یادگیر خارج-سیاست^۶ است و بنابراین ارزش اعمال با این فرض بدست می‌آید که سیاست بهینه دنبال خواهد شد. درحالیکه در روش پیمانهای هیچ‌کدام از پیمانها کنترل عامل را به تنهایی در دست ندارد تا همیشه اعمال پیشنهادی‌اش به اجرا درآید. بنابراین ارزش اعمال به دست آمده از طریق یادگیری-Q، برای مقایسه و انتخاب نامناسب هستند و تصمیم‌گیری براساس آن‌ها باعث کاهش کارایی خواهد شد.

⁷ Sarsa

⁸ On-Policy

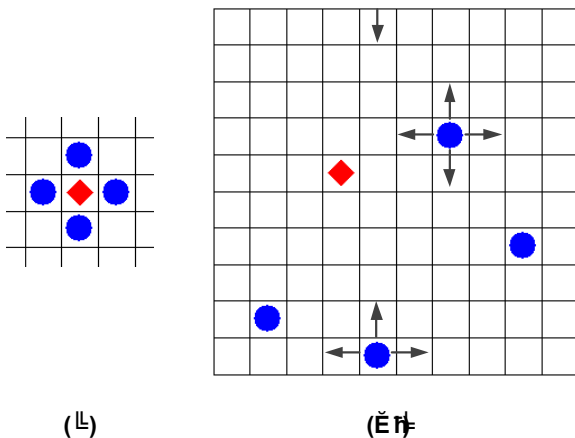
⁵ Greatest Mass

⁶ Off-Policy

وضعیت آن حالت از منظر زیرهدف آن پیمانانه خواهد بود. اگر این ارزش اعمال در یک حالت، از همه ی پیمانانه ها دریافت شده و در کنار هم قرار گیرند، آن گاه وضعیت حالت از منظر همه ی زیراهداف مشخص خواهد شد. با این روش هر حالت از محیط به حالت دیگری نگاشت می شود، و فضای حالت جدیدی به وجود می آید که به نوعی می توان گفت که در آن هر حالت از نظر زیراهداف ارزش گذاری شده است. یادگیر ترکیب کننده در این فضای جدید به یادگیری می پردازد.

۵- پیاده سازی

برای پیاده سازی روش پیشنهادی، مسئله ی تعقیب^۹ انتخاب شده است. در این مسئله که شکار و شکارچی نیز نامیده شده است، هدف چهار شکارچی گرفتن یک شکار در یک محیط گسسته دوبعدی است. شکارچی ها تنها زمانی می توانند شکار را بگیرند که او را از چهار طرف محاصره کنند. شکل ۱.



شکل ۱: (الف) محیط مسئله تعقیب (ب) وضعیت گرفتن شکار

هر موجود (شامل شکارچی ها و شکار) در هر گام می تواند به یکی از چهار خانه ی همسایه ی خود حرکت کرده و یا در جای

وظیفه ای را فرا گرفته اند، می بایست به نحوی نتیجه ی یادگیری آن ها با هم ترکیب شود تا جواب کلی حاصل شود. سازوکاری که از یادگیری برای این ترکیب استفاده کند مطلوب به نظر می رسد، اما این سازوکار چگونه می تواند از پیمانانه ها سود ببرد، و از آن مهمتر چگونه می تواند بر مشکل همیشگی بزرگی فضای حالت غلبه کند؟

روش هایی که در مسائل چند-هدفه ارائه شده اند و از ارزش اعمال برای ترکیب پیمانانه ها استفاده می کنند، در واقع مشکل را با استفاده از ویژگی های خاص مسائل مورد نظرشان دور زده اند. روش یادگیری-Q سلسله مراتبی [۱] نیز در واقع استفاده ی چندانی از نتیجه ی یادگیری پیمانانه ها نمی کند و به همین دلیل پیشرفتی در مورد مشکل بزرگی فضای حالت بدست نمی آورد.

برای استفاده از پیمانانه های یادگیری شده، می بایست به نحوی از دانش و تجربه ی ذخیره شده در آن ها استفاده کرد. اگر از مقادیر ارزش اعمال به عنوان ابزاری برای مقایسه استفاده شود، می بایست زیراهداف دارای اصالت و مقایسه پذیر باشند، که همانطور که پیش از این بیان شد در این صورت روش به مسائل چند-هدفه محدود شده است. بنابراین هم می بایست از ارزش اعمال استفاده کرد و هم نباید آن ها را به عنوان ابزاری برای مقایسه به کار برد.

