



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)  
دانشکده مهندسی کامپیوتر و فناوری اطلاعات

پایان نامه کارشناسی ارشد هوش مصنوعی

عنوان:

یادگیری تقویتی مبتنی بر نقشه خودسازمان‌ده تطبیقی با زمان

توسط:

حسام منتظری

استاد راهنما:

دکتر رضا صفابخش

بهار ۸۵

بسمه تعالی



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)

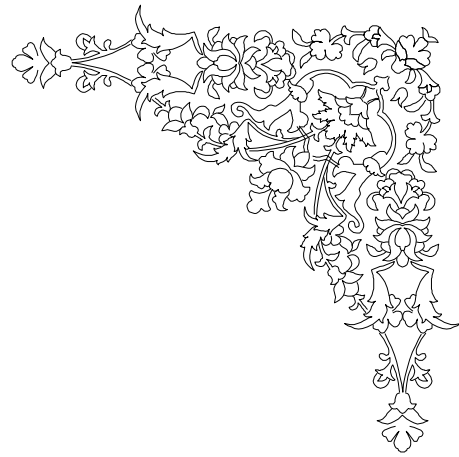
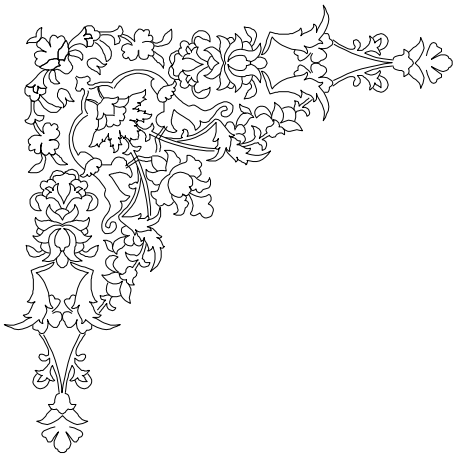
معاونت پژوهشی

فرم اطلاعات پایان نامه  
کارشناسی ارشد و دکترا

تاریخ: .....

پیوست: .....

|   |  |                                 |  |
|---|--|---------------------------------|--|
| نام و نام خانوادگی: حسام منتظری   | دانشجوی آزاد <input checked="" type="checkbox"/> | بورسیه <input type="checkbox"/> | معادل <input type="checkbox"/>                                     |
| شماره دانشجویی: ۸۳۱۳۱۲۳۱  | دانشکده: مهندسی کامپیوتر                         | رشته تحصیلی: هوش مصنوعی         |  |
| نام و نام خانوادگی استاد راهنما: دکتر رضا صفاپخش  |  |                                 |  |
| عنوان پایان نامه به فارسی: یادگیری تقویتی مبتنی بر نقشه خودسازمان‌ده تطبیقی با زمان   |  |                                 |  |
| عنوان پایان نامه به انگلیسی: Reinforcement Learning using Time-Adaptive Self Organizing Map                                       |  |                                 |  |
| نوع پروژه: کارشناسی ارشد <input checked="" type="checkbox"/><br>دکترا <input type="checkbox"/>                                    | کاربردی <input checked="" type="checkbox"/>      | بنیادی <input type="checkbox"/> | توسعه ای <input type="checkbox"/><br>نظری <input type="checkbox"/> |
| تاریخ شروع: ۸۴/۸/۱  | تاریخ خاتمه: ۸۵/۱۰/۳۰                            | تعداد واحد: ۶                   |  |
| سازمان تأمین کننده اعتبار: -  |  |                                 |  |
| واژه های کلیدی به فارسی: یادگیری تقویتی، نقشه خودسازمان‌ده رشدیابنده، نقشه خودسازمان‌ده تطبیقی با زمان، کنترل ترافیک              |  |                                 |  |
| واژه های کلیدی به انگلیسی: Reinforcement Learning, Growing SOM, TASOM, Traffic Light Control                                      |  |                                 |  |
| نظرها و پیشنهادهای به منظور بهبود فعالیت های پژوهشی دانشگاه:<br>استاد راهنما:   |  |                                 |  |
| دانشجو:   |  |                                 |  |
| امضاء استاد راهنما:   | تاریخ:   |                                 |  |
| نسخه ۱: معاونت پژوهشی<br>نسخه ۲: کتابخانه و به انضمام دو جلد پایان نامه به منظور تسویه حساب با کتابخانه و مرکز اسناد و مدارک علمی |  |                                 |  |



تقديم به

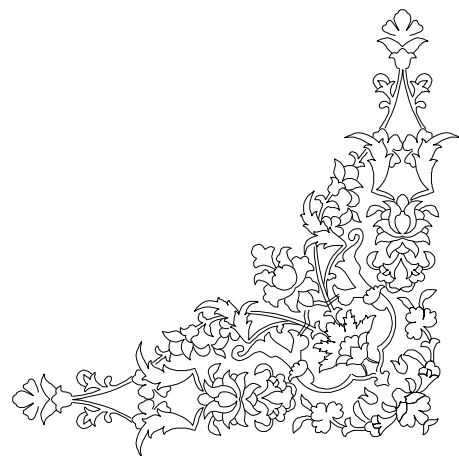
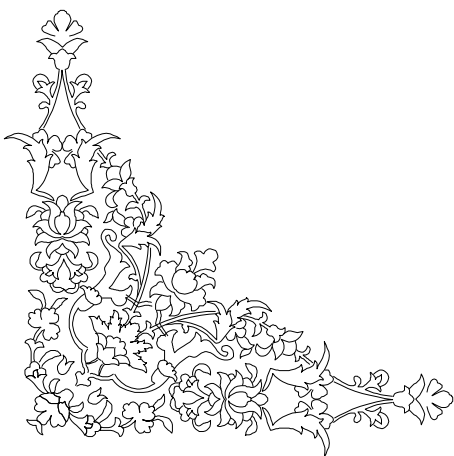
پدر و مادر مهربان

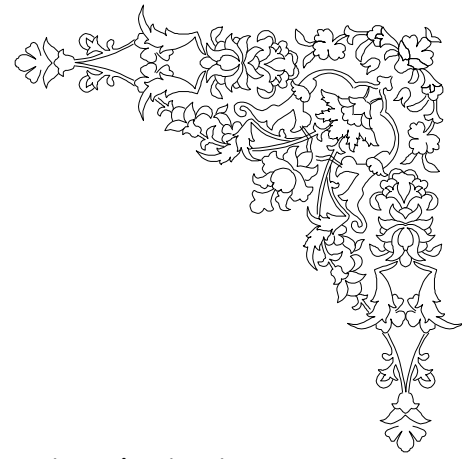
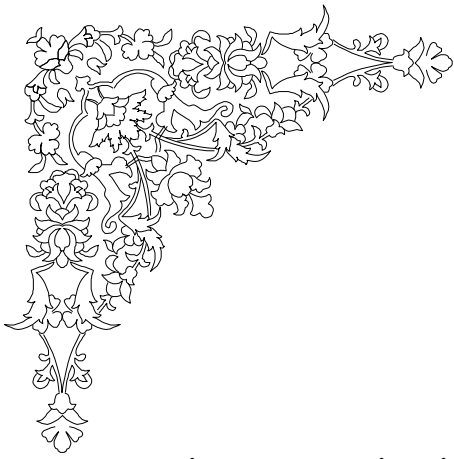
و

همسر عزيزم

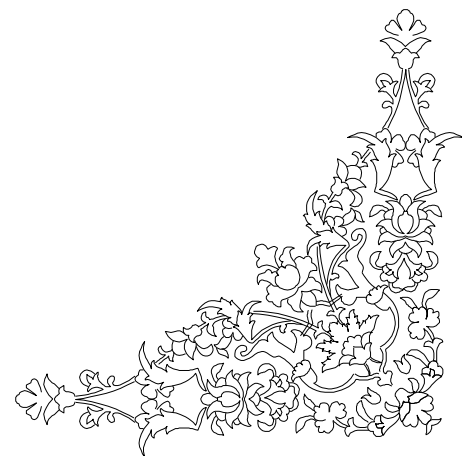
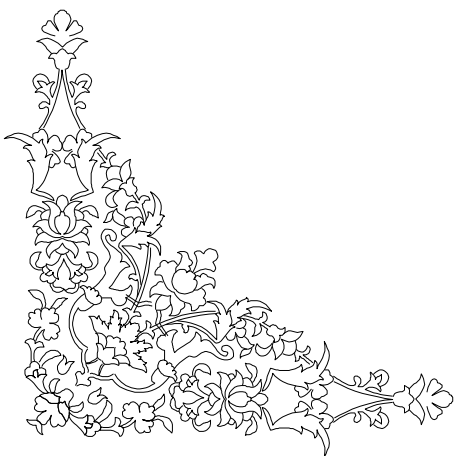
که در همه حال پشتيبان و ياور من

بوده‌اند.





