

استخراج گذرگاهها با استفاده از تشخیص اشیا در یادگیری تقویتی

بهزاد غضنفری، ناصر مزینی و محمدرضا جاهد مطلق

می‌شود، بر حالات و پاداش در زمان $t+1$ تأثیر می‌گذارد. اساساً هیچ مفهومی از شیوه انجام کار که تأکید بر روی بازه متغیری از زمان داشته باشد، وجود ندارد. در نتیجه MDPs به صورت متناول قادر به استفاده از مزایای ساده‌سازی و کارآمدی که در سطوح بالاتر از تحرید زمانی وجود دارد، نیستند. تحرید زمانی می‌تواند درون یادگیری تقویتی به شکل‌های متغیری مطرح شود [۴]. یکی از مهم‌ترین ویژگی‌های انجام تحرید زمانی در چارچوب یادگیری تقویتی، بناهاین تنوری فرایندی‌های تصمیم‌گیری شبیه مارکوف^۱ (SMDPs) است. همان‌گونه که گفته شد، در MDP فرض می‌شود که عامل دارای توانایی دستیابی کامل به حالت محیط است و هر اقدامش را در یک گام زمانی انتخاب می‌کند. در فرض دوم این محدودیت را ندارد و این امکان را به عامل می‌دهد که اقداماتش را در چندین گام زمانی انجام دهد [۴].

خوشنختانه اگرچه تعداد زیادی از مسایل یادگیری تقویتی بسیار پیچیده هستند، آنها اغلب به صورت سلسله مراتبی قابل تجزیه به یک سری از زیروظایف ساده‌تر هستند. انسان‌ها نیز از چنین ساختارهای سلسله مراتبی تجزیه‌شدنی در حل مسایل که پیچیده و دارای مقیاس بزرگی هستند، بهره می‌برند [۳].

همان طور که می‌دانیم، اقدامات توسعه داده شده زمانی در یادگیری تقویتی سلسله مراتبی^۲ (HRL) مورد مطالعه قرار گرفته‌اند. HRL یک چارچوب کلی برای مقیاس پذیر کردن RL به مسایلی با فضای حالات بزرگ با استفاده از ساختار کار (یا اقدام) برای محدودکردن فضای خط مشی‌ها است. اصل کلیدی که HRL در بر گرفته است، توسعه الگوریتم‌های یادگیری است که به یادگیری خط مشی‌ها از ابتداء نیازی ندارند، بلکه در عوض از خط مشی‌های موجود برای زیروظایف ساده‌تر (یا مارکو-اقدامات) دوباره استفاده می‌کنند و استراتژی تقسیم و غلبه را به کار می‌برند.

SMDP به عنوان مدل آماری شناخته‌شده‌ای برای رفتار با اقداماتی با طول‌های متغیر است. در دهه اخیر کارهای زیادی روی مدل SMDP بسط داده شده‌اند که این کارها به صورت مدل‌های کار سلسله مراتبی از زیروظایف سطح پایین‌تر هستند که به صورت جزئی یا کامل مشخص شده‌اند. این کارها منجر به توسعه مدل‌های HRL قدرتمندی مانند سلسله مراتبی از ماشین‌های مجرد^۳ (HAMs)^۴، Option^۵ها^۶ MAXQ^۷ شده است. در مدل Options بررسی شده که چگونه خط مشی‌های سراسری را یا معلوم‌بودن خط مشی‌های کاملاً معنی برای اجرای زیروظایف یادگیریم. در قالب رسمی از HAMs نشان داده شده است که چگونه یادگیری تقویتی می‌تواند حتی زمانی که خط مشی‌ها برای پایین‌ترین سطوح زیروظایف، تنها به صورت جزئی مشخص شده‌اند، انجام وظیفه نماید. مدل MAXQ از اولین روش‌های مطرح شده است که در آن تحرید زمانی با تحرید حالات ترکیب شده است. ویژگی اصلی

نیمه: این مقاله روش جدیدی را مطرح می‌کند که قادر به استخراج گذرگاهها با استفاده از مجموعه اتوماتیک رفتار و مسیریابی جهات الهام گرفته شده است و طله تعاملات عامل با محیط پیرامونی افق عمل می‌کند. عامل با استفاده از بنده و تشخیص اشیا به صورت سلسله مراتبی، نشانه‌هایی را پیدا می‌کند. بن نشانه‌ها در فضای اقدام به هم نزدیک باشند، گذرگاهها با استفاده از های بین آنها استخراج می‌شوند. نتایج آزمایش‌ها بهبود قابل ملاحظه‌ای را ایند یادگیری تقویتی در مقایسه با سایر روش‌های مشابه نشان می‌دهد.

ید و ازه: یادگیری تقویتی، خوشبندی اشیا، یادگیری تقویتی سلسله مراتبی، ات گستریش یافته زمانی.

۱ - مقدمه

یادگیری تقویتی^۸ (RL) در صدد هستیم تنها با استفاده از پاداش و به عامل را برنامه‌ریزی کنیم، بدون این که به عامل بگوییم که به این وظیفه را انجام دهد. عامل پایستی اقداماتی که باعث افزایش دت مجموع ارزش سیگنال‌های تقویتی می‌شوند را انتخاب کند. این کار را می‌تواند در طول زمان از طریق قاعده‌مند کردن آزمایش ملایاد بگیرد [۱].

لی رغم موقعیت روش‌های موجود در یادگیری تقویتی، در مسایلی با بالا این روش‌ها به خوبی مقیاس‌پذیر نیستند. تلاش‌های اخیر برای بر معضل ابعاد^۹ بالا به سمت شیوه‌های مبتنی بر تحرید^{۱۰} در RL برده است، که به صورت طبیعی منجر به معماری کنترلی سلسله مراتبی و یتم‌های یادگیری وابسته به آن می‌شوند [۲] و [۳]. در اکثر موارد راه نای سلسله مراتبی جواب‌های نزدیک به پنهانه‌ای را در کارایی‌شان می‌دهند و همچنین هزینه مناسب‌تری را در زمان اجرا، زمان یادگیری سای مورد نیاز برای حل مسایل در مقابل تکنیک‌های RL مخصوص می‌کنند [۱].

غلب تحقیقات در RL مبتنی بر تحریر چارچوب فضا و اقدام گستره، از فرایند تصمیم‌گیری مارکوف^{۱۱} (MDP) است. چارچوب MDPs ورت مرسوم شامل اقدامات گسترش‌یافته زمانی^{۱۲} نیستند. آنها اصولاً ساس گام‌های گسترش‌زمانی هستند: اقدامی که در زمان t انتخاب

ن مقاله در تاریخ ۱۴ شهریور ماه ۱۳۹۰ دریافت و در تاریخ ۲۰ خرداد ماه ۱۳۹۱ ری شد.

