



دانشگاه صنعتی اصفهان
دانشکده برق و کامپیوتر

**به کارگیری روش یادگیری تقویتی برای مسیریابی پویا در شبکه
به عنوان یک سیستم چندعاملی**

پایان نامه کارشناسی ارشد - هوش مصنوعی

سعید مجیدی

استاد راهنما

دکتر مسعود رضا هاشمی



دانشگاه صنعتی اصفهان

دانشکده برق و کامپیوتر

پایان نامه کارشناسی ارشد مهندسی کامپیوتر - هوش مصنوعی آقای سعید مجیدی

تحت عنوان

به کارگیری روش یادگیری تقویتی برای مسیریابی پویا در شبکه
به عنوان یک سیستم چندعاملی

در تاریخ ۱۳۸۶/۲/۱ توسط کمیته زیر مورد بررسی و تصویب نهائی قرار گرفت .

دکتر مسعودرضا هاشمی

۱- استاد راهنمای پایان نامه

دکتر مازیار پالهننگ

۲- استاد مشاور پایان نامه

دکتر جمال گلستانی

۳- استاد داور

دکتر جواد عسکری

۴- استاد داور

دکتر علی محمد دوست حسینی

سرپرست تحصیلات تکمیلی دانشکده

من لم يشكر المخلوق لم يشكر الخالق

از دست و زبان که بر آید کز عهده شکرش به در آید

حمد و سپاس بی پایان از آن خالق یکتاست که اندک داشته و اندوخته ناچیز خود را مرهون لطف و مرحمت بی کران او می‌باشم. پس از تواضع در برابر آستان حضرت دوست، لازم می‌دانم از خانواده عزیزم، بخصوص پدر و مادر مهربانم به خاطر زحمات بی دریغشان صمیمانه تشکر کنم.

از استاد ارجمند و بزرگوار، جناب آقای دکتر **مسعودرضا هاشمی**، استاد راهنمای پایان نامه، که در طی سال‌های اخیر افتخار شاگردی در مکتب پر خیر و برکت علم و اخلاق ایشان مایه مباهات بنده بوده و همواره از رهنمودهای ارزشمندشان بهره‌مند گردیده‌ام، تشکر و قدردانی می‌نمایم و از ایزد منان سلامتی و بهروزی را برای ایشان خواستارم.

از استاد مشاور پایان نامه، جناب آقای دکتر **مازیار پالهنک**، که بسیار صمیمانه و با صبر و حوصله فراوان تجربیات ارزشمند خود را در اختیار اینجانب قرار دارند تشکر می‌نمایم.

از اساتید محترم، جناب آقای دکتر گلستانی و جناب آقای دکتر عسکری که زحمت داوری این پایان نامه را بر عهده گرفتند و با راهنمایی‌های خود کمک ارزشمندی در تکمیل و بهبود این پایان نامه داشتند تشکر فراوان دارم.

همچنین از جناب آقای دکتر دوست حسینی، سرپرست محترم تحصیلات تکمیلی دانشکده و سرکار خانم نکویی به خاطر زحمات بی‌شائبه ایشان در طی این دوره کمال تشکر و قدردانی را به عمل می‌آورم.

از کلیه دوستان عزیز که دوران خوشی را در کنار آنها گذراندم و افتخار آشنایی با آنها را داشتم، تشکر می‌کنم و از خداوند متعال، سلامتی و بهروزی همه این عزیزان را خواستارم.

سعید مجیدی

بهار ۱۳۸۶

دانشگاه صنعتی اصفهان

کلیه حقوق مادی مترتب بر نتایج مطالعات، ابتکارات و نوآوریهای ناشی از تحقیق موضوع
این پایان نامه متعلق به دانشگاه صنعتی اصفهان است.

این پایان نامه با حمایت مرکز تحقیقات مخابرات ایران (قرارداد
شماره ۵۰۰/۲۹۰۷/ت مورخ ۱۳۸۳/۳/۲۴) به انجام رسیده است.

تقدیم به

پدر و مادر

خوبم،

برادر و

خواهر

عزیزم،

و او که

هرچه داریم

از اوست.

فهرست مطالب

صفحه	عنوان
هشت	فهرست مطالب
۱	چکیده
	فصل اول: مقدمه
۲	۱-۱ مقدمه
۴	۲-۱ روند ارائه پایان نامه
	فصل دوم: یادگیری تقویتی
۶	۱-۲ مقدمه
۷	۲-۲ مدل یادگیری تقویتی
۹	۳-۲ تعریف سیاست بهینه
۱۰	۱-۳-۲ مدل پاداش دراز مدت
۱۰	۲-۳-۲ تابع ارزش بهینه و سیاست بهینه
۱۱	۳-۳-۲ الگوریتم‌های تکرار مقدار و تکرار سیاست
۱۲	۴-۲ راهبردهای یادگیری برای مسائل یادگیری تقویتی
۱۲	۱-۴-۲ راهبردهای مدل-مبنا و مدل-آزاد
۱۳	۲-۴-۲ مزایای روش‌های بر مبنای مدل
۱۳	۵-۲ عمل یادگیری
۱۷	۶-۲ یادگیری Q
۱۸	۱-۶-۲ تابع Q
۱۹	۲-۶-۲ الگوریتمی برای یادگیری Q
۲۱	۷-۲ اعمال و پاداش‌های غیرقطعی
۲۲	۸-۲ مشکلات یادگیری تقویتی
۲۳	۹-۲ روشهای جستجو
۲۳	۱-۹-۲ جستجوی غیر جهت‌دار
۲۴	۲-۹-۲ جستجوی جهت‌دار
۲۶	۱۰-۲ سیستم‌های چندعاملی