از استاد راهنمای خود، جناب آقای دکتر رضا صفا بخش، به پاس  
راهنمایی‌های گرانقدر ایشان در طول مدت تحصیل در دوره کارشناسی ارشد و  
بویژه طی انجام این پژوهش، کمال تشکر را دارم.



## چکیده

یادگیری تقویتی، نگاشت وضعیت‌ها به عمل‌ها با هدف ماکزیمم کردن سیگنال پاداش دریافتی را بررسی می‌کند. در این نوع یادگیری، به عامل گفته نمی‌شود که چه عملی را انتخاب کند، بلکه عامل باید عملی را انتخاب کند که پاداش دریافتی از محیط را بیشینه کند. در چالش برانگیزترین حالات، پاداش عمل‌ها بلافاصله مشخص نمی‌شود. یادگیری تقویتی، از یک سو دارای پشتوانه قوی از قضایا و اثبات‌های ریاضی است؛ و از سوی دیگر، این روش در مسائل مختلفی همچون مسیریابی ربات، اجتناب از مانع، تصمیم‌گیری در بازی‌ها، مسائل مهارت‌ها در روبوکاپ، و کنترل ترافیک به طور موفق عمل کرده است. یکی از مسائل مهمی که در مورد این روش مطرح می‌شود، بسط و توسعه روش به مسائلی با فضای وضعیت پیوسته است. برای حل مسائل با فضای وضعیت پیوسته، روش‌های مختلفی مانند شبکه‌های عصبی پرسپترون چند لایه، کیمک، درخت‌های تصمیم، و نقشه‌های خودسازمان‌ده ارائه شده است. نشان داده شده است که یادگیری تقویتی با استفاده از نقشه‌های خودسازمان‌ده استاندارد در حل بسیاری از مسائل با فضای وضعیت پیوسته و حتی فضای عمل پیوسته موفق بوده‌اند. اما نقشه خودسازمان‌ده استاندارد نمی‌تواند یک تابع هدف متغیر را به خوبی ارائه کند و برای توابع هدفی که با توپولوژی نقشه همخوانی ندارد مورد استفاده واقع نمی‌شود. در این پایان‌نامه، یادگیری تقویتی مبتنی بر نقشه خودسازمان‌ده تطبیقی برای حل مشکل تابع هدف متغیر ارائه شده است. تابع هدف متغیر در یادگیری تقویتی منحصر به داده‌های فضای ورودی نیست، بلکه داده‌هایی که به عنوان ورودی نقشه خروجی داده می‌شود همیشه توزیع چگالی متغیر با زمان دارد. باید توجه داشت در یادگیری تقویتی عامل با گذشت زمان عملکرد خود را بهبود می‌دهد، در نتیجه داده‌های ورودی به نقشه خروجی با گذشت زمان تغییر می‌کنند و توزیع چگالی آن ناپایدار است. روش دیگری که در این پایان‌نامه ارائه شده است روش یادگیری تقویتی با استفاده از نقشه خودسازمان‌ده رشدیابنده است. این روش برای حل هر دو مشکل ذکر شده ارائه شده است. ترکیب یادگیری تقویتی با نقشه‌های خودسازمان‌ده رشدیابنده به سادگی امکان‌پذیر نیست و ترکیب این نوع نقشه با الگوریتم‌هایی که جدول کیو آن در طول زمان ثابت است، میسر نمی‌باشد. در این پایان‌نامه الگوریتم جدیدی مبتنی بر نقشه‌های خودسازمان‌ده رشدیابنده ارائه شده است که جدول کیو آن در طول زمان بزرگ و کوچک می‌شود. نشان داده شده است این الگوریتم در حل مسائل مختلف از بقیه روش‌ها موفق‌تر بوده است.

## فهرست مطالب

|          |   |
|----------|---|
| ..... ۱  | ۱- مقدمه  |
| ..... ۴  | ۲- یادگیری تقویتی   |
| ..... ۴  | ۲-۱-۱ مسأله یادگیری تقویتی  |
| ..... ۴  | ۲-۱-۲ رابط عامل و محیط  |
| ..... ۵  | ۲-۱-۲ هدفها و پاداش‌ها  |
| ..... ۶  | ۲-۱-۲ بازگشت  |
| ..... ۷  | ۲-۱-۲ خصوصیت مارکف  |
| ..... ۷  | ۲-۱-۲ فرآیند تصمیم‌گیری مارکف   |
| ..... ۸  | ۲-۱-۲ توابع ارزش  |
| ..... ۱۰ | ۲-۱-۲ توابع ارزش بهینه  |
| ..... ۱۰ | ۲-۲ تفاضل زمانی   |
| ..... ۱۰ | ۲-۲-۱ پیش بینی یا ارزیابی تفاضل زمانی   |
| ..... ۱۱ | ۲-۲-۲ روش سارسا   |
| ..... ۱۳ | ۲-۲-۲ یادگیری کیو   |
| ..... ۱۳ | ۲-۲-۲ تفاضل زمانی لاندا   |
| ..... ۱۸ | ۲-۲-۲ سارسا لاندا   |
| ..... ۱۸ | ۲-۲-۲ یادگیری کیو لاندا   |
| ..... ۲۱ | ۲-۳ جمع‌بندی  |
| ..... ۲۲ | ۳- نقشه‌های خودسازمان‌ده  |
| ..... ۲۲ | ۳-۱ نقشه خودسازمان‌ده استاندارد   |
| ..... ۲۶ | ۳-۲ نقشه خودسازمان‌ده تطبیقی با زمان  |
| ..... ۲۶ | ۳-۲-۱ نقشه خودسازمان‌ده تطبیقی با نرخ‌های یادگیری جداگانه و پویا                      |
| ..... ۲۹ | ۳-۲-۲ نقشه خودسازمان‌ده تطبیقی با نرخ‌های یادگیری و مجموعه‌های همسایگی جداگانه و پویا |
| ..... ۳۲ | ۳-۲-۳ شبکه خودسازمان‌ده تطبیقی با نرخ‌های یادگیری و تابع‌های همسایگی جداگانه و پویا   |
| ..... ۳۳ | ۳-۳ نقشه‌های خودسازمان‌ده رشدیابنده   |
| ..... ۳۴ | ۳-۳-۱ شبکه گاز عصبی   |
| ..... ۳۵ | ۳-۳-۲ یادگیری هیب رقابتی  |
| ..... ۳۵ | ۳-۳-۳ شبکه گاز عصبی با یادگیری هیب رقابتی   |
| ..... ۳۶ | ۳-۳-۳ شبکه گاز عصبی رشدیابنده   |
| ..... ۳۸ | ۳-۳-۳ ساختار سلولی رشدیابنده  |
| ..... ۳۹ | ۳-۳-۳ توری بزرگ شونده   |
| ..... ۴۲ | ۳-۴ جمع‌بندی  |
| ..... ۴۲ | ۴- یادگیری تقویتی مبتنی بر نقشه‌های خودسازمان‌ده                                      |
| ..... ۴۲ | ۴-۱ روش کیو کوهونن  |
| ..... ۴۴ | ۴-۲ نقشه موتوری   |
| ..... ۴۵ | ۴-۲-۱ تابع چگالی نرمال چندبعدی  |
| ..... ۴۶ | ۴-۲-۲ یادگیری کواریانسی   |
| ..... ۴۸ | ۴-۲-۳ استفاده از یادگیری کواریانسی در پیدا کردن عمل‌ها                                |

|          |   |
|----------|---|
| ..... ۴۹ | ۴-۲-۴ استفاده از تجربه برای بهبود یادگیری                 |
| ..... ۵۰ | ۳-۴ جمع‌بندی  |
| ..... ۵۱ | ۵- مدلهای پیشنهادی  |
| ..... ۵۱ | ۱-۵ مدل تطبیقی  |
| ..... ۵۳ | ۲-۵ مدل رشدیابنده   |
| ..... ۵۴ | ۳-۵ جمع‌بندی  |
| ..... ۵۵ | ۶- نتایج تجربی در مسائلی از رباتیکز                       |
| ..... ۵۵ | ۱-۶ مسأله بازوی ربات دومفصلی - نگاشت یک به یک             |
| ..... ۶۹ | ۲-۶ مسأله بازوی ربات دومفصلی - نگاشت چند به یک            |
| ..... ۸۱ | ۳-۶ جمع‌بندی  |
| ..... ۸۲ | ۷- نتایج تجربی در کنترل چراغ راهنما                       |
| ..... ۸۲ | ۱-۷ شبیه‌ساز ترافیکی                                      |
| ..... ۸۳ | ۲-۷ آزمایش‌های انجام شده                                  |
| ..... ۸۹ | ۳-۷ جمع‌بندی  |
| ..... ۹۰ | ۸- نتیجه‌گیری و پیشنهادات                                 |
| ..... ۹۰ | ۱-۸ نتایج   |
| ..... ۹۱ | ۲-۸ پیشنهادات   |
| ..... ۹۳ | مراجع   |
| ..... ۹۶ | ضمیمه الف- نتایج روشها در وضعیت‌های ترافیکی متوسط و سنگین |