بهزاد غضنفری، دانشکده مهندسی کامپیوتر، دانشگاه علم و صنعت ایران، تهران، (email: be_ghazanfari@ieee.org)

صر مزینی، دانشکده مهندسی کامپیوتر، دانشگاه علم و صنعت ایران، تهران، (email: mozayani@iust.ac.ir)

حمدرضا جاهد مطلق، دانشکده مهندسی کامپیوتر، دانشگاه علم و صنعت ایران، (email: jahedmr@iust.ac.ir)

1. Reinforcement Learning
2. Curse of Dimensionality
3. Abstraction
4. Markov Decision Process
5. Temporally Extended Action

6. Semi-Markov Decision Processes

7. Hierarchical Reinforcement Learning

8. Hierarchies of Abstract Machines

شایستگی آنها برای گذرگاه‌بودن افزایش می‌یابد. این عملیات‌ها دارای پیچیدگی محاسباتی بالایی هستند و همچنین دقت آنها به اندازه خوشه‌ها در فضای حالت بستگی دارد. نوآوری روش پیشنهادی این است که با توجه به تأثیر اقدامات، اشیا و مفاهیم استخراج می‌شوند و با استفاده از سلسله مراتب تحرید نقاط کلیدی برای پیداکردن گذرگاه استخراج می‌شوند.

در بخش ۲ پیش‌زمینه‌ای از شیوه پیشنهادی و در بخش ۳ کارهای مشابهی که در این زمینه انجام شده، ارائه شده است. بخش ۴ به بررسی و توضیح روش پیشنهادی و هر کدام از اجزایش و همچنین به مزایا و پیچیدگی آن اختصاص دارد. ارائه و شرح آزمایش‌ها و نتایج در بخش ۵ آمده است و در نهایت در بخش ۶ جمع‌بندی و نتیجه‌گیری از این نوشتار ارائه شده است.

۲- پیش‌زمینه

اجازه دهید توجیه روش پیشنهادی را بدین ترتیب مطرح نماییم که در ابتداء، یک کودک هیچ ذهنیتی نسبت به دنیای اطراف خود ندارد و تلاش می‌کند تمام اقدامات و حرکاتی را که می‌تواند، انجام دهد. این مهم کاملاً مشابه با آنچه است که برای یک عامل در یادگیری تقویتی برای بهینگی خط مشی به دست آمده‌اش لازم است. ما همچنین مشاهده می‌کنیم که یک کودک بسیار مشتاق است تا از طریق تجربی به ماهیت اجسام پی ببرد، مانند مزه‌کردن هر جسم و به تبع آن نتیجه مثبت و منفی عمل خود را تجربه نماید. باز هم بهطور مشابه برای عامل در یادگیری تقویتی در نتیجه اقداماتش، پاداش‌های مثبت یا منفی را دریافت می‌کند و حالات (اشیا) متفاوت را از طریق پاداشی که از آنها کسب کرده، متمایز می‌کند. عامل برای ماقرئیم کردن مقدار بازگشتی^۳ از محیط باید تأثیرات هر اقدام ممکن در هر حالت را امتحان کند. این یکی از پیش‌شرط‌های بهینگی در همه روش‌های یادگیری تقویتی است. به تدریج مفاهیمی از اشیای متفاوت برای کودک شکل می‌گیرد که وی را قادر می‌سازد تأثیرات اقداماتش را پیش‌بینی کند. کودک رفته رفته قادر است الگوها و اجسام را از هم دیگر تفاکر کند و وی دیگر راغب نخواهد بود که تجارت جدیدی را در برخود با یک دیوار متفاوت کسب کند.

انسان قادر است ویژگی‌های اصلی و متمایز کننده از اشیا را استخراج و سپس الگوهای مشابه را به درستی دسته‌بندی کند. این تحرید در یادگیری وی را قادر می‌سازد که توانایی تشخیص و تضمیم‌گیری دقیق و بالایی داشته باشد. طبق نظر محققان علوم شناختی، انسان دارای ساختار سلسله مراتبی است و داده‌ها در روال‌هایی با چندین مرحله پردازش می‌شوند [۱۵]. ساختارهای سلسله مراتبی برای تشخیص اشیا و استفاده از نشانه‌ها^۴ برای مسیریابی^۵ در انسان و حیوانات، روش‌های دقیق و مقاومی مقاومی در مقایسه با سیستم‌های مصنوعی هستند [۱۶]. مزایای این سیستم‌ها برای سیستم‌های هوش مصنوعی این انگیزه را ایجاد می‌کند که از آنها در فرایند تشخیص گذرگاه‌ها^۶ استفاده شود. در این مقاله نیز سعی شده است که با الهام از سیستم‌های بیولوژیکی و رفتار حیوانات، با در نظر گرفتن محدودیت‌های عامل و تعاریف ممکن در این حوزه یک شیوه کلی مطرح شود که مزیت‌های این سیستم‌ها را تا حد امکان داشته باشد.

MAXQ Option‌ها تحرید زمانی دانش و اقداماتی که در چارچوب یادگیری تقویتی قرار گرفته می‌شوند را در شیوه‌ای طبیعی و عمومی ممکن می‌سازند. در شیوه پیشنهادی از چارچوب Option برای تعریف اقدامات گسترش‌یافته زمانی استفاده شده است. یک سه گانه $\langle I, \pi, B \rangle$ است در جایی که I مجموعه شروع است. $S \subseteq I$ مجموعه‌ای است که هر اقدام گسترش‌یافته زمانی می‌تواند از آن مجموعه فراخوانی شود. $[0, 1] \rightarrow S^* A : \pi$ تابع خط مشی برای هر حالت در مجموعه آغازین از اقدامات گسترش‌یافته زمانی است که این خط مشی نگاشتی به توالی از اقدامات است. $[0, 1] \rightarrow S^+ B : \beta$ شرط خاتمه را مشخص می‌کند که برای هر Option در هر حالتی با چه احتمالی می‌تواند خاتمه بیابد.

می‌توان نشان داد که Option‌ها قادر به استفاده شدن به صورت تبادل‌پذیری با اقدامات ابتدایی^۷ در روش‌های برنامه‌ریزی پویا^۸ (DP) و در در روش‌های یادگیری مانند یادگیری Q هستند. برای ارزیابی ارزش هر اقدام و حالت از جدولی به عنوان جدول Q استفاده می‌شود و برای اطلاعات بیشتر در مورد یادگیری Q به [۱] مراجعه شود.

در این نوشتار همانند سایر مقالاتی که در این زمینه داده شده است، تنها اقدامات گسترش‌یافته زمانی به جدول یادگیری Q که از جدول حالت و اقدامات ابتدایی استفاده می‌کند، اضافه می‌شوند. (اقدامات گسترش‌یافته شده زمانی متناظر با مجموعه حالاتی که در آنها تعریف شده‌اند، به مجموعه اقدامات اولیه آن حالات اضافه می‌شوند). بر اساس [۳] و [۴] مقادیر جدول Q در چارچوب SMDP تنها برای اقدامات گسترش‌یافته زمانی در حالتی که انتخاب شده به روز می‌شوند. در این مقاله همانند [۷] و [۸] بر اساس (۱) بهروز رسانی مقادیر جدول Q انجام می‌گیرد. در واقع اقداماتی به مجموعه اقدامات اضافه می‌گردند و نیاز به ساختار خاص دیگری نیست برای اطلاعات بیشتر به [۳] و [۴] مراجعه شود. هدف، تنها ارائه روشی برای استخراج گذرگاه است تا در فرایند یادگیری تسريع ایجاد کند و زمینه را برای انتقال دانش فراهم نماید. فرایند یادگیری مشابه با [۷] تا [۱۲] است

$$\begin{aligned} Q(s_i, o_i) = & Q(s_i, o_i) + \alpha(n(t, s_i, o_i)) \\ & \times [\gamma^r \max_{a' \in A'(t, +\tau)} Q(s_{i+\tau}, a') - Q(s_i, o_i)] \\ & + r_i + \gamma r_{i+1} + \dots + \gamma^{T-t} r_{T-1} \end{aligned} \quad (1)$$

که (۱) بهروز رسانی یادگیری را انجام می‌دهد، در جایی که T زمان واقعی برای Option، o_i است. $\alpha(n(t, s_i, o_i))$ تابع نرخ یادگیری و آرگومان ورودی آن تعداد دفعاتی است که o_i در حالت s_i تا زمان t تجربه شده است [۳].