فصل سوم: الگوریتم‌های مسیریابی بر مبنای هوش مصنوعی

۲۷ معرفی ۱-۳
۲۸ الگوریتم‌های مسیریابی ۲-۳
۲۸ ویژگی‌های الگوریتم‌های جدید مسیریابی ۱-۲-۳
۲۹ رده‌بندی الگوریتم‌های مسیریابی ۲-۲-۳
۳۲ خط مشی طراحی ۳-۲-۳
۳۴ الگوریتم‌های مسیریابی بر مبنای هوش جمعی ۳-۳
۳۵ ویژگی‌های مهم متاهیوریستیک ACO ۱-۳-۳
۳۷ الگوریتم ABC ۲-۳-۳
۳۸ الگوریتم AntNet ۳-۳-۳
۳۹ الگوریتم‌های تکاملی و مسیریابی ۴-۳
۴۰ عناصر مهم الگوریتم‌های تکاملی در مسیریابی ۱-۴-۳
۴۱ الگوریتم GARA ۲-۴-۳
۴۲ الگوریتم DGA ۳-۴-۳
۴۴ الگوریتم اجتماع زنبورها و مسیریابی ۵-۳
۴۴ الگوریتم BeeHive ۱-۵-۳
۴۸ الگوریتم‌های یادگیری تقویتی و مسیریابی ۶-۳
۴۸ الگوریتم Q-routing ۱-۶-۳
۴۹ الگوریتم PQ-routing ۲-۶-۳
۴۹ نتیجه‌گیری ۷-۳

فصل چهارم: روش پیشنهادی

۵۱ مقدمه ۱-۴
۵۲ مدل سیستم ۲-۴
۵۲ فضای حالت ۳-۴
۵۵ تعیین وضعیت هر گره ۱-۲-۴
۵۷ پیاده‌سازی دستگاه استنتاج فازی ۲-۲-۴
۵۹ یادگیری عامل ۳-۴
۶۰ انتخاب اعمال ۴-۴
۶۱ مسئله اکتشاف-بهره‌برداری ۱-۴-۴

فصل پنجم: شبیه‌سازی و نتایج حاصل از آن

۶۴ ۱-۵ مقدمه
۶۶ ۲-۵ مدل بسته‌ها
۶۶ ۱-۲-۵ بسته‌های داده
۶۷ ۲-۲-۵ بسته‌های پاداش
۶۸ ۳-۵ مدل گره‌های منبع
۶۹ ۱-۳-۵ مدل واحد پردازشگر در گره منبع
۶۹ ۴-۵ مدل گره‌های مقصد
۷۱ ۵-۵ مدل گره‌های مسیر یاب
۷۲ ۱-۵-۵ مدل واحد پردازشگر در واحد مسیر یاب
۷۳ ۶-۵ نتایج شبیه‌سازی
۷۴ ۱-۶-۵ وجود شرایط یکسان در شبکه
۷۸ ۲-۶-۵ وجود شرایط نابرابر در شبکه
۸۲ ۳-۶-۵ به کارگیری ϵ با مقدار متغیر
۸۴ ۴-۶-۵ قطع یک اتصال
۸۷ ۷-۵ نتیجه‌گیری

فصل ششم: نتیجه‌گیری و پیشنهادات

۸۸ ۱-۶ نتیجه‌گیری
۹۰ ۲-۶ پیشنهادات
۹۲ مراجع

چکیده

شبکه‌های کامپیوتری نمونه مهمی از سیستم‌های پویای توزیع شده هستند که در زندگی روزانه نقش زیادی پیدا نموده‌اند. اهمیت استراتژیک و محدودیت‌های ذاتی این گونه سیستم‌ها منجر به نیاز به کنترل توزیع شده، به خصوص برای مسیریابی، می‌شود تا رفتار شبکه را قابل تطبیق با تغییرات در توپولوژی، ترافیک داده‌ها، سرویس‌ها و غیره نماید. اخیراً، محققین بر روی الگوریتم‌های مسیریابی جدیدتری به منظور فراهم آوردن میزان انطباق‌پذیری بیشتر با تغییر حالات شبکه تحقیق کرده‌اند که این الگوریتمها بر پایه روش‌های یادگیری ماشین بنا شده‌اند. یادگیری تقویتی یک روش یادگیری بدون سرپرست است که هدف از آن یادگیری یک سیاست، نگاشتی از مشاهدات به اعمال، بر مبنای بازخوردی است که از محیط دریافت می‌شود. این عمل یادگیری را می‌توان به صورت جستجوی مجموعه‌ای از سیاست‌ها نگریسته که در هر لحظه در تعامل با محیط ارزیابی می‌شوند. در میان الگوریتم‌های مختلفی که برای یادگیری تقویتی تاکنون ارائه شده است الگوریتم یادگیری Q -دارای بیشترین کاربرد است. در این پایان‌نامه به شبکه به عنوان یک سیستم چندعاملی نگریسته شده است که هر گره آن نشان دهنده یک عامل می‌باشد. سپس بر پایه روش یادگیری تقویتی و با استفاده از الگوریتم یادگیری Q -هر یک از عامل‌ها شروع به یادگیری وضعیت شبکه می‌نمایند تا در هر حالت شبکه بهترین عمل ممکن را از خود بروز دهد. در این روش وضعیت هر گره (عامل) به صورت تابعی از وضعیت گره‌های مجاور و خطوط ارتباطی مابین آن‌ها تعریف می‌شود و بدین صورت هر تغییر در وضعیت یک گره (عامل) در وضعیت و نحوه تصمیم‌گیری گره‌های (عاملین) مجاور آن تأثیرگذار است.

