

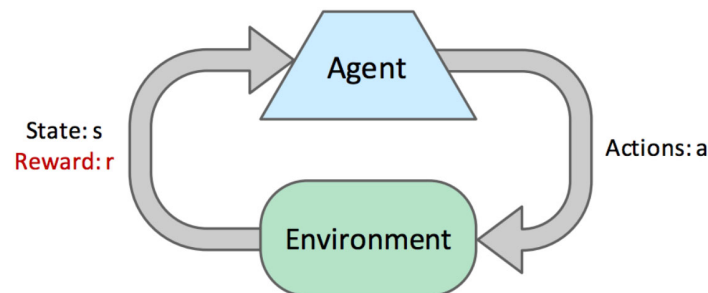
## هوش مصنوعی

## جزوه نهم

یادگیری تقویتی<sup>۱</sup>

[ یادگیری نیروافزوده یا یادگیری تقویتی یا یادگیری پاداش و تانوان یکی از گرایش‌های یادگیری ماشینی است که از روانشناسی رفتارگرایی الهام می‌گیرد. این روش بر رفتارهایی تمرکز دارد که ماشین باید برای بیشینه کردن پاداشش انجام دهد. این مسئله، با توجه به گستردگی‌اش، در زمینه‌های گوناگونی بررسی می‌شود. مانند: نظریه بازی‌ها، نظریه کنترل، تحقیق در عملیات، نظریه اطلاعات، سامانه چندعامله، هوش ازدحامی، آمار، الگوریتم ژنتیک، بهینه‌سازی بر مبنای شبیه‌سازی. یادگیری تقویتی در اقتصاد و نظریه بازیها بیشتر به بررسی تعادل‌های ایجاد شده تحت عقلانیت محدود می‌پردازد. در یادگیری ماشینی با توجه به این که بسیاری از الگوریتم‌های یادگیری نیروافزوده از تکنیک‌های برنامه‌نویسی پویا استفاده می‌کنند معمولاً مسئله تحت عنوان یک فرایند تصمیم‌گیری مارکف مدل می‌شود. تفاوت اصلی بین روش‌های سنتی و الگوریتم‌های یادگیری تقویتی این است که در یادگیری تقویتی نیازی به داشتن اطلاعات راجع به فرایند تصمیم‌گیری ندارد و این که این روش روی فرایندهای مارکف بسیار بزرگی کار می‌کند که روش‌های سنتی در آنجا ناکارآمدند. ]

در یادداشت قبلی، درباره فرآیندهای تصمیم مارکوف بحث کردیم، که با استفاده از تکنیک‌هایی مانند تکرار ارزش و تکرار سیاست برای محاسبه مقادیر بهینه حالت‌ها و استخراج سیاست‌های بهینه، آن‌ها را حل کردیم. حل فرآیندهای تصمیم مارکوف نمونه‌ای از برنامه‌ریزی آفلاین است، که در آن عوامل، اطلاعات کاملی از تابع انتقال و تابع پاداش دارند، تمامی اطلاعاتی که برای پیش محاسبه اقدامات بهینه در جهان نیاز دارند، که توسط MDP کدگذاری شده‌اند، بدون اینکه هیچ اقدامی انجام دهند در دسترس آنهاست. در این یادداشت، درباره برنامه‌ریزی آنلاین بحث می‌کنیم که طی آن یک نماینده هیچ دانش قبلی از پاداش‌ها یا انتقال‌ها در جهان ندارد (که هنوز آنرا به عنوان یک MDP نشان می‌دهیم). در برنامه‌ریزی آنلاین، یک عامل باید جستجویی<sup>۲</sup> را امتحان کند، که در طی آن اقداماتی را انجام می‌دهد و بازخوردی را در قالب حالت‌های جانشینی که وارد می‌شود و پاداش‌های مربوطه دریافت می‌کند، می‌گیرد. عامل از این بازخورد برای تخمین یک خط مشی بهینه، از طریق فرآیندی به نام یادگیری تقویتی قبل از استفاده از این خط مشی تخمین زده شده برای بهره برداری یا به حداکثر رساندن پاداش استفاده می‌کند.