## فهرست اشکال

|          |  |
|----------|--|
| ..... ۵  | شکل ۱-۲ تعامل عامل و محیط در یادگیری تقویتی  |
| ..... ۹  | شکل ۲-۲ نمودار پشتیبان برای $V^{\pi}$  |
| ..... ۱۲ | شکل ۳-۲ تفاضل زمانی جدولی برای تخمین $V^{\pi}$   |
| ..... ۱۲ | شکل ۴-۲ نمودار پشتیبان تفاضل زمانی   |
| ..... ۱۲ | شکل ۵-۲ دنباله وضعیتها و زوجهای عمل-وضعیت در اپیزود  |
| ..... ۱۳ | شکل ۶-۲ شبه کد الگوریتم سارسا  |
| ..... ۱۳ | شکل ۷-۲ شبه کد الگوریتم یادگیری کیو  |
| ..... ۱۴ | شکل ۸-۲ پشتیبان‌های مختلف تفاضل زمانی  |
| ..... ۱۶ | شکل ۹-۲ پشتیبان ترکیبی از پشتیبانهای دو-گام و چهار-گام   |
| ..... ۱۷ | شکل ۱۰-۲ پشتیبان تفاضل زمانی لاندا   |
| ..... ۱۷ | شکل ۱۱-۲ تصویر نشان تجمعی  |
| ..... ۱۸ | شکل ۱۲-۲ شبه کد تفاضل زمانی جدولی لاندا ONLINE   |
| ..... ۱۹ | شکل ۱۳-۲ نمودار پشتیبان سارسا لاندا  |
| ..... ۱۹ | شکل ۱۴-۲ شبه کد سارسا جدولی لاندا  |
| ..... ۲۰ | شکل ۱۵-۲ نمودار پشتیبان یادگیری کیو.   |
| ..... ۲۱ | شکل ۱۶-۲ نسخه جدولی الگوریتم یادگیری کیو WATKIN'S  |
| ..... ۲۳ | شکل ۱-۳ توپولوژی نورون‌های شبکه روی لایه خروجی.  |
| ..... ۲۴ | شکل ۲-۳ نگاهت یک بردار از ابعاد بالاتر به فضایی با ابعاد کوچکتر در نقشه خودسازمان‌ده استاندارد |
| ..... ۲۴ | شکل ۳-۳ تعیین مجموعه نورونهای همسایه با توجه به نورون برنده در زمانهای T و T+1                 |
| ..... ۲۵ | شکل ۴-۳ تابع گاسی همسایگی برای $\sigma(t)=3$   |
| ..... ۲۸ | شکل ۵-۳ تابع نرمالیزه کننده $f(z) = \frac{z}{1+z}$   |
| ..... ۲۸ | شکل ۶-۳ تابع نرمالیزه کننده $f(z) = \frac{1-e^{-z}}{1+e^{-z}}$                                 |
| ..... ۴۳ | شکل ۱-۴ یادگیری کیو با استفاده از نقشه خودسازمان‌ده به روش کیو کوهونن                          |
| ..... ۴۴ | شکل ۲-۴ انتخاب بهترین عمل برای انجام در هر موقعیت در روش کیو کوهونن                            |
| ..... ۴۶ | شکل ۳-۴ پیشنهاد عمل‌های نمونه در نقشه موتوری   |
| ..... ۴۷ | شکل ۴-۴ جهت تقریبی برای بروزرسانی به سمت عمل هدف در روش نقشه موتوری                            |
| ..... ۴۸ | شکل ۵-۴ در این شکل جهت کواریانسی به درستی پیدا شده است ولی اندازه گام بزرگ است.                |
| ..... ۴۹ | شکل ۶-۴ استفاده نامناسب از تجربه در یادگیری کواریانسی  |
| ..... ۵۲ | شکل ۱-۵ مدل یادگیری تقویتی مبتنی بر نقشه خودسازمان‌ده استاندارد                                |



|          |  |
|----------|--|
| ..... ۵۵ | شکل ۶-۱ مسأله بازوی دو مفصلی ربات در حالت محیط ایستا (مسأله ۱-الف)   |
| ..... ۵۶ | شکل ۶-۲ مسأله بازوی دو مفصلی ربات در حالت محیط متغیر با زمان (مسأله ۱-ب)   |
| ..... ۵۶ | شکل ۶-۳ نمودار میانگین پاداش دریافتی بر حسب زمان روش اسمیت در مسأله ۱-الف  |
| ..... ۵۷ | شکل ۶-۴ نقشه خودسازمان‌ده استاندارد استفاده شده برای ارائه فضای ورودی سیستم یادگیر   |
| ..... ۵۸ | شکل ۶-۵ نقشه خودسازمان‌ده استاندارد استفاده شده برای ارائه فضای خروجی سیستم یادگیر   |
| ..... ۵۸ | شکل ۶-۶ نمونه‌ای از تقریب نامناسب فضای خروجی توسط نقشه خودسازمان‌ده استاندارد  |
| ..... ۵۹ | شکل ۶-۷ مکان نوک بازوی ربات در صورت اعمال زوایای نورون‌های نقشه خودسازمان‌ده شکل ۶-۶ به عنوان زوایای مفاصل ربات  |
| ..... ۵۹ | شکل ۶-۸ مقایسه نمودار میانگین پاداش دو روش با استفاده از عمل نویزی (PERTURBED) و عمل پیشنهادی (PROPOSED) یا غیر نویزی (NON-PERTURBED)  |
| ..... ۶۰ | شکل ۶-۹ نمودار میانگین پاداش دریافتی روش تطبیقی بر حسب زمان در مسأله ۱-الف   |
| ..... ۶۱ | شکل ۶-۱۰ نقشه خودسازمان‌ده تطبیقی با زمان استفاده شده برای ارائه فضای ورودی سیستم یادگیر تقویتی  |
| ..... ۶۱ | شکل ۶-۱۱ نقشه خودسازمان‌ده تطبیقی با زمان استفاده شده برای تقریب فضای خروجی عامل یادگیری تقویتی  |
| ..... ۶۲ | شکل ۶-۱۲ مکان نوک بازوی ربات در صورت استفاده از زوایای نورون‌های نقشه خودسازمان‌ده تطبیقی با زمان به عنوان زوایای مفاصل ربات   |
| ..... ۶۲ | شکل ۶-۱۳ مقایسه نمودار میانگین پاداش دریافتی روش تطبیقی در دو حالت استفاده از عمل نویزی (PERTURBED) و عمل پیشنهادی (PROPOSED) یا غیر نویزی (NON-PERTURBED)   |
| ..... ۶۴ | شکل ۶-۱۴ نمودار میانگین پاداش دریافتی روش رشدیابنده بر حسب زمان در مسأله ۱-الف   |
| ..... ۶۴ | شکل ۶-۱۵ مکان هندسی نورون‌های نقشه خودسازمان‌ده رشدیابنده  |
| ..... ۶۵ | شکل ۶-۱۶ مکان نورونها و همسایگی‌های نقشه خودسازمان‌ده رشدیابنده استفاده شده برای تقریب فضای خروجی  |
| ..... ۶۵ | شکل ۶-۱۷ مکان نوک بازوی ربات در صورت اعمال زوایای نورون‌های نقشه خودسازمان‌ده رشدیابنده شکل ۶-۱۶   |
| ..... ۶۶ | شکل ۶-۱۸ نمودارهای میانگین پاداش دریافتی بر حسب زمان برای حالتی که به‌روزرسانی جدول کیو به ازای نورون نویزی باشد (PERTURBED) و حالتی که به‌روزرسانی جدول کیو به ازای نورون غیرنویزی باشد (NON-PERTURBED) |
| ..... ۶۷ | شکل ۶-۱۹ نمودار میانگین پاداش بر حسب زمان برای حالتی که مقداردهی اولیه نورون‌ها به صورت معمولی باشد (NORMAL) و حالتی که مقداردهی اولیه نورون‌ها بر مبنای الگوریتم رشدیابنده باشد (PROPOSED)              |
| ..... ۶۷ | شکل ۶-۲۰ نمودار میانگین پاداش دریافتی روش اسمیت در مسأله ۱-ب   |
| ..... ۶۸ | شکل ۶-۲۱ نمودار میانگین پاداش دریافتی روش تطبیقی در مسأله ۱-ب  |
| ..... ۶۸ | شکل ۶-۲۲ نمودار میانگین پاداش دریافتی روش رشدیابنده در مسأله ۱-ب   |
| ..... ۶۹ | شکل ۶-۲۳ مسأله بازوی دو مفصلی ربات در حالت محیط ایستا (مسأله ۲-الف)  |
| ..... ۷۰ | شکل ۶-۲۴ مسأله بازوی دو مفصلی ربات در حالت محیط متغیر با زمان (مسأله ۲-ب)  |
| ..... ۷۰ | شکل ۶-۲۵ نمودار میانگین پاداش دریافتی روش اسمیت بر حسب زمان برای مسأله ۲-الف   |
| ..... ۷۱ | شکل ۶-۲۶ نمودار میانگین پاداش دریافتی روش اسمیت بر حسب زمان برای مسأله ۲-ب   |
| ..... ۷۲ | شکل ۶-۲۷ نقشه ورودی حاصل از اعمال روش اسمیت در مسأله ۲-الف   |