در روش پیشنهادی این مقاله هر پیمانانه بر اساس ارزشی که به اعمال مختلف در هر حالت داده است، فضای حالت را به فضای حالتی جدید نگاشت می کند. این فضای حالت جدید در واقع فضای حالت از دیدگاه زیرهدف آن پیمانانه است. این فضاهای حالت به وجود آمده از دیدگاه زیراهداف مختلف، در کنار یکدیگر قرار داده شده و یک فضای حالت کلی جدید به وجود می آید. این فضای جدید در اختیار یادگیر دیگری قرار می گیرد تا در آن به یادگیری بپردازد و عمل خروجی عامل را تولید کند. از آنجاییکه این فضای جدید از ترکیب نتایج یادگیری زیراهداف بدست آمده است، انتظار می رود که فضایی کوچکتر، چگال تر از منظر دانش حل مسئله، و آسان تر برای یادگیری باشد. یادگیر مذکور را یادگیر ترکیب کننده می نامیم.

اگر پیمانانه ها یادگیر-Q فرض شوند و هر یادگیر-Q وظیفه مربوط به زیرهدف خود را فرا گرفته باشد، آنگاه ارزش اعمال مختلف یک حالت در جدول-Q یک پیمانانه، نشان دهنده ی

⁹ Pursuit

هنگامیکه عامل به شکار می رسد یک مقدار مثبت، و در گام های دیگر صفر تعیین شده است. یادگیر ترکیب کننده هر عامل نیز هنگامی که شکارچی ها شکار را بگیرند، یک پاداش مثبت بزرگ، و در هر گام دیگر یک پاداش منفی کوچک دریافت می کند.

پیمانها و همچنین یادگیر ترکیب کننده، یادگیر-Q تعیین شده اند. در هر حالت از محیط، ارزش چهار عمل حرکت به چهار جهت در جدول-Q چهار پیمانها، تشکیل یک حالت جدید با ۱۶ متغیر را می دهد. یادگیر ترکیب کننده در این فضای حالت ۱۶ متغیره به یادگیری می پردازد و عملی که عامل در هر حالت انجام می دهد را تعیین می کند، در واقع این یادگیر ترکیب کننده است که کنترل عامل را در دست دارد. شکل ۲.

از آنجاییکه ارزش اعمال در جدول-Q مقادیری پیوسته هستند، این ۱۶ متغیر حالت، مقادیری پیوسته خواهند داشت. اما یادگیر ترکیب کننده که یک یادگیر-Q است، می بایست در فضای گسسته به یادگیری بپردازد. بنابراین یک گسسته سازی ساده ۴-مقداره بر روی ارزش هر عمل انجام می شود، تا فضای حالت یادگیر ترکیب کننده گسسته شود.

با گسسته سازی ۴-مقداره هر متغیر می تواند ۴ مقدار متفاوت داشته باشد، بنابراین فضای حالت یادگیر ترکیب کننده $10^9 \approx 4^{16}$ حالت می تواند داشته باشد. اما در عمل تعداد حالاتی که ترکیب کننده با آن ها برخورد می کند بسیار کمتر از این مقدار و حدود $10^4 \times 2$ حالت است. علت این امر را می توان در وابستگی متغیرهای حالت به یکدیگر دانست، به این معنی که برای هر مقدار از یک متغیر، یک یا تعدادی از متغیرهای دیگر تنها مقادیر خاصی را می توانند اتخاذ کنند. بنابراین همه ی ترکیب های مقادیر متغیرهای حالت قابل دستیابی نیستند و بنابراین فضای حالت بسیار کوچکتر از تعداد کل حالات قابل ساخت با متغیرهای حالت است.

