

به نام خدا

کتابچه

مقدمه‌ای بر یادگیری تقویتی

میثم میرزائی

علی ژاله کریمی

شهریور ۱۴۰۳

## فهرست مطالب

|  |    |
|--|----|
| ۱- مقدمه   | ۵  |
| ۱-۱- یادگیری تقویتی چیست؟                                  | ۵  |
| ۱-۲- تفاوت یادگیری تقویتی با یادگیری با نظارت و بدون نظارت | ۶  |
| ۱-۳- مفاهیم کلیدی یادگیری تقویتی                           | ۷  |
| ۱-۳-۱- عامل (agent)  | ۸  |
| ۱-۳-۲- محیط (environment)                                  | ۸  |
| ۱-۳-۳- حالت و مشاهدات (state & observations)               | ۸  |
| ۱-۳-۴- اقدام (action)                                      | ۹  |
| ۱-۳-۵- خط‌مشی (policy)                                     | ۹  |
| ۱-۳-۶- پاداش (reward)                                      | ۱۰ |
| ۱-۳-۷- تابع ارزش (value function)                          | ۱۰ |
| ۱-۳-۸- مدل (model)   | ۱۱ |
| ۱-۴- چالش‌ها و محدودیت‌های یادگیری تقویتی                  | ۱۲ |
| ۱-۴-۱- کارایی نمونه  | ۱۲ |
| ۱-۴-۲- اکتشاف در مقابل بهره‌برداری                         | ۱۳ |
| ۱-۴-۳- پاداش با تأخیر                                      | ۱۴ |
| ۱-۴-۴- مهندسی پاداش  | ۱۵ |
| ۱-۴-۵- عدم تفسیرپذیری                                      | ۱۵ |
| ۱-۴-۶- ملاحظه‌های اخلاقی                                   | ۱۵ |
| ۱-۴-۷- سایر چالش‌ها  | ۱۵ |
| ۱-۵- کاربردها  | ۱۶ |

|         |  |    |
|---------|--|----|
| ۱-۶-۱-۱ | فرایند مارکوف (MP)                                       | ۱۷ |
| ۱-۶-۱-۲ | فرایند پاداش مارکوف                                      | ۲۱ |
| ۱-۶-۱-۳ | افزودن اقدامات   | ۲۵ |
| ۱-۶-۱-۴ | خطمشی (سیاست)  | ۲۶ |
| ۱-۶-۱-۵ | معادله بهینگی بلمن                                       | ۲۷ |
| ۲-۱     | الگوریتم‌های پایه یادگیری تقویتی                         | ۳۰ |
| ۲-۱-۱   | روش‌های مبتنی بر مدل (برنامه‌نویسی پویا)                 | ۳۰ |
| ۲-۱-۲   | روش‌های مبتنی بر ارزش (مونت کارلو و یادگیری تفاوت زمانی) | ۳۳ |
| ۲-۱-۲-۱ | روش‌های مونت کارلو                                       | ۳۴ |
| ۲-۱-۲-۲ | روش‌های تفاوت زمانی                                      | ۳۶ |
| ۲-۱-۲-۳ | SARSA  | ۳۷ |
| ۲-۱-۲-۴ | Q-learning   | ۳۹ |
| ۲-۱-۳   | روش‌های مبتنی بر خطمشی (گرادیان خطمشی)                   | ۴۲ |
| ۲-۱-۳-۱ | روش‌های گرادیان خطمشی (Policy Gradient)                  | ۴۳ |
| ۲-۱-۴   | تقریب توابع و یادگیری عمیق در یادگیری تقویتی             | ۴۵ |
| ۲-۱-۴-۱ | Deep Q-Networks (DQN)                                    | ۴۸ |
| ۲-۱-۴-۲ | گرادیان خطمشی قطعی عمیق (DDPG)                           | ۵۱ |
| ۳-۱     | مباحث پیشرفته در یادگیری تقویتی                          | ۵۳ |
| ۳-۱-۱   | انواع Deep Q-Networks                                    | ۵۳ |
| ۳-۱-۱-۱ | N-step DQN   | ۵۳ |
| ۳-۱-۱-۲ | Double DQN   | ۵۴ |
| ۳-۱-۱-۳ | Noisy networks   | ۵۵ |

|    |   |
|----|---|
| ۵۵ | ..... Prioritized replay buffer-۳-۱-۴       |
| ۵۶ | ..... Dueling DQN-۳-۱-۵                     |
| ۵۷ | ..... Categorical(Distributional) DQN-۳-۱-۶ |
| ۵۹ | ..... یادگیری انتقالی و فرا یادگیری ۳-۲     |
| ۶۰ | ..... یادگیری تقویتی سلسله‌مراتبی ۳-۳       |
| ۶۱ | ..... یادگیری تقویتی چندعاملی ۳-۴           |
| ۶۳ | ..... منابع ۴                               |
| ۶۳ | ..... پیوست: پیاده‌سازی‌ها                  |

## ۱- مقدمه

این ایده که ما از طریق تعامل با محیط خود یاد می‌گیریم احتمالاً اولین چیزی است که وقتی به ماهیت یادگیری فکر می‌کنیم به ذهنمان خطور می‌کند. وقتی نوزادی بازی می‌کند، دستانش را تکان می‌دهد یا به اطراف نگاه می‌کند، معلم صریحی ندارد، اما ارتباط حسی - حرکتی مستقیمی با محیط خود دارد. اعمال این ارتباط اطلاعات فراوانی در مورد علت و نتیجه، در مورد پیامدهای اعمال و اقداماتی که برای دستیابی به اهداف باید انجام داد، تولید می‌کند. در طول زندگی ما، چنین تعاملاتی بدون شک منبع اصلی دانش در مورد محیط و خودمان است. چه در حال یادگیری رانندگی با ماشین باشیم و چه در حال انجام یک مکالمه، ما کاملاً از واکنش محیط خود به کاری که انجام می‌دهیم آگاه هستیم و به دنبال تأثیرگذاری بر اتفاقات از طریق رفتار خود هستیم. یادگیری از تعامل یک ایده اساسی است که زیربنای تقریباً تمام نظریه‌های یادگیری و هوش است (Richard S. Sutton and Andrew G. Barto).

این کتابچه حاصل یک کار تحقیقاتی است که تمام تلاش در آن ارائه یک دانش حداقلی نسبت به مباحث یادگیری تقویتی و پیشرفت‌ها در این حوزه است. در این کتابچه سعی شده است حجم مطالب ارائه شده حد الامکان به صورت مختصر ارائه شود و تا حد ممکن مباحث اصلی یادگیری تقویتی مورد بررسی قرار گیرد. این کتابچه دریچه‌ای برای ورود به دنیای یادگیری تقویتی است و طبیعتاً تمام نیازهای خواننده برای تسلط بر این حوزه را تأمین نمی‌کند. مخاطب این کتاب دانشجویان و علاقه‌مندان به حوزه هوش مصنوعی و یادگیری تقویتی هستند که می‌خواهند یک دانش اولیه از مباحث یادگیری تقویتی به دست آورند.

در بخش اول این کتابچه ابتدا مفاهیم یادگیری تقویتی پرداخته می‌شود سپس الگوریتم‌های پایه یادگیری تقویتی معرفی می‌شوند. در ادامه مباحث پیشرفته در یادگیری تقویتی بررسی می‌شود و در انتها منابع مورد استفاده معرفی می‌شوند و چند نمونه از پیاده‌سازی‌های الگوریتم‌های یادگیری تقویتی ارائه می‌شود.

## ۱-۱- یادگیری تقویتی چیست؟

یادگیری تقویتی<sup>۱</sup> زیرشاخه‌ای از یادگیری ماشین<sup>۲</sup> است که در آن یک عامل<sup>۳</sup> هوش مصنوعی<sup>۴</sup> با محیط<sup>۵</sup> اطراف خود با کمک روش آزمون و خطا تعامل می‌کند و یک استراتژی رفتاری بهینه را بر اساس سیگنال‌های پاداش<sup>۶</sup> دریافتی از تعاملات قبلی می‌آموزد. به یادگیرنده گفته نمی‌شود که چه اقداماتی را انجام دهد (مانند یادگیری با

<sup>1</sup> Reinforcement learning

<sup>2</sup> Machine learning

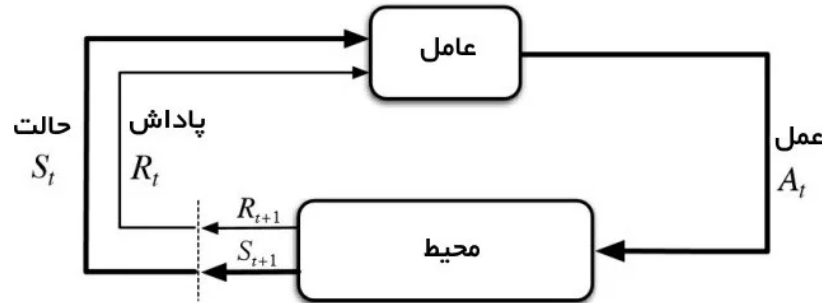
<sup>3</sup> Agent

<sup>4</sup> Artificial intelligence

<sup>5</sup> Environment

<sup>6</sup> Reward

نظارت)، بلکه باید کشف<sup>۱</sup> کند که کدام اقدامات با امتحان کردن آنها بیشترین پاداش را به همراه دارد. دو ویژگی جستجوی آزمون و خطا<sup>۲</sup> و پاداش با تأخیر<sup>۳</sup> دو ویژگی متمایزکننده مهم یادگیری تقویتی هستند. عامل و محیط اجزای اساسی یادگیری تقویتی هستند، همان طور که در شکل ۱ نشان داده شده است. محیط موجودی است که عامل می تواند با آن تعامل داشته باشد.



شکل ۱ - تعامل عامل و محیط

از بین تمام اشکال یادگیری ماشینی، یادگیری تقویتی نزدیک ترین نوع یادگیری است که انسان و سایر حیوانات انجام می دهند، و بسیاری از الگوریتم های اصلی یادگیری تقویتی در اصل از سیستم های یادگیری بیولوژیکی الهام گرفته شده اند. یادگیری تقویتی هم از طریق یک مدل روان شناختی از یادگیری حیوانات که با برخی از داده های تجربی مطابقت دارد و هم از طریق یک مدل تأثیرگذار از بخش هایی از سیستم پاداش مغز، نتیجه داده است.

## ۲-۱- تفاوت یادگیری تقویتی با یادگیری با نظارت و بدون نظارت

یادگیری تقویتی با یادگیری با نظارت<sup>۴</sup> متفاوت است، نوعی یادگیری که در بیشتر تحقیقات فعلی در زمینه یادگیری ماشین مورد مطالعه قرار می گیرد. یادگیری تحت نظارت، یادگیری از مجموعه آموزشی از نمونه های برچسب دار ارائه شده توسط یک ناظر خارجی<sup>۵</sup> آگاه است. هر مثال توصیفی از یک موقعیت همراه با یک مشخصه «برچسب» از اقدام صحیحی است که سیستم باید در آن موقعیت انجام دهد که اغلب برای شناسایی دسته ای است که موقعیت به آن تعلق دارد. هدف این نوع یادگیری این است که سیستم پاسخ های خود را برون یابی یا تعمیم دهد<sup>۶</sup> تا در موقعیت هایی که در مجموعه آموزشی وجود ندارد به درستی عمل کند. این نوع مهمی از

<sup>1</sup> Discover

<sup>2</sup> trial-and-error search

<sup>3</sup> delayed reward

<sup>4</sup> supervised learning

<sup>5</sup> external supervisor

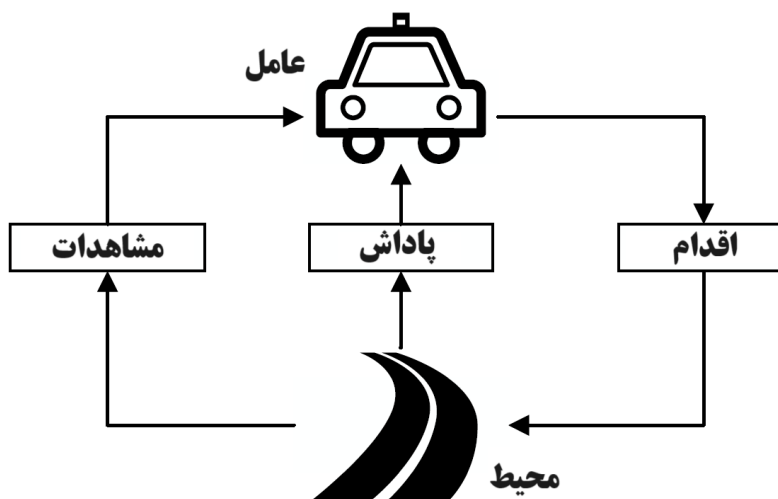
<sup>6</sup> generalize

یادگیری است، اما به‌تنهایی برای یادگیری از تعامل کافی نیست. در مسائل تعاملی، اغلب غیرعملی است که نمونه‌هایی از رفتار مطلوب به دست آوریم که هم درست و هم معرف همه موقعیت‌هایی باشد که عامل باید در آنها عمل کند. در قلمروی ناشناخته (جایی که انتظار می‌رود یادگیری سودمندترین باشد) یک عامل باید بتواند از تجربه خود بیاموزد.

یادگیری تقویتی همچنین با آنچه محققان یادگیری ماشینی آن را **یادگیری بدون نظارت**<sup>۱</sup> می‌نامند، متفاوت است که معمولاً در مورد **یافتن ساختار پنهان** در مجموعه داده‌های **بدون برچسب** است. اگرچه ممکن است کسی وسوسه شود که یادگیری تقویتی را نوعی یادگیری بدون نظارت بداند، زیرا بر نمونه‌هایی از رفتار صحیح تکیه نمی‌کند، اما یادگیری تقویتی به‌جای تلاش برای یافتن ساختار پنهان، سعی در **به حداکثر رساندن سیگنال پاداش** دارد. آشکارسازی ساختار در تجربه یک عامل مطمئناً می‌تواند در یادگیری تقویتی مفید باشد، اما به‌خودی‌خود به مشکل یادگیری تقویتی در به حداکثر رساندن سیگنال پاداش رسیدگی نمی‌کند.

بنابراین، یادگیری تقویتی در کنار یادگیری نظارت شده و یادگیری بدون نظارت و شاید پارادایم‌های دیگر، به‌عنوان سومین پارادایم یادگیری ماشینی شناخته می‌شود.

### ۳-۱- مفاهیم کلیدی یادگیری تقویتی



شکل ۲-۱ اجزاء یادگیری تقویتی در رانندگی در جاده

دو جزء اصلی در یادگیری تقویتی عامل و محیط هستند. عامل تصمیم‌گیرنده است و راه‌حل یک مشکل است. محیط بیانگر یک مشکل است. یکی از تمایزات اساسی یادگیری تقویتی از سایر رویکردهای یادگیری ماشین

<sup>۱</sup> unsupervised learning

این است که عامل و محیط با هم تعامل دارند (شکل ۲). عامل تلاش می‌کند از طریق اعمال بر محیط تأثیر بگذارد و محیط به اقدامات عامل واکنش نشان می‌دهد.

### ۱-۳-۱ - عامل (agent)

عامل **یادگیرنده** یا **تصمیم‌گیرنده** است که با انجام اقدامات خاص، انجام مشاهدات و دریافت پاداش‌های نهایی با محیط تعامل دارد. در اکثر سناریوهای عملی یادگیری تقویتی، عامل **قطعه نرم‌افزار** است که قرار است برخی از مشکلات را به روشی کم‌ویش کارآمد حل کند.

منظور از یک عامل کامل، تعاملی و هدف‌جو همیشه چیزی مانند یک ارگانیسم کامل یا ربات نیست. اینها نمونه‌هایی واضح هستند، اما یک عامل کامل، تعاملی و هدف‌جویی نیز می‌تواند جزء یک سیستم رفتاری بزرگ‌تر باشد. در این حالت، عامل به طور مستقیم با بقیه سیستم بزرگ‌تر و به طور غیرمستقیم با محیط سیستم بزرگ‌تر تعامل دارد. یک مثال ساده، عاملی است که سطح شارژ باتری ربات را کنترل می‌کند و دستوراتی را به معماری کنترل ربات ارسال می‌کند. محیط این عامل بقیه ربات، همراه با محیط ربات است.

### ۱-۳-۲ - محیط (environment)

محیط، محدوده فیزیکی یا مجازی است که عامل در آن قرار می‌گیرد و با آن تعامل دارد. محیط می‌تواند هر چیزی باشد که عامل با آن در تعامل است محیط همه چیز خارج از یک عامل است. در کلی‌ترین مفهوم، بقیه جهان است، اما این کمی بیش از حد است و حتی از ظرفیت رایانه‌های آینده فراتر می‌رود.

ارتباط عامل با محیط به پاداش (به‌دست‌آمده از محیط)، اقدامات (اجرا شده توسط عامل و داده شده به محیط) و مشاهدات (برخی اطلاعات علاوه بر پاداشی است که عامل از محیط دریافت می‌کند) محدود می‌شود.

### ۱-۳-۳ - حالت و مشاهدات (state & observations)

حالت (وضعیت)، وضعیت فعلی محیط را نشان می‌دهد. این وضعیت ممکن است بر اساس داده‌های حسی که عامل دریافت می‌کند، تعریف شود. حالت می‌تواند اطلاعاتی مانند موقعیت، سرعت، شرایط محیطی و غیره را شامل شود. به عبارت دیگر محیط، با مجموعه‌ای از متغیرهای مرتبط با مسئله نمایش داده می‌شود. ترکیبی از تمام مقادیر ممکن که این مجموعه از متغیرها می‌تواند بگیرد، **فضای حالت** نامیده می‌شود. **حالت<sup>۱</sup> مجموعه خاصی از مقادیر است که متغیرها در هر زمان معین می‌گیرند.**

<sup>1</sup> state



عامل‌ها ممکن است به وضعیت واقعی محیط دسترسی داشته باشند یا نداشته باشند. باین‌حال، به‌هرحال، عامل‌ها می‌توانند چیزی را از محیط مشاهده کنند. مجموعه‌ای از متغیرهایی که عامل در هر زمان معین درک می‌کند، **مشاهده**<sup>۱</sup> نامیده می‌شود.

ترکیبی از تمام مقادیر ممکن که این متغیرها می‌توانند بگیرند **فضای مشاهده** است. بدانید که حالت و مشاهده عباراتی هستند که به‌جای یکدیگر در جامعه یادگیری تقویتی استفاده می‌شوند. این به این دلیل است که اغلب عوامل اجازه دیدن وضعیت داخلی محیط را دارند، اما همیشه این‌طور نیست. ممکن است تفاوتی بین حالت‌ها و مشاهدات وجود داشته باشد.

#### ۴-۳-۱ - اقدام (action)

**اقدامات یا کنش‌ها یا عمل‌ها**<sup>۲</sup> کارهایی هستند که یک عامل می‌تواند در محیط انجام دهد. برای مثال، اقدامات می‌توانند حرکاتی باشند که طبق قوانین بازی مجاز است (اگر بازی باشد). در هر حالت، محیط مجموعه‌ای از اقدامات را در دسترس می‌گذارد که عامل می‌تواند از بین آنها انتخاب کند. اغلب مجموعه اقدامات برای همه حالت‌ها یکسان است. مجموعه همه اقدامات در همه حالت‌ها **فضای عمل** نامیده می‌شود.

در زیر نمونه‌ای مختصر از مسئله، عامل، محیط، اقدامات احتمالی و مشاهدات یادگیری تقویتی آمده است:

**مسئله:** شما در حال رانندگی با ماشین خود هستید. **عامل:** بخشی از مغز شما که تصمیم می‌گیرد.

**محیط:** ساخت و مدل ماشین شما، ماشین‌های دیگر، سایر رانندگان، آب و هوا، جاده‌ها، لاستیک‌ها و غیره. **اقدامات (عمل):** هدایت با  $x$ ، شتاب با  $y$ . ترمز توسط  $z$ ، روشن کردن چراغ‌های جلو. مه شکن، پنجره‌ها، پخش موسیقی. **مشاهدات:** شما به مقصد خود نزدیک می‌شوید. یک راه‌بندان در خیابان اصلی وجود دارد. ماشین کنار شما بی احتیاط رانندگی می‌کند. باران شروع به باریدن می‌کند. یک افسر پلیس جلوی شما رانندگی می‌کند.

#### ۵-۳-۱ - خط‌مشی (policy)

یک **خط‌مشی** یا **سیاست**<sup>۳</sup> نحوه رفتار عامل یادگیری را در یک زمان معین تعریف می‌کند. به‌طور کلی، یک خط‌مشی، **نگاشتی**<sup>۴</sup> است از حالت‌های درک شده از محیط به اقداماتی است که باید در آن حالت‌ها انجام شود.

<sup>1</sup> observation

<sup>2</sup> actions

<sup>3</sup> policy

<sup>4</sup> mapping

این با چیزی مطابقت دارد که در روان‌شناسی مجموعه‌ای از قوانین یا ارتباط‌های محرک-پاسخ<sup>۱</sup> نامیده می‌شود. در برخی موارد این خط‌مشی ممکن است یک تابع ساده یا جدول جستجو<sup>۲</sup> باشد، درحالی‌که در موارد دیگر ممکن است شامل محاسبات گسترده مانند فرایند جستجو باشد. خط‌مشی هسته اصلی یک عامل یادگیری تقویتی است به این معنا که به‌تنهایی برای تعیین رفتار کافی است. به‌طور کلی، خط‌مشی‌ها ممکن است تصادفی باشند و احتمالات را برای هر عمل مشخص کنند.

#### ۶-۳-۱- پاداش (reward)

سیگنال پاداش<sup>۳</sup>، هدف یک مسئله یادگیری تقویتی را مشخص می‌کند. در هر مرحله زمانی<sup>۴</sup>، محیط یک عدد واحد به نام پاداش را برای عامل یادگیری تقویتی ارسال می‌کند. تنها هدف عامل به حداکثر رساندن کل پاداشی است که در درازمدت دریافت می‌کند؛ بنابراین سیگنال پاداش، رویدادهای خوب و بد را برای عامل مشخص می‌کند. در یک سیستم بیولوژیکی، ممکن است پاداش‌ها را مشابه تجربه لذت یا درد بدانیم. آنها ویژگی‌های فوری و تعیین‌کننده مسئله‌ای هستند که عامل با آن مواجه است. سیگنال پاداش مبنای اولیه برای تغییر خط‌مشی است. اگر اقدامی که توسط خط‌مشی انتخاب می‌شود با پاداش کم همراه باشد، ممکن است این خط‌مشی برای انتخاب اقدام دیگری در آن موقعیت در آینده تغییر کند. به‌طور کلی، سیگنال‌های پاداش ممکن است توابع تصادفی از وضعیت محیط و اقدامات انجام شده باشد.

