

S	(۱,۱)	(۱,۲)	(۱,۳)	(۲,۱)	(۲,۲)	(۲,۳)
T_0	۰	۰	-۵	۰	۰	+۵
T_1	۰	۰	-۵	۰	۳,۹	+۵
T_2	۰	۲,۱۴	-۵	۲,۵۹	۲,۹۴	+۵

* مقادیر جدول را صورت کمال مستطین بود (جای ۵+۵) صورت کمال مستطین بود (جای ۵+۵) از خانه کار مستطین صورت کمال مستطین بود (جای ۵+۵) از خانه کار مستطین صورت کمال مستطین بود (جای ۵+۵)

$$T_{i,j}(S) = \max_{\alpha} \left(\sum_{S'} T(S, \alpha, S') (R(S, \alpha, S') + \gamma V_i(S')) \right)$$

$$T_1(1,1) = \max_{\alpha} \left[\sum_{S'} T((1,1), \alpha, S') (0 + \gamma T_0(S')) \right]$$

$$= \max_{\alpha} \left[\sum_{S'} 0.9 \times T((1,1), L, S') \times 0 \right] = \max \left[\begin{matrix} 0 \\ 0.9 \times 0.1 \times (-5) \\ 0.9 \times 0.1 \times (-5) \\ 0.9 \times 0.1 \times (-5) \end{matrix} \right] = 0$$

$$\begin{matrix} T((1,1), R, (1,3)) \times T_0(1,3) \\ T((1,1), U, (1,2)) \times T_0(1,2) \\ T((1,1), D, (1,1)) \times T_0(1,1) \end{matrix}$$

$$T_1(2,1) = \max_{\alpha} \left(\sum_{S'} T((2,1), \alpha, S') (0 + \gamma T_0(S')) \right)$$

$$= \max_{S'} \left(\sum_{S'} 0.9 \times T((2,1), L, S') \times 0 \right)$$

$$T((2,1), R, (2,3)) \times T_0(2,3)$$

$$T((2,1), U, (2,2)) \times T_0(2,2)$$

$$T((2,1), D, (2,1)) \times T_0(2,1)$$

$$= \max \left[\begin{matrix} 0 \\ 0.9 \times 0.1 \times 5 \\ 0.9 \times 0.1 \times 5 \\ 0.9 \times 0.1 \times 5 \end{matrix} \right] = 3, 9$$

بجای صفر جدول را با ۵+۵ جایگزین کردیم (مستطین)

$$T_2(1,1) = \max_{\alpha} \left(\sum_{S'} T((1,1), \alpha, S') (0 + \gamma T_1(S')) \right) = \max_{\alpha} \left(0.9 \times \left\{ \begin{matrix} T((1,1), L, (1,1)) T_1(1,1) \\ T((1,1), R, (1,3)) T_1(1,3) \\ T((1,1), D, (1,1)) T_1(1,1) \\ T((1,1), U, (1,2)) T_1(1,2) \end{matrix} \right\} \right)$$

$$= \max \left[\begin{matrix} 0.9 \times [0.1 \times 0 + 0.1 \times 0] \\ 0.9 \times [0.1 \times (-5) + 0.1 \times 3, 9 + 0.1 \times 0] \\ 0.9 \times [0.1 \times 0 + 0.1 \times (-5)] \\ 0.9 \times [0.1 \times 3, 9 + 0.1 \times (-5)] \end{matrix} \right] = \max \left[\begin{matrix} 0 \\ 0.9 \times 0.1 \times 3, 9 \\ -0.9 \times 0.1 \times 5 \\ 0.9 \times 0.1 \times 3, 9 \end{matrix} \right] = 2, 14$$

$$T_2(2,1) = \max_{\alpha} \left(\sum_{S'} T((2,1), \alpha, S') (0 + \gamma T_1(S')) \right) = \max_{\alpha} \left(0.9 \times \left\{ \begin{matrix} T((2,1), L, (2,1)) T_1(2,1) \\ T((2,1), R, (2,3)) T_1(2,3) \\ T((2,1), D, (2,1)) T_1(2,1) \\ T((2,1), U, (2,2)) T_1(2,2) \end{matrix} \right\} \right)$$

$$= \max \left[\begin{matrix} 0.9 \times [0.1 \times 0 + 0.1 \times 0] \\ 0.9 \times [0.1 \times 0 + 0.1 \times 3, 9 + 0.1 \times 0] \\ 0.9 \times [0.1 \times 0 + 0.1 \times (-5)] \\ 0.9 \times [0.1 \times 3, 9 + 0.1 \times (-5)] \end{matrix} \right] = \max \left[\begin{matrix} 0 \\ 0.9 \times 0.1 \times 3, 9 \\ -0.9 \times 0.1 \times 5 \\ 0.9 \times 0.1 \times 3, 9 \end{matrix} \right] = 2, 14$$