فصل اول

مقدمه

۱-۱ مقدمه

شبکه‌های کامپیوتری نمونه مهمی از سیستم‌های پویای توزیع شده هستند که در زندگی روزانه ما نقش زیادی پیدا کرده‌اند. اهمیت استراتژیک و محدودیت‌های ذاتی این گونه سیستم‌ها منجر به نیاز به کنترل توزیع شده، به خصوص برای مسیریابی، می‌شود تا رفتار شبکه را قابل تطبیق با تغییرات در توپولوژی، ترافیک داده‌ها، سرویس‌ها و غیره نماید. بنابراین، انجمن‌های کنترل و یادگیری ماشین همیشه علاقه‌مند به زمینه ارتباطات کامپیوتری بوده‌اند.

در دهه ۱۹۵۰ میلادی، Ford و Bellman روش برنامه‌ریزی پویا^۱ را برای مسئله بهینه‌سازی مسیریابی در شبکه‌ها به کار بردند [۱,۲]. با اینکه الگوریتم مسیریابی Bellman-Ford یک کنترل توزیع شده را پیاده‌سازی می‌کند، با این حال میزان تطبیق پذیری محدودی را ارائه می‌نماید. بنابراین کارآیی آن در حالاتی که وضعیت شبکه با سرعت زیادی تغییر می‌کند کاهش پیدا می‌کند. در طول سالیان متمادی، افراد در جستجوی نسخه‌های تطبیق پذیرتری از الگوریتم اصلی Bellman-Ford بوده‌اند که به هر حال در این راه با مشکلات مختلفی نیز روبرو شدند. اخیراً، محققین بر روی الگوریتم‌های مسیریابی جدیدتری به منظور فراهم آوردن میزان انطباق پذیری بیشتر تحقیق کرده‌اند که این الگوریتمها بر پایه روش‌های یادگیری ماشین [۳] بنا شده‌اند. در این رابطه Nedzelinsk و Narendra راهکاری بر پایه اتوماتای یادگیر تصادفی^۲ تولید کردند [۴]. در سال ۱۹۹۴، Boyan و Littman الگوریتم Q-routing [۵] را پیشنهاد کردند که اقتباسی از الگوریتم Bellman-Ford بود که از الگوریتم Q-Learning [۶]، که در مبحث یادگیری تقویتی^۳ [۷, ۸] مطرح می‌شود، استفاده

^۱ Dynamic Programming

^۲ Stochastic Learning Automata

^۳ Reinforcement Learning

می نمود. در سال ۱۹۹۸، Di Caro و Dorigo الگوریتم AntNet [۹] را ارائه نمودند که از متاهوریستیک بهینه‌سازی اجتماع مورچه (ACO) مشتق می‌شد و یک سیستم نمونه‌برداری توزیع شده مونت کارلو را به منظور یادگیری تصمیم‌های مسیریابی پیاده‌سازی می‌نمود.

علیرغم کارآیی بسیار خوبی که این الگوریتم‌ها و الگوریتم‌های تطبیق‌پذیر مشابه دیگر داشته‌اند، با این حال فن‌آوری‌های شبکه کنونی همچنان از الگوریتم‌های ایستای قبلی استفاده می‌کنند. الگوریتم‌های مسیریابی اینترنت همانند RIP و BGP از الگوریتم پایه Bellman-Ford مشتق شده‌اند. آن‌ها این قابلیت را دارند که با تغییرات توپولوژیکی نادری که ممکن است رخ دهد (همانند از کار افتادگی شبکه) کنار بیایند، اما در زمینه تطبیق‌پذیری ترافیک امکانی را فراهم نمی‌آورند. استراتژی اصلی برای مقابله با نوسانات ترافیکی و فراهم آوردن یک کیفیت سرویس گارانتی شده، محدود ساختن و مقید کردن بیش از اندازه منابع شبکه است که موجب می‌شود تا در عمل دیگر چندان ضرورتی برای وجود الگوریتم‌های مسیریابی کاملاً انطباق‌پذیر وجود نداشته باشد.