|              |   |
|--------------|---|
| .....۷۲..... | شکل ۶-۲۸ نقشه خودسازمان‌ده تطبیقی ورودی حاصل از اعمال روش تطبیقی در مسأله ۲-الف |
| .....۷۳..... | شکل ۶-۲۹ مکان نوک بازوی ربات در صورت بکارگیری روش اسمیت در مسأله ۲-الف          |
| .....۷۳..... | شکل ۶-۳۰ مکان نوک بازوی ربات در صورت بکارگیری روش اسمیت در مسأله ۲-ب            |
| .....۷۴..... | شکل ۶-۳۱ نمودار میانگین پاداش دریافتی روش تطبیقی برحسب زمان برای مسأله ۲-الف    |
| .....۷۵..... | شکل ۶-۳۲ نمودار میانگین پاداش دریافتی روش تطبیقی برحسب زمان برای مسأله ۲-ب      |
| .....۷۵..... | شکل ۶-۳۳ نقشه خودسازمان‌ده تطبیقی ورودی حاصل از اعمال روش تطبیقی در مسأله ۲-الف |
| .....۷۶..... | شکل ۶-۳۴ نقشه خودسازمان‌ده تطبیقی خروجی حاصل از اعمال روش تطبیقی در مسأله ۲-الف |
| .....۷۶..... | شکل ۶-۳۵ مکان نوک بازوی ربات در صورت بکارگیری روش تطبیقی در مسأله ۲-الف         |
| .....۷۷..... | شکل ۶-۳۶ مکان نوک بازوی ربات در صورت بکارگیری روش تطبیقی در مسأله ۲-ب           |
| .....۷۸..... | شکل ۶-۳۷ نمودار میانگین پاداش دریافتی روش رشدیابنده برحسب زمان برای مسأله ۲-الف |
| .....۷۸..... | شکل ۶-۳۸ نمودار میانگین پاداش دریافتی روش رشدیابنده برحسب زمان برای مسأله ۲-ب   |
| .....۷۹..... | شکل ۶-۳۹ نقشه ورودی حاصل از اعمال روش رشدیابنده در مسأله ۲-الف                  |

|              |   |
|--------------|---|
| .....۸۳..... | شکل ۷-۱ تصویری از شبیه‌ساز ترافیکی  |
| .....۸۵..... | شکل ۷-۲ نمودارهای حاصل از اعمال روش اسمیت در مسأله کنترل چراغ راهنما در وضعیت ترافیکی سبک     |
| .....۸۷..... | شکل ۷-۳ نمودارهای حاصل از اعمال مدل تطبیقی در مسأله کنترل چراغ راهنما در وضعیت ترافیکی سبک    |
| .....۸۸..... | شکل ۷-۴ نمودارهای حاصل از اعمال مدل رشدیابنده در مسأله کنترل چراغ راهنما در وضعیت ترافیکی سبک |

|                 |   |
|-----------------|---|
| .....الف=۱..... | شکل الف-۱ نمودارهای حاصل از اعمال روش اسمیت در مسأله کنترل چراغ راهنما برای وضعیت ترافیکی متوسط     |
| .....الف=۲..... | شکل الف-۲ نمودارهای حاصل از اعمال مدل تطبیقی در مسأله کنترل چراغ راهنما برای وضعیت ترافیکی متوسط    |
| .....الف=۳..... | شکل الف-۳ نمودارهای حاصل از اعمال مدل رشدیابنده در مسأله کنترل چراغ راهنما برای وضعیت ترافیکی متوسط |
| .....الف=۴..... | شکل الف-۴ نمودارهای حاصل از اعمال روش اسمیت در مسأله کنترل چراغ راهنما برای وضعیت ترافیکی سنگین     |
| .....الف=۵..... | شکل الف-۵ نمودارهای حاصل از اعمال مدل تطبیقی در مسأله کنترل چراغ راهنما برای وضعیت ترافیکی سنگین    |
| .....الف=۶..... | شکل الف-۶ نمودارهای حاصل از اعمال مدل رشدیابنده در مسأله کنترل چراغ راهنما برای وضعیت ترافیکی سنگین |

### فهرست جدول‌ها

|              |  |
|--------------|--|
| .....۵۷..... | جدول ۱-۶ پارامترهای مورد استفاده روش اسمیت در کنترل بازوی ربات دو مفصلی ایستا                    |
| .....۶.....  | جدول ۲-۶ پارامترهای مورد استفاده روش تطبیقی در کنترل بازوی ربات دو مفصلی ایستا                   |
| .....۶۳..... | جدول ۳-۶ پارامترهای روش رشدیابنده در کنترل بازوی ربات دو مفصلی ایستا                             |
| .....۶۹..... | جدول ۴-۶ مقایسه زمان اجرای روش‌های ارائه شده   |
| .....۷۱..... | جدول ۵-۶ پارامترهای مورد استفاده روش اسمیت در آزمایش کنترل بازوی ربات دو مفصلی متغیر با زمان     |
| .....۷۴..... | جدول ۶-۶ پارامترهای مورد استفاده روش تطبیقی در آزمایش کنترل بازوی ربات دو مفصلی متغیر با زمان    |
| .....۷۷..... | جدول ۷-۶ پارامترهای مورد استفاده روش رشدیابنده در آزمایش کنترل بازوی ربات دو مفصلی متغیر با زمان |
| .....۸۴..... | جدول ۱-۷ پارامترهای مورد استفاده روش اسمیت در مسأله کنترل چراغ راهنما در وضعیت ترافیکی سبک       |
| .....۸۶..... | جدول ۲-۷ پارامترهای مورد استفاده روش تطبیقی در مسأله کنترل چراغ راهنما در وضعیت ترافیکی سبک      |
| .....۸۶..... | جدول ۳-۷ پارامترهای مورد استفاده روش رشدیابنده در مسأله کنترل چراغ راهنما در وضعیت ترافیکی سبک   |

## فصل اول

### مقدمه

تکنیک‌های یادگیری ماشین معمولاً به سه دسته اصلی تقسیم می‌شوند: یادگیری با نظارت، یادگیری بی‌نظارت، و یادگیری تقویتی. یادگیری با نظارت به نمونه‌هایی از ورودی و خروجی تابع هدف (داده‌های آموزشی) نیازمند است. در یادگیری بی‌نظارت داده هدف وجود ندارد و فقط بر روی داده‌های ورودی پردازش صورت می‌گیرد. یادگیری تقویتی از نظر ارزیابی از سوی محیط، در حد فاصل دو روش قبلی است. یادگیری تقویتی مانند روش بانظارت خروجی صحیح را دریافت نمی‌کند، اما سیگنالی را از محیط دریافت می‌کند که صرفاً خوب بودن انجام عمل قبلی را نشان می‌دهد. اثبات شده است که بسیاری از روش‌های یادگیری تقویتی همیشه به سیاست بهینه همگرا خواهند شد [۲۸، ۳]. از جنبه کاربردی نیز روش یادگیری تقویتی در مسائل مختلفی همچون اجتناب از مانع در مسیریابی ربات‌ها [۲۴]، کنترل ترافیک [۴، ۱]، تصمیم‌گیری در بازی‌ها [۳۲]، مهارت‌های عامل‌ها در لیگ‌های ربات‌ها [۲۷]، و غیره به طور گسترده مورد استفاده قرار گرفته است.