$$T_2(2,2) = \max_{\alpha} \left(\sum_{S'} T((2,2), \alpha, S') (0 + \gamma T_1(S')) \right) = \max_{\alpha} \left(0.9 \times \left\{ \begin{matrix} T((2,2), U, (2,3)) T_1(2,3) \\ T((2,2), D, (2,1)) T_1(2,1) \\ T((2,2), L, S') \times 0 \\ T((2,2), R, (2,3)) T_1(2,3) \end{matrix} \right\} \right)$$

$$= \max \left[\begin{matrix} 0.9 \times 0.1 \times 3, 9 \\ 0.9 \times 0.1 \times 3, 9 \\ 0 \\ 0.9 \times 0.1 \times 3, 9 \end{matrix} \right] = \max \left[\begin{matrix} 0.9 \times 0.1 \times 3, 9 \\ 0.9 \times 0.1 \times 3, 9 \\ 0 \\ 0.9 \times 0.1 \times 3, 9 \end{matrix} \right] = 2, 59$$

V Policy

S	(1,1)	(1,2)	(1,3)	(2,1)	(2,2)	(2,3)
$\pi^*(S)$	U	U	-	R	R	-

$$V(1,1) = \frac{-0.5 + 0}{0.5} = -1,9V$$

$$V(2,1) = \frac{0 + 0}{1} = 0$$

$$k) \underline{V_R(s)} = \underline{V_R(s)} + \alpha [R(s, R(s), s') + \gamma \underline{V_R(s')} - \underline{V_R(s)}]$$

$$\hookrightarrow \underline{V_R(1,1)} = \underline{V_R(1,1)} + \alpha [R(1,1, R(1,1)) + \gamma \underline{V_R(1,1)} - \underline{V_R(1,1)}]$$

$$= 0 + 0,1 \times [0 + 0,9 \times 0 - 0] = 0$$

$$\underline{V_R(1,1)} = \underline{V_R(1,1)} + \alpha [R(1,1, R(1,1)) + \gamma \underline{V_R(1,1)} - \underline{V_R(1,1)}]$$

$$= 0 + 0,1 \times [0 + 0,9 \times (-0,1 F \Delta) - 0] = \underline{-0,009 F \Delta}$$

$$\underline{V_R(1,1)} = \underline{V_R(1,1)} + \alpha [R(1,1, R(1,1)) + \gamma \underline{V_R(1,1)} - \underline{V_R(1,1)}]$$

$$= 0 + 0,1 \times [0 + 0,9 \times (-0,009 F \Delta) - 0] = \underline{-0,00081 F \Delta}$$

$$\underline{V_R(1,1)} = \underline{V_R(1,1)} + \alpha [R(1,1, R(1,1)) + \gamma \underline{V_R(1,1)} - \underline{V_R(1,1)}]$$

$$= -0,00081 F \Delta + 0,1 \times [0 + 0,9 \times 0 + 0,00081 F \Delta] = \underline{-0,000729 F \Delta}$$

$$\underline{V_R(1,1)} = \underline{V_R(1,1)} + \alpha [R(1,1, R(1,1)) + \gamma \underline{V_R(1,1)} - \underline{V_R(1,1)}]$$

$$= 0 + 0,1 \times [0 + 0,9 \times 0 + 0] = 0$$

S	(1,1)	(1,1)	(1,1)	(1,1)	(1,1)	(1,1)
V_0	0	0	-0,1	0	0	0
V_1	0	-0,009 F Δ	-0,1	0	0	0
V_2	-0,00081 F Δ	-0,00081 F Δ	-0,1	0	0,00081 F Δ	0

الگوریتم DQN (Deep Q-Network) یکی از الگوریتم‌های یادگیری تقویتی عمیق است که توسط DeepMind توسعه یافته است. این الگوریتم در سال 2015 معرفی شد و توانست به موفقیت‌های قابل توجهی در بازی‌های آتاری دست یابد.

در DQN، از شبکه‌های عصبی عمیق (Deep Neural Networks) برای تقریب تابع Q استفاده می‌شود. تابع Q به ارزیابی ارزش یک اقدام (Action) در یک وضعیت خاص (State) می‌پردازد. هدف DQN این است که یک سیاست (Policy) را بیاموزد که به حداکثر رساندن مجموع پاداش‌ها (Rewards) در طول زمان کمک کند.

اجزای اصلی DQN:

1 تجربه تکراری (Experience Replay): در این روش، تجربه‌های عامل (Agent) در محیط ذخیره می‌شود و به صورت تصادفی از این حافظه برای به‌روزرسانی شبکه عصبی استفاده می‌شود. این کار به کاهش وابستگی زمانی و افزایش کارایی آموزش کمک می‌کند.

2 شبکه هدف DQN (Target Network): از دو شبکه عصبی استفاده می‌کند - شبکه اصلی (Main Network) و شبکه هدف. شبکه هدف هر چند وقت یکبار به‌روزرسانی می‌شود تا مقادیر Q پایدارتر باشند.

کاربردهای DQN:

بازی‌های رایانه‌ای DQN: برای اولین بار در بازی‌های آتاری معرفی شد و توانست عملکردی هم‌سطح یا حتی بهتر از انسان داشته باشد. این نشان داد که الگوریتم می‌تواند سیاست‌های پیچیده‌ای را یاد بگیرد.