#### ۷-۳-۱- تابع ارزش (value function)

درحالی‌که سیگنال پاداش نشان می‌دهد که چه چیزی در معنای فوری خوب است، یک تابع ارزش<sup>۵</sup> مشخص می‌کند که چه چیزی در بلندمدت خوب است. به‌طور کلی، ارزش یک حالت، مقدار کل پاداشی است که یک عامل می‌تواند انتظار داشته باشد در آینده جمع‌آوری کند. درحالی‌که پاداش‌ها، مطلوبیت ذاتی و فوری حالات محیطی را تعیین می‌کنند، ارزش‌ها پس از در نظر گرفتن حالت‌هایی که احتمالاً دنبال می‌شوند و پاداش‌های موجود در آن حالات‌ها، مطلوبیت بلندمدت حالات‌ها را نشان می‌دهند. به‌عنوان مثال، یک حالت ممکن است همیشه یک پاداش فوری کم داشته باشد، اما همچنان ارزش بالایی داشته باشد، زیرا به‌طور منظم توسط سایر حالت‌هایی که پاداش‌های بالایی ارائه می‌کنند، دنبال می‌شود یا برعکس آن نیز می‌تواند درست باشد. برای انجام یک تشبیه انسانی، پاداش‌ها تا حدودی مانند لذت (اگر زیاد) و درد (اگر کم باشد) هستند، درحالی‌که ارزش‌ها با

<sup>1</sup> stimulus-response rules or associations

<sup>2</sup> lookup table

<sup>3</sup> reward signal

<sup>4</sup> time step

<sup>5</sup> value function

یک قضاوت دقیق تر و دوراندیشانه تر از اینکه ما چقدر از محیطمان در حالت خاصی خشنود یا ناراضی هستیم مطابقت دارد.

پاداش‌ها به یک معنا اولیه<sup>۱</sup> هستند، درحالی که ارزش‌ها، به عنوان پیش‌بینی پاداش، ثانویه هستند. بدون پاداش هیچ ارزشی وجود نخواهد داشت و تنها هدف از تخمین ارزش‌ها دستیابی به پاداش بیشتر است. باین وجود، این ارزش‌ها هست که هنگام تصمیم‌گیری و ارزیابی بیشتر به آن‌ها توجه می‌کنیم. انتخاب عمل بر اساس قضاوت‌های ارزشی انجام می‌شود. ما به دنبال اقداماتی هستیم که بالاترین ارزش را به همراه داشته باشد، نه بالاترین پاداش، زیرا این اقدامات در درازمدت بیشترین پاداش را برای ما به همراه دارد. متأسفانه، تعیین ارزش‌ها بسیار دشوارتر از تعیین پاداش است. جوایز اساساً مستقیماً توسط محیط داده می‌شوند، اما ارزش‌ها باید از روی توالی مشاهداتی که یک عامل در طول عمر خود انجام می‌دهد، تخمین زده و مجدداً تخمین زده شوند. در واقع، مهم‌ترین مؤلفه تقریباً همه الگوریتم‌های یادگیری تقویتی که در نظر می‌گیریم، روشی برای تخمین دقیق ارزش‌ها است. نقش اصلی برآورد ارزش مسلماً مهم‌ترین چیزی است که در مورد یادگیری تقویتی طی شش دهه گذشته آموخته شده است.

### ۸-۳-۱ - مدل (model)

برخی از سیستم‌های یادگیری تقویتی از مدلی<sup>۲</sup> از محیط استفاده می‌کنند. چیزی که رفتار محیط را تقلید می‌کند، یا به طور کلی تر، اجازه می‌دهد تا استنتاج‌هایی در مورد نحوه رفتار محیط انجام شود. برای مثال، با توجه به یک حالت و عمل، مدل ممکن است حالت بعدی و پاداش بعدی را پیش‌بینی کند. مدل‌ها برای برنامه‌ریزی استفاده می‌شوند که به معنی استفاده از هر راهی برای تصمیم‌گیری در مورد یک مسیر عملی، با در نظر گرفتن موقعیت‌های احتمالی آینده، قبل از تجربه واقعی است. روش‌هایی برای حل مسائل یادگیری تقویتی که از مدل‌ها و برنامه‌ریزی استفاده می‌کنند، روش‌های مبتنی بر مدل<sup>۳</sup> نامیده می‌شوند، در مقابل روش‌های بدون مدل<sup>۴</sup> که به صراحت یادگیرندگان آزمون و خطا هستند و تقریباً برعکس برنامه‌ریزی تلقی می‌شوند. برخی از سیستم‌های یادگیری تقویتی به طور هم‌زمان هم با آزمون و خطا یاد می‌گیرند، هم مدلی از محیط را یاد می‌گیرند و از مدل برای برنامه‌ریزی استفاده می‌کنند. یادگیری تقویتی مدرن طیفی از یادگیری سطح پایین، آزمون و خطا تا برنامه‌ریزی سطح بالا و مشورتی را در بر می‌گیرد.

<sup>1</sup> primary

<sup>2</sup> model

<sup>3</sup> model-based

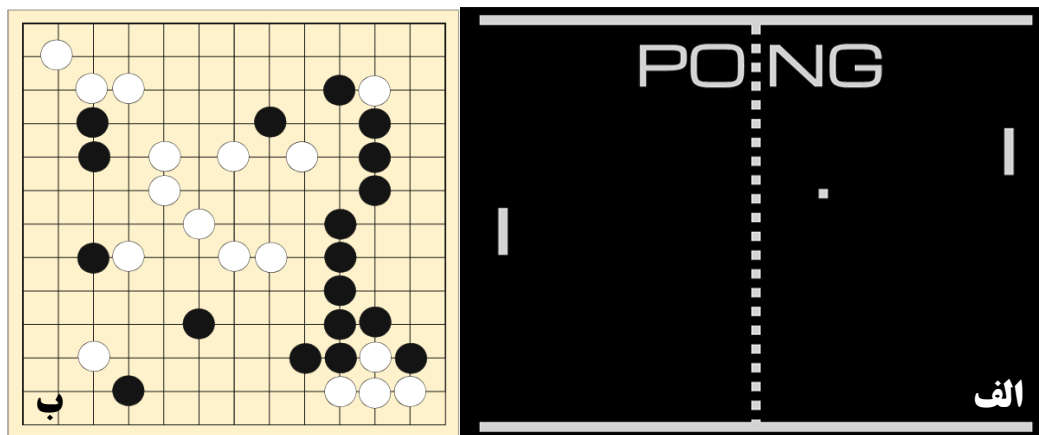
<sup>4</sup> model-free

## ۴-۱- چالش‌ها و محدودیت‌های یادگیری تقویتی

در حالی که یادگیری تقویتی قابلیت‌های قدرتمندی را ارائه می‌کند، با محدودیت‌های خاصی نیز همراه است. در حال حاضر چندین چالش و محدودیت عمده در الگوریتم‌های یادگیری تقویتی وجود دارد:

### ۴-۱-۱- کارایی نمونه

یکی چالش‌های یادگیری تقویتی **کارایی نمونه<sup>۱</sup>** است، به این معنی که یادگیری به داده‌های زیادی نیاز دارد. در کارهایی که شبیه‌سازها در دسترس هستند، یک عامل کامپیوتری می‌تواند تعاملات بی‌حد و حصری با یک محیط ثابت داشته باشد. موفقیت AlphaGo (اولین برنامه رایانه‌ای است که می‌تواند یک بازیکن حرفه‌ای بازی Go (شکل ۳، ب) را شکست دهد) بر اساس تعداد زیادی بازی است که هر بازیکن انسانی نمی‌تواند در طول زندگی آن‌ها را تجربه کند. در محیط‌های فیزیکی واقعی، مانند کنترل ربات یا تعامل انسانی، تجربه واقعی می‌تواند زمان‌بر یا پرهزینه باشد و محیط می‌تواند مدام در حال تغییر باشد به‌طوری‌که یادگیرندگان کند نتوانند از پس آن بر بیایند.



شکل ۳- (الف) محیط بازی Pong و (ب) محیط بازی Go

به‌عنوان مثال بازی Pong را در نظر بگیرید (شکل ۳، الف): یک انسان عادی فقط به ده‌ها آزمایش نیاز دارد تا اساساً در بازی تسلط یابد و به امتیاز نسبتاً خوبی دست یابد. با این حال، برای الگوریتم‌های یادگیری تقویتی موجود (به‌ویژه با روش‌های بدون مدل)، ممکن است حداقل به ده‌ها هزار نمونه نیاز داشته باشد تا به تدریج برخی از سیاست‌های مفید را یاد بگیرد. این یک مشکل اساسی در یادگیری تقویتی ایجاد می‌کند: چگونه می‌توانیم الگوریتم یادگیری تقویتی کارآمدتری طراحی کنیم تا عاملی با مثال‌های کمتر سریع‌تر یاد بگیرد؟

<sup>1</sup> sample efficiency

اهمیت این مشکل بیشتر به دلیل هزینه تعاملات زمان واقعی<sup>۱</sup> یا دنیای واقعی<sup>۲</sup> بین عامل و محیط و یا حتی زمان و انرژی مصرفی تعاملات در محیطهای شبیه‌سازی شده موجود است. بسیاری از الگوریتم‌های یادگیری تقویتی کنونی در یک مسئله در مقیاس بزرگ یا فضای پیوسته آن قدر بازده یادگیری<sup>۳</sup> پایینی دارند که یک فرایند آموزشی معمولی حتی با شبیه‌سازی سریع هنوز به زمان انتظار غیرقابل تحمل با توان محاسباتی فعلی نیاز دارد. این می‌تواند برای تعاملات دنیای واقعی بدتر باشد. مشکلات بالقوه مصرف زمان، فرسودگی تجهیزات، ایمنی در حین اکتشاف یادگیری تقویتی و خطرات موارد شکست، همگی الزامات سخت‌گیرانه‌تری را در مورد بهره‌وری یادگیری روش‌های یادگیری تقویتی در عمل ایجاد می‌کنند.

## ۲-۴-۱- اکتشاف در مقابل بهره‌برداری

یکی از چالش‌هایی که در یادگیری تقویتی مطرح می‌شود، توازن<sup>۴</sup> بین اکتشاف<sup>۵</sup> (انتخاب اقدامات تصادفی تا بتوانیم بیشتر بیاموزیم) و بهره‌برداری<sup>۶</sup> (انتخاب بهترین اقدام بر اساس آنچه تاکنون می‌دانیم) است. برای به‌دست‌آوردن پاداش زیاد، یک عامل یادگیری تقویتی باید اقداماتی را ترجیح دهد که در گذشته امتحان کرده و در تولید پاداش مؤثر بوده است. اما برای کشف چنین اقداماتی، باید اقداماتی را امتحان کند که قبلاً انتخاب نکرده است. عامل باید از آنچه قبلاً تجربه کرده است برای به‌دست‌آوردن پاداش بهره‌برداری کند، اما همچنین باید برای انتخاب اقدامات بهتر در آینده، کاوش کند. معضل این است که نه اکتشاف و نه بهره‌برداری را نمی‌توان منحصرأ بدون شکست در کار دنبال کرد. عامل باید انواع مختلفی از اقدامات را امتحان کند و به تدریج از اقداماتی که به نظر بهترین هستند، حمایت کند. در یک کار تصادفی، هر عمل باید بارها امتحان شود تا تخمین قابل‌اعتمادی از پاداش مورد انتظار آن به دست آید. در ادامه مثالی از این چالش ارائه شده است.

<sup>۱</sup> real-time

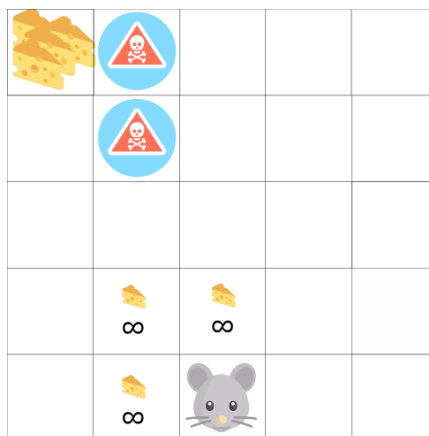
<sup>۲</sup> real-world

<sup>۳</sup> learning efficiency

<sup>۴</sup> trade-off

<sup>۵</sup> exploration

<sup>۶</sup> exploitation



شکل ۴- محیط بازی موش و پنیر

به‌خاطر داشته باشید که هدف عامل در یادگیری تقویتی بیشینه‌سازی پاداش است. بااین‌حال می‌توان در یک تله متداول افتاد. برای مثال در بازی شکل ۴، موش می‌تواند میزان نامتناهی پنیرهای کوچک داشته باشد (۱+ برای هر کدام). اما در بالای هزارتو، حجم زیادی پنیر وجود دارد (۱۰۰۰+). اگر فقط روی پاداش تمرکز شود، عامل هیچ‌وقت به آن کوه عظیم پنیر نمی‌رسد. در عوض، تنها نزدیک‌ترین منابع پاداش را اکتشاف (جستجو) می‌کند، حتی اگر این منابع کوچک باشند. اگر عامل بیشتر جستجو کند و کمتر بهره‌برداری کند، می‌تواند پاداش بزرگی پیدا کند. ایجاد تعادل میان اکتشاف اقدامات جدید و بهره‌برداری از دانش آموخته‌شده چالش‌برانگیز است؛ زیرا اکتشاف بیش از حد ممکن است تصمیم‌گیری بهینه را به تأخیر بیندازد و بهره‌برداری بیش از حد ممکن است به راه‌حل‌های غیربهینه بینجامد. معضل اکتشاف و بهره‌برداری برای چندین دهه توسط ریاضی‌دانان به‌شدت مورد مطالعه قرار گرفته است، اما هنوز حل نشده باقی‌مانده است. به‌هرحال باید قوانینی را تعیین کرد تا به برقراری توازن و مدیریت آن کمک کند.

### ۳-۴-۱- پاداش با تأخیر

یکی از اصول بنیادی یادگیری تقویتی این است که پاداش می‌تواند به‌طور جدی پس از اعمال به تعویق بیفتد<sup>۱</sup>. عامل یادگیری می‌تواند پاداش‌های کوتاه‌مدت را با دستاوردهای بلندمدت معاوضه کند. درحالی‌که این اصل بنیادی یادگیری تقویتی را مفید می‌کند، همچنین کشف سیاست بهینه را برای عامل دشوار می‌کند. این مشکل به‌ویژه در محیط‌هایی که تا زمانی که تعداد زیادی از اقدامات متوالی انجام نشود، نتیجه ناشناخته است، وجود دارد. در این سناریو، اختصاص اعتبار به یک اقدام قبلی برای نتیجه نهایی چالش‌برانگیز است و می‌تواند مغایرت زیادی را در طول آموزش ایجاد کند. بازی شطرنج یک مثال مرتبط در اینجا است، جایی که

<sup>۱</sup> delayed

نتیجه بازی تا زمانی که هر دو بازیکن تمام حرکات خود را انجام ندهند، نامعلوم است به طوری که یک حرکت قوی در میانه بازی می تواند نتایج را تغییر دهد.

#### ۴-۴-۱- مهندسی پاداش

مشکل دیگر یادگیری تقویتی مربوط به **توابع پاداش و درک معنای پاداش** است. طراحی توابع پاداش مناسب که با رفتار مدنظر هماهنگ باشد می تواند پیچیده باشد و تعریف پاداش هایی که به طور دقیق هدف های عامل را نشان می دهند، یک کار غیر ضروری است. اگر یک متخصص انسانی پاداش هایی را که عامل سعی در به حداکثر رساندن آن دارد تعریف کند، آیا به این معنی است که ما تا حدودی بر این عامل "نظارت" داریم؟ و آیا این چیز خوبی است؟ آیا پاداش باید تاحدامکان متراکم باشد که یادگیری را سریع تر می کند، یا تاحدامکان پراکنده که راه حل ها را هیجان انگیز تر و منحصر به فردتر می کند؟

ما به عنوان انسان، به نظر نمی رسد که پاداش های مشخصی داشته باشیم. اغلب، همان فرد می تواند یک رویداد را به سادگی با تغییر دیدگاه خود مثبت یا منفی ببیند. علاوه بر این، طراحی تابع پاداش برای کاری مانند راه رفتن، ساده نیست. آیا این حرکت روبه جلو است که باید هدف قرار دهیم یا سقوط نکردن؟ عملکرد پاداش "کامل" برای پیاده روی انسان چیست؟

#### ۴-۴-۵- عدم تفسیر پذیری

هنگامی که یک عامل یادگیری تقویتی سیاستی بهینه را آموخته و در محیط مستقر می شود، بر اساس تجربه خود اقداماتی را انجام می دهد. برای یک ناظر خارجی، دلیل این اقدامات ممکن است واضح نباشد. این عدم تفسیر پذیری در ایجاد اعتماد بین عامل و ناظر اختلال ایجاد می کند. اگر یک ناظر بتواند اقداماتی را که عامل یادگیری تقویتی انجام می دهد توضیح دهد، به او در درک بهتر مسئله و کشف محدودیت های مدل، به ویژه در محیط های پرخطر کمک می کند.

#### ۴-۴-۶- ملاحظه های اخلاقی

الگوریتم های یادگیری تقویتی می توانند رفتارهای نامطلوب یا مضر را بیاموزند، اگر به دقت طراحی نشده باشند، به طور بالقوه نگرانی های اخلاقی و نیاز به نظارت دقیق را افزایش می دهند.

#### ۴-۴-۷- سایر چالش ها

مشکل دیگری که سر راه یادگیری تقویتی وجود دارد، مدت زمان و منابع محاسباتی ای است که لازم است تا اطمینان حاصل کنیم یادگیری به درستی انجام شده است. از طرفی، هرچه محیط آموزشی بزرگ تر باشد به زمان و منابع بیشتری برای فرایند آموزش الگوریتم نیاز است.

جدای از چالش‌های ذکر شده در بالا برای یادگیری تقویتی چالش‌های دیگری مانند فرا یادگیری و یادگیری بازنمایی برای **تعمیم‌پذیری**<sup>۱</sup> روش‌های یادگیری تقویتی در بین وظایف (این که عامل چگونه می‌تواند بر اساس آنچه که از یک کار قدیمی آموخته است، سریع‌تر در مورد یک کار جدید یاد بگیرد؟)، یادگیری تقویتی چندعاملی<sup>۲</sup> با سایر عوامل به‌عنوان بخشی از محیط، انتقال شبیه‌سازی به واقعیت<sup>۳</sup> برای پر کردن شکاف بین محیط‌های شبیه‌سازی شده و دنیای واقعی، یادگیری تقویتی در مقیاس بزرگ<sup>۴</sup> با چارچوب‌های آموزشی موازی برای کوتاه کردن زمان آموزش و غیره.

## ۵-۱- کاربردها

یادگیری تقویتی کاربردهای متعددی در حوزه‌های مختلف پیدا کرده است، از جمله:

**رباتیک:** یادگیری تقویتی ربات‌ها را قادر می‌کند تا عمل‌ها و حرکات‌های خود را بر اساس آزمون و خطا یاد بگیرند و بهبود بخشند و به آن‌ها اجازه می‌دهد در محیط‌های پیچیده حرکت کنند یا اشیاء را دست‌کاری کنند.

**بازی:** الگوریتم‌های یادگیری تقویتی در انجام دادن بازی‌های پیچیده، مانند شطرنج، Go و بازی‌های ویدئویی، به موفقیت چشمگیری دست یافته‌اند و در برخی موارد از عملکرد انسان پیشی گرفته‌اند.

**وسایل نقلیه خودمختار:** تکنیک‌های یادگیری تقویتی را می‌توان برای آموزش خودروهای خودران برای تصمیم‌گیری بهینه در زمان واقعی به کار برد که به حمل‌ونقل ایمن‌تر و کارآمدتر می‌انجامد.

**مدیریت منابع:** یادگیری تقویتی می‌تواند برای بهینه‌سازی تخصیص منابع، زمان‌بندی و تصمیم‌گیری در حوزه‌هایی مانند مدیریت انرژی، لجستیک و ارتباطات استفاده شود.

**مراقبت‌های بهداشتی:** یکی از راه‌های مهم برای استقرار یادگیری تقویتی، رژیم‌های درمان پویا است. برای تهیه یک رژیم درمانی پویا شخصی باید مجموعه‌ای از مشاهدات بالینی بیمار و ارزیابی‌های پزشکی بیمار را وارد کند و سپس با استفاده از نتایج قبلی و سابقه پزشکی بیمار، عامل پیشنهاداتی درباره نوع درمان، دوز دارو و جدول زمان‌بندی بیمار ارائه می‌دهد. این کار باعث سریع‌تر شدن روند درمان بیمار می‌شود.

<sup>1</sup> generality

<sup>2</sup> multi-agent reinforcement learning

<sup>3</sup> sim-to-real transfer

<sup>4</sup> large-scale



## ۶-۱- فرایند تصمیم‌گیری مارکوف (MDPs)

فرایندهای تصمیم‌گیری مارکوف رسمی‌سازی کلاسیک تصمیم‌گیری متوالی هستند که در آن اقدامات نه تنها بر پاداش‌های فوری، بلکه بر موقعیت‌ها یا حلت‌های بعدی و از طریق آن پاداش‌های آینده تأثیر می‌گذارد؛ بنابراین MDPها شامل پاداش تأخیری و نیازمند توازن پاداش فوری و تسخیری هستند. فرایند تصمیم‌گیری مارکوف یک شکل ریاضی ایده‌آل از مسئله یادگیری تقویتی هستند که می‌توان برای آن گزاره‌هایی نظری دقیقی ارائه داد. فرایند مارکوف و فرایند تصمیم‌گیری مارکوف به طور گسترده در علوم کامپیوتر و سایر زمینه‌های مهندسی استفاده می‌شود.

### ۱-۶-۱- فرایند مارکوف (MP)

فرایند مارکوف ساده‌ترین فرزند از خانواده مارکوف است که به **زنجیره مارکوف** نیز معروف است. تصور کنید که سیستمی در مقابل خود دارید که فقط می‌توانید آن را مشاهده کنید. آنچه مشاهده می‌کنید **حالت‌ها** نامیده می‌شود و سیستم می‌تواند طبق برخی از قوانین دینامیک (حرکت) <sup>۱</sup> بین حالت‌ها جابه‌جا شود. باز هم شما نمی‌توانید سیستم را تحت تأثیر قرار دهید، بلکه فقط می‌توانید شاهد تغییر حالت‌ها باشید.