عموم الگوریتم‌هایی که به استخراج گذرگاه می‌پردازند با استفاده از الگوریتم‌های پارتبیشن‌بندی گراف مانند [۸]، [۱۱] و [۱۲]، یا به بررسی نقش هر حالت در مسیرها مانند [۹]، [۱۳] و [۱۴] به دنبال پیدا کردن گذرگاه هستند. در روش‌های پارتبیشن‌بندی به طرق مختلف عموماً ترکیبات مختلف حالت‌ها بر شماری می‌شود و معیارهایی برای خوشه‌ها چک می‌شود تا بتوانند گراف را خوشه‌بندی کنند. در روشی مانند [۱۴] برای هر مسیری که هر حالت در مسیرهای متفاوت بررسی می‌شوند. مثلاً هزینه‌ترین مسیر بین این دو حالت قرار دارند استخراج می‌شوند و

3. Returned Value

4. Landmarks

5. Navigation

6. Bottlenecks

1. Primitive Actions

2. Dynamic Programming

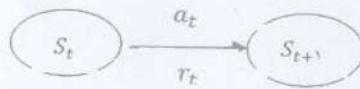
مسایلی با فضاهای حالت فاکتور شده^۲ مطرح کرده‌اند. در [۲۱] تحت مفهومی به عنوان skill تجربید زمانی را با پیدا کردن زیر خط مشی‌هایی که به صورت مکرر در راه حل‌هایی برای یک مجموعه از کارها روی می‌دهند، به دست آورده‌اند.

زیرا هدف [۱۳] حالتی هستند که به صورت مکرر دیده می‌شوند یا آنها که دارای یک گردایان پاداش بالا هستند. در محیط‌هایی که پاداش‌ها با تأخیر همراه هستند این روش ممکن است تواند کارایی خوبی را به نمایش بگذارد. روش [۹] به عنوان زیرا هدف، نواحی‌ای از فضای حالت در نظر می‌گیرد که عامل آنها را به صورت مکرر بر روی خط سیر موفقیت و نه بر روی عدم موفقیت ملاقات می‌کند. این شیوه به تعداد زیادی گام نیاز دارد تا تواند زیرا هدف را استخراج کند.

در شیوه‌هایی مانند [۷]، [۸] و [۱۰] گراف تراکنش حالت‌گرایی مینمایند. در [۱۱] گراف تراکنش حالت MDP را به عنوان یک گراف کامل در نظر می‌گیرند. زیرا هدف را به عنوان max-flow/min-cut روی این گراف، زیرا هدف را به صورت داخلی حالت‌گرایی از نواحی‌ای که در گراف انتقال MDP به صورت داخلي قویاً به هم متصل هستند، پیدا می‌کند. روش ارائه شده در [۸] با استفاده از یک شیوه خوشبندی مبتنی بر چگالی، خوشبندی‌های حالت را پیدا کرده و اقدامات گسترش‌یافته زمانی را تعریف می‌کند که به عنوان یک زیر خط مشی این اجازه را به عامل می‌دهند که به صورت کارایی از یک خوشبندی دیگری انتقال پیدا کند. سر این مقاله از دو مکانیزم خوشبندی متفاوت استفاده شده است که یکی تنها توپولوژی را به کار می‌گیرد و دیگری از ساختار پاداش از مسایل، اضافه بر توپولوژی استفاده می‌کند. زمانی که حالت جدیدی برای T گام زمانی مشاهده نشود، روال خوشبندی فوق فعال می‌شود و تخمین‌های کلی از گذرگاهها می‌زنند و به تدریج دقت خود را بالا می‌برد. بر طبق [۸] شیوه‌های مبتنی بر خوشبندی توانایی تقویت یادگیری را دارند اگر فضای حالت بتواند به خوبی تفکیک شود. در شیوه مطرح شده در [۱۰] که مبتنی بر توزیع گراف طیفی^۳ است، با استفاده از مقایسه در میان یال‌های متصل کننده، دو رأس متصل روش ارائه شده در [۱۲] را ارتقا داده است.

اکثر شیوه‌های مبتنی بر تقسیم‌بندی گراف دارای پیچیدگی زمانی $O(n^2)$ هستند که n تعداد رأس‌های گراف یا حالت‌ها است. پیچیدگی زمانی روش ارائه شده در [۷] $O(m)$ است در جایی که m تعداد لبه‌ها در گراف است. در این شیوه با الهام از ایده‌های مبتنی بر پیدا کردن لبه‌ها در پردازش تصویر، سعی شده است تا گذرگاهها استخراج شوند. برای اعمال این الگوریتم در این شیوه پیش‌فرض‌هایی باید رعایت شود که عملکرد این شیوه را برای محیط‌هایی که در آن اقدامات تصادفی هستند، دشوار می‌سازد. تصادفی عمل کردن در حذف یال‌ها گاهی اوقات باعث نیافتن یک گذرگاه خاص می‌شود و همچنین نوع گذرگاه‌هایی که می‌تواند تشخیص دهد محدود است. در روش‌های [۷]، [۱۰] و [۱۲] راهکاری برای تعیین الگوریتم برای محیط‌هایی که گراف‌های انتقال حالت جهت‌دار و نامتقارن است ارائه نشده است.

در اکثر روش‌های ارائه شده در این زمینه نقاط ضعف مشابهی یافت می‌شود مانند نیاز به کمک طراح برای آنالیز نتایج [۱۱] و [۱۲]، نیاز به ذخیره و پردازش دنیاهای اقدام و حالت [۱۲] و [۹] و همچنین شوابط محدوده کننده‌ای مانند این که خوشبندی‌های حالت‌های محیط باید دارای اندازه‌های مشابهی باشند. این شرایط نقش تعیین کننده‌ای در کیفیت نتایج



شکل ۱: انتقال حالت در RL.

پیدا کردن گذرگاه‌ها شبیه پیدا کردن درب خروج یک اتاق است اما بین توانایی‌های عامل مصنوعی و کودک تفاوت بسیار زیادی در این مورد وجود دارد. کودک به سادگی می‌تواند یک دیوار و مانع را از فضای باز بدون نزدیک شدن به آنها تفکیک کند. می‌توان محدوده دید، برد دید و تجارت قبلى را به عنوان مهم‌ترین ویژگی‌ها برای پیدا کردن و خروج از یک گذرگاه برای انسان‌ها در نظر گرفت. عموماً عامل‌ها در چارچوب MDP و SMDP دارای هیچ محدوده و برد دیدی نیستند یا این که در مقابل وسعت محیط، بسیار جزئی‌اند. در فرایندهای شبیه‌سازی عموماً عامل فاقد این توانایی‌ها است.