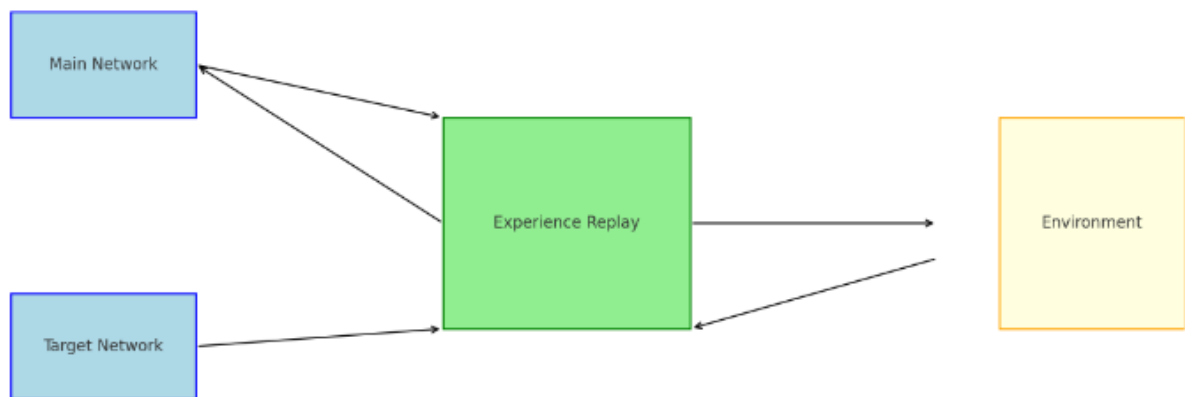
رباتیک DQN: در کنترل و بهینه‌سازی حرکت ربات‌ها مورد استفاده قرار می‌گیرد. مثلاً، در آموزش ربات‌ها برای انجام وظایف پیچیده مانند راه رفتن، پرش و تعامل با اشیاء.

سیستم‌های توصیه‌گر: این الگوریتم می‌تواند در بهینه‌سازی پیشنهادات به کاربران در سیستم‌های توصیه‌گر (Recommender Systems) مورد استفاده قرار گیرد.

مدیریت منابع DQN: در مدیریت منابع شبکه‌های کامپیوتری و تخصیص بهینه منابع برای بهبود عملکرد سیستم‌ها نیز کاربرد دارد.

نتیجه‌گیری:

الگوریتم DQN به عنوان یکی از الگوریتم‌های قدرتمند یادگیری تقویتی عمیق، نقش مهمی در پیشرفت‌های اخیر این حوزه ایفا کرده است. توانایی این الگوریتم در یادگیری سیاست‌های پیچیده از داده‌های خام و اعمال آن‌ها در محیط‌های متنوع، نشان‌دهنده پتانسیل بالای آن برای حل مسائل پیچیده و واقعی است.



این تصویر اجزا و مراحل اصلی الگوریتم DQN را نشان می‌دهد:

1. شبکه عصبی اصلی: (Main Network)

- در سمت چپ بالای تصویر قرار دارد و وظیفه تقریب تابع Q را بر عهده دارد.

2. شبکه عصبی هدف: (Target Network)

- در سمت چپ پایین تصویر قرار دارد و برای محاسبه مقادیر هدف Q استفاده می‌شود.

3. تجربه تکراری: (Experience Replay)

- در مرکز تصویر قرار دارد و تجربیات عامل در این حافظه ذخیره می‌شود.
- این تجربیات به صورت تصادفی برای به‌روزرسانی شبکه‌های عصبی استفاده می‌شوند.

4. محیط: (Environment)

- در سمت راست تصویر قرار دارد و عامل با آن تعامل می‌کند.

مراحل تعامل اجزا:

• انتخاب اقدام و تعامل با محیط:

- عامل یک اقدام را بر اساس وضعیت فعلی و سیاست ϵ حریصانه انتخاب می‌کند.
- اقدام در محیط اعمال می‌شود و پاداش و وضعیت جدید مشاهده می‌شود.
- این تجربه (s, a, r, s') در حافظه تکراری ذخیره می‌شود. (پیکان از "Environment" به "Experience Replay")

- بهروزرسانی شبکه‌های عصبی:

- تجربیات به صورت تصادفی از حافظه تکراری نمونه‌گیری می‌شوند.
 - شبکه عصبی اصلی (Main Network) با استفاده از این نمونه‌ها بهروزرسانی می‌شود (پیکان از "Experience Replay" به "Main Network")
 - شبکه هدف (Target Network) هر چند وقت یک‌بار بهروزرسانی می‌شود تا مقادیر Q پایدارتر باشند (پیکان از "Main Network" به "Target Network")
- این تصویر به خوبی نشان می‌دهد که چگونه الگوریتم DQN با استفاده از شبکه‌های عصبی و حافظه تجربه تکراری به یادگیری سیاست‌های بهینه می‌پردازد .