همان‌طور که گفته شد تمام حالت‌های ممکن برای یک سیستم مجموعه‌ای به نام **فضای حالت** را تشکیل می‌دهند. برای فرایندهای مارکوف، ما نیاز داریم که این مجموعه حالت‌ها محدود باشد. مشاهدات شما دنباله‌ای از حالت‌ها یا **زنجیره‌ای** را تشکیل می‌دهد (به همین دلیل است که فرایند مارکوف، زنجیره مارکوف نیز نامیده می‌شود). به عنوان مثال، با نگاهی به ساده‌ترین مدل آب‌وهوا در برخی از شهرها، می‌توان روز جاری را به صورت آفتابی یا بارانی مشاهده کرد که فضای حالت ماست. دنباله‌ای از مشاهدات در طول زمان زنجیره‌ای از حالت‌ها را تشکیل می‌دهد، مانند [آفتابی، آفتابی، بارانی، آفتابی و... و به این **تاریخ** <sup>۲</sup> می‌گویند.

فرایند مارکوف، باید **ویژگی مارکوف** <sup>۳</sup> را برآورده کند، به این معنی که دینامیک سیستم آینده از هر حالتی، باید فقط به این حالت بستگی داشته باشد. نکته اصلی ویژگی مارکوف این است که هر حالت قابل مشاهده را برای توصیف آینده سیستم خودکفا کند. به عبارت دیگر، **ویژگی مارکوف مستلزم آن است که حالات سیستم از یکدیگر متمایز و منحصر به فرد باشند**. در این حالت، تنها یک حالت برای مدل‌سازی حرکت‌های آینده سیستم مورد نیاز است و نه کل تاریخچه یا مثلاً آخرین  $N$  حالت.

<sup>1</sup> dynamics

<sup>2</sup> history

<sup>3</sup> markov property

در مورد مثال آب‌وهوا، ویژگی مارکوف مدل ما را محدود می‌کند تا فقط مواردی را نشان دهد که یک روز آفتابی را می‌توان با یک روز بارانی با احتمال یکسان دنبال کرد، صرف‌نظر از تعداد روزهای آفتابی که در گذشته دیده‌ایم. این یک مدل خیلی واقعی<sup>۱</sup> نیست، زیرا از روی عقل سلیم می‌دانیم که احتمال بارندگی فردا نه تنها به شرایط فعلی بلکه به تعداد زیادی از عوامل دیگر مانند فصل، عرض جغرافیایی و وجود کوه‌ها و دریا بستگی دارد؛ بنابراین، مثال ما واقعاً ساده‌لوحانه است، اما درک محدودیت‌ها و تصمیم‌گیری آگاهانه در مورد آنها مهم است.

البته، اگر بخواهیم مدل خود را پیچیده‌تر کنیم، همیشه می‌توانیم این کار را با گسترش فضای حالت خود انجام دهیم که به ما امکان می‌دهد وابستگی‌های بیشتری را در مدل به قیمت یک فضای حالت بزرگ‌تر ثبت کنیم. به عنوان مثال، اگر می‌خواهید به طور جداگانه احتمال روزهای بارانی در تابستان و زمستان را ثبت کنید، می‌توانید فصل را در حالت‌های خود قرار دهید. در این صورت فضای حالت شما [آفتابی+تابستان، آفتابی+زمستان، بارانی+تابستان، بارانی+زمستان] خواهد بود.

از آنجایی که مدل سیستم شما با ویژگی Markov مطابقت دارد، می‌توانید احتمالات انتقال را با یک ماتریس انتقال<sup>۲</sup> که یک ماتریس مربع به اندازه  $N \times N$  است که  $N$  تعداد حالت‌های مدل ما است، ثبت کنید. هر سلول در یک ردیف،  $i$ ، و یک ستون،  $j$ ، در ماتریس حاوی احتمال انتقال سیستم از حالت  $i$  به حالت  $j$  است. به عنوان مثال، در مثال آفتابی/بارانی ما، ماتریس انتقال می‌تواند به صورت زیر باشد:

جدول ۱ - ماتریس انتقال فضای حالت (آفتابی/بارانی)

|        | آفتابی | بارانی |
|--------|--------|--------|
| آفتابی | ۰.۸    | ۰.۲    |
| بارانی | ۰.۱    | ۰.۹    |

در این صورت، اگر روز آفتابی داشته باشیم، ۸۰ درصد احتمال آفتابی بودن روز بعد و ۲۰ درصد احتمال بارانی بودن روز بعد وجود دارد. اگر یک روز بارانی را مشاهده کنیم، به احتمال ۱۰٪ هوا بهتر می‌شود و به احتمال ۹۰٪ روز بعد بارانی است.

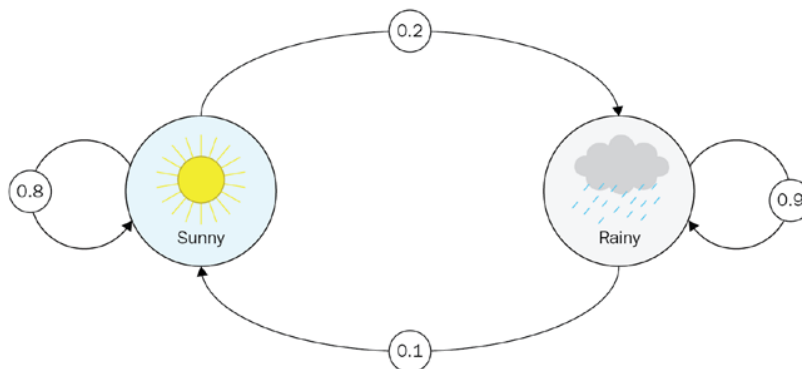
بنابراین، تعریف رسمی فرایند مارکوف به شرح زیر است:

- مجموعه‌ای از حالات (S) که یک سیستم می‌تواند در آن باشد
- یک ماتریس انتقال (T)، با احتمالات انتقال که دینامیک سیستم را تعریف می‌کند

<sup>1</sup> realistic

<sup>2</sup> transition matrix

شکل ۵ یک نمایش بصری مفید از فرایند مارکوف یک نمودار با گره‌های مربوط به حالت‌ها و یال‌های سیستم است که با احتمالاتی برچسب‌گذاری شده است که انتقال احتمالی از یک حالت به حالت دیگر را نشان می‌دهد. اگر احتمال انتقال ۰ باشد، یال نمی‌کشیم (راهی برای رفتن از یک حالت به حالت دیگر وجود ندارد).



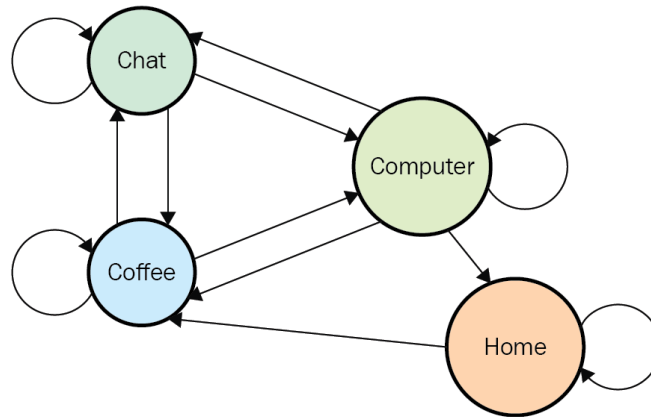
شکل ۵ - مدل آب و هوای آفتابی/بارانی

باز هم ما فقط در مورد مشاهده صحبت می‌کنیم. هیچ راهی برای تأثیرگذاری بر آب‌وهوا وجود ندارد، بنابراین فقط آن را مشاهده می‌کنیم و مشاهدات خود را ثبت می‌کنیم.

برای مثال پیچیده‌تر، بیایید مدل دیگری به نام کارمند اداری را در نظر بگیریم. فضای حالت او در مثال ما دارای حالات زیر است:

- خانه (Home): او در اداره نیست
- رایانه (Computer): او در اداره کار روی رایانه خود دارد
- قهوه (Coffee): او در اداره مشغول نوشیدن قهوه است
- گفتگو کردن (Chat): او در حال گفتگو با همکارانش در اداره است

نمودار انتقال حالت در شکل صفحه بعد نشان داده شده است:



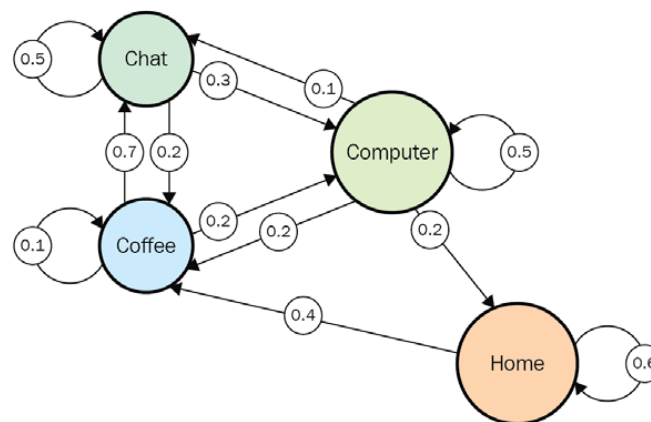
شکل ۶- نمودار انتقال حالت برای کارمند اداره

ما فرض می‌کنیم که روز کاری کارمند اداری ما معمولاً از حالت خانه شروع می‌شود و او بدون استثنا روز خود را با قهوه شروع می‌کند (بدون یال خانه ← کامپیوتر و بدون یال خانه ← گفتگو). نمودار بالا همچنین نشان می‌دهد که روزهای کاری همیشه از حالت رایانه به پایان می‌رسد (رفتن به حالت خانه) ماتریس انتقال برای نمودار قبل به شرح زیر است:

جدول ۲- ماتریس انتقال برای کارمند اداره

|          | خانه | قهوه | گفتگو | کامپیوتر |
|----------|------|------|-------|----------|
| خانه     | ۶۰٪  | ۴۰٪  | ۰٪    | ۰٪       |
| قهوه     | ۰٪   | ۱۰٪  | ۷۰٪   | ۲۰٪      |
| گفتگو    | ۰٪   | ۲۰٪  | ۵۰٪   | ۳۰٪      |
| کامپیوتر | ۲۰٪  | ۲۰٪  | ۱۰٪   | ۵۰٪      |

همان‌طور که در اینجا نشان داده شده است، احتمالات انتقال را می‌توان مستقیماً در نمودار انتقال حالت قرار داد:



شکل ۷- نمودار انتقال حالت با احتمالات انتقال

در عمل، ما به‌ندرت ماتریس انتقال دقیق را می‌دانیم. یک موقعیت بسیار واقعی‌تر در دنیای واقعی، زمانی است که ما فقط مشاهداتی از حالات سیستم خود داریم که به آنها **قسمت (اپیزود)**<sup>۱</sup> نیز می‌گویند:

- خانه ← قهوه ← قهوه ← گفتگو ← گفتگو ← قهوه ← کامپیوتر ← کامپیوتر ← خانه
- رایانه ← رایانه ← گفتگو ← گفتگو ← قهوه ← رایانه ← رایانه ← رایانه
- خانه ← خانه ← قهوه ← گفتگو ← رایانه ← قهوه ← قهوه

تخمین ماتریس انتقال از روی مشاهداتمان پیچیده نیست - ما فقط همه انتقال‌ها را از هر حالتی می‌شماریم و آنها را تا مجموع ۱ نرمال می‌کنیم. هر چه داده‌های مشاهدات بیشتری داشته باشیم، تخمین ما به مدل اصلی واقعی نزدیک‌تر خواهد بود.

همچنین شایان ذکر است که ویژگی مارکوف بر ایستایی<sup>۲</sup> دلالت دارد (توزیع انتقال برای هیچ حالتی در طول زمان تغییر نمی‌کند). غیرایستایی به این معنی است که یک عامل<sup>۳</sup> پنهان وجود دارد که بر دینامیک سیستم ما تأثیر می‌گذارد و این عامل در مشاهدات گنجانده نمی‌شود. باین‌حال، این در تضاد با ویژگی مارکوف است که مستلزم آن است که توزیع احتمال اساسی برای همان حالت بدون در نظر گرفتن تاریخچه انتقال یکسان باشد.

درک تفاوت بین انتقال‌های واقعی مشاهده شده در یک قسمت و توزیع اساسی ارائه شده در ماتریس انتقال بسیار مهم است. قسمت‌های مشخصی که مشاهده می‌کنیم به طور تصادفی از توزیع مدل نمونه‌برداری می‌شوند، بنابراین می‌توانند از قسمتی به قسمت دیگر متفاوت باشند. باین‌حال، احتمال انتقال واقعی برای نمونه‌برداری یکسان باقی می‌ماند. اگر این‌طور نباشد، اصالت شکلی<sup>۴</sup> زنجیره مارکوف غیرقابل اجرا می‌شود.

اکنون می‌توانیم جلوتر برویم و مدل فرایند مارکوف را گسترش دهیم تا آن را به مشکلات یادگیری تقویتی نزدیک‌تر کنیم. بیایید به نمودار جایزه پاداش کنیم!

## ۲-۶-۱- فرایند پاداش مارکوف

برای معرفی پاداش، باید مدل فرایند مارکوف خود را کمی گسترش دهیم. اول، ما باید برای انتقال خود از حالتی به حالت دیگر ارزش بیفزاییم. ما در حال حاضر احتمال داریم، اما از احتمال برای ثبت دینامیک سیستم استفاده می‌شود، بنابراین اکنون یک عدد اسکالر اضافی بدون بار اضافی داریم.

<sup>1</sup> episodes

<sup>2</sup> stationarity

<sup>3</sup> factor

<sup>4</sup> formalism

پاداش را می‌توان به اشکال مختلف نشان داد. کلی‌ترین راه این است که یک ماتریس مربع دیگر، شبیه به ماتریس انتقال، با پاداش برای انتقال از حالت  $i$  به حالت  $j$  که در ردیف  $i$  و ستون  $j$  قرار دارد، داشته باشید.

همان‌طور که گفته شد، پاداش می‌تواند مثبت یا منفی، بزرگ یا کوچک باشد. در برخی موارد، این نمایش اضافی است و می‌توان آن را ساده کرد. به‌عنوان مثال، اگر برای رسیدن به حالتی بدون توجه به حالت قبلی، پاداش داده شود، می‌توانیم فقط جفت‌های **حالت** ← **پاداش** را نگه داریم که نمایش فشرده‌تری هستند. باین‌حال، این فقط در صورتی قابل‌اعمال است که ارزش پاداش صرفاً به حالت هدف بستگی داشته باشد که همیشه این‌طور نیست.

دومین موردی که ما به مدل اضافه می‌کنیم ضریب تخفیف  $\gamma$  (گاما) است که یک عدد واحد از ۰ تا ۱ (شامل) است. معنای این پس از تعریف ویژگی‌های اضافی فرایند پاداش مارکوف توضیح داده می‌شود.

همان‌طور که به یاد دارید، ما زنجیره‌ای از انتقال حالت را در یک فرایند مارکوف مشاهده می‌کنیم. این مورد همچنان در مورد فرایند پاداش مارکوف صادق است، اما برای هر انتقال، ما مقدار اضافی خود یعنی پاداش را داریم؛ بنابراین اکنون، همه مشاهدات ما دارای ارزش پاداشی هستند که به هر انتقال سیستم متصل است.

برای هر قسمت، **بازده**<sup>۲</sup> در زمان  $t$  را به‌عنوان کمیت تعریف می‌کنیم:

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

برای هر نقطه زمانی، **بازده** را به‌عنوان مجموع پاداش‌های بعدی محاسبه می‌کنیم، اما پاداش‌های دورتر در ضریب تخفیف  $\gamma$  افزایش‌یافته به توان تعداد قدم‌آهایی که از نقطه شروع در  $t$  فاصله داریم (همان  $k$ ) ضرب می‌شوند. ضریب تخفیف به معنای آینده‌نگری عامل است. اگر گاما ( $\gamma$ ) برابر با ۱ باشد، بازده،  $G_t$ ، فقط برابر است با مجموع تمام پاداش‌های بعدی و مربوط به عاملی است که هر پاداش بعدی را کاملاً مشاهده می‌کند. اگر گاما برابر با ۰ باشد،  $G_t$  فقط یک پاداش فوری و بدون هیچ حالت بعدی خواهد بود و با کوتاه بینی<sup>۴</sup> (عدم مشاهده) مطلق مطابقت دارد.

این مقادیر مغرط<sup>۵</sup> فقط در موارد گوشه‌ای مفید هستند و اغلب اوقات، گاما روی چیزی در این بین تنظیم می‌شود، مانند ۰.۹ یا ۰.۹۹. در این مقادیر، ما به پاداش‌های آینده نگاه خواهیم کرد، اما نه خیلی دور. مقدار  $\gamma = 1$  ممکن است در موقعیت‌های قسمت‌های محدود کوتاه، قابل استفاده باشد. پارامتر گاما در یادگیری تقویتی مهم است و

<sup>1</sup> discount factor

<sup>2</sup> return

<sup>3</sup> steps

<sup>4</sup> short-sightedness

<sup>5</sup> extreme

در حال حاضر، در مورد آن به عنوان معیاری برای تخمین بازده آینده در نظر بگیرید. هر چه به عدد ۱ نزدیکتر باشد، قدم های بیشتری را در پیش خواهیم داشت.

این مقدار بازده در عمل چندان مفید نیست، زیرا برای هر زنجیره خاصی که از فرایند پاداش مارکوف خود مشاهده کردیم، تعریف شده است، بنابراین می تواند به طور گسترده ای متفاوت باشد، حتی برای همان حالت. باین حال، اگر ما امید ریاضی<sup>۱</sup> بازده را برای هر حالت محاسبه کنیم (با میانگین تعداد زیادی زنجیره)، مقدار بسیار مفیدتری به دست خواهیم آورد که به آن مقدار (ارزش) حالت می گویند.

$$V(s) = \mathbb{E}[G | S_t = s]$$

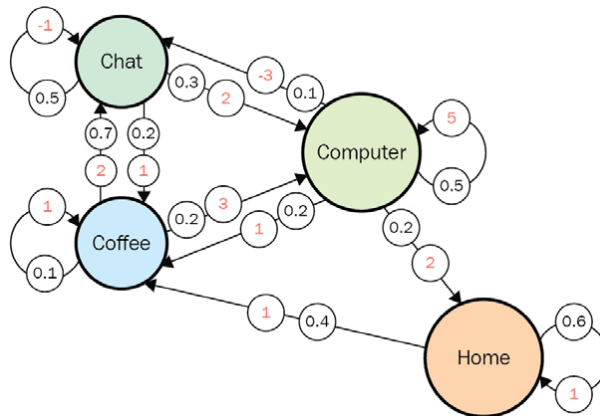
برای هر حالت،  $s$ ، مقدار  $V(s)$ ، میانگین (مورد انتظار) بازدهی است که با دنبال کردن فرایند پاداش مارکوف به دست می آوریم.

برای نشان دادن این موارد نظری در عمل، بیایید فرایند کارمند اداره را با پاداش گسترش دهیم و آن را به فرایند پاداش تبدیل کنیم. مقادیر پاداش ما به شرح زیر خواهد بود:

- خانه ← خانه: ۱ (چون خانه بودن خوب است)
- خانه ← قهوه: ۱
- کامپیوتر ← کامپیوتر: ۵ (سخت کار کردن چیز خوبی است)
- کامپیوتر ← گفتگو: ۳- (پرت شدن حواس خوب نیست)
- گفتگو ← کامپیوتر: ۲
- کامپیوتر ← قهوه: ۱
- قهوه ← کامپیوتر: ۳
- قهوه ← قهوه: ۱
- قهوه ← گفتگو: ۲
- گفتگو ← قهوه: ۱
- گفتگو ← گفتگو: ۱- (مکالمات طولانی خسته کننده می شود)

<sup>1</sup> mathematical expectation

نموداری از این مقادیر در اینجا نشان داده شده است:



شکل ۸- نمودار انتقال حالت با احتمالات انتقال (مشکی) و پاداش (روشن)

بیایید به پارامتر گامای خود برگردیم و در مورد مقادیر (ارزش) حالت‌هایی با مقادیر مختلف گاما فکر کنیم. ما با یک حالت ساده شروع می‌کنیم:  $\gamma = 0$ . چگونه مقادیر حالت‌ها را در اینجا محاسبه می‌کنید؟ برای پاسخ به این سوال، بیایید حالت خود را در **گفتگو** درست کنیم. انتقال بعدی چه می‌تواند باشد؟ پاسخ این است که بستگی به شانس دارد. با توجه به ماتریس انتقال ما برای فرایند مارکوف، احتمال ۵۰ درصد این است که حالت بعدی دوباره **گفتگو**، ۲۰ درصد **قهوه** و ۳۰ درصد **رایانه** باشد. وقتی  $\gamma = 0$  باشد، بازده ما فقط با مقداری از حالت فوری بعدی برابر است؛ بنابراین، اگر بخواهیم مقدار حالت **گفتگو** را محاسبه کنیم، باید همه مقادیر انتقال را جمع کنیم و آن را در احتمالات آنها ضرب کنیم:

$$V(\text{گفتگو}) = -1 \times 0.5 + 2 \times 0.3 + 1 \times 0.2 = 0.3$$

$$V(\text{قهوه}) = 2 \times 0.7 + 1 \times 0.1 + 3 \times 0.2 = 2.1$$

$$V(\text{خانه}) = 1 \times 0.6 + 1 \times 0.4 = 1.0$$

$$V(\text{رایانه}) = 5 \times 0.5 + (-3) \times 0.1 + 1 \times 0.2 + 2 \times 0.2 = 2.8$$

بنابراین، ریلنه با ارزش‌ترین حالتی است که می‌توان در آن قرار گرفت (اگر فقط به پاداش فوری اهمیت دهیم) که جای تعجب نیست؛ زیرا رایانه ← رایانه مکرر است، پاداش زیادی دارد و نسبت وقفه‌ها خیلی زیاد نیست.