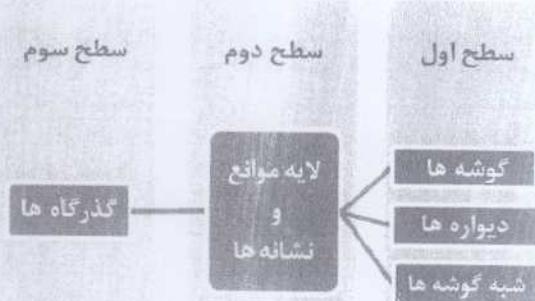
گروه‌بندی اشیا مبتنی بر یک سری از ویژگی‌ها، شیوه متدالوی است. در انتقال‌های حالت که در RL انجام می‌گیرد تنها چهار ویژگی برای یک عامل وجود دارد: حالت فعلی s ، اقدام انجام‌شده a ، پاداش r ، و حالت بعدی s' که در شکل ۱ نمایش داده شده است. این مورد می‌تواند این شکل دیده شود که حالت به عنوان اشیا به حساب بیانند و نتیجه هر اقدام بر روی هر حالت - حالت بعدی - به عنوان ویژگی‌ها در نظر گرفته شوند. به این ترتیب مانع یا دیوار، اشیایی هستند که اگر عامل به سمت آنها برود حالت بعدی و فعلی اش یکسان خواهد بود و در محیط‌های نویزی و تصادفی می‌توان گفت که در اکثر موقعیت این رویداد رخ می‌دهد. استفاده از نشانه‌ها برای تشخیص گذرگاه‌ها و شیوه استخراج گذرگاه‌ها را می‌توان از نوآوری‌های روش پیشنهادشده در این حوزه قلمداد کرد. در این مقاله، نشانه به عنوان یک ویژگی ادراکی بر جسته در نظر گرفته شده است و به عبارت دیگر بعضی از حالات در فضای حالت، نشانه‌هایی برای یافتن گذرگاه‌ها هستند. در این مقاله آنها به وسیله یک مکانیزم تشخیص اشیا به صورت سلسله مرتبی از عناصر محیطی که به ندرت تغییر می‌کنند، استخراج می‌شوند.

عامل ویژگی‌های سطح بالایی را مبتنی بر ترکیب هوشمندانه‌ای از تجارت و مشاهدات سطح پایین می‌تواند استخراج کند. از یک شیوه خوشبندی، مبتنی بر اقدامات عامل برای تشخیص اشیایی مانند دیوارها و گوششها استفاده شده است. این اشیا در یک شیوه سلسله مرتبی برای شکل دادن اشیای سطح بالاتر بررسی و ترکیب می‌شوند.

۳- کارهای مشابه

روش ارائه شده در [۱۷]، بعضی از مفاهیم تقریباً مشابه با آنچه در این مقاله به کار گرفته شده را داراست مانند خطوط بحرانی^۴. خطوط بحرانی در این مقاله متناظر با گذرگاه‌ها یا حالاتی است که در بین نشانه‌ها قرار گرفته‌اند. اما این تفاوت وجود دارد که او از یک سیستم سنسوری (شامل چندین سنسور) برای ریات استفاده کرده است که به عامل این توانایی را می‌دهد که برد و محدوده دید داشته باشد و با ترکیب اطلاعات سنسوری از هر کدام از اجزا می‌تواند یک نگاشت از دنیای اطراف خودش را مهیا کند. در حالی که در روش پیشنهادی نحوه به دست آوردن گذرگاه‌ها کاملاً متفاوت است.

روش‌های [۱۸] تا [۲۰] ساخت به سلسله مرتبی از تجربیدها را در



(ب)

-

(الف)

-

(ج)

۱) تشخیص اشیا به صورت سلسله مراتبی

الف) محاسبه اقدامات موثر برای هر حالت.

ب) خوشه بندی حالتها بر اساس اقدامات موثر.

ج) استخراج خوشه‌های حالت‌های گوش، دیوار و شبکه گوش‌ها.

۲) استخراج نشانه‌ها

ترکیب اشیا سطح پایین پیداکردن لایه‌های مانع و نشانه‌ها

۳) پیدا کردن گذرگاه‌ها

استخراج گذرگاه‌ها بر اساس جستجوی محلی

۴) شکل گیری option

ایجاد اقدامات گسترش یافته زمانی براساس گذرگاه‌های استخراج شده.

(الف)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

(ب)

-

(ج)

-

(د)

-

(ه)

-

(ز)

-

شرایط بین آنها برقرار باشد، آنها با یکدیگر ترکیب شده و اشیای سطح بالاتری را شکل می‌دهند. فروانند ترکیب عموماً از لبه‌های خوش‌های مربوط به حالت‌های مانع به عنوان مهم‌ترین بخش از خوش‌های اشیا از حالات برای ترکیب و نمایش آنها استفاده می‌کند. برای حرکت کردن یا مرتبط کردن اشیا در فضای بدهطور ضمته، ما به نزدیک‌ترین بخش از موانع یا اجسام توجه می‌کنیم. در فضاهای دو بعدی، لبه‌های دیوارها، دو حالت هستند که در انتهای دو سر آن قرار گرفته‌اند (شکل ۳-الف). در فضای سه بعدی، گوشه‌های یک صفحه، لبه‌ها هستند. در گام نهایی، لبه‌های از اشیای سطح بالاتر به عنوان نشانه‌ها شناخته می‌شوند (شکل ۳-ب). در حقیقت ما به دنبال خطوط بحرانی در فضای حالت هستیم. خطوط بحرانی همان گونه که در [۱۷] اشاره شده است منتظر با گذرگاه‌های باریک هستند. به عبارت دیگر اگر تعداد اقداماتی که مورد نیاز برای رفتن از یک نشانه به دیگری کمتر از یک آستانه باشد، این مقدار آستانه می‌تواند مناسب با اندازه خوش‌های متصل شده باشد. می‌توان حالت‌هایی که بین نشانه‌ها قرار دارند را با در نظر گرفتن شرایطی، به عنوان حالت‌های گذرگاه در نظر بگیریم. این شرایط بسیار ساده هستند، مانند این که جسم خوش‌های متصل شده نسبت به اندازه حالت‌های قرارگرفته در خطوط بحرانی مقدار متناسبی باشد.

۴-۳-۳ پیداکردن گذرگاه‌ها

از آنجایی که ما گذرگاه‌ها را بر اساس نشانه‌ها (لبه‌های حالت‌های موانع) تشخیص می‌دهیم، عامل در واقع خطوط بحرانی را در نظر می‌گیرد. اگر یک نشانه پیدا شد، از آنجایی که عامل دارای محدوده دید نیست و نمی‌تواند رابطه نشانه را با یقین عناصر به صورت درجا در نظر بگیرد، او باید آن نشانه را مورد بررسی قرار دهد که آیا هیچ نشانه یا مانع دیگری وجود دارد به گونه‌ای که آنها دارای فاصله نزدیکی در فضای اقدام باشند (شکل ۳-ج).