با برخی از اصطلاحات اولیه شروع می‌کنیم. در هر مرحله زمانی در طول برنامه ریزی آنلاین، یک عامل در یک حالت  $s$  شروع می‌کند، سپس یک عمل  $a$  را انجام می‌دهد و در حالت جانشین  $s'$  تمام می‌کند و پاداش  $r$  را دریافت می‌کند. هر  $(s, a, s', r)$  به عنوان نمونه<sup>۳</sup> شناخته می‌شود. اغلب، یک عامل به

<sup>1</sup> Reinforcement Learning<sup>2</sup> exploration<sup>3</sup> sample

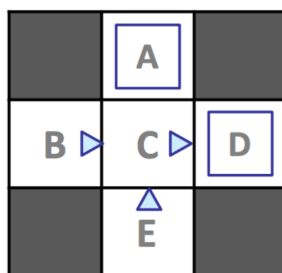
انجام اقدامات و جمع آوری نمونه های متوالی تا رسیدن به حالت پایانی ادامه می دهد. چنین مجموعه ای از نمونه ها به عنوان یک قسمت<sup>4</sup> شناخته می شود. عوامل معمولاً در حین اکتشاف قسمت های زیادی را پشت سر می گذارند تا داده های کافی مورد نیاز برای یادگیری را جمع آوری کنند.

دو نوع یادگیری تقویتی داریم: یادگیری مبتنی بر مدل<sup>5</sup> و یادگیری بدون مدل<sup>6</sup>. یادگیری مبتنی بر مدل تلاش می کند تا توابع انتقال و پاداش را با نمونه های به دست آمده در طول اکتشاف، قبل از استفاده از این تخمین ها برای حل MDP به طور معمول با تکرار ارزش یا تکرار خط مشی، تخمین بزند. از سوی دیگر، یادگیری بدون مدل، تلاش می کند تا مقادیر یا مقادیر Q حالت ها را مستقیماً تخمین بزند، بدون اینکه از هیچ حافظه ای برای ساخت مدلی از پاداش ها و انتقال ها در MDP استفاده کند.

### یادگیری مبتنی بر مدل<sup>7</sup>

[در یادگیری تقویتی مبتنی بر مدل، عامل یادگیری تقویتی که مدلی از محیط را یاد می گیرند و از آن مدل برای ایجاد یک قانون کنترلی بدون آزمون و خطای مستقیم استفاده می نمایند.]

در یادگیری مبتنی بر مدل، عامل تقریبی از تابع انتقال  $\hat{T}(s,a,s')$  را با نگه داشتن تعداد دفعاتی که پس از رسیدن به هر حالت  $Q(s,a)$  به حالت  $s'$  می رسد محاسبه می کند. سپس عامل می تواند تابع انتقال تقریبی  $\hat{T}$  را با عادی سازی شمارش هایی که جمع کرده است محاسبه کند، سپس تعداد هر  $(s,a,s')$  مشاهده شده را بر مجموع شمارش ها برای همه مواردی که عامل در حالت  $Q(s,a)$  قرار داشته است تقسیم می کند. عادی سازی شمارش ها آنها را به گونه ای مقیاس می کند که مجموع آنها به یک برسد و بتوان آنها را به عنوان احتمالات تفسیر کرد. مثال MDP زیر را با حالت های  $S = \{A,B,C,D,E,x\}$  در نظر بگیرید که  $x$  نشان دهنده حالت پایانی است و ضریب کاهش  $\gamma = 1$  است:



فرض کنید به عامل اجازه بدهیم تا MDP را برای چهار قسمت تحت خط مشی  $\pi$  explore که در بالا مشخص شده است کاوش کند (مثلاً جهت دار حرکت در جهتی را که مثلث نشان می دهد و مربع های آبی نشان دهنده خروج به عنوان اقدام انتخابی است) را نشان می دهد و نتایج زیر را به دست می آید:

<sup>4</sup> episode

<sup>5</sup> model-based learning

<sup>6</sup> model-free learning

<sup>7</sup> Model-Based Learning

### Episode 1

B, east, C, -1  
C, east, D, -1  
D, exit, x, +10

### Episode 2

B, east, C, -1  
C, east, D, -1  
D, exit, x, +10

### Episode 3

E, north, C, -1  
C, east, D, -1  
D, exit, x, +10

### Episode 4

E, north, C, -1  
C, east, A, -1  
A, exit, x, -10

اکنون ۱۲ نمونه جمعی داریم، ۳ نمونه از هر بخش با شمارش به شرح زیر است:

s	a	s'	count
A	exit	x	1
B	east	C	2
C	east	A	1
C	east	D	3
D	exit	x	3
E	north	C	2

می دانیم که  $T(s, a, s') = P(s'|a, s)$  پس می توانیم تابع انتقال را با این شمارش ها با تقسیم تعداد هر  $(s, a, s')$  بر تعداد کل مواقعی که در حالت  $Q(s, a)$  قرار داشتیم و تابع پاداش مستقیماً از پاداش هایی که در طول اکتشاف به دست آوردیم، تخمین بزنیم:

#### • Transition Function: $\hat{T}(s, a, s')$

$$\begin{aligned} - \hat{T}(A, \text{exit}, x) &= \frac{\#(A, \text{exit}, x)}{\#(A, \text{exit})} = \frac{1}{1} = 1 \\ - \hat{T}(B, \text{east}, C) &= \frac{\#(B, \text{east}, C)}{\#(B, \text{east})} = \frac{2}{2} = 1 \\ - \hat{T}(C, \text{east}, A) &= \frac{\#(C, \text{east}, A)}{\#(C, \text{east})} = \frac{1}{4} = 0.25 \\ - \hat{T}(C, \text{east}, D) &= \frac{\#(C, \text{east}, D)}{\#(C, \text{east})} = \frac{3}{4} = 0.75 \\ - \hat{T}(D, \text{exit}, x) &= \frac{\#(D, \text{exit}, x)}{\#(D, \text{exit})} = \frac{3}{3} = 1 \\ - \hat{T}(E, \text{north}, C) &= \frac{\#(E, \text{north}, C)}{\#(E, \text{north})} = \frac{2}{2} = 1 \end{aligned}$$

#### • Reward Function: $\hat{R}(s, a, s')$

$$\begin{aligned} - \hat{R}(A, \text{exit}, x) &= -10 \\ - \hat{R}(B, \text{east}, C) &= -1 \\ - \hat{R}(C, \text{east}, A) &= -1 \\ - \hat{R}(C, \text{east}, D) &= -1 \\ - \hat{R}(D, \text{exit}, x) &= +10 \\ - \hat{R}(E, \text{north}, C) &= -1 \end{aligned}$$

طبق قانون اعداد بزرگ<sup>۸</sup>، هرچقدر که عامل نمونه های بیشتری را با تجربه کردن قسمت های بیشتر جمع آوری کند، مدل های  $\hat{T}$  و  $\hat{R}$  ما بهبود می یابند، و به سمت  $T$  و  $R$  همگرا می شود، و با کشف  $(s, a, s')$  های جدید، دانش پاداش هایی قبلاً کشف نشده اند به دست می آیند. هر زمان که بخواهیم، می توانیم آموزش عامل خود را برای ایجاد یک سیاست  $\pi$  exploit با اجرای روش تکرار ارزش یا تکرار خط مشی با مدل های فعلی خود برای  $\hat{T}$  و  $\hat{R}$  پایان دهیم و از  $\pi$  exploit برای بهره برداری استفاده کنیم، و از عامل بخواهیم تا MDP را طی کند و اقداماتی را برای به حداکثر رساندن پاداش انجام دهد به جای اینکه به دنبال یادگیری باشد. به زودی روش هایی را برای چگونگی تخصیص زمان موثر بین اکتشاف و بهره برداری را ارائه می کنیم. یادگیری مبتنی بر مدل بسیار

<sup>8</sup> law of large numbers

ساده و شهودی و در عین حال بسیار مؤثر است و  $\hat{T}$  و  $\hat{R}$  را با شمارش و عادی سازی تولید می کند. با این حال، حفظ شمارش برای هر  $(s, a, s')$  که بررسی می شود می تواند هزینه بر باشد، و بنابراین در بخش بعدی در مورد یادگیری بدون مدل، روش هایی را برای دور زدن کل حفظ شمارش و جلوگیری از سربار حافظه مورد نیاز روش یادگیری مبتنی بر مدل شرح خواهیم داد.

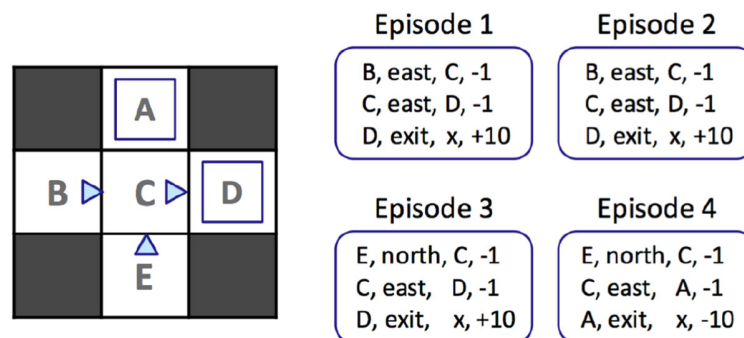
### یادگیری بدون مدل<sup>۹</sup>

[در یادگیری تقویتی بدون مدل، عامل یادگیری تقویتی که به مدلی از محیط متکی نیستند. آنها از تجربه مستقیم در تعامل با محیط یادگیری را انجام می دهند.]