اکنون یک سوال پیچیده‌تر، ارزش زمانی که  $\gamma = 1$  است چیست؟ در این مورد با دقت فکر کنید. پاسخ این است که مقدار برای همه حالت‌ها بی‌نهایت است. نمودار ما شامل حالات فرورفتگی (حالت‌های بدون انتقال خروجی) نیست، و وقتی تخفیف ما برابر با ۱ باشد، ما به تعداد بالقوه بی‌نهایتی از انتقال‌ها در آینده اهمیت



می‌دهیم. همان‌طور که در مورد گاما  $\gamma = 0$  مشاهده کردید، تمام مقادیر ما در کوتاه مدت مثبت هستند، بنابراین مجموع تعداد نامتناهی مقادیر مثبت، صرف‌نظر از حالت شروع، یک مقدار بی‌نهایت به ما می‌دهد.

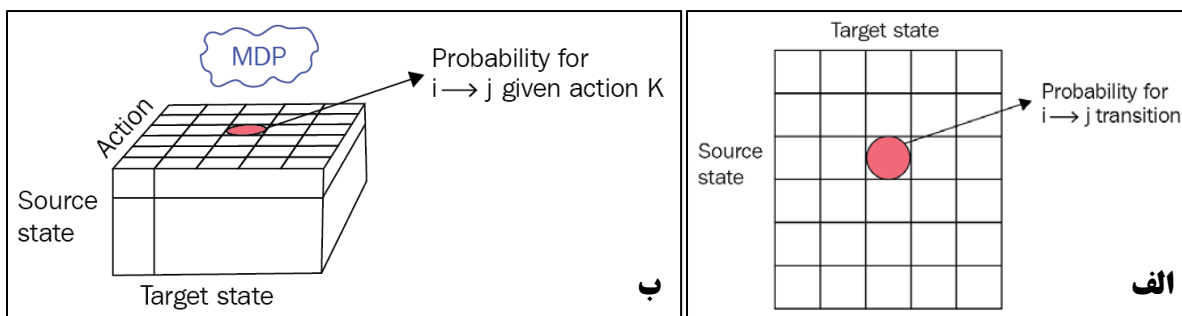
این نتیجه بی‌نهایت یکی از دلایلی را به ما نشان می‌دهد که به‌جای جمع کردن تمام پاداش‌های آینده، گاما را در فرایند پاداش مارکوف وارد کنیم. در بیشتر موارد، فرایند می‌تواند بی‌نهایت (زیاد) انتقال داشته باشد. از آنجایی که پرداختن به مقادیر نامتناهی چندان عملی نیست، می‌خواهیم افقی را که مقادیر را برای آن محاسبه می‌کنیم محدود کنیم. گامای با مقدار کمتر از ۱ چنین محدودیتی را فراهم می‌کند. از طرف دیگر، اگر با محیط‌های افق محدود سر و کار دارید (مثلاً بازی tic-tac-toe که حداکثر با ۹ مرحله محدود شده است)، استفاده از گاما  $\gamma = 1$  خوب است.

در حال حاضر، اجازه دهید لایه دیگری از پیچیدگی را در اطراف فرایندهای پاداش مارکوف خود قرار دهیم و آخرین قطعه گمشده را معرفی کنیم: اقدامات.

### ۳-۶-۱- افزودن اقدامات

ممکن است در حال حاضر ایده‌هایی در مورد چگونگی گسترش فرایند پاداش مارکوف برای شامل اقدامات داشته باشید. ابتدا باید مجموعه‌ای از اعمال (A) را اضافه کنیم که باید متناهی باشد. این **فضای اقدامات** عامل ماست. ثانیاً، ما باید ماتریس انتقال خود را با اقدامات شرطی کنیم که اساساً به این معنی است که ماتریس ما به یک بعد اقدام اضافی نیاز دارد که آن را به یک مکعب تبدیل می‌کند.

اگر به‌خاطر داشته باشید، در مورد فرایند مارکوف و فرایندهای مارکوف، ماتریس انتقال شکل مربعی داشت (شکل ۹) که حالت مبدأ در ردیف‌ها و حالت هدف در ستون‌ها قرار داشت؛ بنابراین، هر ردیف  $i$  حاوی لیستی از احتمالات برای پرش به هر حالت بود.



شکل ۹- (الف) ماتریس انتقال به شکل مربع و (ب) احتمالات انتقال برای فرایند تصمیم‌گیری مارکوف

اکنون عامل دیگر به صورت غیرفعال انتقال حالت را مشاهده نمی‌کند، بلکه می‌تواند به طور فعال اقدامی را برای انجام در هر انتقال حالت انتخاب کند؛ بنابراین، برای هر حالت منبع، ما فهرستی از اعداد نداریم، اما ماتریسی داریم که در آن بعد عمق شامل اقداماتی است که عامل می‌تواند انجام دهد و بعد دیگر آن چیزی است که سیستم حالت هدف پس از انجام اقدامات توسط عامل به آن خواهد پدید. نمودار زیر جدول انتقال جدید ما را نشان می‌دهد که به یک مکعب با حالت مبدأ به عنوان بعد ارتفاع (i) و حالت هدف به عنوان عرض (j) تبدیل شد و عملی که عامل می‌تواند به عنوان عمق (k) جدول انتقال انجام دهد.

بنابراین، به طور کلی با انتخاب یک عمل عامل می‌تواند بر احتمالات حالت‌های هدف تأثیر بگذارد که این یک توانایی مفید است.

در نهایت، برای تبدیل فرایند پاداش مارکوف به یک فرایند تصمیم‌گیری مارکوف، باید اقداماتی را به ماتریس پاداش خود اضافه کنیم، همان طور که با ماتریس انتقال انجام دادیم. ماتریس پاداش ما نه تنها به حالت بلکه به عمل نیز بستگی دارد. به عبارت دیگر، پاداشی که عامل به دست می‌آورد، نه تنها به حالتی که در نهایت به آن می‌رسد، بلکه به عملی که منجر به این حالت می‌شود نیز بستگی دارد.

#### ۴-۶-۱ - خطمشی (سیاست)

ممکن است به یاد داشته باشید که هدف اصلی عامل در یادگیری تقویتی این است که تاحدامکان بازدهی بیشتری را جمع‌آوری کند؛ بنابراین، دوباره، سیاست‌های مختلف می‌توانند مقادیر متفاوتی بازدهی را به ما بدهند که پیدا کردن یک سیاست خوب را مهم می‌کند. به همین دلیل است که مفهوم سیاست مهم است.

به طور رسمی، خطمشی به عنوان توزیع احتمال بر روی اقدامات برای هر حالت ممکن تعریف می‌شود:

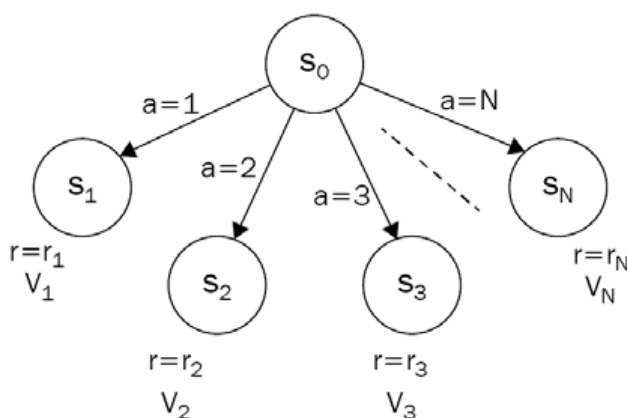
$$\pi(a|s) = P[A_t = a | S_t = s]$$

این به عنوان احتمال تعریف می‌شود و نه به عنوان یک اقدام مشخص برای وارد کردن تصادفی بودن به رفتار (رفتار تصادفی) یک عامل. خطمشی قطعی یک مورد خاص از احتمال است که احتمال عمل مورد نیاز برابر با ۱ است.

مفهوم مفید دیگر این است که اگر خطمشی ما ثابت باشد و تغییر نکند، فرایند تصمیم‌گیری مارکوف ما به یک فرایند پاداش مارکوف تبدیل می‌شود، زیرا می‌توانیم ماتریس‌های انتقال و پاداش را با احتمالات یک سیاست کاهش دهیم و از ابعاد عمل خلاص شویم.

## ۵-۶-۱- معادله بهینگی بلمن

برای توضیح معادله بلمن<sup>۱</sup> بهتر است کمی انتزاعی پیش برویم (شکل ۱۰). بیایید با یک مورد قطعی<sup>۲</sup> شروع کنیم، زمانی که همه اقدامات ما یک نتیجه ۱۰۰٪ تضمین شده دارند. تصور کنید که عامل ما حالت  $s_0$  را مشاهده می‌کند و  $N$  عمل در دسترس دارد. هر عملی به حالت دیگری منتهی می‌شود،  $s_1 \dots s_N$ ، با پاداش مربوطه،  $r_1 \dots r_N$ . همچنین، فرض کنید که مقادیر  $V_i$ ، همه حالت‌های متصل به حالت  $s_0$  را می‌دانیم. بهترین اقدامی که عامل می‌تواند در چنین حالتی انجام دهد خواهد بود؟



شکل ۱۰- یک محیط انتزاعی با  $N$  حالت قابل دسترسی از حالت اولیه

اگر عمل مشخص،  $a_i$  را انتخاب کنیم و مقدار داده شده به این عمل را محاسبه کنیم، مقدار  $V_0(a = a_i) = r_i + V_i$  خواهد بود. بنابراین، برای انتخاب بهترین اقدام ممکن، عامل باید مقادیر حاصل از هر اقدام را محاسبه کند و حداکثر نتیجه ممکن را انتخاب کند. به عبارت دیگر،  $V_0 = \max_{a \in 1 \dots N} (r_a + V_a)$ . اگر از ضریب تخفیف  $\gamma$  استفاده می‌کنیم، باید مقدار حالت بعدی را در گاما ضرب کنیم:

$$V_0 = \max_{a \in 1 \dots N} (r_a + \gamma V_a)$$

این ممکن است بسیار شبیه به روش‌های حریصانه<sup>۳</sup> باشد و در واقع همین‌طور است. با این حال، یک تفاوت وجود دارد: در این مورد وقتی حریصانه عمل می‌کنیم، فقط به پاداش آنی عمل نگاه نمی‌کنیم، بلکه به پاداش فوری به اضافه ارزش بلندمدت دولت نگاه می‌کنیم.

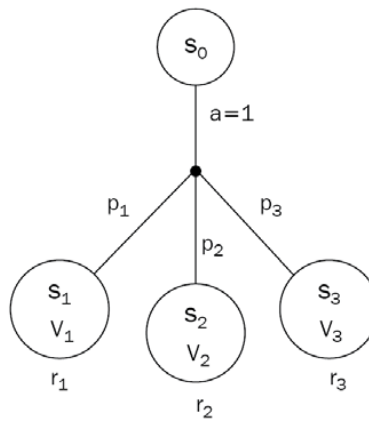
<sup>۱</sup> Bellman equation

<sup>۲</sup> deterministic

<sup>۳</sup> greedy

بلمن ثلثت کرد که با این بسط، رفتار ما بهترین نتیجه ممکن را خواهد گرفت. به عبارت دیگر بهینه خواهد بود؛ بنابراین، معادله قبلی معادله ارزش بلمن نامیده می‌شود (برای یک حالت قطعی).

گسترش این ایده برای یک مورد تصادفی<sup>۱</sup>، زمانی که اقدامات ما این شانس را دارند که به حالت‌های مختلف ختم شوند، چندان پیچیده نیست. کاری که ما باید انجام دهیم این است که به جای اینکه فقط مقدار حالت بعدی را بگیریم، مقدار مورد انتظار را برای هر اقدام محاسبه کنیم. برای نشان دادن این موضوع، بیایید یک اقدام واحد را که از حالت  $s_0$  در دسترس است، با سه نتیجه ممکن در نظر بگیریم.



شکل ۱۱- مثالی از انتقال از حالت در یک مورد تصادفی

در اینجا ما یک عمل داریم که می‌تواند منجر به سه حالت مختلف با احتمالات متفاوت شود. با احتمال  $p_1$ ، عمل می‌تواند به حالت  $s_1$ ، با  $p_2$  در حالت  $s_2$ ، و با  $p_3$  در حالت  $s_3$  ختم شود (البته  $p_1 + p_2 + p_3 = 1$ ). هر حالت هدف پاداش مخصوص به خود را دارد ( $r_1$  یا  $r_2$  یا  $r_3$ ). برای محاسبه مقدار مورد انتظار پس از صدور عمل ۱، باید همه مقادیر را ضرب در احتمالات آنها جمع کنیم:

$$V_0(a=1) = p_1(r_1 + \gamma V_1) + p_2(r_2 + \gamma V_2) + p_3(r_3 + \gamma V_3)$$

یا به صورت رسمی‌تر:

$$V_0(a) = \mathbb{E}_{s \sim S}[r_{s,a} + \gamma V_s] = \sum_{s \in S} p_{a,0 \rightarrow s}(r_{s,a} + \gamma V_s)$$

<sup>1</sup> stochastic

با ترکیب معادله بلمن، برای یک حالت قطعی، با مقداری برای اقدامات تصادفی، معادله بهینگی بلمن برای یک حالت کلی به دست می‌آید:

$$V_0 = \max_{a \in A} \mathbb{E}_{s \sim S} [r_{s,a} + \gamma V_s] = \max_{a \in A} \sum_{s \in S} p_{a,0 \rightarrow s} (r_{s,a} + \gamma V_s)$$

توجه داشته باشید که  $p_{a,i \rightarrow j}$  به معنای احتمال عمل  $a$  است که در حالت  $i$  صادر شده و به حالت  $j$  ختم می‌شود. مقدار بهینه حالت برابر است با عمل که حداکثر پاداش احتمالی فوری مورد انتظار را به ما می‌دهد، به علاوه پاداش بلندمدت با تخفیف برای حالت بعدی. همچنین ممکن است متوجه شوید که این تعریف بازگشتی است: مقدار (ارزش) حالت، از طریق مقادیر حالت‌های بلافاصله قابل دسترسی تعریف می‌شود.

این ارزش‌ها نه تنها بهترین پاداشی را که می‌توانیم به دست آوریم، به ما می‌دهند، بلکه اساساً خط‌مشی بهینه را برای به دست آوردن آن پاداش به ما می‌دهند: اگر عامل ما ارزش هر حالت را بداند، به طور خودکار می‌داند چگونه همه این پاداش را جمع‌آوری کند. به لطف اثبات بهینه‌سازی بلمن، در هر حالتی که عامل به آن می‌رسد، باید اقدامی را با حداکثر پاداش مورد انتظار انتخاب کند که مجموع پاداش فوری و پاداش بلندمدت با تخفیف یک مرحله‌ای است؛ بنابراین، دانستن این مقادیر واقعاً مفید است.

## ۲- الگوریتم‌های پایه یادگیری تقویتی

در این بخش به بررسی الگوریتم‌های پایه در یادگیری تقویتی خواهیم پرداخت که محیط‌های مختلف بسیار پر کاربرد هستند و الگوریتم‌های پیشرفته بر پایه آن‌ها شکل گرفته‌اند. چندین الگوریتم و روش تصمیم‌گیری در یادگیری تقویتی وجود دارد که تفاوت آن‌ها عمدتاً به دلیل استراتژی‌های مختلفی است که برای کشف محیط خود استفاده می‌کنند. در ادامه به برخی از پرکاربردترین روش‌های تصمیم‌گیری در یادگیری تقویتی اشاره خواهیم کرد.

### روش‌های مبتنی بر مدل:

- ساخت مدلی از محیط
- برنامه‌ریزی<sup>۱</sup> با استفاده از مدل (مثلاً با نگاه پیش رو<sup>۲</sup>).

### روش‌های مبتنی بر ارزش:

- تخمین تابع مقدار بهینه  $Q^*(s, a)$
- این حداکثر مقدار قابل دستیابی تحت هر سیاستی است

### روش‌های مبتنی بر سیاست:

- مستقیماً سیاست بهینه  $\pi^*$  را جستجو می‌کند
- این سیاست دستیابی به حداکثر پاداش آینده است

### روش‌های تقریب توابع:

- استفاده از شبکه‌های عصبی عمیق برای تقریب توابع ارزش
- استفاده از توابع پایه برای تقریب توابع

## ۱-۲- روش‌های مبتنی بر مدل (برنامه‌نویسی پویا)

یادگیری تقویتی مبتنی بر مدل مانند پسرعموی باهوش خانواده یادگیری تقویتی است؛ جایی که عامل به شکل کورکورانه به محیط واکنش نشان نمی‌دهد. در عوض با ایجاد نسخه کوچکی از خود تلاش دارد رویدادهای آینده را پیش‌بینی کند. این مدل حرکات مختلف را امتحان می‌کند تا ببیند در آینده چه اتفاقی می‌افتد. در این شیوه،

<sup>1</sup> Plan

<sup>2</sup> lookahead

عامل می‌تواند سناریوهایی را در مدل خود اجرا کند، بدون آنکه نتیجه منفی بگیرد؛ درست مانند چیدن استراتژی حرکت در بازی شطرنج پیش از لمس مهره.

اصطلاح برنامه‌نویسی پویا<sup>۱</sup> به مجموعه‌ای از الگوریتم‌ها اشاره دارد که می‌توانند برای محاسبه سیاست‌های بهینه باتوجه به یک مدل کامل از محیط به‌عنوان فرایند تصمیم مارکوف استفاده شوند. الگوریتم‌های برنامه‌نویسی پویا کلاسیک هم به دلیل فرض یک مدل کامل و هم به دلیل هزینه محاسباتی زیادشان در یادگیری تقویتی کاربرد محدودی دارند، اما از نظر تئوری همچنان مهم هستند. برنامه‌نویسی پویا یک پایه اساسی برای درک اکثر روش‌های ارائه شده در یادگیری تقویتی را فراهم می‌کند. در واقع، همه این روش‌ها را می‌توان به‌عنوان تلاش‌هایی برای دستیابی به نتایج مشابه برنامه‌نویسی پویا، تنها با محاسبات کمتر و بدون فرض یک مدل کامل از محیط مشاهده کرد.

- ایده کلیدی DP استفاده از توابع ارزش برای سازماندهی و ساختار جستجوی سیاست‌های خوب است
- هنگامی که توابع مقدار بهینه را یافتیم که معادلات بهینه بلمن را برآورده می‌کند، به راحتی می‌توانیم سیاست‌های بهینه را به دست آوریم:

$$V^*(s) = \max_a \mathbb{E} [r_{t+1} + \gamma V^*(s_{t+1}) | s_t = s, a_t = a] = \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^*(s')].$$

- ارزش اقدام a در حالت s تحت یک خط‌مشی  $\pi$ ، بازده مورد انتظار هنگام شروع در s، انجام اقدام a و دنبال کردن  $\pi$  پس از آن است.

$$Q^*(s, a) = \mathbb{E} [r_{t+1} + \gamma \max_{a'} Q^*(s_{t+1}, a') | s_t = s, a_t = a] = \sum_{s'} P_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma \max_{a'} Q^*(s', a')].$$

که در آن  $V^*(s)$ : تابع مقدار بهینه در حالت s. s: متغیر حالت. a: متغیر اقدام/کنترل.  $x_{t+1}$ : پاداش در زمان t+1.  $\gamma$ : ضریب تخفیف.  $s_t$ : حالت در زمان t.  $a_t$ : اقدام در زمان t. E: عملگر امید ریاضی.  $P_{ss'}^a$ : احتمال انتقال از حالت s به حالت s تحت عمل a.  $R_{ss'}^a$ : پاداش مرتبط با انتقال از حالت s به حالت s تحت عمل a.

### تکرار سیاست:

- تکرار خط‌مشی یک فرایند تکراری است

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi^* \xrightarrow{E} V^*$$

<sup>۱</sup> dynamic programming

- تکرار خطمشی دو مرحله دارد: ارزیابی و بهبود سیاست.
- در ارزیابی خطمشی، ما توابع مقدار حالت یا حالت-عمل را محاسبه می‌کنیم:

$$V^\pi(s) = \mathbb{E}_\pi[R_t | s_t = s] = \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s\right] = \sum_{\pi} \pi(s, a) \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^\pi(s')]$$

- در بهبود خطمشی، سیاست را تغییر می‌دهیم تا خطمشی بهتری به دست آوریم

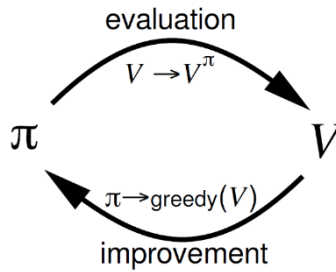
$$\pi'(s) = \operatorname{argmax}_a Q^\pi(s, a) = \operatorname{argmax}_a \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V^\pi(s')].$$

### ارزش و تکرار خطمشی تعمیم‌یافته:

- در تکرار ارزش داریم:

$$V_{k+1}(s) = \max_a \mathbb{E}[r_{t+1} + \gamma V_k(s_{t+1}) | s_t = s, a_t = a] = \max_a \sum_{s'} P_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V_k(s')]$$

- تکرار سیاست تعمیم‌یافته

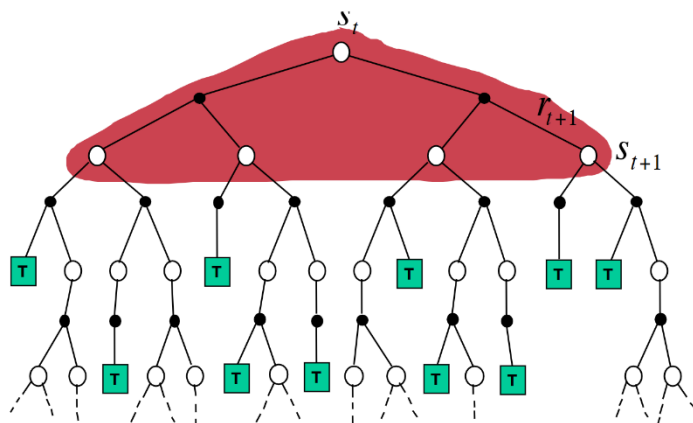


شکل ۱۲ - چرخه ارزیابی و بهبود در برنامه‌نویسی پویا



نمودار پشتیبان<sup>۱</sup>:

$$V(S_t) \leq \mathbb{E}_{\pi} [R_{t+1} + \gamma V(S_{t+1})]$$



شکل ۱۳- نمودار پشتیبان روش برنامه‌نویسی پویا

## ۲-۲- روش‌های مبتنی بر ارزش (مونت کارلو و یادگیری تفاوت زمانی)

در این روش، تمرکز بر یادگیری تابع مقدار بهینه است که با  $Q^*$  یا  $V^*$  نشان داده می‌شود. تابع ارزش پاداش‌های تجمعی آینده مورد انتظار را برای قرارگرفتن در یک حالت معین و پیروی از سیاست فعلی پس از آن تخمین می‌زند.