یادگیری در دوره های^{۱۳} متوالی انجام می شود. هر دوره با محاصره شدن شکار توسط شکارچی ها، و یا با رسیدن تعداد گام

خود ثابت بماند. محیط چنبره ای^{۱۰} است، یعنی حرکت به سمت بیرون در یک طرف محیط، باعث ورود از طرف دیگر می شود. بنابراین محیط دیوار و گوشه ندارد که باعث شود کمتر از چهار شکارچی بتوانند شکار را محاصره کنند. موجودات محیط در هر گام به طور همزمان یک عمل (شامل حرکت به چهار جهت و ثابت ایستادن) انجام می دهند. هیچ دو موجودی نمی توانند در یک خانه قرار بگیرند، و تصادم بین آن ها از طریق اولویت در حرکت حل می شود که این اولویت در هر گام بصورت تصادفی به موجودات محیط تعلق می گیرد.

شکار در هر گام یکی از پنج عمل را به صورت تصادفی انتخاب می کند. عامل های یادگیر شکارچی ها هستند که می بایست گرفتن شکار را فرا بگیرند. این مسئله یک مسئله ی چند-عامله^{۱۱} است، اما چون عامل های آن به دیگر عامل ها تنها به عنوان یک جزء از محیط نگاه می کنند و به مدل کردن یکدیگر نمی پردازند، مسئله از دید هر عامل تبدیل به مسئله ای تک-عامله با محیطی بسیار پویا می شود. بنابراین بر طبق دسته بندی [۷]، این مسئله یک مسئله ی چند-عامله متجانس فاقد ارتباطات^{۱۲} است، به علاوه در آن عامل ها یکدیگر را مدل نمی کنند.

اگر ابعاد محیط 10×10 باشد، فضای حالت برای عامل 10^8 حالت خواهد داشت. این فضای حالت بسیار بزرگتر از آن است که یک یادگیر تقویتی ساده بتواند در زمان معقول آن را فرا بگیرد.

هریک از زیراهداف عامل در رابطه با یکی از موجودات محیط تعیین شده است، بنابراین عامل چهار زیرهدف دارد که سه تای آن مربوط به سه شکارچی دیگر، و یکی مربوط به شکار است. پاداش زیرهدف مربوط به یک شکارچی، هنگامیکه عامل باعث تصادم با آن می شود مقداری منفی، و در گام های دیگر صفر تعیین شده است. پاداش زیرهدف مربوط به شکار نیز