با این حال وضع موجود اکنون به سرعت در حال تغییر است. پیشرفت‌های حاصل در تکنولوژی‌های بی‌سیم، همانند WiFi و Bluetooth، درجه آزادی بیشتری را در هنگام ایجاد و یا تغییر یک شبکه داده‌ای فراهم آورده‌اند. همچنین معرفی مدل‌های ارتباطی و سرویس‌های کاربری جدید، همانند شبکه‌سازی Peer-to-Peer و Voice-over-IP، منجر به نیازهای جدید و در حال تغییری در زمینه ترافیک داده‌ها شده است. شبکه‌ها به سرعت در حال پویا و ناهمگون شدن هستند و از آنجاییکه این شبکه‌های جدید بیشتر بر مبنای کاربران هستند، یعنی ویژگی‌های آن‌ها توسط کاربران تعیین می‌شود و نه توسط یک قدرت مرکزی، لذا محدودسازی‌هایی که قبلاً ذکر شد نمی‌توانند چندان گزینه موثری باشند. انتظار می‌رود که این تحول و تکامل سرعت بیشتری به خود بگیرد و بنابراین نیاز به الگوریتم‌های کنترلی پویا و کارآمد در حال افزایش است. این الگوریتم‌ها می‌بایست وضعیت فعلی شبکه و کاربران را فرا گرفته و سیاست‌های تصمیم‌گیری خود را مطابق با آن سازند و حتی بتوانند پارامترهای درونی خود را تنظیم نمایند. حرکت گره‌های شبکه و تغییرات در الگوهای ترافیک داده‌ها که به خاطر پدید آمدن سرویس‌های جدید رخ می‌دهد، منجر به حالت‌های مختلف شبکه شده است که با خصوصیات همانند پهنای باند، اتصال و غیره تعریف می‌شوند. حالت شبکه ممکن است در طول زمان تغییر کند و یا حالات مختلفی همزمان در یک شبکه ناهمگون وجود داشته باشند. همچنین محدودیت‌ها و قیودی که توسط فن‌آوری‌های شبکه اعمال می‌شوند، پیچیدگی بیشتری را اضافه می‌کنند.

افزایش سرعت سخت‌افزار و دسترس بودن آن، اجازه می‌دهد تا کنترل‌کننده‌هایی را ایجاد کنیم که قابلیت توانایی یادگیری داشته باشند. این کنترل‌کننده‌ها می‌توانند این عمل یادگیری را با استفاده از روش‌های

⁴ Ant Colony Optimization

هوش مصنوعی همانند یادگیری تقویتی، شبکه‌های عصبی، الگوریتم‌های ژنتیک و غیره انجام دهند. این کنترل کننده‌های مبتنی بر هوش مصنوعی این توانایی را دارند که تعدادی از توابع و اعمال کنترلی موجود در شبکه‌های کامپیوتری و مخابراتی را پیاده‌سازی نمایند. این روش‌های جدید وفقی می‌بایست منجر به قابلیت انطباق بالا شوند که امکان برآورده‌سازی نیازهای سرویس، چه از منظر کاربران و چه از منظر متصدیان شبکه، را داشته باشند.

یادگیری تقویتی یک روش یادگیری بدون سرپرست است که هدف از آن یادگیری یک سیاست، نگاشتی از مشاهدات به اعمال، بر مبنای بازخوردی است که از محیط دریافت می‌شود. این عمل یادگیری را می‌توان به صورت جستجوی مجموعه‌ای از سیاست‌ها نگریست که در هر لحظه در تعامل با محیط ارزیابی می‌شوند. در میان الگوریتم‌های مختلفی که برای یادگیری تقویتی تاکنون ارائه شده است الگوریتم یادگیری Q -دارای بیشترین کاربرد است. با به کارگیری الگوریتم یادگیری Q می‌توان یک راهکار مسیریابی خودمختار به صورت وفقی بدست آورد [۵]. با این حال یادگیری Q یک روش متمرکز است که برای یادگیری یک سیستم تک عامل به کار می‌رود. در این پایان‌نامه به شبکه به عنوان یک سیستم چندعاملی^۵ نگریسته شده است که هر گره آن نشان دهنده یک عامل می‌باشد. سپس بر پایه روش یادگیری تقویتی و با استفاده از الگوریتم یادگیری Q هر یک از عامل‌ها شروع به یادگیری وضعیت شبکه می‌نماید تا در هر حالت شبکه بهترین عمل ممکن را از خود بروز دهد. در این روش وضعیت هر گره (عامل) به صورت تابعی از وضعیت گره‌های مجاور و خطوط ارتباطی مابین آن‌ها تعریف می‌شود و بدین صورت هر تغییر در وضعیت یک گره (عامل) در وضعیت و نحوه تصمیم‌گیری گره‌های (عاملین) مجاور آن تأثیرگذار است.