یکی از مسائل مهمی که در مورد یادگیری تقویتی مطرح می‌شود، حل مسائل با فضای وضعیت و عمل پیوسته است. بیشتر نظریه‌هایی که در مورد یادگیری تقویتی است، مبتنی بر فرض گسسته بودن فضای وضعیت است. برای حل مسائل پیوسته یادگیری تقویتی، روش‌های مختلفی همچون شبکه‌های عصبی پرسپترون چند لایه [۱۴، ۳]، کیمک [۲۶، ۲۸، ۲]، و نقشه‌های خودسازمانده مطرح می‌شود [۲۴، ۳۱].

این پایان‌نامه استفاده از نقشه‌های خودسازمانده در یادگیری تقویتی را مورد توجه قرار داده است. کارهای مختلفی در این زمینه صورت گرفته است. Q-Kohon روش تقویتی مبتنی بر این نقشه است که توسط توژنت ارائه شده است [۱۸]. در Q-Kohon جدول وضعیت-عمل با نقشه خودسازمانده جایگزین شده است. هر نورون شامل سه گانه وضعیت، عمل، و ارزش کیو است. در این روش، انتخاب نورون برای اجرا مبتنی بر مقدار وضعیت و ارزش کیو انجام می‌شود، اما انتخاب نورون برای یادگیری بر اساس مقدار وضعیت و عمل است. روش دیگری که برای حل مسأله یادگیری تقویتی از نقشه خودسازمانده استفاده کرده است، روش نقشه موتوری (Motoric Map) است [۳۱]. در این روش، یک نقشه خودسازمانده در ورودی قرار می‌گیرد. هر نورون این نقشه شامل وزن عمل و مقدار کیو است. در نقشه موتوری تخمین مقدار کیو با استفاده از روش سارسا (sarsa) انجام می‌شود و عمل بهینه متناظر هر نورون با استفاده از روشی به نام یادگیری کواریانسی محاسبه می‌گردد.

مهمترین کاری که در این زمینه صورت گرفته است، روش اسمیت است که از یادگیری کیو به عنوان تکنیک یادگیری تقویتی استفاده می‌کند و فضای ورودی و خروجی را به طور مجزا با استفاده از نقشه‌های خودسازمانده استاندارد گسسته می‌کند [۲۴، ۲۵]. جدول کیو بر روی فضای ورودی و خروجی گسسته شده تشکیل می‌شود و مطابق با الگوریتم استاندارد یادگیری کیو (Q-Learning) بروز می‌شود. نقشه ورودی و ورودی‌هایی که از محیط می‌گیرد، به روز می‌شود. اما برای بروزرسانی نقشه خروجی داده مستقیمی وجود ندارد. اسمیت با مطرح کردن ایده نوپزی کردن عمل اولیه (نورون برنده در فضای خروجی)، مشکل کاوش در فضای عمل را حل نموده است. نشان داده شده است که عملکرد این روش در کاربردهایی از رباتیک موفق بوده است [۲۵].

کاری که در این پروژه صورت گرفته است، توسعه‌ای بر روش اسمیت است. همانگونه که گفته شد روش اسمیت از دو نقشه خودسازمانده استاندارد برای ارائه فضای ورودی و خروجی استفاده می‌کند. اما نقشه خودسازمانده استاندارد از سویی نمی‌تواند تابع هدف متغیر را به خوبی ارائه کند و از سویی قادر نیست توابع هدفی که با توپولوژی نقشه همخوانی ندارد را ارائه کند. در این پایان‌نامه، یادگیری تقویتی مبتنی بر نقشه خودسازمانده تطبیقی برای حل مشکل تابع هدف متغیر ارائه شده است. تابع هدف متغیر در یادگیری تقویتی منحصر به داده‌های فضای ورودی نیست، بلکه داده‌های ورودی به نقشه خروجی همیشه توزیع چگالی متغیر با زمان دارد. باید توجه داشت در یادگیری تقویتی عامل با گذشت زمان عملکرد خود را بهبود می‌دهد، در نتیجه داده‌های ورودی به نقشه خروجی با گذشت زمان تغییر می‌کنند و توزیع چگالی آن ناپایستا است. روش یادگیری تقویتی با استفاده از نقشه خودسازمانده رشدیابنده برای حل هر دو مشکل ذکر شده ارائه شده است. ترکیب یادگیری تقویتی با نقشه‌های خودسازمانده رشدیابنده به سادگی امکان‌پذیر نیست و ترکیب این نوع نقشه با الگوریتم اسمیت که جدول کیو آن در طول زمان ثابت است، امکان‌پذیر نیست. در این پایان‌نامه الگوریتم جدیدی مبتنی بر نقشه‌های خودسازمانده رشدیابنده ارائه شده است که جدول کیو آن در طول زمان بزرگ و کوچک می‌شود. نشان داده شده است این الگوریتم در حل مسائل مختلف از بقیه روش‌ها موفق‌تر بوده است.

نتایج حاصل از اعمال روش‌های پیشنهادی بر روی مسائل کلاسیکی از رباتیکز بررسی شده است. همه روش‌ها در مسائل ایستا عملکرد تقریباً مناسبی دارند. اما روش رشدیابنده نسبت به دیگر روش‌ها عملکرد بهتری دارد. از سویی میانگین پاداش دریافتی این روش بیشتر از دیگر روش‌ها است و از سوی دیگر روش نسبت به دیگر روش‌ها پایدارتر است. نقشه‌های ورودی و خروجی در روش رشدیابنده بهتر از روش‌های دیگر ارائه‌کننده توزیع ورودی هستند. در مسائل ناپایستا، یادگیری تقویتی با استفاده از نقشه خودسازمانده استاندارد نسبت به دیگر روش‌ها عملکرد بدتری دارد. این روش پایداری کمی در برابر تغییرات محیط دارد. اما دو روش یادگیری تقویتی با نقشه خودسازمانده تطبیقی و رشدیابنده نسبت به تغییرات محیط پایدارتر هستند و عملکرد قابل قبولی را ارائه می‌کنند.

مسئله کنترل چراغ راهنما در قلمرو ترافیک شهری کاربرد دیگری است که برای ارزیابی مدل‌های پیشنهادی استفاده شده است. مشکل ترافیک از مسائل پیچیده شهری است که به عوامل مختلفی همچون تردد بیش از حد خودرو از جاده‌ها، عدم کنترل صحیح چراغ‌های راهنما، تصادفات، عدم پیروی رانندگان و عابران پیاده از قوانین راهنمایی و رانندگی، ورودی‌های خیابان‌های فرعی به اصلی، و دوربرگردان‌ها بستگی دارد. جهت انجام آزمایش‌ها یک شبیه‌ساز ترافیکی پیاده‌سازی شد که وضعیت‌های مختلف ترافیکی همانند ترافیک سبک و سنگین را ایجاد کند. استفاده از روش‌های پیشنهادی این پایان‌نامه در حل مسئله کنترل چراغ راهنما نشان داد که از ترکیب چندین عامل یادگیری تقویتی می‌توان برای حل مسائل پیچیده استفاده نمود. در بین روش‌های بررسی شده، عملکرد مدل رشدیابنده در سه وضعیت ترافیکی سبک، متوسط، و سنگین از روش‌های دیگر بهتر است. لازم به ذکر است عملکرد مدل تطبیقی نیز در این مسئله قابل قبول است و کارایی و قابلیت اطمینان این روش در وضعیت‌های مختلف ترافیکی از روش اسمیت بهتر است.

این پایان‌نامه از هشت فصل تشکیل شده است. فصل دوم به معرفی روش یادگیری تقویتی می‌پردازد. در این فصل مسئله یادگیری تقویتی و روش‌های کلاسیک حل آن بررسی می‌شود. نقشه‌های خودسازمانده و انواع آن در فصل سوم شرح داده می‌شوند. در این فصل نقشه‌های خودسازمانده استاندارد، تطبیقی با زمان، و رشدیابنده ارائه می‌شوند. فصل چهارم پیشینه پژوهش در ترکیب یادگیری تقویتی با نقشه‌های خودسازمانده را بررسی می‌کند. مدل‌های پیشنهادی ترکیب یادگیری تقویتی با نقشه خودسازمانده تطبیقی و نقشه خودسازمانده رشدیابنده در فصل پنجم مورد بررسی واقع می‌شوند. فصل‌های ششم و هفتم نتایج تجربی مدل‌ها را بررسی می‌کند. فصل ششم دو کاربرد از رباتیکز که در کارهای قبلی مورد توجه قرار گرفته است را بررسی می‌کند. در فصل هفتم مسئله ترافیکی

معرفی و نتایج حاصل از مدل‌های پیشنهادی ارائه و مقایسه می‌شود. این پایان‌نامه در فصل هشتم با بیان نتیجه‌گیری و پیشنهادات به پایان می‌رسد.