برای یافتن نشانه‌ها یک برد محدود چک می‌شود که آیا نشانه دیگری وجود دارد یا نه. اگر وجود داشته باشد، حالت‌های بین آنها به عنوان گذرگاه شناخته می‌شوند. چندین استراتژی را برای پیداکردن گذرگاه می‌توان در نظر داشت. اول این که اگر هر حالت محیط با یک مختصات (x, y) مشخص بشود، عامل می‌تواند از این مختصات و معیار فاصله‌ای مبتنی بر آن به عنوان یک راهنمای برای تشخیص نشانه دیگر استفاده کند. راه دیگر استفاده از فاصله حقیقی است که برای اهداف دیگری (تعیین حالت‌هایی که در مجموعه اولیه باید قرار بگیرند) در [۱۲] استفاده می‌شود. اگر فاصله نشانه با مانع یا نشانه دیگری کمتر از مقدار مشخصی باشد، این حالت‌ها یک ورودی از گذرگاه را تشکیل می‌دهند. این که تا چه بردی از نشانه یافت شده مورد بررسی قرار بگیرد تنها پارامتر موجود در روش پیشنهادی است.

۴-۴ شکل‌گیری Option

بعد از این که عامل یک گذرگاه را پیدا کند، می‌تواند یک Option را برای این خوش‌های شکل دهد. در واقع، زمانی که خوش‌های و گذرگاه‌ها بین آنها شکل گرفتند، یک Option برای هر مجموعه گذرگاه در هر خوش‌های از حالات شکل می‌گیرد. بر طبق تعریف Option برای همه حالت‌ها در هر یک از خوش‌های مجموعه ورودی Option، I ، به مقدار یک مقداردهی می‌شود. شرط خاتمه برای این Option برای هر حالت در این خوش‌های به مقدار صفر نسبت دهی می‌شود و برای حالت‌های دیگر خارج از این خوش‌های یک مقداردهی می‌شود. خط مشی Option، π ، بر طبق روش پیشنهادی در [۴] می‌تواند شکل بگیرد. ما زیر مسایلی داریم و باید یک یادگیری

۴-۱ تشخیص اشیا به صورت سلسله مراتب

عامل حالات را مبتنی بر ویژگی‌ها دسته‌بندی می‌کند. اگر اقدامات عامل تصادفی هستند، عامل می‌تواند از استراتژی رأی اکثریت برای پیداکردن اقدامات حاکم اصلی در هر حالت استفاده کند. در ابتدا هیچ خوش‌های وجود ندارد و حالت اول، یک خوش‌های جداگانه را شکل می‌دهد. حالات با خوش‌های موجود بر اساس شرایطی که در ادامه توضیح داده شده‌اند چک می‌شود. اگر این شرایط بین حالت و خوش‌های برقرار باشد، این حالت به آن خوش‌های اضافه و جستجو متوقف می‌شود، در غیر این صورت یک خوش‌های جدید شکل خواهد گرفت.

دو شرط فوق به شکل یک سلسله مراتب تعریف شده‌اند. اگر اولین شرط رد بشود، این حالت با خوش‌های دیگری چک می‌شود و در غیر این صورت شرط دیگری باید چک شود. در پایین ترین سطح از سلسله مراتب، هر کدام از حالات دیوار برای عامل به عنوان یک خوش‌های از حالات تفسیر می‌شود که دارای دو ویژگی هستند. در پایین ترین سطح سلسله مراتب، اولین ویژگی برای اعضای یک خوش‌های این است که اقداماتی که به ازای آنها حالت فعلی با حالت نتیجه برابر است، برای اعضای یک خوش‌های یکسان هستند. شرط دوم در پایین ترین سطح عبارت است از این که اگر خوش‌هایی بیشتر از یک عضو داشته باشد، برای هر کدام از اعضای باید حداقل عضو دیگری وجود داشته باشد به گونه‌ای که آنها بایستی با یکدیگر مجاور باشند. در سطح دوم، حالت‌هایی که در یک خوش‌های قرار دارند دارای این خاصیت هستند که به ازای هر دو حالت که با یکدیگر مجاورند، حداقل باید همسایه‌ای از این حالت‌ها وجود داشته باشد که آنها تیز با یکدیگر مجاور باشند. به عبارت دیگر شرط دوم در صورتی برقرار است که اگر یک حالت برای مثال 'S' تحت یک اقدام به حالت دیگری مانند 'S' برود، همسایه‌های 'S' تحت حداقل یک اقدام با 'S' مجاور باشند. اگر این شرط برقرار نباشد ما می‌توانیم 'S' و 'S' را به عنوان یک ورودی گذرگاه در نظر بگیریم. این ویژگی گذرگاه را در مواردی مانند مثال می‌بینیم در [۱۲] و مثال تاکسی که در [۶] معرفی شده است، تشخیص می‌دهد. با این شرایط، خوش‌هایی از حالت‌ها شکل می‌گیرند که با یک دیوار مجاور هستند (التبه بدون گوشه‌ها). هر کدام از گوشه‌ها دارای خوش‌هایی متعلق به خودشان هستند. حالت‌های مانع، ترکیبی از حالات گوشه و دیوارها هستند که با یکدیگر مجاور بوده‌اند. حالت‌هایی به عنوان شیه گوشه‌ها نیز استخراج می‌شوند که به عنوان عناصر متصل کننده گوشه‌ها و دیوارها استفاده می‌شوند. این عناصر در شکل ۳-الف نشان داده شده‌اند و به ترتیبی که در شکل ۲-ب آمده است عناصر سطح بالاتری از آنها استخراج شده‌اند.

۴-۲ استخراج نشانه‌ها

یک حالت شیه گوشه ممکن است یک نشانه باشد، در غیر این صورت در سطوح بالاتر برای ترکیب حالت‌های مانع و پیداکردن نشانه‌های دقیق تر استفاده می‌شود. این روال نیز تا حدی شبیه فرایند تشخیص اشیا به صورت سلسله مراتبی است. با استفاده از شرایط تعریف شده، خوش‌هایی از عناصر حالت‌های مانع، مانند دیوارها و گوشه‌ها استخراج می‌شوند که بعد از این رویداد آنها با یکدیگر ترکیب می‌شوند.

با استفاده از مفهوم گوشه‌ها و شبیه گوشه‌ها، می‌توانیم اشیا را در سطوح پایین تر ترکیب کنیم و آنها را در یک خوش‌های قرار دهیم. در سطح دوم از سلسله مراتب، شبیه گوشه‌ها بررسی می‌شوند که آیا آنها نشانه هستند و همچنین بعضی از اشیای سطح پایین تر شبیه دیوارها، گوشه‌ها و شبیه گوشه‌ها برای ترکیب شدن با یکدیگر مورد بررسی قرار می‌گیرند. اگر

۵- نتایج آزمایش

از آنجایی که ایده پیشنهادی مقاله مربوط به تئوری الگوریتم‌های RL است، بنابراین Data Set خاصی وجود ندارد و از بسترهاست استاندارد آزمایشی متداولی همانند [۷] تا [۱۴] و [۲۲] که در این زمینه داده می‌شوند استفاده می‌کنیم.