چندین الگوریتم یادگیری بدون مدل وجود دارد، و ما به سه مورد از آنها خواهیم پرداخت: ارزیابی مستقیم (direct evaluation)، یادگیری تفاوت زمانی (temporal difference learning) و یادگیری Q (Q-learning). ارزیابی مستقیم و یادگیری تفاوت زمانی تحت دسته ای از الگوریتم ها قرار می گیرند که به عنوان یادگیری تقویتی غیرفعال<sup>۱۰</sup> شناخته می شوند. در یادگیری تقویتی غیرفعال، به یک عامل خط مشی داده می شود تا دنبال کند و ارزش ارزش حالت های تحت آن خط مشی را در حین تجربه قسمت ها بیاموزد، این دقیقاً همان کاری است که با ارزیابی خط مشی برای MDP ها زمانی که  $T$  و  $R$  را داشتیم، انجام می شد. یادگیری Q تحت دسته دومی از الگوریتم های یادگیری بدون مدل قرار می گیرد که به عنوان یادگیری تقویتی فعال<sup>۱۱</sup> شناخته می شوند، که طی آن عامل می تواند از بازخوردی که دریافت می کند برای به روزرسانی مکرر خط مشی خود در حین یادگیری استفاده کند تا سرانجام پس از کاوش کافی، خط مشی بهینه را تعیین کند.

### ارزیابی مستقیم

اولین تکنیک یادگیری تقویتی غیرفعال که به آن می پردازیم، ارزیابی مستقیم است، روشی که مانند نامش خسته کننده و ساده است. تنها کاری که ارزیابی مستقیم انجام می دهد این است که برخی از خط مشی های  $\pi$  را اصلاح می کند و نماینده چندین قسمت را حین پیروی از  $\pi$  تجربه می کند. همانطور که عامل از طریق این قسمت ها نمونه ها را جمع آوری می کند، تعداد کل سودمندی به دست آمده از هر حالت و تعداد دفعاتی که از هر حالت بازدید کرده است را نیز حفظ می کند. در هر نقطه، می توانیم ارزش تخمینی هر حالت  $S$  را با تقسیم کل سودمندی به دست آمده از حالت  $S$  بر تعداد دفعاتی که حالت  $S$  بازدید شده است، محاسبه کنیم. اکنون ارزیابی مستقیم را بر روی مثال قبلی خود اجرا می کنیم، به یاد داشته باشید که  $\gamma = 1$  است.



با گذر از قسمت اول، می بینیم که از حالت D تا پایان، پاداش کلی ۱۰، از وضعیت C پاداش کلی  $9 = 10 + (-1)$ ، و از حالت B پاداش کلی  $8 = 10 + (-1) + (-1)$  را به دست می آوریم. با تکمیل این فرآیند، پاداش کل در هر قسمت برای هر حالت و مقادیر تخمینی حاصل به شرح زیر به دست می آیند:

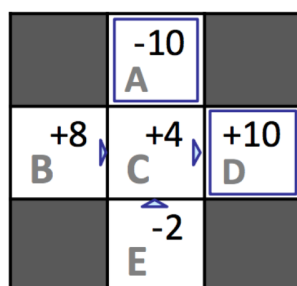
<sup>۹</sup> Model-Free Learning

<sup>۱۰</sup> passive reinforcement learning

<sup>۱۱</sup> active reinforcement learning

s	Total Reward	Times Visited	$V^\pi(s)$
A	-10	1	-10
B	16	2	8
C	16	4	4
D	30	3	10
E	-4	2	-2

اگرچه ارزیابی مستقیم در نهایت مقادیر حالت را برای هر حالت بررسی می‌کند، اما اغلب به کندی همگرا می‌شود زیرا اطلاعات مربوط به انتقال بین حالت‌ها را هدر می‌دهد.