## فصل دوم

### یادگیری تقویتی

در این فصل به معرفی و بیان روش یادگیری تقویتی می‌پردازیم. در ابتدای این فصل، مسأله یادگیری تقویتی<sup>۱</sup> را تعریف می‌کنیم. در ادامه، روش تفاضل زمانی<sup>۲</sup> معرفی می‌شود. روش تفاضل زمانی یکی از مهمترین روش‌های حل مسأله یادگیری تقویتی است که در قسمت دوم این فصل مورد بررسی قرار می‌گیرد.

#### ۱-۲ مسأله یادگیری تقویتی

در این قسمت به بیان مسأله یادگیری تقویتی می‌پردازیم. از آنجایی که همه روش‌های حل این مسأله از جمله روش‌های یادگیری تقویتی به شمار می‌روند، این مسأله تعریف کننده شاخه یادگیری تقویتی است. در این قسمت همچنین حالت ایده‌آل ریاضی مسأله یادگیری تقویتی را بیان می‌کنیم. در انتهای این قسمت اشاره‌ای به روابط عناصر اصلی آن مانند توابع ارزش و معادلات بلمن می‌شود.

#### ۱-۱-۲ رابط عامل و محیط

مسأله یادگیری تقویتی، چارچوب مناسبی را برای یادگیری در حین تعامل (با محیط) و رسیدن به هدف ایجاد می‌کند. یادگیرنده یا تصمیم گیرنده را عامل گویند و آن چیزی که با آن تعامل می‌شود، شامل هر چیزی خارج از عامل، محیط نامیده می‌شود. این تعامل به طور پیوسته صورت می‌پذیرد و عامل عملی را انتخاب می‌کند و محیط پاسخ آن را می‌دهد و عامل را در وضعیت جدیدی قرار می‌دهد. محیط همچنین پاداش<sup>۴</sup> عمل‌های عامل را می‌دهد. پاداش مقدار عددی است که عامل باید آن را در طول زمان ماکزیمم کند [۲۸].

عامل و محیط در بازه‌های زمانی گسسته  $t = 0, 1, 2, \dots$  با یکدیگر تعامل می‌کنند و در هر بازه زمانی  $t$ ، عامل وضعیتی<sup>۵</sup>  $(s_t)$  را از محیط دریافت می‌کند که  $s_t \in S$  و  $S$  همه وضعیت‌های ممکن محیط است و بر مبنای آن عملی<sup>۶</sup>  $(a_t)$  را انتخاب می‌کند که  $a_t \in A(s_t)$  و  $A(s_t)$  مجموعه عمل‌هایی است که در وضعیت  $s_t$  موجود می‌باشند. در بازه زمانی بعد، و در نتیجه عملش، عامل مقدار عددی  $r_{t+1} \in R$  را دریافت می‌کند و در وضعیت جدید  $s_{t+1}$  قرار می‌گیرد. در شکل ۱-۲ چگونگی تعامل عامل و محیط ترسیم شده است.

در هر بازه زمانی، عامل نگاشتی از وضعیت‌ها به احتمال انتخاب هر عمل ممکن را انجام می‌دهد. این نگاشت سیاست<sup>۷</sup> عامل نامیده و با  $\pi_t$  داده می‌شود که  $\pi(s, a)$  احتمال انتخاب عمل  $a$  در وضعیت  $s$  است و یا احتمال  $a_t = a$  در وضعیت  $s_t = s$  است. روش‌های یادگیری تقویتی چگونگی تغییر سیاست عامل در اثر کسب تجربه را بررسی می‌کنند.

<sup>1</sup> Reinforcement Learning Problem

<sup>2</sup> Temporal Difference

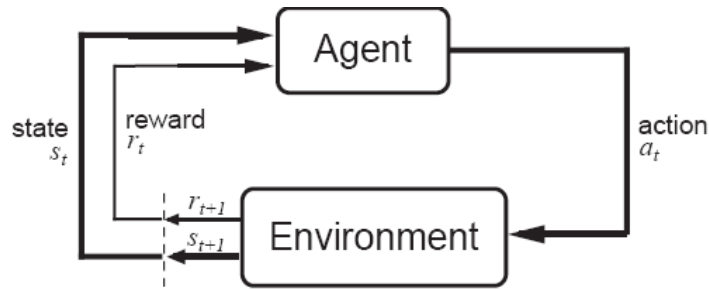
<sup>3</sup> Interface

<sup>4</sup> Reward

<sup>5</sup> state

<sup>6</sup> action

<sup>7</sup> policy



شکل ۱-۲ تعامل عامل و محیط در یادگیری تقویتی [۲۸]

این چارچوب کاملاً مجرد و انعطاف‌پذیر است و می‌تواند به مسائل مختلف و با روش‌های گوناگون اعمال شود. برای مثال، نیازی به ثابت بودن طول بازه‌های زمانی نیست و بازه‌ها می‌توانند مراحل متوالی و دلخواه تصمیم‌گیری و عمل کردن باشند. اعمال می‌توانند کنترل سطح پایین، مانند ولتاژ عملی به بازوی ربات، یا تصمیم‌گیری سطح بالا مانند صرف نهار یا رفتن به دانشگاه باشند. به طور مشابه، وضعیت‌ها می‌توانند اشکال گوناگونی داشته باشند. وضعیت‌ها می‌توانند با ورودی‌های سطح پایین، همچون دریافت‌های مستقیم حس‌گر<sup>۱</sup>، و یا سطح بالای مجرد مانند توصیف سمبلیک یک شیء در اتاق توصیف شوند. وضعیت‌ها می‌توانند بر مبنای حافظه‌ای از دریافت‌های گذشته و یا حتی کاملاً ذهنی باشند. برای مثال، عامل می‌تواند در وضعیتی باشد که مکان شیء خاصی را نداند. به همین ترتیب بعضی از عمل‌ها کاملاً ذهنی یا محاسباتی هستند. مثلاً بعضی از عمل‌ها کنترل می‌کنند که عامل درباره چه چیزی فکر کند. به طور کلی، عمل‌ها می‌توانند هر تصمیمی باشند که می‌خواهیم چگونگی اتخاذ آن را یاد بگیریم، و وضعیت‌ها هر چیزی که در اتخاذ چنین تصمیمی مفید هستند می‌باشند.

وضعیت‌ها و عمل‌ها از کاربردی تا کاربرد دیگر تغییر می‌کنند و چگونگی ارائه<sup>۲</sup> آن‌ها تاثیر زیادی در کارایی دارد. در یادگیری تقویتی همانند دیگر روش‌های یادگیر، انتخاب نوع ارائه بیش از آن که علم باشد، هنر است.

## ۲-۱-۲ هدف‌ها و پاداش‌ها

در یادگیری تقویتی، هدف<sup>۳</sup> عامل براساس سیگنال پاداش که از محیط دریافت می‌کند فرموله شده است. در هر بازه زمانی، پاداش، عددی حقیقی ( $r_t \in \mathbb{R}$ ) می‌باشد و هدف عامل ماکزیمم کردن مجموع پاداش‌ها در طول زمان است که به این معنی است که پاداش بلاواسطه را ماکزیمم نمی‌کند، بلکه مجموع پاداش دریافتی در طول زمان<sup>۴</sup> را ماکزیمم می‌کند [۱۱].

استفاده از سیگنال تقویتی، ویژگی متمایز یادگیری تقویتی است. در ابتدا ممکن است این روش فرموله کردن، محدود کننده به نظر برسد، ولی در عمل اثبات شده است که کاملاً انعطاف‌پذیر و عملی است. برای اینکه چگونگی استفاده از سیگنال تقویتی کاملاً مشخص شود به مثال‌های زیر توجه شود. برای مثال، رباتی می‌خواهد فرار از ماز<sup>۵</sup> را یاد بگیرد. به این منظور، پاداش همیشه مقدار صفر دارد مگر اینکه عامل از ماز خارج شود که پاداش +۱ می‌گیرد. روش دیگر پاداش دهی مسئله ماز عبارت است از پاداش -۱ در هر بازه زمانی و پاداش صفر در صورت فرار از ماز. که این نحوه پاداش دهی باعث فرار ربات از ماز در کمترین زمان ممکن می‌شود. برای عاملی که

<sup>1</sup> sensor  
<sup>2</sup> representation  
<sup>3</sup> goal  
<sup>4</sup> long run  
<sup>5</sup> maze



می‌خواهد شطرنج یاد بگیرد. پاداش‌دهی معمول عبارت است از: +۱ برای برد، -۱ برای باخت، و پاداش صفر در صورت تساوی و همه حالت‌های غیر پایانی دیگر.