در این بخش ما نتایج الگوریتم پیشنهادی را بر روی یک محیط آزمایش ارائه می‌کنیم و آنها را با روش‌های تجزیه متوازن که در [۱۰] ارائه شده است، مقایسه می‌کنیم. روش پیشنهادی بر اساس کل گراف تراکنش‌های حالت و اقدام کار می‌کند (حالت غیر همزمان). پارامتر ارزیابی که در این مقاله مشابه مقالات دیگر در این حوزه مورد بررسی قرار می‌گیرد عبارت است از سرعت یادگیری (نمودار تعداد گام‌های مورد نیاز برای رسیدن به هدف در تعداد trial)، در این نوشتار نیز با استفاده از این معیار به مقایسه روش پیشنهادی با سایر روش‌های دیگر پرداخته شده است که در نمودارهای شکل ۴ مشاهده می‌شود. همچنین روش‌های مورد نظر از منظر نیاز به طراحی برای مقداردهی پارامترهای مورد نیاز برای استخراج مورد ارزیابی قرار می‌گیرند که در بخش کارهای مشابه به آن پرداخته شده است. تقریباً در تمامی مقالاتی که در این زمینه داده شده است معيارهای ارزیابی که مورد بررسی قرار می‌گیرند عبارتند از سرعت یادگیری و پارامترهای مورد نیاز که توسط روش مورد استفاده قرار می‌گیرند و در این نوشتار نیز از این دو معیار استفاده شده است.

فرض می‌کنیم نتیجه حرکت عامل با احتمال α مطابق با اقدامی است که انجام داده است و در غیر این صورت، به شکل تصادفی با احتمال $1 - \alpha$ در یکی از جهت‌های دیگر حرکت می‌کند. حالت‌های شروع و هدف به صورت تصادفی در هر اجرا انتخاب می‌شوند. پاداش - برای هر اقدام در نظر گرفته می‌شود مگر اقداماتی که عامل را به هدف پرساند که برای آنها پاداش $+10$ قابل می‌شود. نرخ کاهش با مقدار $\gamma = 0.9$ مقداردهی شده است و نرخ یادگیری با مقدار $\alpha = 0.1$ ثابت نگه داشته می‌شود. نرخ ϵ در خط مشی greedy (مکانیزم انتخاب اقدام) برابر با 0.1 در نظر گرفته می‌شود و مقادیر اولیه Q با صفر مقداردهی می‌شوند. برای هر مقایسه تعداد پنج حالت تصادفی برای شروع و پنج حالت هدف در خوش‌هایی که گذرگاه آنها تشخیص داده نشده‌اند تعیین می‌شوند. اگر خوش‌هایی حالت به درستی تشخیص داده نشوند آنگاه یادگیری تقویتی با استفاده از Optionها ممکن است یادگیری را نه تنها بهبود ندهد بلکه آن را بدتر نیز نماید. به این منظور نقاط هدف برای نشان دادن این مهم در این خوش‌های انتخاب شده‌اند. برای هر کدام از این جفت‌ها تعداد اجرایها برابر ۵ در نظر گرفته می‌شود. در هر دور زمانی که عامل به نقطه هدف برسد، دور جاری خاتمه می‌یابد. در آزمایشاتی که در ادامه به آنها پرداخته می‌شود، فرض می‌کنیم عامل در هر گام تنها چهار اقدام ممکن دارد: بالا، پایین، چپ و راست.

به طور مشابه با چارچوب MDP در هر گام زمانی t عامل در یک حالت که با a_t مشخص می‌شود، قرار دارد و می‌تواند یک اقدام a_{t+1} را از مجموعه اقدامات ممکن خود انتخاب کند. در برایند این اقدام، یک پاداش r_{t+1} از محیط او دریافت می‌کند و حالت محیط به s_{t+1} تغییر می‌کند. اگر مانع در جهت اقدامی که عامل آن را انتخاب کرده است وجود داشته باشد، مکان عامل تغییر نمی‌کند.

از یک ماز ساده که در شکل ۴-الف آمده است، به عنوان محیط آزمایش استفاده شده است. این شکل با اعمال تغییراتی از محیط آزمایشی در [۸] برگرفته شده است. در این آزمایش روش ارائه شده در [۱۰] نمی‌تواند همه گذرگاه‌ها را تشخیص دهد (شکل ۴-ب). روش ارائه شده

جدید برای هر کدام از آنها انجام بدھیم. برای مثال اقداماتی که باعث رسیدن به گذرگاه می‌شوند یک پاداش مثبت و بقیه اقدامات در حالت عادی یک پاداش صفر دریافت می‌کنند. عامل بدین صورت زیر خط مشی‌هایی را برای خروج از هر خوشی یاد می‌گیرد. در این مقاله همانند سایر مقالاتی که در زمینه استخراج گذرگاه داده شده است، بررسی چارچوب Option به صورت مفصل‌تر را به [۳] و [۴] ارجاع می‌دهیم.

۴- مزایا

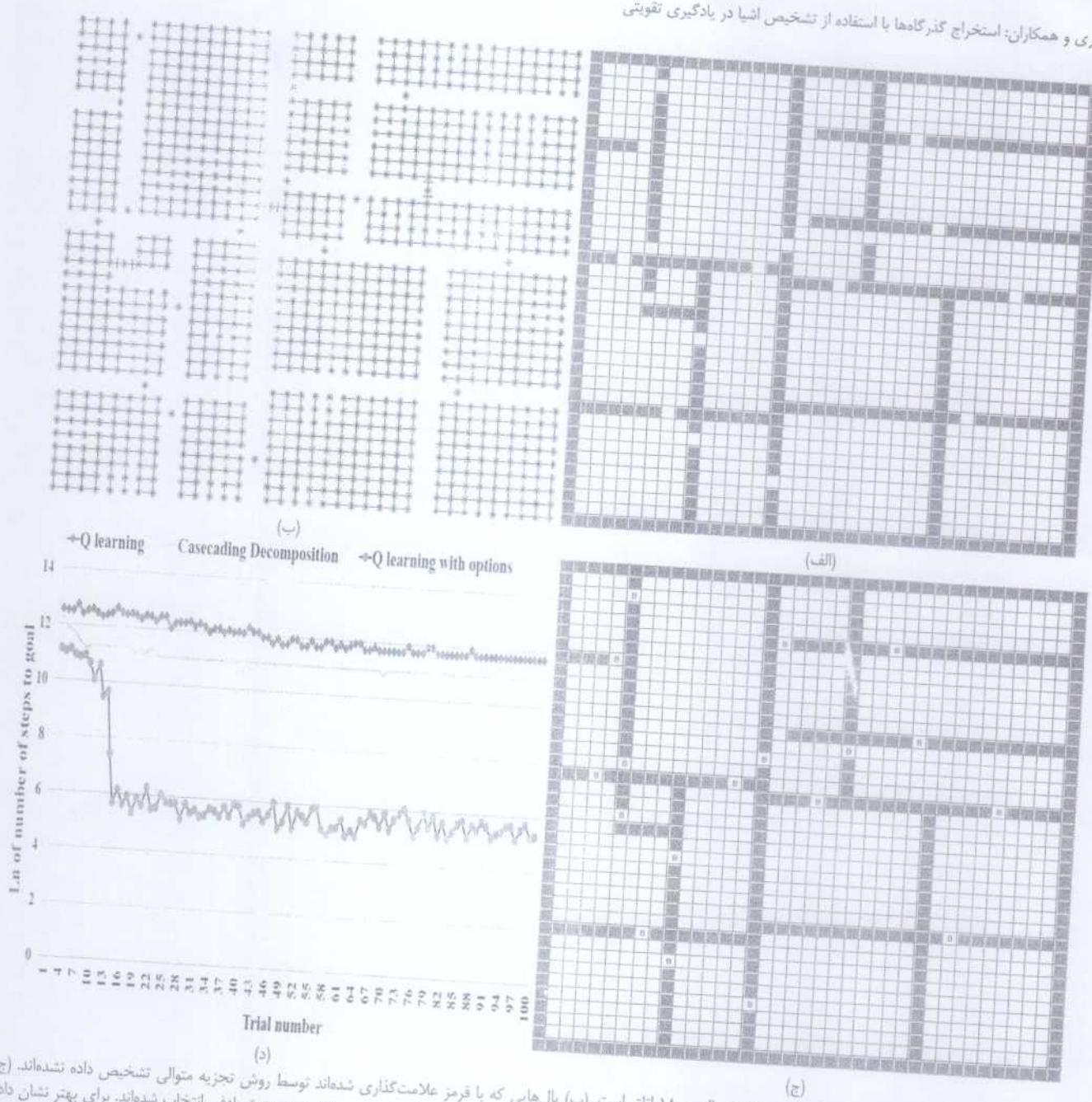
ما ادعا می‌کنیم که روش ارائه شده می‌تواند گذرگاه‌ها را با دقت بالای تشخیص دهد چرا که او نیازمندی‌های اصلی - ویژگی‌های لازم - از امکان گذرگابودن را تشخیص می‌دهد. برخلاف روش‌هایی مانند [۸]، [۱۰]، [۱۲] و [۱۴] که دقت آنها به اندازه خوشی‌ها در مقایسه با یکدیگر وابسته است یا فقط توانایی تشخیص گذرگاه‌های خاصی را دارند، روش پیشنهادی گذرگاه‌ها را با دقت بالایی استخراج می‌کند. همچنین هزینه محاسباتی روش ارائه شده در مقایسه با سایر روش‌ها از نکات قوت آن به حساب می‌آید.