در رویکرد مبتنی بر ارزش، هدف بهینه‌سازی تابع ارزش  $V(s)$  است. تابع ارزش، تابعی است که پاداش بیشینه آینده را مشخص می‌کند که عامل در هر حالت دریافت می‌کند. ارزش هر حالت برابر است با ارزش کل پاداشی که عامل می‌تواند انتظار داشته باشد در آینده با آغاز از آن حالت جمع‌آوری کند.

$$v_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s]$$

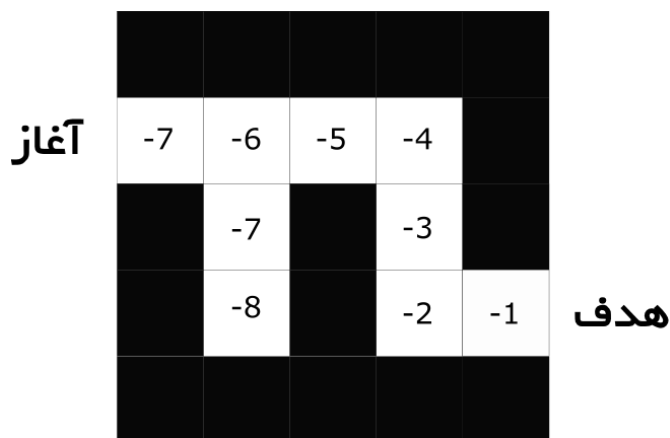
مورد انتظار

پاداش  
تنزیل یافته

با توجه به حالت

عامل از این تابع ارزش برای انتخاب آنکه کدام حالت در هر گام انتخاب شود، استفاده می‌کند. عامل حالتی با بیشترین ارزش را انتخاب می‌کند.

<sup>۱</sup> نمودار پشتیبان (backup diagram) یک نمایش بصری است که در یادگیری تقویتی برای نشان دادن نحوه انتشار یا "پشتیبان گیری" مقادیر از حالت‌های آینده به حالت فعلی استفاده می‌شود.



شکل ۱۴- محاسبه ارزش همه حالت‌ها در مثال هزارتو

در مثال هزارتو، در هر گام بیشترین ارزش، یعنی ۷-، ۶- و سپس ۵- (و به همین ترتیب) برای رسیدن به هدف اتخاذ می‌شود.

- این روش‌ها به طور ضمنی بر تابع سیاست تکیه دارند.
- این روش‌ها ابتدا تابع مقدار  $Q(s, a)$  را یاد می‌گیرند.
- سپس خط‌مشی  $\pi(s, a)$  را از  $Q(s, a)$  استنتاج کنید.
- نمونه‌ها:

۱. روش‌های مونت کارلو<sup>۱</sup>

۲. روش‌های تفاوت زمانی<sup>۲</sup>

۳. Q-learning

۴. SARSA

## ۱-۲-۲- روش‌های مونت کارلو

روش مونت کارلو برخلاف روش برنامه‌ریزی پویا که نیازمند اطلاعات کامل در مورد توزیع احتمالات همه انتقال‌های ممکن بود؛ تنها به نمونه‌ای از توزیع‌های احتمال بسنده می‌کند. به بیانی دیگر، در روش مونت کارلو شناخت کامل محیط لازم نیست و با برقراری تعامل واقعی یا شبیه‌سازی شده با یک محیط می‌توان به توالی نمونه‌ای از حالت‌ها، اقدامات و پاداش‌های دست‌یافت. به همین دلیل است که در این روش، پاداش‌ها در انتهای دوره حساب می‌شود تا بتوان از دانش کسب شده برای دوره جدید استفاده نمود.

<sup>1</sup> Monte Carlo

<sup>2</sup> Temporal-difference

- روش‌های مونت کارلو به طور مستقیم از قسمت (اپیزود) های تجربه یاد می‌گیرند.
- مونت کارلو بدون مدل است: بدون آگاهی از انتقال / پاداش MDP
- مونت کارلو از قسمت‌های کامل درس می‌گیرد
- مونت کارلو از ساده‌ترین ایده ممکن استفاده می‌کند: ارزش = میانگین بازده
- هدف مونت کارلو: یادگیری  $V_\pi$  از قسمت‌های تجربه تحت خط‌مشی  $\pi$

$$[S_1 \xrightarrow{R_1} S_2 \xrightarrow{R_2} S_3 \xrightarrow{R_3} S_4 \cdots \xrightarrow{R_{k-1}} S_k]$$

- بازده کل پاداش با تخفیف است:

$$G_t = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{T-1} R_T$$

- تابع مقدار بازده مورد انتظار است:

$$V_\pi(s) = \mathbb{E}_\pi [G_t | S_t = s]$$

- ارزیابی خط‌مشی مونت کارلو از میانگین بازده تجربی به جای بازده مورد انتظار استفاده می‌کند.

### ارزیابی سیاست هر بازدید از مونت کارلو:

- برای ارزیابی حالت  $s$
- هر مرحله زمانی  $t$  که حالت  $s$  را در یک قسمت مشاهده می‌کند، شمارنده افزایش می‌یابد

$$N(s) \leftarrow N(s) + 1$$

- بازده کل را افزایش می‌دهد:

$$S(s) \leftarrow S(s) + G_t$$

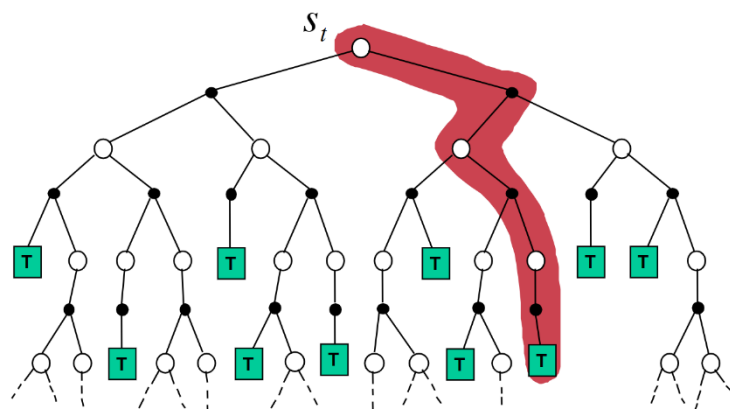
- ارزش با بازده متوسط تخمین زده می‌شود

$$V(s) = \frac{S(s)}{N(s)}$$

- طبق قانون اعداد زیاد،  $V(s) \rightarrow V_\pi(s)$  به عنوان  $N(s) \rightarrow \infty$

## نمودار پشتیبان:

$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$$



شکل ۱۵- نمودار پشتیبان روش مونت کارلو

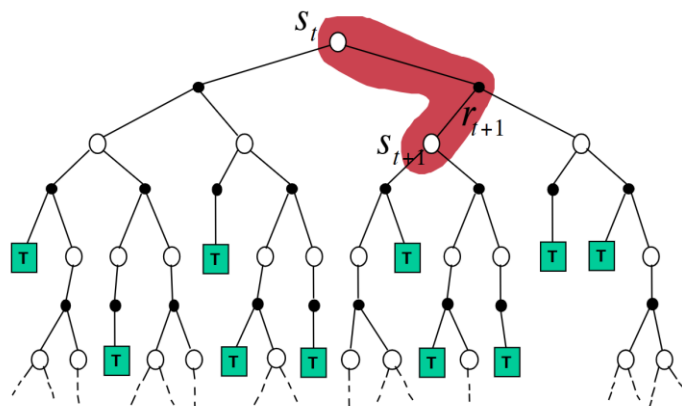
## ۲-۲-۲- روش‌های تفاوت زمانی

روش‌های تفاوت زمانی (TD) یک نوع روش پیش‌بینی است که از ایده‌های مونت کارلو و برنامه‌نویسی پویا برای حل مسائل تقویتی استفاده می‌کند. روش‌های تفاوت زمانی مانند روش‌های مونت کارلو، مستقیماً از تجربه خام بدون نیاز به مدلی پویا از محیط، می‌آموزند و مانند برنامه‌ریزی پویا، تخمین‌ها را تا حدی بر اساس سایر تخمین‌های آموخته شده به‌روز می‌کنند، بدون اینکه منتظر نتیجه نهایی باشند.

- تفاوت زمانی ترکیبی از ایده‌های مونت کارلو و ایده‌های برنامه‌نویسی پویا است.
- مانند روش‌های مونت کارلو، روش‌های تفاوت زمانی می‌توانند مستقیماً از تجربه خام بدون مدلی از پویایی محیط بیاموزند.
- مانند برنامه‌نویسی پویا، روش‌های تفاوت زمانی تخمین‌ها را تا حدی بر اساس سایر تخمین‌های آموخته شده به‌روزرسانی می‌کنند، بدون اینکه منتظر نتیجه نهایی باشند (آنها راه‌اندازی<sup>۱</sup> می‌کنند).
- روش‌های مونت کارلو منتظر می‌مانند تا بازده پس از بازدید مشخص شود، سپس از آن بازده به‌عنوان یک هدف برای  $V(s_t)$  استفاده می‌کند درحالی‌که روش‌های تفاوت زمانی فقط تا مرحله زمانی بعدی باید منتظر بمانند.
- ساده‌ترین روش تفاوت زمانی که با نام TD (۰) شناخته می‌شود، است

$$V(s_t) \leftarrow V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

<sup>۱</sup> bootstrap



شکل ۱۶- بردار پشتیبان روش تفاوت زمانی

الگوریتم TD(0):

```

Initialize V(s) arbitrarily,  $\pi$  to the policy to be evaluated

Repeat (for each episode):
  Initialize s

  Repeat (for each step of episode):
    a  $\leftarrow$  action given by  $\pi$  for s
    Take action a; observe reward, r, and next state, s'
     $V(s) \leftarrow V(s) + \alpha [r + \gamma V(s') - V(s)]$ 
    s  $\leftarrow$  s'
  until s is terminal
  
```

### ۳-۲-۲- SARSA

اکنون به استفاده از روش‌های پیش‌بینی TD برای مسئله کنترل می‌پردازیم. طبق معمول، ما از الگوی تکرار خط‌مشی تعمیم‌یافته<sup>۱</sup> پیروی می‌کنیم، فقط این بار از روش‌های TD برای بخش ارزیابی یا پیش‌بینی استفاده می‌کنیم. همانند روش‌های مونت کارلو، ما با نیاز به موازنه اکتشاف و بهره‌برداری مواجه هستیم، و دوباره رویکردها به دو دسته اصلی تقسیم می‌شوند: مبتنی بر سیاست<sup>۲</sup> و خارج از سیاست<sup>۳</sup>. در این بخش ما یک روش کنترل TD مبتنی بر سیاست را ارائه می‌کنیم.

اولین گام این است که به‌جای یک تابع ارزش-حالت، یک تابع ارزش-عمل را یاد بگیریم. به طور خاص، برای یک روش مبتنی بر سیاست، باید  $q_{\pi}(s, a)$  را برای خط‌مشی رفتار فعلی  $\pi$  و برای همه حالت‌ها  $s$  و اقدامات  $a$

<sup>۱</sup> generalized policy iteration

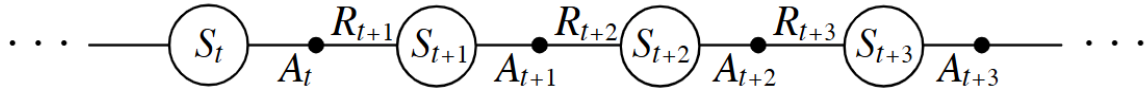
<sup>۲</sup> on-policy

<sup>۳</sup> off-policy

تخمین بزنیم. این را می‌توان با استفاده از همان روش TD که در زیر برای یادگیری  $v_\pi$  توضیح داده شد انجام داد:

$$v_\pi(s) \stackrel{\text{def}}{=} \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi[R_t + 1 + \gamma G_{t+1} | S_t = s] = \mathbb{E}_\pi[R_t + 1 + \gamma v_\pi(S_{t+1}) | S_t = s]$$

به یاد بیاورید که یک قسمت (اپیزود) از یک توالی متناوب از حالت‌ها و جفت‌های حالت-عمل تشکیل شده است:



در بخش قبل انتقال از حالتی به حالت دیگر را در نظر گرفتیم و مقادیر حالت‌ها را یاد گرفتیم. اکنون انتقال از جفت حالت-عمل به جفت حالت-عمل را در نظر می‌گیریم و مقادیر جفت حالت-عمل را یاد می‌گیریم. به طور رسمی، این موارد یکسان هستند: هر دو زنجیره مارکوف با فرایند پاداش هستند. قضایایی که همگرایی مقادیر حالت را تحت TD(0) تضمین می‌کنند، در الگوریتم مربوطه برای مقادیر عمل نیز اعمال می‌شوند:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)].$$



این به‌روزرسانی پس از هر انتقال از حالت غیر پایانی  $S_t$  انجام می‌شود. اگر  $S_{t+1}$  ترمینال باشد،  $Q(S_{t+1}, A_{t+1})$  به‌عنوان صفر تعریف می‌شود. این قانون از هر عنصر از پنج‌گانه رویدادها  $(A_t, S_t)$  استفاده می‌کند که انتقال از یک جفت حالت-عمل به جفت بعدی را تشکیل می‌دهند. این پنج‌گانه باعث ایجاد نام Sarsa(State-Action-Reward-State-Action) برای

الگوریتم می‌شود. نمودار پشتیبان برای Sarsa در سمت چپ نشان داده شده است. Sarsa

طراحی یک الگوریتم کنترل روی خط‌مشی بر اساس روش پیش‌بینی Sarsa ساده است. مانند همه روش‌های مبتنی بر سیاست، ما به طور مداوم  $q_\pi$  را برای خط‌مشی رفتاری  $\pi$  تخمین می‌زنیم، و درعین حال  $\pi$  را نسبت به  $q_\pi$  به سمت حریصانه<sup>۱</sup> تغییر می‌دهیم. شکل کلی الگوریتم کنترل Sarsa در کادر صفحه بعد آورده شده است.

ویژگی‌های همگرایی<sup>۲</sup> الگوریتم Sarsa به ماهیت وابستگی خط‌مشی به  $Q$  بستگی دارد. برای مثال، می‌توان از سیاست‌های «طمع آمیز» ( $\epsilon$ -greedy) یا «نرم» ( $\epsilon$ -soft) استفاده کرد. Sarsa با احتمال ۱، به یک خط‌مشی بهینه و تابع ارزش-عمل همگرا می‌شود تا زمانی که همه جفت‌های حالت-عمل بی‌نهایت بار بازدید شوند و خط‌مشی در حد خط‌مشی حریصانه همگرا شود.

<sup>۱</sup> greediness

<sup>۲</sup> convergence properties

الگوریتم Sarsa (کنترل TD مبتنی بر سیاست) برای تخمین  $Q \approx q^*$ :

```

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\epsilon > 0$ 
Initialize  $Q(s, a)$  for all  $s \in S^+$ ,  $a \in A(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$ 

Loop for each episode:
  Initialize  $S$ 
  Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
  Loop for each step of episode:
    Take action  $A$ , observe  $R, S'$ 
    Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$ 
     $S \leftarrow S'; A \leftarrow A'$ 
  until  $S$  is terminal
  
```

#### ۴-۲-۲ - Q-learning

یکی از پیشرفت‌های اولیه در یادگیری تقویتی، توسعه یک الگوریتم کنترل TD خارج از سیاست (off-policy) معروف به یادگیری  $Q$  بود که به صورت زیر تعریف می‌شود:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

الگوریتم یادگیری  $Q$  شباهت زیادی با الگوریتم یادگیری SARSA دارد. تفاوت کلیدی این دو الگوریتم در این است که SARSA برخلاف الگوریتم یادگیری  $Q$ ، در دسته الگوریتم‌های مبتنی بر سیاست قرار می‌گیرد؛ بنابراین، الگوریتم SARSA مقدار  $Q$ -value را با توجه به اقدامی که ناشی از سیاست فعلی است محاسبه می‌کند نه اقدام ناشی از سیاست حریصانه.

در این مورد، تابع ارزش-عمل آموخته شده،  $Q$ ، به طور مستقیم به  $q^*$ ، تابع ارزش-عمل بهینه، مستقل از سیاست مورد استفاده، تقریب می‌زند. این به طور چشمگیری تجزیه و تحلیل الگوریتم را ساده می‌کند و اثبات همگرایی<sup>۱</sup> اولیه را فعال می‌کند. این خط‌مشی همچنان دارای اثری است که تعیین می‌کند کدام جفت‌های حالت-عمل بازدید و به روزرسانی می‌شوند. با این حال، تنها چیزی که برای همگرایی صحیح مورد نیاز است این است که همه جفت‌ها به روز شوند. این یک حداقل نیاز است به این معنا که هر روشی که تضمین شده است رفتار بهینه را در حالت کلی پیدا کند، باید آن را داشته باشد. با این فرض و نوعی از شرایط تقریب تصادفی معمول در دنباله پارامترهای اندازه-گام، نشان داده شده است که  $Q$  با احتمال ۱ به  $q^*$  همگرا می‌شود.

<sup>1</sup> convergence proofs

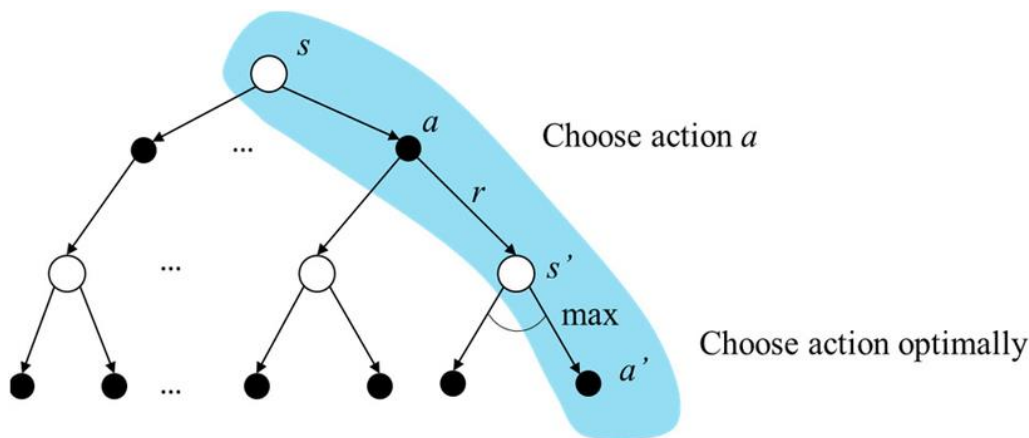
الگوریتم یادگیری Q (کنترل TD خارج از سیاست) برای تخمین  $\pi \approx \pi^*$ :

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\epsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in S^+$ ,  $a \in A(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

**Loop** for each episode:  
 Initialize  $S$   
**Loop** for each step of episode:  
   Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)  
   Take action  $A$ , observe  $R, S'$   
    $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$   
    $S \leftarrow S'$   
 until  $S$  is terminal

در مورد نمودار پشتیبان برای یادگیری Q، قانون یادگیری Q یک جفت حالت-عمل را به روز می‌کند، بنابراین گره بالایی، ریشه به روزرسانی، باید یک گره اقدام کوچک و پر باشد. به روزرسانی نیز از گره‌های اقدام است که بیش از همه آن اقدامات ممکن در حالت بعدی به حداکثر می‌رسد؛ بنابراین گره‌های پایینی نمودار پشتیبان باید تمام این گره‌های عمل باشند (شکل ۱۷).



شکل ۱۷-بردار پشتیبان روش تفاوت زمانی

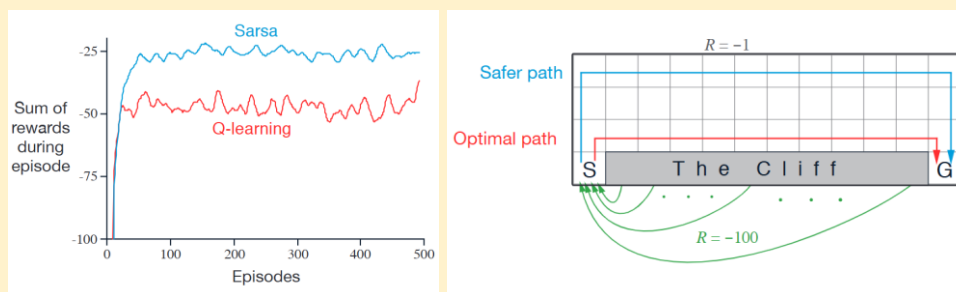
یکی از مزایای یادگیری خارج از سیاست این است که تابع مقدار بهینه با یک خط‌مشی قطعی را می‌توان در حین پیروی از یک خط‌مشی اکتشافی تصادفی یاد گرفت. اشکالات یادگیری خارج از سیاست این است که عملکرد در طول یادگیری می‌تواند با نادیده گرفتن تأثیر اکتشاف به خطر بیفتد و یادگیری می‌تواند در صورت ترکیب با یک تقریب تابع، ناپایدار باشد.



مثال پیاده روی صخره (Cliff Walking): این مثال دنیای شبکه‌ای Sarsa (gridworld) و یادگیری Q را مقایسه می‌کند و تفاوت بین روش‌های مبتنی بر سیاست (Sarsa) و روش‌های خارج از سیاست (یادگیری Q) را برجسته می‌کند.

دنیای شبکه‌ای را که در زیر نشان داده شده است در نظر بگیرید. این یک کار اپیزودیک و بدون تخفیف استاندارد است، با حالت‌های شروع (start) و هدف (goal)، و اعمال معمولی که باعث حرکت به سمت بالا، پایین، راست و چپ می‌شود. پاداش در همه انتقال‌ها ۱- است به جز مواردی که در منطقه با علامت "The Cliff" مشخص شده است. قدم گذاشتن در این منطقه دارای پاداش ۱۰۰- است و عامل را فوراً به شروع باز می‌گرداند.

نمودار زیر عملکرد روش‌های Sarsa و Q-learning را با انتخاب عمل  $\epsilon$ -greedy نشان می‌دهد،  $\epsilon = 0.1$ . پس از یک دوره گذرا اولیه، Q-learning مقادیر سیاست بهینه را می‌آموزد، چیزی که درست در امتداد لبه صخره حرکت می‌کند. متأسفانه، به دلیل انتخاب اکشن  $\epsilon$ -greedy، گاهی اوقات از صخره سقوط می‌کند. از سوی دیگر، سارسا، انتخاب عمل را در نظر می‌گیرد و مسیر طولانی‌تر اما امن‌تر را از طریق قسمت بالای شبکه یاد می‌گیرد. اگرچه Q-learning در واقع مقادیر خط مشی بهینه را می‌آموزد، عملکرد آن‌لاین آن بدتر از Sarsa است که خط مشی دورگرد را می‌آموزد. البته، اگر  $\epsilon$  به تدریج کاهش یابد، هر دو روش به طور مجانبی به خط مشی بهینه همگرا می‌شوند.