در مثال‌های بالا، عامل ماکزیمم کردن مقادیر پاداش را یاد می‌گیرد. اگر بخواهیم که عامل مسأله‌ای را حل کند باید پاداش‌دهی را به گونه‌ای تعریف کنیم که ماکزیمم شدن آن در طول زمان، باعث حل مسأله شود بنابراین تعیین نوع پاداش‌دهی مناسب برای حل مسأله کاملاً ضروری است. یکی از نکات دیگر پاداش‌دهی این است که سیگنال‌های تقویتی نباید دانش ما را به عامل منتقل کنند و همچنین نباید چگونگی رسیدن به هدف را مشخص نمایند، بلکه تنها باید هدف و آن چیزی را که عامل می‌خواهد به آن برسد را نشان دهند.

نکته دیگری که وجود دارد این است که محاسبه پاداش در خارج از عامل و در محیط صورت می‌گیرد. به عبارت دیگر پاداش‌دهی خارج از کنترل عامل است و عامل نباید قادر به تغییرات این سیگنال باشد. البته عامل می‌تواند پاداش‌های داخلی برای خود تعریف کند (بسیاری از روش‌های یادگیری تقویتی از چنین پاداش‌هایی استفاده می‌کنند).

## ۲-۱-۳ بازگشت<sup>۱</sup>

گفته شد هدف عامل، ماکزیمم کردن مجموع پاداش دریافتی در طول زمان است، اما چگونه آن را به صورت فرمول نشان می‌دهیم؟ اگر دنباله پاداش‌های دریافتی بعد از بازه زمان  $t$  را  $r_{t+1}, r_{t+2}, r_{t+3}, \dots$  و ... تعریف کنیم، بطور کلی می‌خواهیم مقدار بازگشت مورد انتظار<sup>۲</sup> را ماکزیمم کنیم. بازگشت دارای تعاریف مختلفی است که در ساده‌ترین حالت عبارت است از:

$$R_t = r_{t+1} + r_{t+2} + r_{t+3} + \dots + r_T$$

که  $T$  بازه زمانی پایانی اپیزود است. این تعریف در کاربردهای عملی، بسیار مورد استفاده قرار می‌گیرد. زیرا معمولاً ارتباط عامل و محیط به زیردنباله‌هایی از بازه‌ها تقسیم می‌شوند که اپیزود<sup>۳</sup> نام دارند. مثالی از آن انجام بازی فرار از ماز، و یا هر تعامل تکراری دیگر است. هر اپیزود در وضعیت خاصی که وضعیت پایانی نامیده می‌شود، پایان می‌یابد و ابتدای اپیزود بعدی، در وضعیت شروع استاندارد قرار می‌گیرد. مسائل با ویژگی‌های فوق را مسائل اپیزودیک گویند. بعضی اوقات در مسائل اپیزودیک مجموعه وضعیت‌های غیر پایانی را با  $S$  و مجموعه وضعیت‌های غیر پایانی و پایانی را با  $S^+$  نشان می‌دهند. غیر از تعامل اپیزودیک، تعامل غیر اپیزودیک عامل و محیط نیز وجود دارد که در [۱۱، ۲۸] شرح داده شده است و ارتباط این دو حالت تعامل با محیط مورد بررسی قرار گرفته است.

نرخ کاهندگی<sup>۴</sup> مفهوم دیگری است که بایستی تعریف شود. اگر از نرخ کاهندگی در تعریف بازگشت استفاده شود، عامل عمل‌هایی را انتخاب می‌کند که مجموع پاداش‌های کاهش یافته<sup>۵</sup> که در آینده دریافت می‌کند را ماکزیمم کند. در این حالت رابطه بازگشت عامل به صورت زیر تعریف می‌شود:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{k=0}^T \gamma^k r_{t+k+1}$$

<sup>۱</sup> return

<sup>۲</sup> expected return

<sup>۳</sup> episode

<sup>۴</sup> discounting

<sup>۵</sup> discounted

که پارامتر  $\gamma$  ( $0 \leq \gamma \leq 1$ ) نرخ کاهندگی است. نرخ کاهندگی، مقدار فعلی پاداش‌های آینده را مشخص می‌کند. پاداشی که در  $k$  بازه زمانی بعد دریافت می‌شود  $\gamma^{k-1}$  برابر پاداشی است که مستقیماً دریافت می‌شود.

## ۴-۱-۲ خصوصیت مارکف

مناسب است که سیگنال وضعیتی، همه دریافت‌های گذشته را خلاصه کند و همه اطلاعات مربوط به آنها را نگهداری کند. به صورت عادی برای تحقق این مورد، نیاز به چیزی بیش از یک دریافت مستقیم است. البته این نیاز فراتر از کل دریافت‌های گذشته نمی‌شود. سیگنال وضعیتی که موفق در نگهداری همه اطلاعات گذشته می‌باشد، دارای خصوصیت مارکف است. برای مثال، مکان مهره‌های شطرنج دارای خصوصیت مارکف است. زیرا دنباله کل حرکاتی که به این وضعیت رسیده است را خلاصه می‌کند. اگرچه اکثر اطلاعاتی که منجر به این وضعیت شده است از دست رفته‌اند، ولی همه نکات مهم برای ادامه بازی نگهداری شده است. به این خصوصیت گاهی خصوصیت مستقل از مسیر<sup>۱</sup> هم گویند که به این معنا است که وضعیت فعلی، مستقل از مسیر و تاریخچه سیگنال‌های گذشته است. با فرض تعداد محدود وضعیت و مقادیر پاداشها، تعریف ریاضی خصوصیت مارکف را بیان می‌کنیم. توجه کنید که یک محیط بطور کلی در زمان  $t+1$  به عمل زمان  $t$  پاسخ می‌دهد و در کلی‌ترین حالت، این پاسخ ممکن است به همه پاسخ‌های گذشته ارتباط داشته باشد. در این حالت دینامیک محیط با توزیع احتمال کامل زیر مشخص می‌کنیم:

$$\Pr\{s_{t+1}=s', r_{t+1}=r | s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_1, s_0, a_0\} \quad (1-2)$$

برای همه  $r, s$  ها و همه مقادیر ممکن رویدادهای گذشته:  $s_0, a_0, r_1, s_1, a_1, \dots, r_t, a_t, s_t$ . ولی در حالتی که سیگنال وضعیتی، دارای خصوصیت مارکف باشد، پاسخ محیط در زمان  $t+1$  فقط بستگی به وضعیت و عمل در زمان  $t$  دارد که در این حالت دینامیک محیط فقط با رابطه زیر بیان می‌شود:

$$\Pr\{s_{t+1}=s', r_{t+1}=r | s_t, a_t\} \quad (2-2)$$

برای همه مقادیر  $s', r, s_t, a_t$ . به عبارت دیگر می‌گوییم سیگنال وضعیتی دارای خصوصیت مارکف است، در صورتی که دو رابطه ۱-۲ و ۲-۲ برای همه مقادیر  $r, s$  و گذشته‌های  $s_0, a_0, r_1, s_1, a_1, \dots, r_t, a_t, s_t$  معادل یکدیگر باشند. در این حالت اصطلاحاً گفته می‌شود که محیط دارای خصوصیت مارکف است. محیط‌هایی که دارای خصوصیت مارکف هستند باعث سادگی حل مسأله می‌شوند. حتی برای سیگنال وضعیتی که خصوصیت مارکف ندارد می‌توانیم فرض کنیم که سیگنال تقریبی از خصوصیت مارکف را دارد.

## ۵-۱-۲ فرآیند تصمیم‌گیری مارکف

نمونه‌ای از مسأله یادگیری تقویتی که دارای خصوصیت مارکف است را فرآیند تصمیم‌گیری مارکف<sup>۲</sup> گویند. اگر فضای وضعیت‌ها و اعمال محدود باشند، آنگاه فرآیند تصمیم‌گیری مارکف متناهی<sup>۳</sup> است. فرآیند تصمیم‌گیری مارکف متناهی برای درک روش یادگیری تقویتی بسیار مهم می‌باشد. فرآیند تصمیم‌گیری مارکف متناهی توسط مجموعه وضعیت‌ها و اعمالش در دینامیک یک مرحله بعد محیط تعریف می‌شود. برای هر وضعیت و عمل داده شده  $r, s$ ، احتمال انتقال به وضعیت بعدی  $s'$  برابر است با:

$$P_{ss'}^a = \Pr\{s_{t+1}=s' | s_t=s, a_t=a\}$$

<sup>1</sup> independence of path

<sup>2</sup> Markov Decision Process (MDP)

<sup>3</sup> finite MDP

این مقادیر را احتمال‌های انتقال گویند. بطور مشابه برای هر وضعیت و عمل فعلی  $s, r$  و با وضعیت بعدی  $s'$ ، امید ریاضی مقدار حالت بعدی عبارت است:

$$R_{ss'}^a = E\{r_{t+1} | s_t = s, a_t = a, s_{t+1} = s'\}$$

این مقادیر  $P_{ss'}^a$  و  $R_{ss'}^a$ ، کاملاً ویژگی‌های مهم دینامیک فرآیند تصمیم‌گیری مارکف متناهی را بیان می‌کنند.