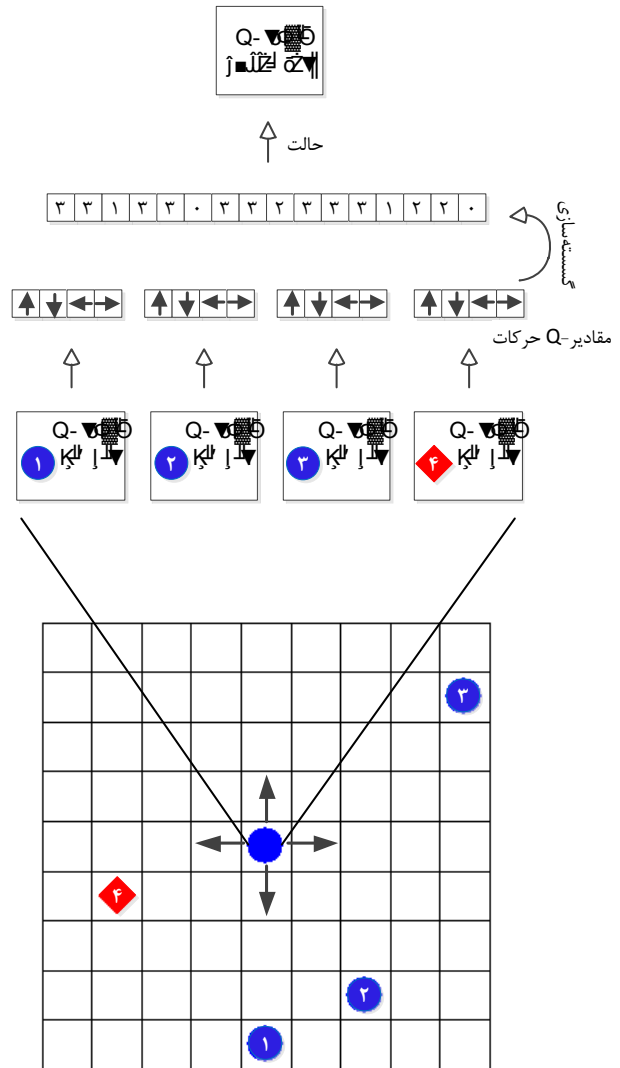
¹⁰ Toroidal

¹¹ Multi-Agent

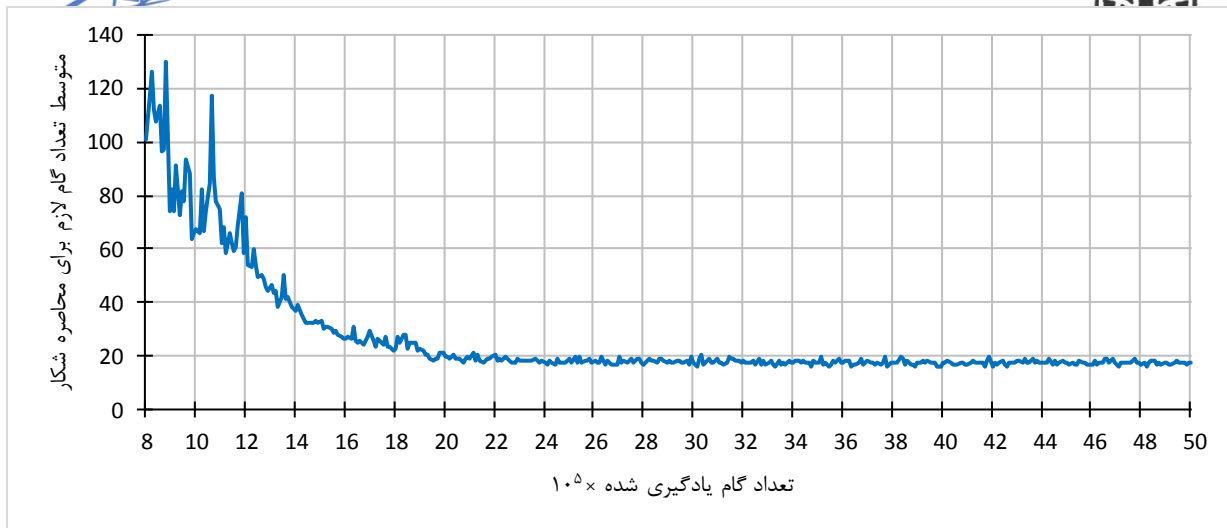
¹² Homogeneous Non-Communicating

شکل ۳ نمودار کارایی روش پیشنهادی را نشان می دهد. با گذشت هر ۱۰ هزار گام از یادگیری، یادگیری بطور موقت متوقف شده و عامل ها هزار گام مورد تست قرار گرفته اند تا متوسط تعداد گام های لازم برای محاصره کردن شکار بدست آید. برای رسم نمودار، ۱۰ دفعه یادگیری و تست با روش ذکر شده به انجام رسیده و بین آن ها میانگین گیری انجام شده است. پس از گذشت ۲ میلیون گام از یادگیری، عامل ها تقریبا به حداکثر کارایی خود نزدیک شده اند. برای مسئله ای با فضای حالت به اندازه 10^8 ، این تعداد گام برای یادگیری بسیار مناسب است. به علاوه پس از یادگیری، شکارچی ها در طی کمتر از ۲۰ گام شکار را محاصره می کنند، که در چنین محیطی، تقریبا نزدیک به بیشینه کارایی هوشمندترین عامل ها است.

دوره به ۱۰۰ پایان می یابد. در ابتدای هر دوره موجودات به صورت تصادفی در خانه های محیط پخش می شوند. نرخ یادگیری ثابت و به مقدار ۰/۱، و ضریب کاهش نیز ثابت و به مقدار ۰/۹ تنظیم شده است. همه ی یادگیرها از یک سیاست e-greedy استفاده می کنند که مقدار پارامتر e در ۵۰ گام اول هر دوره از ۰/۴ به ۰ به صورت خطی تنزل پیدا می کند.



شکل ۲: نمایش پیاده سازی روش پیشنهادی در مسئله تعقیب



شکل ۳: نمودار کارایی عامل‌های شکارچی - متوسط تعداد گام‌های لازم برای محاصره‌ی شکار برحسب تعداد گام‌های یادگیری عامل‌های شکارچی

با این سازوکار انتظار می‌رود که این مشکل مرتفع گردد.

