

(1) MDP

S	(1,1)	(1,2)	(1,3)	(2,1)	(2,2)	(2,3)
v_0	0	0	-5	0	0	5
v_1	0	0	-5	0	7.9	5
v_2	0	8.552	-5	5.472	11.184	5

S	(1,1)	(1,2)	(1,3)	(2,1)	(2,2)	(2,3)
$\pi^*(s)$	↑	↑	-	→	→	-

(2)

(3) Monte Carlo Model-free و Reinforcement Learning (TD) learning با action-value (نمایه اقدام) این روش با Temporal Difference (TD) learning (نمایه تفاوت زمانی) تفاوت دارد. bootstrapping (نمایه پیاپی) این روش با TD learning (نمایه تفاوت زمانی) تفاوت دارد.

$$\frac{0.9 \times (-5) + (0.9)^2 \times 5}{1.9} = 1.2 \quad \text{برای حالت (1,1) :}$$

$$\frac{1 \times 5 + 1 \times 5}{2} = 3 \quad \text{برای حالت (1,2) :}$$

$$(1,1) \xrightarrow{\text{iter 1}} (1,2) \xrightarrow{\text{iter 2}} (1,3) \quad (4)$$

$$\text{iter 1: sample} = R + \gamma v^\pi((1,2)) = 0 + \frac{9}{11} \times 0 = 0 \Rightarrow$$

$$v^\pi((1,1)) = (1-\alpha) v^\pi((1,1)) + \alpha \times \text{sample} = \frac{9}{11} \times 0 + \frac{1}{11} \times 0 = 0$$

$$\text{iter 2: sample} = R + \gamma v^\pi((1,3)) = (-5) + \frac{9}{11} \times (-5) = -9.5 \Rightarrow$$

$$v^\pi((1,2)) = (1-\alpha) v^\pi((1,2)) + \alpha \times \text{sample} = \frac{9}{11} \times 0 + \frac{1}{11} \times (-9.5) = -0.86$$

DQN از یک شبکه عصبی برای تقریب تابع Q استفاده می‌کند که پارامترهای خود را انتظار آینده را برای انجام اقدامات مختلف در حالت‌های مختلف تعیین می‌کند. این اولین و ساده‌ترین نوع یادگیری تقویتی است که با استفاده از شبکه‌های عصبی برای یادگیری انجام می‌گیرد.

در DQN یک مدل شبکه عصبی است که شبکه Q اصلی از یک شبکه Q هدف است. شبکه اصلی مقادیر Q را بر اساس وضعیت فعلی و اقدامات خود در محیط می‌کند، با توجه به شبکه هدف که ارزش‌های Q برای شبکه اصلی را هم می‌کند تا از آن یاد بگیرد.

DQN از یک تابع هزینه (loss function) برای بهینه‌سازی مدل استفاده می‌کند. در بین مقادیر Q پیش‌بینی شده، مقدار Q هدف استفاده می‌کند که با استفاده از شبکه هدف می‌تواند بهینه شود. این نوع از شبکه‌ها بیشتر به عنوان یک شبکه عصبی برای یادگیری تقویتی استفاده می‌شود. شبکه‌های Q را می‌توان به دو دسته تقسیم کرد: شبکه‌های Q که به عنوان Q هدف استفاده می‌شوند و شبکه‌های Q که به عنوان Q اصلی استفاده می‌شوند.

DQN قادر به یادگیری از یک محیط است که پارامترهای آن به تدریج در حال تغییر است.