## ۲-۱-۶ توابع ارزش

همه الگوریتم‌های یادگیری تقویتی، بر مبنای تخمین توابع ارزش<sup>۱</sup> هستند. توابع ارزش، توابعی از وضعیت‌ها (یا زوج‌های وضعیت عمل) هستند که خوب بودن عامل در وضعیتی خاص را تخمین می‌زنند (یا خوب بودن انجام عمل داده شده در وضعیتی خاص را تخمین می‌زنند). خوب بودن را با پاداش‌های دریافتی در آینده تعریف می‌کنیم. واضح است پاداش‌هایی که عامل، انتظار دریافت آن را در آینده دارد وابسته به عمل فعلی است. بنابراین، توابع ارزش با توجه به سیاست انتخابی تعریف می‌شوند. همانگونه که ذکر شد  $\pi$  نگاشتی از وضعیت  $s \in \mathcal{S}$  و عمل  $a \in A(s)$  به احتمال انتخاب عمل  $a$  در وضعیت  $s$  یا  $\pi(s, a)$  است. ارزش وضعیت  $s$  در سیاست  $\pi$ ،  $V^\pi(s)$ ، امید ریاضی بازگشت است در صورتی که حالت اولیه  $s$  باشد و سیاست  $\pi$  در ادامهٔ پیروز انتخاب شود. تابع ارزش وضعیت برای فرآیند تصمیم‌گیری مارکف به صورت زیر تعریف شده است:

$$V^\pi(s) = E_\pi\{R_t | s_t = s\} = E_\pi\left\{\sum_{k=0}^T \gamma^k r_{t+k+1} | s_t = s\right\}$$

که  $E_\pi\{\cdot\}$  نشان‌دهندهٔ امید ریاضی بازگشت در صورت اتخاذ سیاست  $\pi$  است. توجه کنید که ارزش وضعیت‌های پایانی در صورت وجود، صفر می‌باشد.  $V^\pi$  را تابع ارزش وضعیت برای سیاست  $\pi$  می‌نامیم.

بطور مشابه،  $Q^\pi(s, a)$  ارزش انتخاب عمل  $a$  در وضعیت  $s$  در صورت اتخاذ سیاست  $\pi$  است. در واقع  $Q^\pi(s, a)$  امید ریاضی بازگشت است در صورتی که از وضعیت  $s$  شروع و عمل  $a$  را انتخاب کنیم و پس از آن سیاست  $\pi$  را پیش بگیریم.

$$Q^\pi(s, a) = E_\pi\{R_t | s_t = s, a_t = a\} = E_\pi\left\{\sum_{k=0}^T \gamma^k r_{t+k+1} | s_t = s, a_t = a\right\}$$

$Q^\pi$  را تابع ارزش عمل برای سیاست  $\pi$  می‌نامیم.

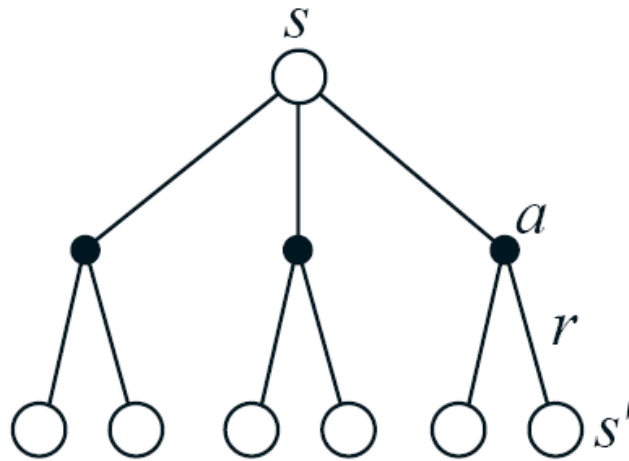
می‌توان توابع ارزش  $V^\pi$  و  $Q^\pi$  را با استفاده از تجربه، تخمین زد. برای مثال، فرض کنید عامل سیاست  $\pi$  را اتخاذ کرده باشد و برای هر وضعیتی که مواجه می‌شود میانگینی از بازگشت واقعی که پس از آن وضعیت دریافت می‌کند را نگه دارد. در صورت زیاد شدن تعداد دفعات مشاهدهٔ وضعیت، این میانگین به ارزش وضعیت  $V^\pi(s)$  همگرا خواهد شد. اگر میانگین‌های مجزایی برای هر زوج وضعیت-عمل نگه‌داری شود، این میانگین‌ها به مقادیر عمل  $Q^\pi(s, a)$  همگرا خواهند شد. به این گونه روش‌های تخمین، روش‌های مونت کارلو [۲۸، ۱۱] گویند. البته، اگر تعداد وضعیت‌ها زیاد شود، آنگاه عملی نیست که برای هر وضعیت مستقلاً، میانگینی نگه‌داری شود. در مقابل  $V^\pi$  و  $Q^\pi$  را به صورت توابعی پارامتری بیان می‌کنیم و پارامترها را به گونه‌ای که بهتر با بازگشت‌های مشاهده شده منطبق باشند، تعیین می‌کنیم.

دارا بودن رابطهٔ بازگشتی، یکی از خصوصیت اساسی توابع ارزش می‌باشد. برای هر سیاست  $\pi$  و هر وضعیت  $s$  رابطهٔ زیر بین ارزش وضعیت  $s$  و ارزش وضعیت‌های ممکن بعدی آن وجود دارد:

<sup>1</sup> value function

$$\begin{aligned}
V^\pi(s) &= E\{R_t | s_t = s\} = \\
&= E_\pi \left\{ \sum_{k=0}^T \gamma^k r_{t+k+1} | s_t = s \right\} \\
&= E_\pi \left\{ r_{t+1} + \gamma \sum_{k=0}^T \gamma^k r_{t+k+2} | s_t = s \right\} \\
&= \sum_a \pi(s, a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma E_\pi \left\{ \sum_{k=0}^T \gamma^k r_{t+k+2} | s_{t+1} = s' \right\}] \\
&= \sum_a \pi(s, a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^\pi(s')] \tag{۳-۲}
\end{aligned}$$

که  $s' \in S^+$  و  $a \in A(s)$  است. معادله (۳-۲) را معادله بلمن<sup>۱</sup> برای  $V^\pi$  گوئیم که رابطه بین ارزش یک وضعیت و ارزش وضعیتهای بعدی آن را بیان می‌کند. تغییر از یک وضعیت به وضعیت‌های ممکن بعدی در شکل ۲-۲ نشان داده شده است. در این شکل هر دایره سفید نمایانگر یک وضعیت و هر دایره سیاه نشان دهنده زوج وضعیت-عمل است. با شروع از وضعیت  $s$  گره ریشه در شکل ۲-۲، عامل می‌تواند مجموعه‌ای از عمل‌ها را انتخاب کند (سه عمل در شکل نشان داده شده است). با انتخاب هر کدام از عمل‌ها، محیط به یکی از وضعیت‌های بعدی،  $s'$ ، می‌رود و پاداش  $r$  را به عامل می‌دهد. معادله بلمن (۳-۲) میانگین وزن‌دار همه حالات ممکن و احتمال وقوع را به صورت وزن دار بیان می‌کند. ارزش وضعیت اولیه برابر امید ریاضی (کاهش‌یافته) وضعیت بعدی، به اضافه امید ریاضی پاداش دریافتی تا زمان بعدی است. این نمودار را نمودار پشتیبان<sup>۲</sup> می‌گویند و برای شرح تصویری الگوریتم‌ها استفاده می‌شود.



شکل ۲-۲ نمودار پشتیبان برای  $V^\pi$

این نمودار را از آن جهت پشتیبان گویند که شامل عملیات پشتیبان یا به روز رسانی می‌باشد که هسته اصلی روش‌های یادگیری تقویتی است. عمل پشتیبان ارزش وضعیت‌های بعدی را به وضعیت قبلی منتقل می‌کند.

<sup>۱</sup> Bellman equation

<sup>۲</sup> backup