شکل ۱۸- دنیای شبکه‌ای پیاده روی صخره (راست) و عملکرد روش‌های Sarsa و Q-learning

### ۳-۲- روش‌های مبتنی بر خط‌مشی (گرادیان خط‌مشی)

روش‌های یادگیری تقویتی مبتنی بر خط‌مشی<sup>۱</sup> (سیاست) به طور مستقیم خط‌مشی مدل را بهینه می‌کنند؛ درحالی‌که تکنیک‌های مبتنی بر ارزش به دنبال یافتن تابع ارزش بهینه هستند که به نوبه خود بهترین خط‌مشی را ارائه می‌دهد.

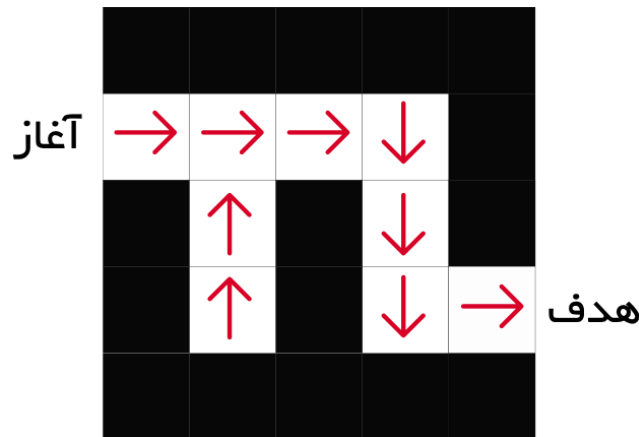
در روش‌های مبتنی بر خط‌مشی، قصد بهینه‌سازی تابع سیاست  $\pi(s)$  بدون استفاده از تابع ارزش است. سیاست چیزی است که رفتار عامل را در یک زمان داده شده، تعیین می‌کند. عامل یک تابع سیاست را می‌آموزد. این امر به او کمک می‌کند تا هر حالت را به بهترین عمل ممکن نگاشت کند.

دو دسته از سیاست‌ها وجود دارند.

- قطعی: سیاست برای یک حالت داده شده همیشه عمل مشابهی را باز می‌گرداند.
- تصادفی: برای هر یک از اعمال یک توزیع احتمالی در نظر می‌گیرد.

$$\pi(a|s) = \mathbb{P}[A_t = a | S_t = s] \text{ : سیاست تصادفی}$$

انجام یک عمل مشخص  
مشروط شده به یک حالت



شکل ۱۹- محاسبه بهترین عمل بر اساس سیاست در مثال هزارتو

همان‌طور که مشهود است، سیاست مستقیماً بهترین عمل برای هر حالت را بیان می‌کند.

<sup>۱</sup> Policy-based

در این نوع الگوریتم‌ها، عامل مستقیماً یک سیاست (تابعی که اقدام مناسب در هر حالت را مشخص می‌کند) را یاد می‌گیرد، بدون اینکه نیاز به تخمین تابع ارزش داشته باشد. این الگوریتم‌ها به‌ویژه در مسائلی که فضای عمل پیوسته است، مؤثر هستند.

### ۱-۳-۲- روش‌های گرادیان خط‌مشی (Policy Gradient)

نوعی دیگر از الگوریتم یادگیری تقویتی که یک خط‌مشی را یاد می‌گیرند. خط‌مشی تابعی است که از حالت‌ها به اقدامات صحیح می‌رسد. این شیوه با تخمین پاداش مورد انتظار با توجه به هر خط‌مشی، گرادیان را به‌روز می‌کند. به طور مثال، در آموزش سگ، تصور کنید که به او یاد می‌دهید که بنشیند، اما گاهی اوقات او خیلی آهسته می‌نشیند یا اصلاً نمی‌نشیند. با استفاده از گرادیان خط‌مشی<sup>۱</sup>، می‌توانید بر اساس میزان پیروی از خط‌مشی، به سگ پاداش‌های مختلفی بدهید. به عنوان مثال، اگر سگ سریع بنشیند، می‌توانید غذای زیادی به او بدهید، اما اگر آرام بنشیند، می‌توانید غذای کمی در اختیارش بگذارید. این تفاوت در پاداش‌ها به سگ کمک می‌کند تا بفهمد کدام اقدامات بهتر است و او را تشویق می‌کند تا رفتار مورد انتظار را به طور مؤثرتری انجام دهد. با تکرار مکرر این فرایند، سگ به تدریج یاد می‌گیرد که این سیاست را بهتر و بیشتر دنبال کند. به طور مشابه، در هوش مصنوعی، ما مدل را با مثال‌های زیادی آموزش می‌دهیم و زمانی که تصمیم‌های خوب می‌گیرد به آن پاداش می‌دهیم و وقتی تصمیمات بد می‌گیرد، جریمه می‌کنیم. این کار به هوش مصنوعی اجازه می‌دهد تا سیاست خود را بهبود بخشد و در موقعیت‌های مختلف تصمیمات بهتری بگیرد.

هر کدام از این روش‌ها از نظر پیچیدگی متفاوت‌اند و برای انواع مختلفی از مشکلات بر اساس ماهیت محیط و اطلاعات موجود مناسب هستند. علاوه بر روش‌های فوق، تعدادی تکنیک دیگر وجود دارد که می‌توان از آن‌ها برای بهبود فرایند تصمیم‌گیری در یادگیری تقویتی استفاده کرد. این تکنیک‌ها شامل موارد زیر می‌شوند:

**اکتشاف حریصانه اپسیلون:** تکنیکی است که به عامل اجازه می‌دهد در عین حال که پاداش‌های بهینه را یاد می‌گیرد، محیط را هم کاوش کند. عامل با کاوش تصادفی محیط یادگیری را شروع می‌کند؛ اما با یادگیری بیشتر در مورد محیط حریص‌تر می‌شود و اقدامی را با بالاترین ارزش تخمینی انتخاب می‌کند.

**پاسخ بر اساس تجربه:** تکنیکی که تجربیات عامل را در یک بافر ذخیره می‌کند. این تجربیات ذخیره‌شده به عامل اجازه می‌دهد تا با یادگیری از تعداد بیشتری از تجربیات به طور کارآمدتری فرایند آموزش را طی کند.

<sup>1</sup> Policy Gradient

**شکل‌دهی پاداش:** این تکنیک پاداش‌های موجود در محیط را تغییر می‌دهد تا مشکل را آسان‌تر حل کند. این کار را می‌توان با افزودن پاداش برای رفتارهای دلخواه یا با حذف پاداش برای رفتارهای ناخواسته انجام داد.

- در یادگیری مبتنی بر سیاست، هیچ تابع ارزشی وجود ندارد.
- خط‌مشی  $\pi(s, a)$  توسط بردار  $\pi(s, a; \theta)$  پارامتری می‌شود.
- به طور صریح خط‌مشی  $\pi(s, a; \theta)$  را یاد بگیرید که به طور ضمنی پاداش را در تمام خط‌مشی‌ها به حداکثر می‌رساند.
- باتوجه به سیاست  $\pi(s, a; \theta)$  با پارامترهای  $\theta$ ، بهترین  $\theta$  را پیدا کنید.
- چگونه کیفیت یک خط‌مشی  $\pi(s, a; \theta)$  را اندازه‌گیری کنیم؟
- فرض کنید تابع هدف  $J(\theta)$  باشد.
- پارامترهای خط‌مشی  $\theta$  را که  $J(\theta)$  را به حداکثر می‌رسانند، پیدا کنید.
- یکی از الگوریتم‌های معروف در این دسته، الگوریتم REINFORCE است که با استفاده از نمونه‌گیری از توزیع احتمال اقدامات، سعی در بهینه‌سازی سیاست عامل دارد.

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_{\pi}(s, a) \nabla \pi(a|s, \theta) = \mathbb{E}_{\pi} \left[ \sum_a q_{\pi}(S_t, a) \nabla \pi(a|S_t, \theta) \right]$$

$$\theta_{t+1} \doteq \theta_t + \alpha \sum_a \hat{q}(S_t, a, \mathbf{w}) \nabla \pi(a|S_t, \theta)$$

### الگوریتم نمونه REINFORCE:

```

Input: a differentiable policy parametrization  $\pi(a|s, \theta)$ 
Algorithm parameter: step size  $\alpha > 0$ 
Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  (e.g., to  $\theta$ )

Loop forever (for each episode):
    Generate an episode  $S_0, A_0, R_1, \dots, S_{t-1}, A_{t-1}, R_t$ , following  $\pi(\cdot|\cdot, \theta)$ 
    Loop for each step of the episode  $t = 0, 1, \dots, T - 1$ :
         $G \leftarrow \sum_{k=t+1}^{\infty} \gamma^{k-t-1} R_k$ 
         $\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t|S_t, \theta)$ 
    
```

### روش‌های عملگر-منتقد

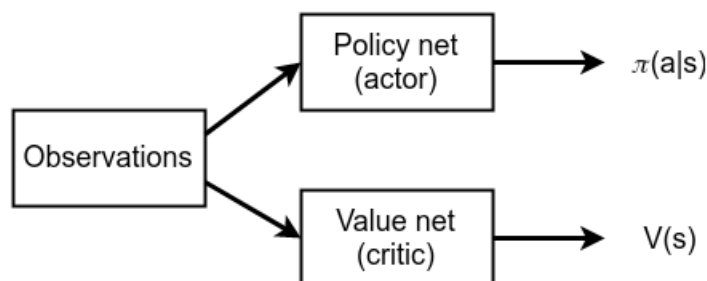
روش عملگر-منتقد<sup>۱</sup> به‌طور کلی به یادگیری هم‌زمان یک خط‌مشی و یک تابع ارزش اشاره دارد، جایی که تابع ارزش برای ارزیابی خط‌مشی استفاده می‌شود. عملگر مسئول ایجاد خط‌مشی‌ها، انتخاب کنش‌ها و تعامل با

<sup>۱</sup> Actor-Critic

محیط است. عملگر با گرادینانهای محاسبه شده از توابع زیر به روز می شود. منتقد کارکرد ارزشی سیاست عملگر را در هر مرحله زمانی ارزیابی می کند. معیارهای مختلفی را می توان برای ارزیابی خطمشی عملگر استفاده کرد، مانند تابع ارزش عمل  $Q(s, a)$ ، تابع ارزش حالت  $V(s)$  یا تابع مزیت  $A(s, a)$ .

$$\nabla_{\theta} J(\theta) = E_{\pi_{\theta}(s)} \left[ \sum_a Q(s, a) \nabla_{\theta} \pi_{\theta}(a|s) \right]$$

$$\theta = \theta + \alpha \nabla_{\theta} J$$



شکل ۲۰ - عملکرد عملگر-منتقد

الگوریتم عملگر-منتقد

```
Hyperparameters: step size  $\eta_{\theta}$  and  $\eta_{\psi}$ , reward discount factor  $\gamma$ 

Input: initial policy parameters  $\theta_0$ , initial value function parameters  $\psi_0$ 

Initialize  $\theta = \theta_0$  and  $\psi = \psi_0$ 

for  $t = 0, 1, 2, \dots$  do
    Run policy  $\pi_{\theta}$  for one step, collection  $\{S_t, A_t, R_t, S_{t+1}\}$ 
    Estimate advantages  $\hat{A}_t = R_t + \gamma V_{\psi}^{\pi_{\theta}}(S_{t+1}) - V_{\psi}^{\pi_{\theta}}(S_t)$ 
     $J(\theta) = \sum_t \log \pi_{\theta}(A_t|S_t) \hat{A}_t$ 
     $J_{V^{\pi_{\theta}}}(\psi) = \sum_t \hat{A}_t^2$ 
     $\psi = \psi + \eta_{\psi} \nabla J_{V^{\pi_{\theta}}}(\psi)$ ,  $\theta = \theta + \eta_{\theta} \nabla J(\theta)$ 
end for
Return  $(\theta, \psi)$ 
```

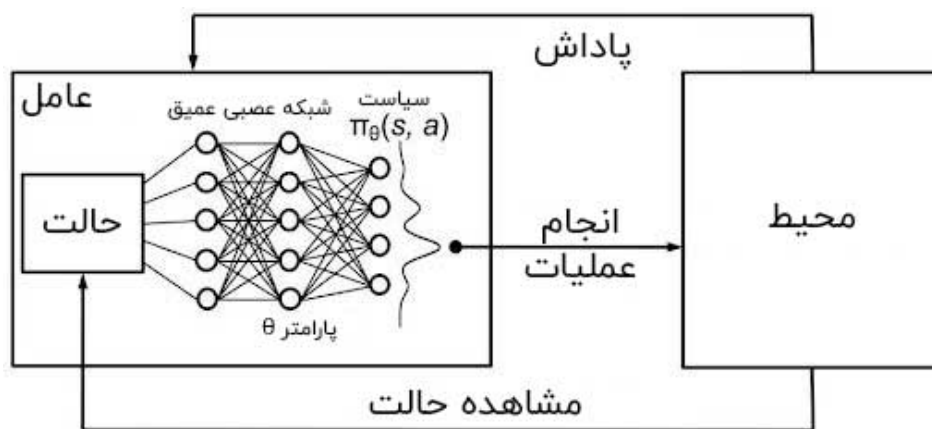
## ۴-۲- تقریب توابع و یادگیری عمیق در یادگیری تقویتی

عامل یادگیرنده در فرایند یادگیری تقویتی نیاز دارد تا تابع ارزش را محاسبه و ذخیره نماید. اما در مواقعی که ابعاد مسئله بزرگ باشد، ذخیره تمامی ارزشهای حالتهای مختلف سیستم عملاً غیرممکن خواهد بود. در این گونه مواقع به سراغ روشهای تخمین تابع ارزش<sup>۱</sup> خواهیم رفت. روشهای تقریب مختلفی وجود دارد که برخی از آنها عبارتاند از:

<sup>۱</sup> Approximate value function

- تقریب زننده خطی (Linear Approximator)
- تقریب زننده غیرخطی (Non-linear approximator)
- شبکه عصبی (Neural Network)
- ...

در صورتی که در روش یادگیری تقویتی، برای تخمین تابع ارزش از شبکه عصبی استفاده شود، شاخه جدیدی به نام یادگیری تقویتی عمیق<sup>۱</sup> شکل می گیرد که در سال های اخیر کاربردهای یادگیری تقویتی در طیف وسیعی از حوزه ها آزمایش شده است.



شکل ۲۱- تعامل عامل و محیط در یادگیری تقویتی عمیق

یکی از چالش های موجود در پیاده سازی روش های یادگیری تقویتی در دنیای واقعی نیاز به حجم زیادی از داده جهت دستیابی به نتایج مطلوب است. در روش های کلاسیک یادگیری تقویتی، این داده ها که به واسطه تعامل با محیط و انجام آزمایش های سعی و خطا تولید می گردند، قابلیت تعمیم دهی کمی دارند. به بیان دیگر برای مسئله جدید مجدداً باید فرایند جمع آوری داده طی شود که می تواند بسیار دشوار و زمان بر باشد، به ویژه در مواردی که تعداد ترکیبی از اقدامات و حالت ها بسیار زیاد است و یا جایی که محیط غیرقطعی بوده و می تواند حالت های تقریباً نامحدودی داشته باشد.

بدین منظور مسائل یادگیری تقویتی عمیق که استفاده از یادگیری عمیق در بستر یادگیری تقویتی را فراهم کرده، مطرح شده است تا بتواند بسیاری از چالش های موجود در این زمینه را رفع نماید. باتوجه به ساختار و تابع بهینه سازی مناسب، یک شبکه عصبی عمیق می تواند یک خط مشی بهینه را بدون گذر از تمام حالت های ممکن

<sup>1</sup> Deep Reinforcement Learning

یک سیستم، یاد بگیرد. هرچند عوامل یادگیری تقویتی عمیق همچنان به حجم عظیمی از داده‌ها نیاز دارند، اما می‌توانند مشکلاتی را که حل آنها با سیستم‌های کلاسیک یادگیری تقویتی غیرممکن بود، حل کنند.

یادگیری عمیق، شامل استفاده از تقریب توابع غیرخطی چندلایه، معمولاً شبکه‌های عصبی است. یادگیری عمیق کار متفاوتی با آنچه قبلاً توضیح داده شد نیست. یادگیری عمیق مجموعه‌ای از تکنیک‌ها و روش‌ها برای استفاده از شبکه‌های عصبی برای حل وظایف یادگیری ماشین، اعم از یادگیری با نظارت، یادگیری بدون نظارت، یا یادگیری تقویتی است. یادگیری تقویتی عمیق استفاده از یادگیری عمیق برای حل وظایف یادگیری تقویتی است.

منظور از تقریب توابع<sup>۱</sup> چیست؟ در واقع منظور از تقریب توابع غیرخطی همان یادگیری عمیق است. در این روش ما روش‌های ارائه‌شده را گسترش می‌دهیم تا برای مسائل با فضاهای حالت بزرگ دلخواه اعمال شود. در بسیاری از وظایفی که ما می‌خواهیم یادگیری تقویتی را برای آنها اعمال کنیم، فضای حالت ترکیبی و عظیم است. در چنین مواردی ما نمی‌توانیم انتظار داشته باشیم که یک خط‌مشی بهینه یا تابع مقدار بهینه را حتی در محدوده زمان و داده بی‌نهایت پیدا کنیم. هدف ما در عوض یافتن یک راه‌حل تقریبی خوب با استفاده از منابع محاسباتی محدود است.

مشکل فضاهای حالت بزرگ فقط حافظه موردنیاز برای جداول بزرگ نیست، بلکه زمان و داده‌های موردنیاز برای پر کردن دقیق آنهاست. در بسیاری از وظایف هدف ما، تقریباً هر حالتی که با آن مواجه می‌شویم هرگز قبلاً دیده نشده است. برای تصمیم‌گیری معقول در چنین حالت‌هایی، لازم است از رویارویی‌های قبلی با حالت‌های متفاوتی که به‌نوعی شبیه به حالت فعلی هستند، تعمیم دهیم. به عبارت دیگر، موضوع کلیدی بحث تعمیم است.

مزایای استفاده از تقریب تابع مقدار نه تنها شامل مقیاس‌پذیری برای وظایف در مقیاس بزرگ، بلکه سهولت تعمیم به حالت‌های نادیده از حالت‌های دیده شده با فضاهای حالت پیوسته است. علاوه بر این، تقریب تابع مبتنی بر شبکه‌های عصبی نیز نیاز به طراحی دستی ویژگی‌ها برای نمایش حالت‌ها را کاهش می‌دهد یا حذف می‌کند.

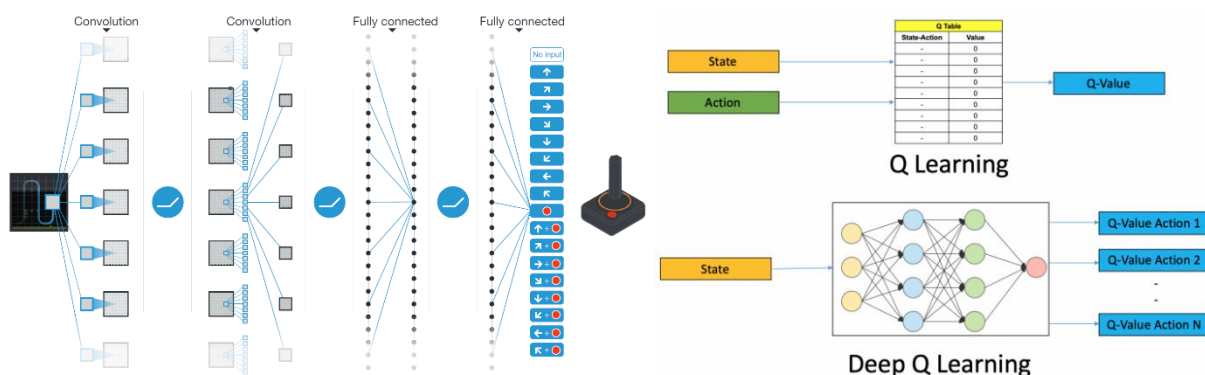
انگیزه ما برای استفاده از تقریب تابع نه تنها حل مسائلی است که در غیر این صورت قابل حل نیستند، بلکه برای حل کارآمدتر مسائل است.

<sup>1</sup> function approximation

## ۱-۴-۲- Deep Q-Networks (DQN)

الگوریتم یادگیری Q الگوریتم قدرتمندی است، اما قابلیت تعمیم‌پذیری ندارد و همین مسئله را می‌توان بزرگ‌ترین نقطه ضعف آن دانست. اگر الگوریتم یادگیری Q را به‌روزرسانی اعداد موجود در یک آرایه دوبعدی (شامل: فضای اقدام×فضای حالت) در نظر بگیرید، متوجه شباهت آن با برنامه‌نویسی پویا خواهید شد. این موضوع برای ما روشن می‌سازد که وقتی عامل تصمیم‌گیرنده در الگوریتم یادگیری Q با وضعیتی کاملاً جدید روبه‌رو شود، هیچ راهی برای شناسایی و انتخاب اقدام مناسب نخواهد داشت. به عبارت دیگر، عامل تصمیم‌گیرنده الگوریتم یادگیری Q توانایی تخمین ارزش حالت‌های ناشناخته را ندارد. برای حل این مشکل، شبکه DQN آرایه دوبعدی را حذف و شبکه عصبی را جایگزین آن می‌کند.

شبکه DQN به کمک یک شبکه عصبی، تابع Q-value را تخمین می‌زند. حالت فعلی به‌عنوان ورودی به این شبکه داده می‌شود، سپس مقدار Q-value متناظر با هر اقدام به‌عنوان خروجی از شبکه دریافت خواهد شد.



شکل ۲۲- تفاوت یادگیری Q با یادگیری Q عمیق (راست) و مثالی از شبکه DQN در آتاری (چپ)

شرکت دیپ مایند در سال ۲۰۱۳، شبکه DQN را همان‌طور که در تصویر بالا ملاحظه می‌کنید، در بازی آتاری به کار گرفت. ورودی که به این شبکه داده می‌شد یک تصویر خام از حالت جاری بازی بود. این ورودی از چندین لایه مختلف از جمله لایه‌های پیچشی و تماماً متصل عبور می‌کند و خروجی نهایی شامل مقادیر Q-value های مربوط به تمام اقدامات احتمالی عامل تصمیم‌گیرنده است.