۱-۲ روند ارائه پایان‌نامه

در ادامه مطالب پایان‌نامه و در فصل دوم به معرفی روش یادگیری تقویتی می‌پردازیم. در این فصل مسئله یادگیری تقویتی و مدل آن معرفی می‌شود. سپس راهبردهای موجود برای این مسئله بیان می‌گردند و مزایا و مشکلات که این روش دارد بررسی می‌شود. همچنین روش‌های مختلف جستجو و اکتشاف که از آن‌ها در یک سیستم یادگیر تقویتی می‌توان بهره برد ارائه می‌شوند. در نهایت از بین الگوریتم‌های مورد استفاده در مسائل یادگیری تقویتی، از الگوریتم یادگیری Q یاد می‌کنیم و به بیان دقیق عملکرد آن و نحوه یادگیری توسط این الگوریتم می‌پردازیم.

در فصل سوم مسئله مسیریابی در شبکه را مطرح می‌کنیم و ویژگی‌های مختلفی که انواع الگوریتم‌های مسیریابی را از هم متمایز می‌نماید معرفی می‌کنیم. سپس به معرفی الگوریتم‌های ارائه شده برای مسیریابی که

⁵ Multi Agent

از روش‌های هوش مصنوعی در آن‌ها استفاده شده است می‌پردازیم و ویژگی‌ها و خصوصیات هر یک را به تفکیک بیان می‌کنیم.

در فصل چهارم به معرفی روش ارائه شده می‌پردازیم. در این روش شبکه را به عنوان یک سیستم چند عاملی مورد بررسی قرار داده‌ایم که در آن هر گره معادل یک عامل می‌باشد. هر عامل توسط الگوریتم یادگیری- Q آموزش می‌بیند تا رفتار شبکه را یاد گرفته و بتواند یک رفتار هوشمندانه مطابق با وضعیت شبکه از خود بروز دهد. در این فصل همچنین به چگونگی شبیه‌سازی این سیستم پرداخته می‌شود و نتایج حاصل از آن ارائه می‌گردد.

در نهایت در فصل آخر نتیجه‌گیری و پیشنهاداتی که برای ادامه کار مناسب به نظر می‌رسند ارائه می‌شوند.

فصل دوم

یادگیری تقویتی

۲-۱ مقدمه

یادگیری تحت سرپرستی^۱ یک روش عمومی در یادگیری ماشین است که در آن به یک سیستم مجموعه جفت‌های ورودی - خروجی ارائه شده و سیستم تلاش می‌کند تا تابعی از ورودی به خروجی را فرا گیرد. پس از یک دوره یادگیری، می‌توان داده ورودی را به سیستم ارائه نمود و سیستم تلاش می‌نماید تا خروجی متناسب با آن را تولید کند.

یادگیری تحت سرپرستی نیازمند تعدادی داده ورودی به منظور آموزش^۲ سیستم است. با این حال رده‌ای از مسائل وجود دارند که خروجی مناسب که یک سیستم یادگیری تحت سرپرستی نیازمند آن است، برای آن‌ها موجود نیست. برای نمونه، در مسائل کنترل پویا (برای مثال کنترل ترافیک خطوط هوایی) ممکن است که پاسخ‌های صحیح زیادی وجود داشته باشد، اما این پاسخ‌ها به راحتی قابل دسترسی و بدست آوردن نباشند. این نوع از مسائل چندان قابل جوابگویی با استفاده از یادگیری تحت سرپرستی نیستند. در یادگیری تقویتی^۳ (RL)، سیستم تلاش می‌کند تا تقابلات خود با یک محیط پویا را از طریق خطا و آزمایش بهینه نماید. یادگیری تقویتی مدلی برای مسائلی از این قبیل فراهم می‌آورد.

سابقه یادگیری تقویتی (RL) به روزهای اولیه سایبرنتیک^۴ و کار در زمینه‌های آمار، روانشناسی، علوم عصبی^۵ و علوم کامپیوتر برمی‌گردد. در ۱۰ تا ۱۵ سال گذشته، این زمینه توجه بسیاری را در انجمن‌های هوش مصنوعی و یادگیری ماشین^۶ با سرعت زیاد به خود جلب نموده است [۷]. هدف اولیه یادگیری تقویتی

^۱ Supervised Learning

^۴ Cybernetic

^۲ Train

^۵ Neuroscience

^۳ Reinforcement Learning

^۶ Machine Learning

برنامه‌ریزی عامل‌ها^۷ با استفاده از تنبیه و تشویق است بدون آنکه ذکر از چگونگی انجام وظیفه آن‌ها شود. با این حال موانع محاسباتی دشواری برای رسیدن به این هدف وجود دارد. یادگیری تقویتی از بسیاری از جهات با مسئله یادگیری تحت سرپرستی متفاوت است. مهم‌ترین تفاوت، آن است که در یادگیری تقویتی هیچ نوع زوج ورودی/خروجی ارائه نمی‌شود. به جای آن، پس از اتخاذ یک عمل، حالت بعدی و پاداش بلافاصله به عامل ارائه می‌شود. برای عامل لازم است که تجربیات مفیدی در مورد حالات ممکن سیستم، اعمال، انتقال‌ها و پاداش‌ها جمع‌آوری نماید تا بتواند به بهترین حالت عمل نماید. تفاوت دیگر با یادگیری تحت سرپرستی آن است که عملکرد همزمان در اینجا بسیار مهم است: ارزیابی سیستم معمولاً همزمان با یادگیری صورت می‌گیرد.