در این مثال،  $V^\pi(B) = 8$  و  $V^\pi(E) = -2$  را محاسبه کردیم، اگرچه بر اساس بازخوردی که دریافت کردیم، هر دو حالت فقط C را به عنوان حالت جانشین دارند و در هنگام انتقال به C، پاداش یکسان 1- را متحمل می‌شوند. با توجه به معادله بلمن، این بدان معنی است که B و E باید هر دو مقدار یکسانی را تحت خط مشی  $\pi$  داشته باشند. با این حال، از 4 باری که عامل ما در وضعیت C بود، به D منتقل شد و سه بار پاداش 10 را کسب کرد و به A منتقل شد و یک بار پاداش 10- را کسب کرد. این کاملاً تصادفی بوده است که هنگامی که پاداش 10- را دریافت کرد از حالت E به جای B شروع کرده بود، اما این به شدت ارزش تخمینی را برای E تغییر داده است. با تعداد قسمت‌های کافی، مقادیر B و E به مقادیر واقعی خود همگرا می‌شوند، اما مواردی مانند این باعث می‌شوند که فرآیند بیش از آنچه ما می‌خواهیم طول بکشد. این مشکل را می‌توان با استفاده از دومین الگوریتم یادگیری تقویتی غیرفعال، یعنی یادگیری تفاوت زمانی، کاهش داد.

### یادگیری تفاوت زمانی

[ یادگیری تفاوت زمانی (Temporal difference learning) یک روش پیش‌بینی است. این روش به صورت عمده برای حل مسائل یادگیری تقویتی مورد استفاده بود است. روش تفاوت زمانی ترکیبی از ایده‌های مونت کارلو و برنامه‌ریزی پویا است. این روش مشابه روش مونت کارلو است چرا که یادگیری در آن با استفاده از نمونه برداری از محیط با توجه به یک یا چند سیاست خاص انجام می‌شود. روش تفاوت زمانی به این دلیل به تکنیک‌های برنامه‌ریزی پویا شباهت دارد که این روش تخمین کنونی را بر اساس تخمین‌های یادگیری شده به دست می‌آورد. به عنوان یک روش پیش‌بینی، یادگیری تفاوت زمانی این واقعیت را در نظر می‌گیرد که پیش‌بینی‌های آینده نیز معمولاً از جهاتی دارای همبستگی هستند. در روش‌های یادگیری مبتنی بر پیش‌بینی نظارتی، عامل تنها از مقادیر دقیقاً مشاهده شده یاد می‌گیرد: یک پیش‌بینی انجام می‌شود، و زمانی که مشاهده ممکن باشد، پیش‌بینی به تطابق بهتری با مشاهده خواهد رسید. ایده اساسی یادگیری تفاوت زمانی این است که پیش‌بینی‌ها با پیش‌بینی‌هایی دقیق‌تر دیگری از آینده تنظیم شوند. ]

یادگیری تفاوت زمانی (TD Learning) از ایده یادگیری از هر تجربه، به جای اینکه صرفاً کل پاداش‌ها و تعداد دفعات بازدید از حالات و یادگیری در پایان را مانند روش ارزیابی مستقیم پیگیری کند، استفاده می‌کند. در ارزیابی خط مشی، ما از سیستم معادلات تولید شده توسط خط مشی ثابت خود و معادله بلمن برای تعیین مقادیر حالت‌های تحت آن خط مشی استفاده کردیم (یا از به روز رسانی‌های تکراری مانند روش تکرار ارزش استفاده کردیم).

$$V^\pi(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

هر یک از این معادلات، ارزش یک حالت را برابر با میانگین وزنی نسبت به مقادیر کاهش یافته جانشینان آن حالت به علاوه پاداش های به دست آمده در انتقال به آنها قرار می دهد. یادگیری تفاوت زمانی سعی می کند به این سوال پاسخ دهد که چگونه می توان این میانگین وزنی را بدون وزن محاسبه کرد و این کار را هوشمندانه با میانگین متحرک نمایی<sup>12</sup> انجام می دهد. با مقداردهی اولیه  $\forall s, V^\pi(s) = 0$  شروع می کنیم. در هر گام زمانی، یک عامل یک عمل  $\pi(s)$  را از حالت  $s$  انجام می دهد، به حالت  $s'$  منتقل می شود و پاداش  $R(s, \pi(s), s')$  را کسب می کند. می توانیم با جمع کردن پاداش دریافتی با مقدار فعلی کاهش یافته  $s'$  تحت  $\pi$ ، یک ارزش نمونه<sup>13</sup> بدست بیاوریم:

$$\text{sample} = R(s, \pi(s), s') + \gamma V^\pi(s')$$

این ارزش نمونه، تخمین جدیدی برای  $V^\pi(s)$  است. گام بعدی این است که این تخمین نمونه را در مدل موجود خود برای  $V^\pi(s)$  با میانگین متحرک نمایی، که به قانون به روز رسانی زیر پایبند است، بگنجانیم:

$$V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + \alpha \cdot \text{sample}$$