**دو تکنیک دیگر نیز برای آموزش شبکه DQN ضروری است:**

۱. **تکرار تجربه:** در سازوکار متداول یادگیری تقویتی، نمونه‌های آموزشی همبستگی بالایی دارند و از نظر مقدار داده موردنیاز نیز کارآمد نیستند. به همین دلیل، سخت است که مدل به مرحله همگرایی برسد. یک راه برای حل مسئله توزیع نمونه، به کارگیری تکنیک تکرار تجربه است. در این روش، تابع انتقال



نمونه‌ها ذخیره می‌شود و سپس الگوریتم این تجربیات را از محلی به نام «مجموعه انتقال» به طور تصادفی انتخاب می‌کند تا دانش خود را بر اساس آن به‌روزرسانی نماید.

۲. **شبکه هدف مجزا:** ساختار شبکه هدف<sup>۱</sup> Q مشابه ساختار شبکه‌ای است که ارزش را تخمین می‌زند. همان‌طور در نمونه کد بالا ملاحظه کردید، بعد از C مرحله، در شبکه هدف جدیدی تعریف می‌شود تا از شدت نوسانات کاسته شود و فرایند آموزش ثبات بیشتری داشته باشد.

- شبکه Q عمیق شامل
  - شبکه Q که Q-value را پیش‌بینی می‌کند
  - شبکه هدف که ساختاری مشابه شبکه Q دارد
  - مؤلفه تکرار تجربه
- تکرار تجربه یک عمل حریصانه ( $\epsilon$ -greedy) را از حالت فعلی انتخاب می‌کند، آن را در محیط اجرا می‌کند و یک پاداش و حالت بعدی را پس می‌گیرد. این مشاهده را به‌عنوان نمونه‌ای از داده‌های آموزشی ذخیره می‌کند
- دسته‌ای از داده‌های آموزشی به هر دو شبکه داده می‌شود.
- شبکه Q حالت فعلی و عمل را از هر نمونه داده می‌گیرد و Q-value را برای آن عمل خاص پیش‌بینی می‌کند.
- شبکه Target(Q') حالت بعدی را از هر نمونه داده می‌گیرد و بهترین Q-value را از بین تمام اقداماتی که می‌توان از آن حالت انجام داد، پیش‌بینی می‌کند.
- تابع ضرر<sup>۲</sup> در تکرار i به‌صورت تعریف شده است:

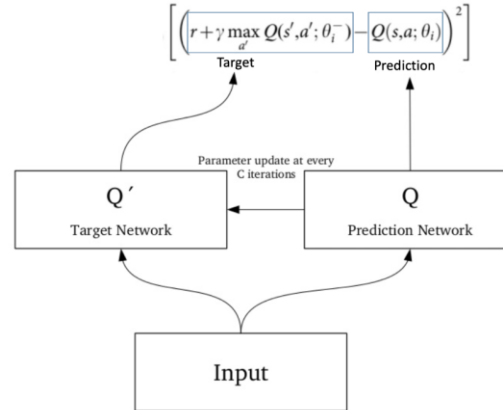
$$J_i(\theta_i) = \mathbb{E}_{(s,a,r,s') \sim U(S)} \left[ \left( (r + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta_i)) \right)^2 \right]$$

که در آن  $U(S)$  توزیع یکنواخت از مجموعه آموزشی S و  $\theta_i^-$  پارامترهای شبکه هدف است.

- فقط شبکه Q آموزش داده‌شده و شبکه هدف ثابت است.
- در هر مرحله C، وزن شبکه Q در شبکه هدف کپی می‌شود.

<sup>۱</sup> Target Network

<sup>۲</sup> loss function



شکل ۲۳- عملکرد DQN

الگوریتم DQN همراه با تکرار تجربه:

```

Initialize replay memory D to capacity N
Initialize action-value function Q with random weights  $\theta$ 
Initialize target action-value function  $Q'$  with weights  $\theta^- = \theta$ 
For episode = 1, M do
    Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$ 
    For t = 1, T do
        With probability  $\epsilon$ , select a random action  $a_t$ 
        Otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocesses  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(\phi_{t+1}, a_t, r_t, \phi_{t+1})$  in D
        Sample random minibatch of transitions  $(\phi_j, a_j, r_j; \phi_{j+1})$  from D
        Set  $y_j = \{$ 
             $r_j$  if episode terminates at step j+1
             $r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta^-)$  otherwise
         $\}$ 
        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to
        the network parameters  $\theta$ 
        Every C steps reset  $Q' = Q$ 
    End For
End For

```

کد ساختار شبکه DQN:

```

class DQN(nn.Module):

    def __init__(self, n_observations, n_actions):
        super(DQN, self).__init__()
        self.layer1 = nn.Linear(n_observations, 128)
        self.layer2 = nn.Linear(128, 128)
        self.layer3 = nn.Linear(128, n_actions)

    # Called with either one element to determine next action, or a batch
    # during optimization. Returns tensor([[left0exp,right0exp]...]).
    def forward(self, x):
        x = F.relu(self.layer1(x))
        x = F.relu(self.layer2(x))
        return self.layer3(x)

```

## ۲-۴-۲- گرادیان خط‌مشی قطعی عمیق (DDPG)

اگرچه شبکه DQN در مسائلی که ابعاد زیادی دارند (همچون بازی آتاری)، بسیار موفق عمل کرده است، اما فضای اقدام در این شبکه گسسته است و از آنجاکه بسیاری از مسائل موردعلاقه و حائز اهمیت برای ما از جمله مسائل مربوط به کنترل فیزیکی، دارای فضای اقدام پیوسته هستند، گسسته بودن فضا در شبکه DQN یک نقطه‌ضعف به شمار می‌آید.

ما الگوریتمی به نام گرادیان خط‌مشی قطعی عمیق<sup>۱</sup> (DDPG) را بررسی می‌کنیم. DDPG را می‌توان به‌عنوان یک DQN برای فضاهای عمل پیوسته در نظر گرفت. DDPG از بسیاری از تکنیک‌های مشابه موجود در DQN استفاده می‌کند. با این حال، DDPG همچنین سیاستی را آموزش می‌دهد که عمل بهینه را تقریب می‌کند. به همین دلیل، DDPG یک روش گرادیان خط‌مشی قطعی است که به فضاهای کنش پیوسته محدود شده است.

با تبدیل یک فضای اقدام پیوسته به یک فضای گسسته به‌صورت دقیق و جزئی، یک فضای اقدام بسیار گسترده به دست می‌آید. برای مثال، فرض کنید درجه آزادی سیستم تصادفی ۱۰ باشد. به‌ازای هر درجه، باید فضا را به ۴ قسمت تقسیم کنیم. به این ترتیب، در آخر  $4^{10} = 1,048,576$  عدد اقدام خواهیم داشت. سخت است که در چنین فضای اقدام بزرگی به همگرایی برسید.

الگوریتم DDPG مبتنی بر معماری عملگر منتقد عمل می‌کند. این معماری دو عنصر اصلی دارد: عملگر و منتقد. عنصر عملگر، تنظیم پارامتر با تابع سیاست را برعهده دارد، به این ترتیب، عملگر تصمیم می‌گیرد که بهترین اقدام ممکن برای هر حالت چیست.

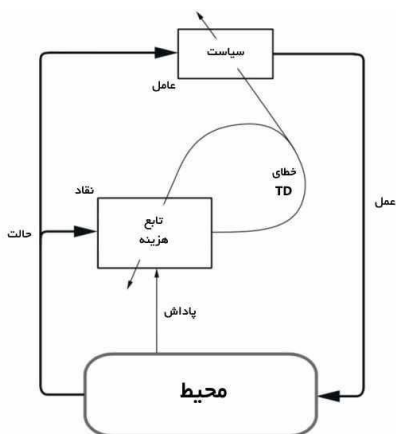
$$\pi_{\theta}(s, a) = \mathbb{P}[a | s, \theta]$$

وظیفه منتقد نیز ارزیابی تابع سیاستی است که عملگر بر اساس تابع خطای تفاوت موقتی تخمین زده است.

$$r_{t+1} + \gamma V^{\nu}(s_{t+1}) - V^{\nu}(s_t)$$

در این تابع حرف  $\nu$  کوچک نمایانگر سیاستی است که عملگر برمی‌گزیند. این فرمول کمی آشنا به نظر نمی‌رسد؟ بله، درست است. این فرمول دقیقاً مشابه معادله به‌روزرسانی الگوریتم یادگیری Q است. یادگیری تفاوت زمانی روشی است که الگوریتم به کمک آن می‌تواند نحوه پیش‌بینی یک تابع ارزش را بر پایه ارزش‌های آتی یک حالت معین بیاموزد. الگوریتم یادگیری Q نوعی خاصی از یادگیری تفاوت زمانی در حوزه یادگیری Q-value است.

<sup>1</sup> deep deterministic policy gradient



شکل ۲۴- معماری عملگر-منتقد

الگوریتم DDPG تکنیک‌های تکرار تجربه و شبکه هدف مجزا در شبکه DQN را نیز به کار می‌گیرد. اما یکی از مشکلات DDPG این است که به‌ندرت اقدامات را جست‌وجو می‌کند. یک راه‌حل برای این مشکل، **ایجاد اختلال** در فضای پارامترها یا فضای اقدام است.

البته محققین OpenAI در مقاله خود ادعا کرده‌اند که ایجاد اختلال در فضای پارامترها بهتر از ایجاد اختلال در فضای اقدام است. یکی از رایج‌ترین و پرکاربردترین اختلالات در این زمینه فرایند تصادفی **اورنستین-یولنبرگ** نام دارد.

الگوریتم DDPG:

```
Randomly initialize critic network  $Q(s, a | \theta^Q)$  and actor  $\mu(s | \theta^\mu)$  with weights  $\theta^Q$  and  $\theta^\mu$ 
Initialize target network  $Q'$  and  $\mu'$  with weights  $\theta^{Q'} \leftarrow \theta^Q$ ,  $\theta^{\mu'} \leftarrow \theta^\mu$ 
Initialize replay buffer  $R$ 
for episode = 1, M do
    Initialize a random process  $N$  for action exploration
    Receive initial observation state  $s_1$ 
    for t = 1, T do
        Select action  $a_t = \mu(s_t | \theta^\mu) + N_t$  according to the current policy and exploration noise
        Execute action  $a_t$  and observe reward  $r_t$  and observe new state  $s_{t+1}$ 
        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $R$ 
        Sample a random minibatch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $R$ 
        Set  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1} | \theta^{\mu'}) | \theta^{Q'})$ 
        Update critic by minimizing the loss:  $L = 1/N \sum_i (y_i - Q(s_i, a_i | \theta^Q))^2$ 
        Update the actor policy using the sampled policy gradient:
         $\nabla_{\theta^\mu} J \approx 1/N \sum_i \nabla_a Q(s, a | \theta^Q) |_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s | \theta^\mu) |_{s_i}$ 
        Update the target networks:
         $\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$ 
         $\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$ 
    end for
end for
```

کد ساختار شبکه DDPG:

```
class DDPGActor(nn.Module):
    def __init__(self, obs_size, act_size):
        super(DDPGActor, self).__init__()
        self.net = nn.Sequential(
            nn.Linear(obs_size, 400), nn.ReLU(), nn.Linear(400, 300),
            nn.ReLU(), nn.Linear(300, act_size), nn.Tanh()
        )
    def forward(self, x):
        return self.net(x)

class DDPGCritic(nn.Module):
    def __init__(self, obs_size, act_size):
        super(DDPGCritic, self).__init__()
        self.obs_net = nn.Sequential(
            nn.Linear(obs_size, 400), nn.ReLU(),
        )
        self.out_net = nn.Sequential(
            nn.Linear(400 + act_size, 300), nn.ReLU(),
            nn.Linear(300, 1)
        )
    def forward(self, x, a):
        obs = self.obs_net(x)
        return self.out_net(torch.cat([obs, a], dim=1))
```

### ۳- مباحث پیشرفته در یادگیری تقویتی

پس از بررسی الگوریتم‌های یادگیری تقویتی در بخش قبل، در این بخش از دنیای شبیه‌سازی و بازی فاصله گرفته و به دنیای واقعی و محیط‌های پیچیده‌تر پرداخته می‌شود و انواع تکنیک‌های پیشرفته‌تر و نحوه مقابله با محیط‌های پیچیده‌تر بررسی می‌شود.

#### ۳-۱- انواع Deep Q-Networks

از زمانی که مدل شبکه Q عمیق (DQN) در سال ۲۰۱۵ معرفی شد، پیشرفت‌های زیادی، همراه با تغییراتی در معماری پلایه، ارائه شده است که به طور قابل توجهی همگرایی، پایداری و کارایی نمونه پلایه DQN را بهبود بخشیده است. در این قسمت نگاهی به برخی از این ایده‌ها خواهیم داشت.

##### ۳-۱-۱ N-step DQN

اولین پیشرفتی که بررسی خواهیم کرد، یک پیشرفت کاملاً قدیمی است. اولین بار در سال ۱۹۸۸ معرفی شد. برای درک ایده، اجازه دهید یک‌بار دیگر به آپدیت بلمن مورد استفاده در یادگیری Q نگاه کنیم:

$$Q(s_t, a_t) = r_t + \gamma \max_a Q(s_{t+1}, a_{t+1})$$

این معادله بازگشتی است، به این معنی که می‌توانیم  $Q(s_{t+1}, a_{t+1})$  را بر حسب خودش بیان کنیم که این نتیجه را به ما می‌دهد:

$$Q(s_t, a_t) = r_t + \gamma \max_a \left[ r_{a,t+1} + \gamma \max_{a'} Q(s_{t+2}, a') \right]$$

مقدار  $r_{a,t+1}$  به معنای پاداش محلی در زمان  $t+1$ ، پس از صدور اقدام  $a$  است. باین حال، اگر فرض کنیم که عمل  $a$  در مرحله  $t+1$  به طور بهینه یا نزدیک به بهینه انتخاب شده است، می‌توانیم عملیات  $\max_a$  را حذف کنیم و این را به دست آوریم:

$$Q(s_t, a_t) = r_t + \gamma r_{t+1} + \gamma^2 \max_a Q(s_{t+2}, a')$$

برای سرعت بخشیدن به همگرایی (convergence) این مقدار می‌تواند بارها و بارها باز شود. این باز کردن<sup>۱</sup> را می‌توان به راحتی با جایگزین کردن نمونه انتقال یک مرحله‌ای (1-Step) با دنباله‌های انتقال طولانی‌تر از  $n$  مرحله (n-Step)، به به روزرسانی DQN ما اعمال کرد. این باز کردن به ما کمک می‌کند تا سرعت آموزش را افزایش دهیم.

<sup>1</sup> unrolling

الگوریتم یادگیری N-step Q به جز تغییرات زیر به روشی مشابه DQN عمل می‌کند:

- هیچ بافر اجرا مجدد استفاده نمی‌شود. به جای نمونه‌برداری از دسته‌های تصادفی انتقال، شبکه هر N مرحله با استفاده از آخرین N مرحله‌ای که توسط عامل اجرا می‌شود، آموزش داده می‌شود.
- به منظور تثبیت یادگیری، چندین کارگر با هم برای به‌روزرسانی شبکه کار می‌کنند. این همان اثر عدم همبستگی نمونه‌های مورد استفاده برای آموزش را ایجاد می‌کند.
- به جای استفاده از اهداف Q تک مرحله‌ای برای شبکه، پاداش‌های حاصل از مراحل بعدی  $\$N\$$  جمع‌آوری می‌شوند تا اهداف N مرحله Q را تشکیل دهند.

## ۲-۱-۳ Double DQN

ایده ثمربخش بعدی در مورد چگونگی بهبود DQN پایه توسط محققان DeepMind در مقاله‌ای ارائه شد. در این مقاله، نویسندگان نشان دادند که DQN اصلی تمایل دارد مقادیر Q را بیش از حد تخمین بزند که ممکن است برای عملکرد آموزشی مضر باشد و گاهی اوقات می‌تواند منجر به سیاست‌های غیربهبه‌شود. علت اصلی آن، عملیات max در معادله بلمن است. به عنوان راه‌حلی برای این مشکل، نویسندگان پیشنهاد کردند که آپدیت بلمن را کمی تغییر یابد:

$$Q(s_t, a_t) = r_t + \gamma \max_a Q'(s_{t+1}, a)$$

$Q'(s_{t+1}, a)$  مقادیر Q بود که با استفاده از شبکه هدف ما محاسبه شد، بنابراین ما با شبکه آموزش دیده هر n مرحله را به روز می‌کنیم. نویسندگان مقاله پیشنهاد کردند اقداماتی را برای حالت بعدی با استفاده از شبکه آموزش دیده انتخاب کنند، اما مقادیر Q را از شبکه هدف بگیرند؛ بنابراین، عبارت جدید برای Q-values هدف به شکل زیر خواهد بود:

$$Q(s_t, a_t) = r_t + \gamma \max_a Q'(s_{t+1}, \arg \max_a Q(s_{t+1}, a))$$

نویسندگان ثابت کردند که این پیچ‌دادن ساده<sup>۱</sup> تخمین بیش از حد را کاملاً برطرف می‌کند و آنها این معماری جدید را DQN دوگانه (Double DQN) نامیدند.

پیاده‌سازی آن بسیار ساده است. کاری که ما باید انجام دهیم این است که تابع ضرر خود را کمی اصلاح کنیم.

<sup>1</sup> simple tweak fixes

### ۳-۱-۳ Noisy networks

پیشرفت بعدی که می‌خواهیم به آن بپردازیم به یکی دیگر از مشکلات یادگیری تقویتی می‌پردازد: اکتشاف در محیط. شبکه‌های نویزی<sup>۱</sup> یک ایده بسیار ساده برای یادگیری ویژگی‌های اکتشاف در طول آموزش به‌جای داشتن یک برنامه زمانی جداگانه مرتبط با اکتشاف است.

DQN کلاسیک با انتخاب اقدامات تصادفی با اپسیلون هاپرپارامتر ویژه تعریف شده به اکتشاف دست می‌یابد که به آرامی در طول زمان از ۱ (عملکردهای کاملاً تصادفی) به نسبت کوچکی ۰.۱ یا ۰.۰۲ کاهش می‌یابد. این فرایند برای محیط‌های ساده به خوبی کار می‌کند. اما حتی در چنین موارد ساده‌ای، برای کارآمد کردن فرایندهای آموزشی نیاز به تنظیم دارد.

در روش Noisy Networks نویز را به وزن لایه‌های کاملاً متصل شبکه اضافه می‌کنند و پارامترهای این نویز را در حین تمرین با استفاده از پس انتشار تنظیم می‌کنند. نویز توسط یک جریان نویز اضافی به لایه خطی  $y = (Wx + b)$  اضافه می‌شود:

$$y = (Wx + b) + ((W_{noisy} \odot \epsilon_w)x + b_{noisy} \odot \epsilon_b)$$

که در آن  $\odot$  به ضرب از نظر عنصر اشاره دارد، هر دو  $W_{noisy}$  و  $b_{noisy}$  پارامترهای قابل آموزش هستند، درحالی‌که  $\epsilon_w$  و  $\epsilon_b$  مقیاس‌های تصادفی هستند که به صفر می‌رسند. آزمایش‌ها نشان می‌دهد که شبکه‌های نویزی برای طیف وسیعی از بازی‌های آتاری در چندین خط پایه امتیازات بالاتری به دست می‌آورند.

### ۳-۱-۴ Prioritized replay buffer

یکی از زمینه‌های بهبود در DQN پایه، استراتژی نمونه‌گیری بهتر برای تکرار تجربه است. تکرار تجربه اولویت‌دار<sup>۲</sup> (PER) تکنیکی برای اولویت‌بندی تجربه است، به‌گونه‌ای که انتقال‌های مهم را بیشتر تکرار کنیم. ایده اصلی تکرار تجربه اولویت‌دار در نظر گرفتن اهمیت انتقال با خطای تفاوت زمانی  $\delta$  است که می‌تواند به‌عنوان یک اندازه‌گیری شگفت‌انگیز در نظر گرفته شود.

دلیل این که تکرار تجربه اولویت‌دار می‌تواند کمک‌کننده باشد این است که برخی از تجربیات ممکن است حاوی اطلاعات بیشتری برای یادگیری در مقایسه با دیگران باشد. دادن شانس بیشتری برای بازپخش شدن به آن تجربه غنی از اطلاعات بیشتر به سریع‌تر و کارآمدتر شدن کل فرایند یادگیری کمک می‌کند.

<sup>1</sup> Noisy Networks

<sup>2</sup> Prioritized experience replay

تکرار تجربه اولویت‌دار از خطای تفاوت زمانی (TD) هر انتقال به عنوان معیار ارزیابی اولویت استفاده می‌کند. شکل خطای TD به صورت زیر است:

$$\delta_t = R_t + \gamma \max_a Q(s_t, a) - Q(s_{t-1}, a_{t-1})$$

هر چه قدر مطلق  $\delta_t$  بزرگ‌تر باشد و احتمال اینکه انتقال مربوطه انتخاب شود بیشتر باشد، سهم بیشتری در بهبود سیاست خواهد داشت. در فرایند نمونه‌گیری، روش‌های اولویت‌بندی تصادفی<sup>۱</sup> و نمونه‌گیری اهمیت<sup>۲</sup> به کار گرفته می‌شود. عملیات اولویت‌بندی تصادفی نه تنها می‌تواند از انتقال‌ها استفاده کامل کند، بلکه تنوع را نیز تضمین می‌کند. نمونه‌برداری اهمیت سرعت به روزرسانی پارامتر را کاهش می‌دهد و ثبات یادگیری را تضمین می‌کند.

برای پیاده‌سازی کارآمد، تابع چگالی تجمعی<sup>۳</sup> احتمال نمونه‌برداری با یک تابع خطی تکه‌ای<sup>۴</sup> با  $k$  بخش تقریب زده می‌شود. به طور دقیق‌تر، اولویت‌ها در یک ساختار داده کارآمد به نام درخت قطعه<sup>۵</sup> ذخیره می‌شوند. در طول زمان اجرا، تکرار تجربه اولویت‌دار ابتدا از یک قطعه نمونه‌برداری می‌کند، و سپس به طور یکنواخت از بین انتقال‌های درون آن نمونه‌برداری می‌کند.