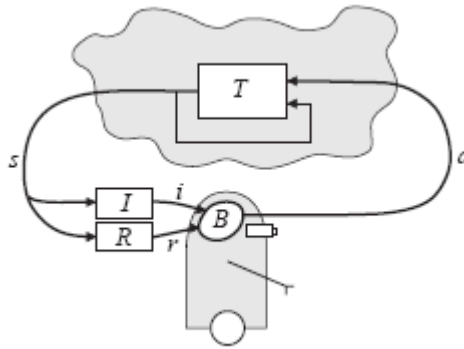
یادگیری تقویتی مسئله‌ای است که یک عامل که می‌بایست رفتار خود را از طریق تعاملات آزمایش و خطا با یک محیط پویا فرا گیرد، با آن مواجه است. دو راهکار اصلی برای حل مسائل یادگیری تقویتی وجود دارد. روش اول آن است که در فضای رفتارها به منظور یافتن بهترین رفتار در محیط مورد نظر جستجو نماییم. این راهکار در الگوریتم‌های ژنتیک و برنامه‌ریزی ژنتیک به کار گرفته شده است [۱۰]. روش دوم به کارگیری تکنیک‌های آماری و روش‌های برنامه‌ریزی پویا^۸ به منظور تخمین منفعت اتخاذ هر عمل در حالات مختلف محیط یادگیری است. این فصل به معرفی تکنیک‌های روش دوم می‌پردازد. چرا که این روش‌ها از مزایای ساختار ویژه مسائل یادگیری تقویتی استفاده می‌کنند که در سایر مسائل بهینه‌سازی به طور کل موجود نیست.

۲-۲ مدل یادگیری تقویتی

در مدل استاندارد یادگیری تقویتی، یک عامل از طریق ادراکات و اعمال خود با محیط در ارتباط است (شکل ۲-۱). در هر مرحله عامل به عنوان ورودی، i ، نشانه‌ای از حالت فعلی محیط، s ، را دریافت می‌کند. سپس عامل عمل خود، a ، را انتخاب نموده و به عنوان خروجی تولید می‌کند. این عمل حالت محیط را تغییر داده و مقدار این تغییر حالت از طریق اسکالر سیگنال تقویت، r ، به عامل منعکس می‌گردد. رفتار عامل (نحوه انتخاب یک عمل از میان اعمال ممکن)، B ، باید به گونه‌ای اعمال خود را انتخاب نماید که منجر به افزایش جمع مقادیر سیگنال تقویتی دریافتی در طول زمان طولانی اجرا شود. عامل می‌تواند این رفتار را در طول زمان از طریق آزمایش و خطای سیستماتیک که با الگوریتم‌های مختلفی قابل هدایت است، فرا گیرد.

⁷ Agent

⁸ Dynamic Programming



شکل 1-2: مدل استاندارد یادگیری تقویتی [7]

مدل تشکیل شده است از :

- یک مجموعه گسسته از حالت‌های محیط، S .
 - یک مجموعه گسسته از اعمال عامل، A .
 - یک مجموعه از سیگنال‌های تقویت اسکالر؛ به طور معمول $\{0, 1\}$ ، و یا اعداد حقیقی.
- دو تابع R و T به ترتیب تعیین کننده پاداش دریافتی عامل از انجام عمل a در حالت s و حالت بعدی که عامل در اثر این زوج حالت-عمل به آن منتقل می‌شود، هستند. شکل ۱ همچنین شامل یک تابع ورودی I نیز می‌شود که تعیین می‌کند که عامل چگونه حالات محیط را ببیند. فرض می‌شود که I یک تابع همانی باشد، بدین معنی که عامل دقیقاً چیزی را درک می‌کند که حالت محیط است.

یک روش شهودی برای فهمیدن ارتباط ما بین عامل و محیط آن در مثال زیر آمده است:

محیط: در حالت ۶۵ هستید. ۴ عمل ممکن دارید.

عامل: عمل ۲ را انجام می‌دهم.

محیط: شما ۷ واحد پاداش دریافت می‌کنید. اکنون در حالت ۱۵ هستید. ۲ عمل ممکن دارید.

عامل: عمل ۱ را انجام می‌دهم.

محیط: شما ۴- واحد پاداش دریافت می‌کنید. اکنون در حالت ۱۵ هستید. ۲ عمل ممکن دارید.

عامل: عمل ۲ را انجام می‌دهم.

محیط: شما ۵ واحد پاداش دریافت می‌کنید. اکنون در حالت ۱۵ هستید. ۲ عمل ممکن دارید.

⋮

این مجموعه تراکنش‌های مابین عامل و محیط تا رسیدن عامل به حالت هدف و دستیابی به پاداش مناسب ادامه می‌یابد.