نتایج ما نشان می‌دهد که نشانه‌های استخراج شده توسط روال تشخیص اشیاء، نشانه‌های مطمئنی برای تشخیص گذرگاه هستند. همچنین کمتر نیازمند و حساس به تنظیم پارامتر ورودی‌اند، بنابراین به داشتن اولیه در مورد خصوصیات محیط یا اطلاعاتی که عامل عموماً هیچ ذهنیتی نسبت به آنها ندارد، کمتر نیازمند هستند و بنابراین خودمختاری عامل بهتر حفظ می‌شود و همچنین پارامتر ورودی آن قابل درک است. گذرگاه‌های استخراج شده دارای قطبیت هستند به این معنی که در همه اجرایها، حالت‌های مشخصی به عنوان گذرگاه شناخته می‌شوند. زمانی که عامل به دنبال گذرگاه تنها روی مسیرهای موقیت‌آمیز است (زیرا هدف)، ممکن است که گذرگاه‌های دیگری نادیده گرفته شوند. بدون شک همه گذرگاه‌ها روی مسیرهای موقیت‌آمیز نیستند. با تشخیص ویژگی‌های اصلی برای وجود گذرگاه می‌توان از آنها یک شک شده‌اند، فارغ از هدف جاری استفاده کرد. آنها برای کاربرکردن اکتشاف در سایر بخش‌های فضای حالت و یا برای تعمیم داشت فعلی برای محیط‌هایی دیگر مشابه (که تنها تابع پاداش آنها متفاوت است)، می‌توانند مفید باشند. عامل خواهد توانست با تلاش سیار کمتر روی کارهای متفاوت نتیجه بگیرد. بعضی از مزیت‌های دیگر استفاده از گذرگاه در [۱۲] شرح داده شده است.

۴- پیچیدگی الگوریتم

روش ارائه شده نیازی به ذخیره با پردازش اطلاعات بر روی مسیرها مانند عملیات‌های پیش‌پردازشی شبیه حذف حلقه‌ها ندارد. پیچیدگی زمانی الگوریتم برای یک دنبای توری $O(n!)$ حالت برابر با $O(n^t)$ است. فرض کنید فضای حالت داری k خوشی باشد. حالت‌ها به ترتیج قرار است در خوشی‌ها قرار بگیرند و خوشی‌ها نیز به ترتیج شکل می‌گیرند. برای این که یک عنصر در یک خوشی قرار بگیرد باید بررسی شود که آیا حالت مورد بررسی با حالتی در این خوشی از طریق اقسامی قابل دسترسی است یا نه. بنابراین ارتباط هر عنصر با عناصر فعلی موجود در خوشی‌هایی که تاکنون شکل گرفته‌اند باید بررسی شود. از آنجایی که ارتباط هر عنصر باید با عناصر یک خوشی چک شود و خوشی‌ها به ترتیج شکل می‌گیرند پس تعداد مقایسه‌ها برابر است با

$$(2) \quad 1 + 2 + 3 + \dots + (n-1) = O(n^t)$$



(ج) شکل ۴: (الف) یک دنیای Maze ساده که دارای ۱۰۹۰ حالت و ۱۸ اتاق است. (ب) مالهایی که با قرمت علامت گذاری شده‌اند توسط روش تشخیص متولی تشخیص داده شده‌اند. (ج) گذرگاه‌ها در درسی توسط روش پیشنهادی تشخیص داده شده‌اند. (د) برای ارزیابی سرعت یادگیری تعداد ۵ حالت شروع و هدف به صورت تصادفی انتخاب شده‌اند. برای بهتر نشان داده شدن تفاوت‌ها از لگاریتم مقادیر استفاده شده است (مقایسه‌ای بین روش‌هایی که گذرگاه را به درستی تشخیص داده‌اند مانند روش پیشنهادی، یادگیری Q و روش تشخیص متولی).

ست و همچنین تیازی به ذخیره مسیرها ندارد.
تصور کنید که در یک سبد ۲ مهره قرمز و ۱۰ مهره آبی وجود دارد و
ما می‌خواهیم مهره‌ها را تنها از همدیگر تفکیک کنیم. ما ۲ مهره قرمز را
بیندا می‌کنیم و آنها را بر می‌داریم تا این که بخواهیم ۱۰ مهره آبی را
برداریم و در واقع با انتخاب شیء اقلیت، تفکیک هوشمندانه‌ای را انجام
می‌دهیم. از آنجایی که عموماً در محیط‌هایی که عامل‌ها با آنها سر و
کار دارند، فضای حالت دارای تعداد بسیار زیادی حالت است، ما معتقد
همستیم استفاده از نشانه‌ها برای پیداکردن گذرگاه کاری معقول و با توجیه
بیولوژیک است.

مراجع

- [1] L. Kaelbling, M. Littman, and A. Moore, "Reinforcement learning: a survey," *J. of Artificial Intelligence Research*, vol. 4, pp. 237-285, 1996.