مسئله‌ی دیگر از بین رفتن بخشی از اطلاعاتی که توسط پیمان‌ها یادگرفته شده، به دلیل گسسته‌سازی مقادیر ارزش اعمال است. برای مقابله با این مشکل می‌توان از یک یادگیر پیوسته برای ترکیب‌کننده استفاده کرد. بنابراین در صورت استفاده از یک یادگیر پیوسته مناسب، استفاده کامل از اطلاعات موجود در ارزش اعمال، می‌تواند باعث افزایش کارایی عامل شود. با توجه به نتایج موفق بدست آمده، روش ارائه شده می‌تواند آغاز راهی باشد که روش‌های پیمان‌های را به عنوان راه‌کاری مناسب برای غلبه بر مشکل نفرین ابعاد در یادگیری تقویتی مطرح کند.

مراجع

- [1] L. J. Lin, "Scaling up Reinforcement Learning for robot control," *Proc. 10th. Int. Conf. on Machine Learning*, pp. 182-196, 1993.
- [2] C. J. C. H. Watkins, *Learning from Delayed Rewards*. PhD Thesis, University of Cambridge, Psychology Department, 1989.
- [3] J. Karlsson, *Learning to Solve Multiple Goals*. PhD Thesis, University of Rochester, 1997.
- [4] M. Humphrys, "Action Selection Methods using Reinforcement Learning," *From Animals to Animats 4: Proc. 4th. Int. Conf. on Simulation of Adaptive Behavior*. Cambridge, MA, pp. 134-144, 1996.
- [5] S. Whitehead, J. Karlsson and J. Tenenbergs, *Learning Multiple Goal Behavior via Task Decomposition and Dynamic Policy Merging*. Kluwer Academic Press, pp. 45-78, 1993.
- [6] N. Sprague and D. H. Ballard, "Multiple-Goal Reinforcement Learning with Modular Sarsa(0)," *Proc. Int. Joint Conf. on Artificial Intelligence*, pp. 1445-1447, 2003.
- [7] P. Stone, M. Veloso, "Multiagent systems: A survey from a machine learning perspective" *Autonomous Robots*, Vol. 8, pp. 348-383, 2000.

۶- جمع‌بندی و نتیجه‌گیری

در این مقاله روشی برای حل مشکل ترکیب وظایف، در روش‌های یادگیری تقویتی پیمان‌های ارائه شد. مزیت روش ارائه شده نه تنها در استفاده از یادگیری برای ترکیب پیمان‌ها، بلکه در استفاده‌ی مناسب از تجربه‌ی یادگیری انجام شده در پیمان‌ها برای ترکیب آن‌ها نیز هست؛ و به همین علت به نتایج قابل توجهی دست پیدا می‌کند.

مزیت دیگر روش پیشنهادی را می‌توان در عدم وابستگی آن به اندازه‌ی محیط مسئله دانست، چراکه پیمان‌ها تنها ارزش اعمال را به یادگیر ترکیب‌کننده عرضه می‌کنند، و با افزایش اندازه محیط، تغییر در اندازه فضای حالت یادگیر ترکیب‌کننده ایجاد نخواهد شد. خود پیمان‌ها نیز چون تنها یک یا تعدادی کمی از اجزاء محیط را در نظر دارند، حساسیت زیادی به اندازه‌ی محیط ندارند.

روش ارائه شده نیز مانند هر روشی علاوه بر فواید، با مسائلی نیز روبروست که می‌بایست برای آن‌ها چاره‌ای اندیشیده شود. یکی از این مشکلات افزایش نمایی اندازه فضای حالت یادگیر ترکیب‌کننده، با افزایش تعداد وظایف است. بنابراین در مسائلی که تعداد وظایف زیاد است، روش کارایی خود را از دست می‌دهد. برای رفع این مشکل می‌توان ابتدا وظایف را در دسته‌های کوچک‌تر با یکدیگر ترکیب کرد، سپس نتایج بدست آمده از دسته‌ها را نیز برای رسیدن به جواب نهایی ترکیب کرد.