وظیفه عامل پیدا نمودن یک سیاست π است که حالات را به گونه‌ای به اعمال نگاشت می‌کند که یک مقیاس تقویت طولانی مدت را بیشینه نماید. به طور کلی انتظار می‌رود که محیط غیر قطعی^۹ باشد، بدین

^۹ Non Deterministic

معنی که اتخاذ یک عمل در یک حالت مشابه در دو موقعیت متفاوت ممکن است منجر به حالت بعدی و مقدار تقویت متفاوتی شود. این حالت در مثال فوق نیز اتفاق افتاده است: در حالت ۶۵ به کارگیری عمل ۲، در دو موقعیت متفاوت سیگنال تقویت و حالت بعدی متفاوتی را ایجاد کرده است. با این حال فرض می‌شود که محیط پایدار^{۱۰} باشد، بدین معنی که احتمالات انتقال حالت و یا دریافت یک سیگنال تقویت خاص در طول زمان تغییر نمی‌کند.

مسائل یادگیری تقویتی معمولاً به صورت فرآیندهای تصمیم‌گیری مارکف (MDP)^{۱۱} مدل می‌شوند. یک MDP تشکیل شده است از:

- یک مجموعه از حالات، S
- یک مجموعه از اعمال، A
- یک تابع پاداش $R: S \times A \rightarrow \mathcal{R}$. $r = R(s, a)$ پاداش بلافاصله مورد انتظار از انجام عمل a در حالت s است.
- یک تابع توزیع انتقال حالت: $T: S \times A \rightarrow \Pi(S)$ که در آن $\Pi(S)$ مجموعه‌ای از توزیع‌های تصادفی بر روی مجموعه S است. $T(s, a, s')$ به عنوان احتمال انتقال از حالت s به حالت s' با اتخاذ عمل a در نظر گرفته می‌شود.

سیستم ممکن است شامل یک سری حالت نهایی^{۱۲} باشد. یک حالت نهایی F در یک MDP به صورت حالتی بیان می‌شود که در آن هر عمل انتقال به حالت F دارای پاداش ۰ است. یک MDP جذب‌کننده^{۱۳} گونه‌ای است که در آن از هر حالت غیرنهایی بتوان سرانجام به یک حالت نهایی رسید. یک مدل مارکف است اگر انتقال حالات تنها به حالت فعلی عامل بستگی داشته و مستقل از حالات و اعمال قبلی وی باشد.

۳-۲ تعریف سیاست بهینه

پاداش بلافاصله یک عمل نتایج آتی آن عمل را در نظر نمی‌گیرد، اگرچه که ممکن است آن عمل اثر بزرگی بر روی عملکرد کلی سیستم داشته باشد. مسئله ارزیابی اعمال از طریق پاداش تأخیر یافته مربوط به مسئله تخصیص اعتبار زمانی^{۱۴} می‌شود. سیاست بهینه نیازمند آن است که علاوه بر نتیجه بلافاصله عمل، نتایج آتی آن نیز در نظر گرفته شود.

¹⁰ Stationary

¹³ Absorbing MDP

¹¹ Markov Decision Process

¹⁴ Temporal Credit Assignment Problem

¹² Terminal State

۲-۳-۱ مدل پاداش دراز مدت

لازم است که رفتار بهینه هر عامل تعریف شود. تعریف سیاست بهینه از ارزیابی عملکرد درازمدت عامل تشکیل می‌شود. تعدادی مدل استاندارد وجود دارد که برای تعریف این عملکرد به کار برده می‌شوند. این مدل‌ها تعیین می‌کنند که اثرات آتی عمل تا چه میزان باید در نظر گرفته شود، و آیا پاداش‌هایی که زودتر دریافت می‌شوند می‌بایست از ارزش بیشتری نسبت به پاداش‌هایی که دیرتر دریافت می‌شوند برخوردار باشند.

مدل *افق محدود*^{۱۵} به سادگی، میزان عملکرد یک عامل را به صورت جمع پاداش مورد انتظار در h مرحله بعد ارزش دهی می‌کند:

$$E\left(\sum_{t=0}^h r_t\right) \quad (1-2)$$

که در آن r_t بیانگر پاداش دریافت شده در t مرحله بعد در آینده است. این مدل پاداشی را که در $h+1$ مرحله بعد و بعد از آن دریافت می‌شود در نظر نمی‌گیرد.

مدل *افق نامحدود کاهش یافته*^{۱۶} پاداش‌های دراز مدت یک عامل را در نظر می‌گیرد، اما این پاداش‌هایی که در آینده دریافت می‌شوند طبق یک فاکتور کاهش $0 < \gamma \leq 1$ به صورت هندسی کاهش می‌یابند:

$$E\left(\sum_{t=0}^{\infty} \gamma^t r_t\right) \quad (2-2)$$

این مدل به یک حاصل جمع همگرا شونده با توجه به پاداش نامحدود در آینده منجر می‌شود، در عین این که به پاداش‌های دریافت شده در آینده نزدیک تاکید بیشتری دارد.