### ۵-۱-۳- Dueling DQN

برای برخی از حالات‌ها، اقدامات مختلف به مقدار مورد انتظار مربوط نمی‌شود و ما نیازی به یادگیری تأثیر هر عمل برای چنین حالت‌هایی نداریم. برای مثال تصور کنید روی کوه ایستاده‌اید و طلوع خورشید را تماشا می‌کنید. منظره دلپذیر به شما آرامش زیادی می‌دهد که پاداش بالایی را به همراه دارد. شما می‌توانید اینجا بمانید و Q-value اقدامات مختلف مهم نیست؛ بنابراین جداکردن ارزش مستقل از عمل حالت و Q-value ممکن است به یادگیری قوی‌تری منجر شود.

مدل شبکه‌های عصبی عمیق مورد استفاده در DQN دارای سه لایه کانولوشن و به دنبال آن دو لایه کاملاً متصل است و خروجی مقادیر  $Q$  هر عمل است. متفاوت از معماری DQN، معماری دوئل ویژگی‌های استخراج شده از لایه‌های کانولوشن را به دو جریان تقسیم می‌کند. به طور دقیق‌تر، Q-value را می‌توان به بخش مقدار حالت و بخش مزیت عمل به صورت زیر تقسیم کرد:

$$Q^\pi(s, a) = V^\pi(s) + A^\pi(s, a)$$

<sup>1</sup> stochastic prioritization

<sup>2</sup> Importance-sampling

<sup>3</sup> cumulative density function

<sup>4</sup> piece-wise linear function

<sup>5</sup> segment tree



و دوئل DQN نمایش (representations) های این دو بخش را به صورت زیر جدا می کند:

$$Q(s, a; \theta, \theta_v, \theta_a) = V(s; \theta, \theta_v) + (A(s, a'; \theta, \theta_a) - \max_{a'} A(s, a'; \theta, \theta_a))$$

که در آن  $\theta_v$  و  $\theta_a$  پارامترهای دو جریان لایه های کاملاً متصل<sup>۱</sup> هستند،  $\theta$  پارامترهای لایه های کانولوشن را نشان می دهد. توجه داشته باشید که عملگر max در معادله قابلیت شناسایی<sup>۲</sup> را تضمین می کند که Q-value مقدار حالت و مزیت عمل را به طور منحصر به فرد بازایی کند. در غیر این صورت، آموزش ممکن است عبارت مقدار حالت را نادیده بگیرد و تابع مزیت را فقط به Q-value همگرا کند. پیشنهاد شده است که برای پایداری بهتر، max با میانگین به صورت زیر جایگزین شود:

$$Q(s, a; \theta, \theta_v, \theta_a) = V(s; \theta, \theta_v) + (A(s, a; \theta, \theta_a) - \frac{1}{|A|} \sum_{a'} A(s, a'; \theta, \theta_a))$$

که توسط آن تابع مزیت به جای دنبال کردن مزیت بهینه، فقط باید با مزیت مستقیم میانگین تطبیق یابد. آموزش معماری های دوئل، مانند DQN استاندارد، تنها به لایه های بیشتری نیاز دارد. آزمایش ها نشان می دهند که معماری های دوئل منجر به ارزیابی بهتر سیاست در حضور بسیاری از اقدامات بالارزش مشابه می شوند.

### ۳-۱-۶ - Categorical(Distributional) DQN

آخرین و پیچیده ترین روش در جعبه ابزار بهبود DQN مربوط به یک مقاله است که توسط DeepMind در ژوئن ۲۰۱۷ منتشر شده است، به نام دیدگاه توزیعی<sup>۳</sup> در یادگیری تقویتی. در مقاله، نویسندگان بخش اساسی یادگیری Q که همان Q-values است را زیر سوال بردند و سعی کردند آنها را با توزیع احتمال Q-value عمومی تر جایگزین کنند. بیایید سعی کنیم ایده را درک کنیم. هر دو روش یادگیری Q و تکرار ارزش با مقادیر اعمال یا حالت هایی که به صورت اعداد ساده نشان داده شده اند کار می کنند و نشان می دهند که ما چقدر می توانیم از یک حالت، یا یک عمل و یک حالت به پاداش کلی برسیم. باین حال، آیا فشرده کردن همه پاداش های احتمالی آینده در یک عدد عملی است؟ در محیط های پیچیده، آینده می تواند تصادفی باشد و مقادیر متفاوتی با احتمالات متفاوت به ما بدهد.

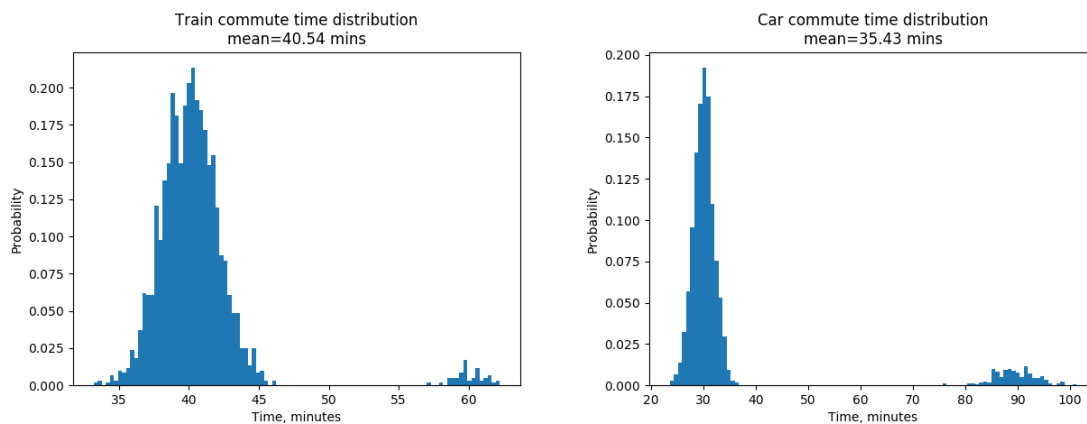
به عنوان مثال، سناریوی رفت و آمد را زمانی که به طور منظم از خانه به محل کار رانندگی می کنید، تصور کنید. بیشتر اوقات، ترافیک آن چنان سنگین نیست و حدود ۳۰ دقیقه طول می کشد تا به مقصد برسید. دقیقاً ۳۰ دقیقه نیست، اما به طور متوسط ۳۰ است. هرازگاهی اتفاقی می افتد، مانند تعمیر جاده یا تصادف، و به دلیل

<sup>1</sup> fully connected layers

<sup>2</sup> identifiability

<sup>3</sup> Distributional

ترافیک، سه برابر بیشتر طول می کشد تا به سر کار برسید. احتمال زمان رفت و آمد شما را می توان به عنوان توزیعی از متغیر تصادفی "زمان رفت و آمد" نشان داد و در نمودار زیر نشان داده شده است.



شکل ۲۵- توزیع احتمال زمان رفت و آمد با ماشین (راست) و قطار (چپ)

حالا تصور کنید که یک راه جایگزین برای رسیدن به محل کار دارید: قطار. کمی بیشتر طول می کشد، زیرا باید از خانه به ایستگاه قطار و از ایستگاه به دفتر بروید، اما بسیار قابل اعتمادتر است. مثلاً بگویید که زمان رفت و آمد قطار به طور متوسط ۴۰ دقیقه است، با احتمال کمی اختلال در قطار که ۲۰ دقیقه زمان اضافی به سفر اضافه می کند. توزیع رفت و آمد قطار در نمودار بالا نشان داده شده است.

تصور کنید که اکنون می خواهیم در مورد نحوه رفت و آمد تصمیم گیری کنیم. اگر فقط میانگین زمان ماشین و قطار را بدانیم، ماشین جذاب تر به نظر می رسد، زیرا به طور متوسط ۳۵:۴۳ دقیقه طول می کشد که بهتر از قطار ۴۰:۵۴ دقیقه ای است.

با این حال، اگر به توزیع کامل نگاه کنیم، ممکن است تصمیم بگیریم با قطار برویم، زیرا حتی در بدترین سناریو، یک ساعت رفت و آمد در مقابل یک ساعت و ۳۰ دقیقه خواهد بود. با تغییر به زبان آماری، توزیع خودرو واریانس بسیار بالاتری دارد، بنابراین در شرایطی که واقعاً باید حداکثر در ۶۰ دقیقه در دفتر باشید، قطار بهتر است.

وضعیت در سناریوی فرایند تصمیم مارکوف پیچیده تر می شود، زمانی که توالی تصمیم ها باید گرفته شود و هر تصمیمی ممکن است بر وضعیت آینده تأثیر بگذارد. در مثال رفت و آمد، ممکن است زمان یک جلسه مهم باشد که با توجه به راهی که می خواهید رفت و آمد کنید، باید آن را ترتیب دهید. در این صورت، کار با مقادیر میانگین پاداش ممکن است به معنای از دست دادن اطلاعات زیادی در مورد دینامیک های اساسی باشد.

دقیقاً همین ایده توسط نویسندگان دیدگاه توزیعی در یادگیری تقویتی ارائه شد. چرا با تلاش برای پیش‌بینی یک مقدار متوسط برای یک اقدام، خودمان را محدود می‌کنیم، درحالی‌که ارزش اساسی ممکن است توزیع زیربنایی پیچیده‌ای داشته باشد؟ شاید به ما کمک کند که مستقیماً با توزیع‌ها کار کنیم.

نتایج ارائه شده در مقاله نشان می‌دهد که در واقع، این ایده می‌تواند مفید باشد، اما به قیمت معرفی یک روش پیچیده‌تر. ایده کلی این است که توزیع ارزش را برای هر عمل، شبیه به توزیع‌های مثال ماشین/قطار پیش‌بینی کنیم. در مرحله بعد، نویسندگان نشان دادند که معادله بلمن را می‌توان برای یک حالت توزیع تعمیم داد، و به شکل  $Z(x, a) \stackrel{D}{=} R(x, a) + \gamma Z(x', a')$  که بسیار شبیه به معادله بلمن آشنا است، اما اکنون  $R(x, a)$  و  $Z(x, a)$  توزیع‌های احتمال هستند و نه اعداد.

توزیع به‌دست‌آمده می‌تواند برای آموزش شبکه ما استفاده شود تا پیش‌بینی‌های بهتری از توزیع ارزش برای هر عمل در حالت داده شده ارائه دهد، دقیقاً به همان روشی که با یادگیری Q انجام می‌شود. تنها تفاوت در تابع ضرر خواهد بود که اکنون باید با چیزی مناسب برای مقایسه توزیع جایگزین شود. چندین گزینه در دسترس وجود دارد، به‌عنوان مثال، واگرایی (Kullback-Leibler (KL) (از دست دادن آنتروپی متقابل) که در مسائل طبقه‌بندی استفاده می‌شود، یا معیار Wasserstein. در مقاله، نویسندگان برای معیار Wasserstein توجیه نظری ارائه کردند، اما زمانی که سعی کردند آن را در عمل اعمال کنند، با محدودیت‌هایی مواجه شدند. بنابراین، در نهایت از واگرایی KL استفاده شد.

## ۲-۳- یادگیری انتقالی و فرا یادگیری

در اکثر الگوریتم‌های یادگیری تقویتی، تغییرات کوچک حالت باعث شکست استراتژی‌های آموزش‌دیده قبلی می‌شود و آموزش از ابتدا گران است. یادگیری انتقالی<sup>۱</sup> (TL) می‌تواند با به‌کارگیری تجربه‌ای که از کار منبع به دست می‌آید، کارایی آموزشی وظیفه هدف را بهبود بخشد؛ بنابراین، برای بهبود سرعت یادگیری در وظایف جدید، یادگیری انتقالی برای یادگیری تقویتی تبدیل به یک جهت تحقیقاتی جذاب شده است. با توجه به این که چه دانش منتقل می‌شود، یادگیری انتقالی می‌تواند به الگوریتم یادگیری تقویتی عمیق در جنبه‌های زیر کمک کند: شکل دادن به پاداش<sup>۲</sup>، یادگیری از نمایش‌ها<sup>۳</sup>، انتقال خط‌مشی<sup>۴</sup>، نگاشت بین وظایف<sup>۵</sup>، و انتقال نمایش<sup>۶</sup>. لان و همکاران یک رویکرد یادگیری تقویتی انتقال جدید را از طریق استخراج فرا دانش با استفاده از درختان تصمیم

<sup>1</sup> Transfer learning

<sup>2</sup> reward shaping

<sup>3</sup> learning from demonstrations

<sup>4</sup> policy transfer

<sup>5</sup> intertask mapping

<sup>6</sup> representation transfer

هرس شده پیشنهاد کرده‌اند. ژانگ و همکاران روش  $DQDR^1$  را پیشنهاد کردند که دانش دامنه را از انسان استخراج می‌کند و آن را به‌عنوان مجموعه‌ای از قوانین بیان می‌کند و آن را با یک DQN جفت می‌کند تا قابلیت انتقال الگوریتم‌های یادگیری تقویتی را بهبود بخشد.

هدف فرا یادگیری تعمیم دانش به وظایف جدید با یادگیری تجربیات از چندین کار دیگر است. ایده اصلی یادگیری فرا-تقویتی<sup>۲</sup> این است که دانش قبلی آموخته شده از تعداد زیادی از وظایف یادگیری تقویتی را در وظایف یادگیری تقویتی جدید به کار گیرد و سرعت یادگیری و توانایی تعمیم عامل را بهبود بخشد. یادگیری فرا تقویتی عمده‌تاً کاستی‌های زیر یادگیری تقویتی عمیق را حل و بهینه می‌کند: استفاده کم از نمونه، مشکل در طراحی توابع پاداش، استراتژی‌های کاوش در وظایف ناشناخته، و فقدان توانایی تعمیم.

### ۳-۳- یادگیری تقویتی سلسله‌مراتبی

اغلب، ما خود را در حال توسعه محیط‌هایی می‌بینیم که با افق‌های متعدد مشکل دارند. به‌عنوان مثال، اگر می‌خواهیم یک عامل بهترین استراتژی سطح بالا را پیدا کند، اما فقط دستورات کنترل سطح پایین را برای اقدامات به او بدهد، در این صورت عامل باید یاد بگیرد که از فضای عمل سطح پایین به سطح بالا برود. به طور شهودی، یک سلسله‌مراتب در سیاست‌ها برای اکثر عامل‌ها وجود دارد. وقتی برنامه‌ریزی می‌کنم این کار را در فضای عمل سطح بالاتر انجام می‌دهم من به رفتن به فروشگاه فکر می‌کنم نه اینکه دستانم را برای رسیدن به فروشگاه تکان دهم. یادگیری تقویتی سلسله‌مراتبی، عوامل را قادر می‌سازد تا سلسله‌مراتبی از اقدامات را در داخل برای مقابله با مشکلات افق بلند ایجاد کنند. عوامل، دیگر در مورد دستورات چپ-راست استدلال نمی‌کنند، بلکه بیشتر در مورد رفتن به اینجا یا آنجا استدلال می‌کنند.

با کمک یادگیری عمیق، HRL می‌تواند مسائل مربوط به فضاهای حالت بزرگ و فضاهای عمل را حل کند و می‌تواند مستقیماً ورودی‌های با ابعاد بالا را پردازش کند. علاوه بر این، توانایی بازنمایی یادگیری عمیق، انتخاب‌های بیشتری را برای شکل اهداف فرعی فراهم می‌کند. یادگیری تقویتی عمیق سلسله‌مراتبی (HDRL) به یک جهت تحقیقاتی مهم یادگیری تقویتی عمیق تبدیل شده است.

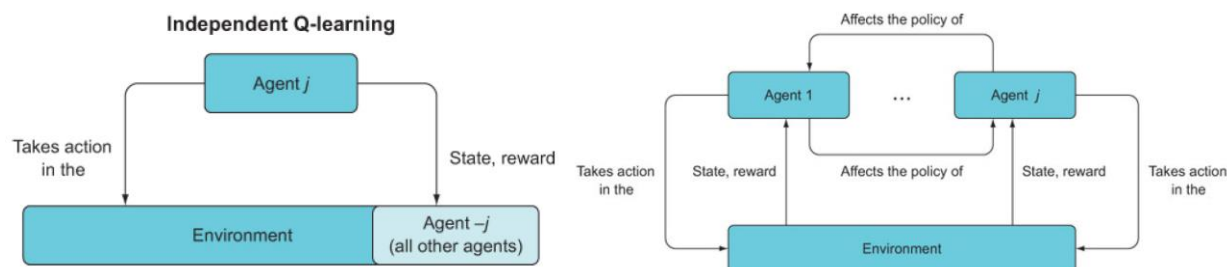
الگوریتم‌های مطرح: HIRO و STRAW

<sup>1</sup> Deep Q-learning with transferable Domain Rules

<sup>2</sup> Meta-Reinforcement Learning

### ۴-۳- یادگیری تقویتی چندعاملی

یادگیری تقویتی چندعاملی<sup>۱</sup> (MARL) یکی از مهم‌ترین زمینه‌های تحقیقاتی در یادگیری تقویتی است زیرا می‌تواند بسیاری از مشکلات تصمیم‌گیری مشارکتی پیچیده دنیای واقعی را برطرف کند. در بحث‌های قبلی، تنها الگوریتم‌های تک عامل یادگیری تقویتی مورد بحث قرار گرفته‌اند. یک عامل یادگیری تقویتی منفرد که با محیط تعامل دارد نمی‌تواند تعداد زیادی از مشکلات کنترل را به طور کامل حل کند. با این حال، چندین عامل یادگیری تقویتی که هم‌زمان با یک محیط در تعامل هستند، می‌توانند به یادگیری کارآمدتری دست یابند. MARL بر تجزیه و تحلیل و کنترل رفتارهای چندین عامل یادگیری تقویتی که در یک محیط مشترک وجود دارند متمرکز است. هر عاملی از عملکرد پاداش خود می‌آموزد و بر اساس منافع خود عمل می‌کند؛ این منافع فردی ممکن است با منافع برخی از عوامل دیگر در تضاد باشد و در نتیجه تعامل گروهی پیچیده شود. MARL یک پیچیدگی دیگر به مسئله یادگیری تقویتی اضافه می‌کند و آن معرفی غیریستایی به دلیل تغییر در رفتار سایر عوامل در محیط پس از یادگیری است. یک شبکه ارتباطی کارآمد بین عوامل می‌تواند یک رویکرد مؤثر برای بهبود همکاری و یادگیری در MARL باشد.



شکل ۲۶- تعامل در محیط‌های چندعاملی (راست) و تعامل در یادگیری Q مستقل (چپ)

MARL یک موضوع مطالعاتی جذاب بوده است که برای برنامه‌های کاربردی در بازی‌های رایانه‌ای و رباتیک مورد توجه زیادی قرار گرفته است. محققان الگوریتم‌های MARL فوق‌العاده‌ای را در سال‌های اخیر پیشنهاد کرده‌اند که کاربرد خود را در زمینه‌های مختلف کاربردی یافته‌اند. تیم تحقیقاتی OpenAI از طریق یک مشکل پنهان و جستجوی چندعاملی و الگوریتم‌های یادگیری استاندارد RL دریافته‌اند که عوامل می‌توانند یک استراتژی اقدام در حال تکامل خود نظارتی ایجاد کنند (محتوای پیشنهادی-[ویدیو](#)). در محیط، آنها شش مرحله در حال تکامل را در استراتژی مقابله با حریفان یافتند. عوامل در بهره‌برداری از کاستی‌های محیط برای دستیابی به پاداش‌های بهتر کاملاً موفق بودند.

<sup>1</sup> Multi agent reinforcement learning

یکی دیگر از افکار جذاب در هنگام بررسی یادگیری تقویتی چندعاملی این است که یادگیری تقویتی سلسله‌مراتبی را می‌توان به‌عنوان مورد دیگری از یادگیری تقویتی چندعاملی در نظر گرفت. چطور؟ به چندین عامل فکر کنید که در افق‌های مختلف تصمیم می‌گیرند. ساختار چند افق مشابه روشی است که اکثر شرکت‌ها در تجارت انجام می‌دهند. افرادی که در بالا هستند، اهداف سطح بالاتری را برای چند سال آینده برنامه‌ریزی می‌کنند و افراد دیگر تصمیم می‌گیرند که چگونه به‌صورت ماهانه و روزانه به این هدف برسند. آنهایی که در بالا هستند اهداف را برای کسانی که در پایین هستند تعیین می‌کنند. کل سیستم برای عملکرد همه عوامل پاداش می‌گیرد.

الگوریتم‌های مطرح: DDPG، MADDPG و Neighborhood Q-learning

## ۴- منابع

R. S. Sutton and A. Barto, *Reinforcement learning : an introduction*(2nd ed.). Cambridge, Ma ; Lodon: The Mit Press, 2018.

Maxim Lapan, *Deep reinforcement learning hands-on : apply modern RL methods to practical problems of chatbots, robotics, discrete optimization, web automation, and more*. Birmingham ; Mumbai Packt January, 2020.

H. Dong, Z. Ding, and S. Zhang, *Deep Reinforcement Learning Fundamentals, Research and Applications*. Singapore: Springer Singapore : Imprint: Springer, 2020.

K. Doya, "Reinforcement Learning," in *The Cambridge Handbook of Computational Cognitive Sciences*, R. Sun, Ed. Cambridge: Cambridge University Press, 2023, pp. 350–370

X. Wang *et al.*, "Deep Reinforcement Learning: A Survey," in *IEEE Transactions on Neural Networks and Learning Systems*, vol. 35, no. 4, pp. 5064-5078, April 2024, doi: 10.1109/TNNLS.2022.3207346

M. Morales and C. Isbell, *Grokking deep reinforcement learning*. Shelter Island, Ny: Manning Publications, 2020.

Alexander Alexander Zai, *Deep Reinforcement Learning in Action*. Manning Publications Company, 2020.

اسلایدهای درس یادگیری تقویتی دانشگاه صنعتی شریف

<https://blog.faradars.org/an-introduction-to-reinforcement-learning/>

<https://parsinfotech.com/reinforcement-learning/>

<https://hooshio.com/?p=5198>

<https://intellabs.github.io/coach/components/agents/index.html>

[یادگیری تقویتی یا Reinforcement Learning و کاربردهایش چیست؟ | کافه تدریس \(cafetadris.com\)](https://cafetadris.com)

محتوای پیشنهادی:

<https://www.youtube.com/@aiwarehouse>

<https://www.youtube.com/watch?v=kopoLzvH5jY>

<https://cs.stanford.edu/people/karpathy/reinforcejs/>

## پیوست: پیاده سازی ها

در این بخش چند نمونه از پیاده سازی های الگوریتم های یادگیری تقویتی ارائه می شود (به صورت فایل).