در معادله بالا،  $\alpha$  پارامتری است با محدوده  $0 \leq \alpha \leq 1$  که به عنوان نرخ یادگیری شناخته می شود، که وزنی را که می خواهیم مدل موجود خود را برای  $V^\pi(s)$ ،  $1 - \alpha$  اختصاص دهیم، و وزنی را که می خواهیم تخمین نمونه جدید خود یعنی  $\alpha$  را به آن اختصاص دهیم، مشخص می کند. به طور معمول با نرخ یادگیری  $\alpha = 1$  شروع می کنیم، بر این اساس  $V^\pi(s)$  را به هر نمونه اولیه نسبت می دهیم، و به آرامی آن را به سمت  $\cdot$  کاهش می دهیم، در این مرحله تمام نمونه های بعدی صفر می شوند و دیگر تاثیری بر مدل ما یعنی  $V^\pi(s)$  ندارند.

اکنون قانون به روز رسانی را تجزیه و تحلیل می کنیم. با تشریح وضعیت مدل خود در مقاطع مختلف زمانی با تعریف  $V_k^\pi(s)$  و  $\text{sample}_k$  به عنوان مقدار تخمینی حالت  $s$  پس از به روزرسانی  $k$ ام و  $k$ امین نمونه، می توانیم قانون به روزرسانی خود را دوباره بیان کنیم:

$$V_k^\pi(s) \leftarrow (1 - \alpha)V_{k-1}^\pi(s) + \alpha \cdot \text{sample}_k$$

این تعریف بازگشتی را برای  $V_k^\pi(s)$  بسط می دهیم:

$$\begin{aligned} V_k^\pi(s) &\leftarrow (1 - \alpha)V_{k-1}^\pi(s) + \alpha \cdot \text{sample}_k \\ V_k^\pi(s) &\leftarrow (1 - \alpha)[(1 - \alpha)V_{k-2}^\pi(s) + \alpha \cdot \text{sample}_{k-1}] + \alpha \cdot \text{sample}_k \\ V_k^\pi(s) &\leftarrow (1 - \alpha)^2 V_{k-2}^\pi(s) + (1 - \alpha) \cdot \alpha \cdot \text{sample}_{k-1} + \alpha \cdot \text{sample}_k \\ &\vdots \\ V_k^\pi(s) &\leftarrow (1 - \alpha)^k V_0^\pi(s) + \alpha \cdot [(1 - \alpha)^{k-1} \cdot \text{sample}_1 + \dots + (1 - \alpha) \cdot \text{sample}_{k-1} + \text{sample}_k] \\ V_k^\pi(s) &\leftarrow \alpha \cdot [(1 - \alpha)^{k-1} \cdot \text{sample}_1 + \dots + (1 - \alpha) \cdot \text{sample}_{k-1} + \text{sample}_k] \end{aligned}$$

از آنجایی که  $0 \leq (1 - \alpha) \leq 1$  است، با افزایش مقدار  $(1 - \alpha)$  به توان های بزرگتر، این مقدار به  $\cdot$  نزدیک تر و نزدیک تر می شود. با بسط قانون به روز رسانی ای که به دست آوردیم، این بدان معناست که به نمونه های قدیمی به طور تصاعدی وزن کمتری داده می شود، دقیقاً همان چیزی که می خواهیم زیرا این نمونه های قدیمی با استفاده از نسخه های قدیمی تر (و در نتیجه بدتر) مدل ما برای  $V_k^\pi(s)$  محاسبه می شوند! این خوبی روش یادگیری تفاوت های زمانی است، با یک قانون به روز رسانی ساده، ما می توانیم:

- در هر مرحله زمانی یاد می گیریم، از این رو از اطلاعات مربوط به انتقال حالت به هنگام دریافت آنها استفاده می کنیم، زیرا ما از نسخه های به روز شده مکرر  $V^\pi(s')$  در نمونه های خود استفاده می کنیم به جای اینکه برای انجام هر محاسباتی تا پایان منتظر بمانیم.
- به نمونه های قدیمی تر و بالقوه کمتر، وزن کمتری می دهیم.

<sup>12</sup> Exponential moving average

<sup>13</sup> sample value

- با قسمت های کمتر نسبت به ارزیابی مستقیم، یادگیری مقادیر حالت واقعی خیلی سریعتر همگرا می شوند.