۲-۳-۲ تابع ارزش بهینه و سیاست بهینه

در یک مدل با عملکرد بهینه‌ی از پیش تعریف شده، منظور از *مقدار بهینه*^{۱۷} یک حالت، عملکرد مورد انتظار یک عامل است در صورتی که از آن حالت شروع کرده و یک سیاست بهینه را اجرا کند. به عنوان مثال، با استفاده از مدل کاهش یافته و نشان دادن یک سیاست تصمیم‌گیری کامل با π ، تعریف می‌کنیم:

$$V^*(s) = \max_{\pi} E\left(\sum_{t=0}^{\infty} \gamma^t r_t\right) \quad (3-2)$$

$V^*(s)$ مقدار بهینه s است. این تابع یکتا بوده و راه‌حل معادلات *Bellman* است:

$$V^*(s) = \max_a \left(R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V^*(s') \right) \quad (4-2)$$

که در آن $T(s, a, s')$ و $R(s, a)$ ، مدل سیستم هستند.

¹⁵ Finite Horizon¹⁶ Infinite Horizon Discounted¹⁷ Optimal Value

با توجه به رابطه (۲-۴)، سیاست بهینه $\pi^*(s)$ به صورت زیر می تواند بیان شود:

$$\pi^*(s) = \arg \max_a \left(R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V^*(s') \right) \quad (۲-۵)$$

۳-۳-۲ الگوریتم های تکرار مقدار و تکرار سیاست

برای یافتن تابع مقدار بهینه و پس از آن سیاست بهینه، می توان از یک الگوریتم ساده تکراری که تکرار مقدار^{۱۸} نام دارد استفاده نمود. در [۱۱، ۱۲] نشان داده شده است که الگوریتم تکرار مقدار به مقادیر صحیح $V^*(s)$ همگرا خواهد شد. الگوریتم تکرار مقدار در جدول (۲-۱) نشان داده شده است.

جدول ۱-۲: الگوریتم تکرار مقدار

Value Iteration algorithm

initialize $V(s) \forall s \in S$ arbitrary

repeat

for all $s \in S$ **do**

for all $a \in A$ **do**

$$Q(s, a) \leftarrow R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V(s')$$

end for

$$V(s) \leftarrow \max_a Q(s, a)$$

end for

until policy is good enough

در برابر تکرار مقدار، که یک سیاست بهینه را به صورت غیر مستقیم با استفاده از تابع مقدار بهینه پیدا می کند، الگوریتم تکرار سیاست^{۱۹} وجود دارد که به صورت مستقیم بر روی یک سیاست کار می کند. Puterman در [۱۳] نشان داده است که الگوریتم تکرار سیاست در حداکثر تعداد نمایی تکرار خاتمه می یابد. الگوریتم تکرار سیاست در جدول (۲-۲) آورده شده است.

¹⁸ Value Iteration

¹⁹ Policy Iteration

Policy Iteration algorithm

choose an arbitrary policy

repeat**for all** $s \in S$ **do**

$$V(s) \leftarrow R(s, \pi(s)) + \gamma \sum_{s' \in S} T(s, \pi(s), s') V(s')$$

end for**for all** $s \in S$ **do**

$$\pi'(s) \leftarrow \arg \max (R(s, a) + \gamma \sum_{s' \in S} T(s, a, s') V(s'))$$

end for

$$\pi \leftarrow \pi'$$

until no further improvement is possible

۴-۲ راهبردهای یادگیری برای مسائل یادگیری تقویتی

همانگونه که در قسمت ۲-۳-۳ نشان داده شد، سیاست بهینه یک عامل را می توان با حل معادلات Bellman بدست آورد. اگر چه حل کامل این معادلات دستیابی به سیاست بهینه را تضمین می کند، با این حال روش های تخمینی وجود دارند که می توانند رفتارهای نزدیک به بهینه را بدون نیاز به محاسبه کامل تابع ارزش بهینه بدست آورند.

از این روش های تخمین به عنوان راهبردهای یادگیری یاد می شود. این راهبردها تعیین می کنند که هنگامی که عامل در حال یادگیری است، چه تجربیاتی در سیستم نگهداری شوند و اطلاعات جمع آوری شده از این تجربیات چگونه برای هدایت یادگیری و عملکرد بهتر به کار برده شوند.

۴-۲-۱ راهبردهای مدل-مبنا و مدل-آزاد

الگوریتم های یادگیری تقویتی را می توان به دو رده کلی تقسیم نمود: آنهایی که تلاش می کنند تا یک مدل از سیستم را فرا گیرند و آنهایی که نیاز به یادگیری مدل سیستم ندارند. هر دو رده روش های بر مبنای مدل و بی نیاز از مدل قادر به پیدا نمودن سیاست بهینه هستند. روش های بی نیاز از مدل معمولاً نیاز به زمان محاسباتی کمتری در هر تکرار الگوریتم دارند، اما در این روش ها تعداد تکرار بیشتری برای رسیدن (نزدیک شدن) به سیاست بهینه لازم است.

در یک راهکار بر مبنای مدل، سیستم تلاش می کند تا احتمالات انتقال حالت و پاداش ها را فرا گیرد. احتمال رفتن به حالت s' هنگامی که در حالت s هستیم و عمل a را انجام می دهیم را با $T(s, a, s')$ نشان می دهیم. پاداش انجام عمل a و رفتن به حالت s' در هنگامی که در حالت s هستیم با $R(s, a, s')$ نشان داده می شود. به زوج توابع T و R ، مدل تخمین زده گفته می شود.