در این مقاله و در [۱۰] نمی‌تواند همه گذرگاهها را تشخیص دهد (شکل ۴-ب). روش ارائه شده در این مقاله و [۷] همه گذرگاهها را به درستی تشخیص می‌دهند (شکل ۴-ج). در شکل ۴-د برای نشان دادن تأثیر مناسب تشخیص درست همه گذرگاهها، مقایسه‌ای بین روش ارائه شده یا [۷] که توانسته‌اند گذرگاهها را به درستی تشخیص دهنده و یادگیری $O(n)$ و تجزیه متولی که بعضی از گذرگاهها را تشخیص نمی‌دهد، ارائه شده است.

۶- نتیجه‌گیری

- [17] S. Thrun, "Learning metric-topological maps for indoor mobile robot navigation," *Artificial Intelligence*, vol. 99, no. 1, pp. 21-71, 1998.
- [18] N. Mehta, S. Ray, P. Tadepalli, and T. Dietterich, "Automatic discovery and transfer of task hierarchies in reinforcement learning," *AI Magazine*, vol. 32, no. 1, p. 35, 2011.
- [19] A. Jonsson, *A Causal Approach to Hierarchical Decomposition in Reinforcement Learning*, Ph. D. Thesis, University of Massachusetts Amherst, Feb. 2006.
- [20] B. Hengst, *Discovering Hierarchy in Reinforcement Learning*, Ph. D. Thesis, University of New South Wales, Australia, Dec. 2003.
- [21] S. Thrun and A. Schwartz, "Finding structure in reinforcement learning," *Proc. 5th Annual Conf. on Advances in Neural Information Processing Systems, NIPS'95*, pp. 385-392, 1995.
- [22] C. C. Chiu, "Subgoal identification for reinforcement learning and planning in multiagent problem solving," in *Proc. of 5th German Conf. on Multiagent System Technologies*, pp. 37-48, 2007.
- بهزاد غضنفری تحصیلات خود را در مقطع کارشناسی در دانشگاه فردوسی مشهد و کارشناسی ارشد در دانشگاه علم و صنعت ایران بهترتبی در سال‌های ۱۳۸۸ و ۱۳۹۰ به پایان رسانده است. زمینه‌های تحقیقاتی ایشان عبارتند از: بردازش نرم، یادگیری ماشین، یادگیری تقویتی، سیستم‌های چند عاملی و الگوریتم‌های تعابق الگو و رشته.
- ناصر مژیینی در سال ۱۳۶۹ مدرک کارشناسی خود را از دانشگاه صنعتی شریف در مهندسی برق گرایش کامپ - سخت افزار اخذ نمود و سپس در سال ۱۳۷۲ مدرک کارشناسی ارشد را در رشته سیستم‌های اطلاعاتی و تله‌ماینیک از سویلک فرانسه و همچنین در سال ۱۳۷۷ از دانشگاه رن یک فرمانه در رشته انفورماتیک دریافت نمود. وی از سال ۱۳۷۹ در دانشکده مهندسی کامپیوتر دانشگاه علم و صنعت ایران مشغول به فعالیت گردید و اینک نیز عضو هیأت علمی این دانشکده می‌باشد. زمینه‌های علمی مورد علاقه نامبرده عمدتاً در زمینه رایانش نرم، فناوری اطلاعات و شبکه‌های کامپیوتری است.
- محمد رضا جاهد مطلق در سال ۱۳۵۷ مدرک کارشناسی مهندسی برق خود را از دانشگاه صنعتی شریف دریافت نمود و پس از ۶ سال در فعالیت‌های صنعتی برای ادامه تحصیلات خود به انگلستان رفت و مدارک کارشناسی ارشد و دکتری خود را در سال‌های ۱۳۶۶ و ۱۳۷۰ از دانشگاه برادفورد انگلستان در زمینه مهندسی کنترل اخذ نمود. پس از بازگشت از سال ۱۳۷۰ تاکنون در دانشگاه علم و صنعت ایران مشغول به فعالیت می‌باشد. زمینه‌های علمی مورد علاقه نامبرده متعدد بوده و شامل موضوعاتی مانند سیستم‌های پیچیده محاسبات آشوب گونه سیستم‌های هایبرید رباتیک و کنترل سیستم‌های پیچیده می‌باشد.
- [2] M. Ghavamzadeh, S. Mahadevan, and R. Makar, "Hierarchical multi-agent reinforcement learning," *Autonomous Agents and Multi-Agent Systems*, vol. 13, no. 2, pp. 197-229, Sep. 2006.
- [3] A. Barto and S. Mahadevan, "Recent advances in hierarchical reinforcement learning markov and semi-markov decision processes," *Discrete Event Dynamic Systems*, vol. 13, pp. 41-77, 2003.
- [4] R. S. Sutton, D. Precup, and S. Singh, "Between MDPs and semi-MDPs: a framework for temporal abstraction in reinforcement learning," *Artificial Intelligence*, vol. 112, no. 1-2, pp. 181-211, Aug. 1999.
- [5] R. Parr and S. Russell, "Reinforcement learning with hierarchies of machines," in *Proc. Conf. on Advances in Neural Information Processing Systems*, pp. 1043-1049, 1997.
- [6] T. G. Dietterich, "Hierarchical reinforcement learning with the MAXQ value function decomposition," *J. of Artificial Intelligence Research*, vol. 13, pp. 227-303, 2000.
- [7] G. Kheradmandian and M. Rahmati, "Automatic abstraction in reinforcement learning using data mining techniques," *Robotics and Autonomous Systems*, vol. 57, no. 11, pp. 1119-1128, Nov. 2009.
- [8] S. Mannor, I. Menache, A. Hoze, and U. Klein, "Dynamic abstraction in reinforcement learning via clustering," in *Proc. 21st Int. Conf. on Machine learning, ICML'04*, p. 560-567, 2004.
- [9] E. A. McGovern, *Autonomous Discovery of Temporal Abstractions from Interaction with an Environment*, Citeseer, 2002.
- [10] C. Chiu and V. W. Soo, "Automatic complexity reduction in reinforcement learning," *Computational Intelligence*, vol. 26, no. 1, pp. 1-25, Feb. 2010.
- [11] I. Menache, S. Mannor, and N. Shimkin, "Q-cut - dynamic discovery of sub-goals in reinforcement learning," in *Proc. of the 13th European Conf. on Machine Learning*, pp. 295-306/2002.
- [12] O. Simsek, A. P. Wolfe, and A. G. Barto, "Identifying useful subgoals in reinforcement learning by local graph partitioning," in *Proc. of the 22nd Int. Conf. on Machine Learning, ICML'05*, pp. 816-823, 2005.
- [13] B. Digney, "Learning hierarchical control structures for multiple tasks and changing environments," in: *Proc. of 5th Int. Conf. on Simulation of Adaptive Behavior: From Animals to Animats 5*, pp. 321-330, 1998.
- [14] O. Simsek and A. Barto, "Skill characterization based on betweenness," in *Proc. 22nd Annual Conf. on Advances in Neural Information Processing Systems, NIPS'08*, pp. 1497-1504, 2008.
- [15] M. Riesenhuber and T. Poggio, "Hierarchical models of object recognition in cortex," *Nature Neuroscience*, vol. 2, no. 11, pp. 1019-25, Nov. 1999.
- [16] T. S. Collett and P. Graham, "Animal navigation: path integration, visual landmarks, and cognitive maps," *Current Biology*, vol. 14, no. 12, pp. 475-457, Jun. 2004.