## یادگیری Q

[ یادگیری Q، یک تکنیک یادگیری تقویتی است که با یادگیری یک تابع اقدام/مقدار، سیاست مشخصی را برای انجام حرکات مختلف در وضعیت های مختلف دنبال می کند. یکی از نقاط قوت این روش، توانایی یادگیری تابع مذکور بدون داشتن مدل معینی از محیط می باشد. یادگیری Q در تلاش است با توجه به شرایط فعلی، بهترین اقدامات را انجام دهد. این الگوریتم خارج از خط مشی در نظر گرفته می شود. چراکه، تابع یادگیری Q از اقداماتی خارج از خط مشی فعلی یاد می گیرد. به طور کلی می توان گفت که، یادگیری Q به دنبال یادگیری خط مشی است تا مجموع پاداش را بیشینه کند.]

هم ارزیابی مستقیم و هم یادگیری TD در نهایت ارزش واقعی همه حالت ها را تحت سیاستی که دنبال می کنند یاد خواهند گرفت. با این حال، هر دوی آنها یک مسئله ذاتی اصلی دارند، ما می خواهیم یک خط مشی بهینه برای عامل خود پیدا کنیم که نیاز به آگاهی از مقادیر Q حالت ها دارد. برای محاسبه مقادیر Q از مقادیری که داریم، به یک تابع انتقال و تابع پاداش که توسط معادله بلمن ایجاد می شود، نیاز داریم.

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

در نتیجه، یادگیری TD یا ارزیابی مستقیم معمولاً همراه با برخی از یادگیری های مبتنی بر مدل برای بدست آوردن تخمین T و R به منظور به روز رسانی موثر خط مشی دنبال شده توسط عامل در حال یادگیری استفاده می شوند. این امر با ایده انقلابی جدید به نام یادگیری Q، که یادگیری مستقیم مقادیر Q حالت ها را پیشنهاد می کرد، اجتناب پذیر شد و نیاز به دانستن هر گونه ارزش، توابع انتقال یا توابع پاداش را دور زد. در نتیجه، یادگیری Q کاملاً یک یادگیری بدون مدل است. یادگیری Q از قانون به روز رسانی زیر برای انجام آنچه به عنوان تکرار ارزش Q<sup>۱۴</sup> شناخته می شود استفاده می کند:

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

توجه کنید که این به روز رسانی فقط یک تغییر جزئی در قانون به روز رسانی برای تکرار ارزش است. در واقع، تنها تفاوت این است که موقعیت حداکثر عملگر بر روی اقدامات تغییر کرده است، زیرا ما یک عمل را قبل از انتقال، زمانی که در یک حالت هستیم انتخاب می کنیم، اما قبل از انتخاب یک عمل جدید، زمانی که در یک حالت Q هستیم، انتقال را انجام می دهیم.

با این قانون به روز رسانی جدید تحت، یادگیری Q اساساً به همان روش یادگیری TD، با به دست آوردن نمونه های ارزش Q<sup>۱۵</sup> مشتق می شود:

$$\text{sample} = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$

و ترکیب آنها در یک میانگین متحرک نمایی به صورت زیر خواهد شد:

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha \cdot \text{sample}$$

تا زمانی که زمان کافی را صرف کاوش کنیم و نرخ یادگیری  $\alpha$  را با سرعت مناسب کاهش دهیم، یادگیری Q ارزش های Q بهینه را برای هر حالت Q می آموزد. این همان چیزی است که یادگیری Q را بسیار انقلابی می کند، در حالی که یادگیری TD و ارزیابی مستقیم ارزش های حالت های تحت یک خط مشی را با پیروی از خط مشی، قبل از تعیین بهینه بودن خط مشی از طریق تکنیک های دیگر یاد می گیرند، یادگیری Q می تواند مستقیماً خط مشی

<sup>14</sup> Q-value iteration

<sup>15</sup> Q-value samples

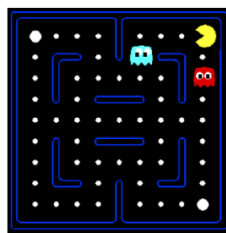
بهینه را حتی با استفاده از روش‌های زیر بهینه یا اقدامات تصادفی یاد بگیرد. این را یادگیری خارج از سیاست<sup>۱۶</sup> می‌نامند (برخلاف ارزیابی مستقیم و یادگیری TD که نمونه‌هایی از یادگیری درون سیاست<sup>۱۷</sup> هستند).

### یادگیری Q تقریبی

یادگیری Q یک تکنیک یادگیری باورنکردنی است که همچنان در مرکز تحولات در زمینه یادگیری تقویتی قرار دارد. با این حال، هنوز هم جای پیشرفت دارد. یادگیری Q فقط تمام ارزش‌های Q را برای حالت‌ها به شکل جدول ذخیره می‌کند، که با توجه به اینکه اکثر برنامه‌های یادگیری تقویتی چندین هزار یا حتی میلیون‌ها حالت دارند، کارایی خاصی ندارد. این بدان معناست که ما نمی‌توانیم در طول آموزش از همه حالت‌ها بازدید کنیم و نمی‌توانیم تمام ارزش‌های Q را به دلیل کمبود حافظه ذخیره کنیم.



شکل ۱



شکل ۲



شکل ۳

در بالا، اگر یک من یاد گرفته بود که شکل ۱ پس از اجرای یادگیری Q نامطلوب است، باز هم متوجه نمی‌شد که شکل ۲ یا حتی شکل ۳ نیز نامطلوب هستند. یادگیری Q تقریبی تلاش می‌کند این مشکل را با یادگیری در مورد چند موقعیت کلی و برون‌یابی به بسیاری از موقعیت‌های مشابه توضیح دهد. کلید تعمیم تجارب یادگیری، نمایش حالت‌های مبتنی بر ویژگی<sup>۱۸</sup> است که هر حالت را به عنوان یک بردار نشان می‌دهد که به عنوان بردار ویژگی<sup>۱۹</sup> شناخته می‌شود. به عنوان مثال، یک بردار ویژگی برای یک من ممکن است به صورت‌های زیر رمزگذاری شود:

- فاصله تا نزدیکترین روح
- فاصله تا نزدیکترین گلوله غذا
- تعداد ارواح
- آیا یکمن به دام افتاده است؟ 0 یا ۱

با بردارهای ویژگی، می‌توانیم ارزش حالت‌ها و حالت‌های Q را به عنوان توابع ارزش خطی<sup>۲۰</sup> در نظر بگیریم:

$$V(s) = w_1 \cdot f_1(s) + w_2 \cdot f_2(s) + \dots + w_n \cdot f_n(s) = \vec{w} \cdot \vec{f}(s)$$

$$Q(s, a) = w_1 \cdot f_1(s, a) + w_2 \cdot f_2(s, a) + \dots + w_n \cdot f_n(s, a) = \vec{w} \cdot \vec{f}(s, a)$$

که در آن

$$\vec{f}(s) = [f_1(s) \quad f_2(s) \quad \dots \quad f_n(s)]^T$$

<sup>16</sup> off-policy learning

<sup>17</sup> on-policy learning

<sup>18</sup> feature-based representation

<sup>19</sup> feature vector

<sup>20</sup> linear value functions



$$\vec{f}(s,a) = [f_1(s,a) \quad f_2(s,a) \quad \dots \quad f_n(s,a)]^T$$

به ترتیب نشان دهنده بردارهای ویژگی برای حالت  $s$  و حالت  $Q(s,a)$  و

$$\vec{w} = [w_1 \quad w_2 \quad \dots \quad w_n]$$

یک بردار وزن را نشان می دهد. تفاوت به صورت زیر بیان می کنیم:

$$\text{difference} = [R(s,a,s') + \gamma \max_{a'} Q(s',a')] - Q(s,a)$$

یادگیری تقریبی  $Q$ ، تقریباً مشابه یادگیری  $Q$  با استفاده از قانون به روز رسانی زیر عمل می کند:

$$w_i \leftarrow w_i + \alpha \cdot \text{difference} \cdot f_i(s,a)$$

به جای ذخیره کردن ارزش های  $Q$  برای هر حالت، با یادگیری تقریبی  $Q$  فقط به ذخیره یک بردار وزنی نیاز داریم و می توانیم ارزش های  $Q$  را در صورت نیاز محاسبه کنیم. در نتیجه، این نه تنها نسخه تعمیم یافته تری از یادگیری  $Q$  را در اختیار ما قرار می دهد، بلکه نسخه ای از نظر حافظه کارآمدتر نیز به ما می دهد.

به عنوان آخرین نکته در مورد یادگیری  $Q$ ، می توانیم قانون به روز رسانی را برای یادگیری دقیق  $Q$  با استفاده از تفاوت به صورت زیر بیان کنیم:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot \text{difference}$$

این نماد دوم تفسیر کمی متفاوت اما به همان اندازه ارزشمند از به روز رسانی را به ما می دهد: این محاسبه تفاوت بین مدل تخمینی نمونه و مدل فعلی  $Q(s,a)$  است، و مدل را در جهت تخمین جابجا می کند تا بزرگی تغییر، متناسب با بزرگی تفاوت باشد.