



⑤

0	1/2	+1
0	0	-1
1	2	2

⑥ (1,1) → ↑ ← ↓  
0 0 0 max = 0

(2,1) → ↑ ← ↓  
1/2 1/2 0 max = 1/2

→ :  $\frac{1}{2} \times \frac{1}{2} \times 1/2 + \frac{1}{2} \times 0 + \frac{1}{2} \times 0 = 1/4$

↑ :  $\frac{1}{2} \times \frac{1}{2} \times 1/2 = 1/8$

← : 0

↓ : 0 > 1/8 > 0 (↑) = 1/8

(1,2) → ↑ ← ↓  
-1/2 1/2 1/2 max → 1/2

→ :  $\frac{1}{2} \times \frac{1}{2} \times -1 + \frac{1}{2} \times \frac{1}{2} \times 1/2 = -1/4$

↑ :  $\frac{1}{2} \times \frac{1}{2} \times 1/2 + \frac{1}{2} \times \frac{1}{2} \times -1 + 0 = -1/4$

← :  $\frac{1}{2} \times \frac{1}{2} \times 1/2 = 1/8$

↓ :  $\frac{1}{2} \times \frac{1}{2} \times -1 = -1/4$

(2,2) → ↑ ← ↓  
1/2 1/2 1/2 max = 1/2

→ :  $\frac{1}{2} \times \frac{1}{2} \times 1 + \frac{1}{2} \times \frac{1}{2} \times 1/2 = 3/8$

↑ :  $\frac{1}{2} \times \frac{1}{2} \times 1/2 + \frac{1}{2} \times \frac{1}{2} \times 1 = 3/8$

← :  $\frac{1}{2} \times \frac{1}{2} \times 1/2 = 1/8$

↓ : 0

1/2 → 1/2 1/2 0
↑ 1/2 1/2 -1

3	1,1	1,2	1,3	2,1	2,2	2,3	
1/2	↑	↑	-	→	→	-	

Monte Carlo

(1,1) (1,2) (1,3)

(2,1) (2,2) (2,3)

S1	
S2	
S3	
S4	
S5	
S6	

②

I) (1,1), (1,2), (1,3)

$$\text{reward } (1,1) \rightarrow (1,2) = 0$$

$$,, (1,2) \rightarrow (1,3) = -\omega$$

②) (1,1), (1,2), (2,2), (2,3)

$$,, (1,1) \rightarrow (1,2) = 0$$

$$,, (1,2) \rightarrow (2,2) = 0$$

$$,, (2,2) \rightarrow (2,3) = \omega$$

③) (1,1), (2,1), (2,2), (2,3)

$$\text{reward } (1,1) \rightarrow (2,1) = 0$$

$$,, (2,1) \rightarrow (2,2) = 0$$

$$,, (2,2) \rightarrow (2,3) = -\omega$$

$$\nabla^* (1,1) = \frac{1}{2} [(1 \times 0 + 0.9 \times \omega) + (1 \times 0 + 0.9 \times 0 + (0.9)^2 \times \omega) + (1 \times 0 + 0.9 \times 0 + (0.9)^2 \times \omega)] = 1/2$$

$$\nabla^* (1,2) = \frac{1}{2} [(1 \times \omega) + (1 \times \omega)] = \omega$$

الگوریتم (DQN) یکی از الگوریتم‌های سرورس که حوزن یادگیری تقویتی است این الگوریتم با ترکیب یادگیری تقویتی و شبکه عصبی عمیق توانست موفقیت‌های قابل توجهی داشته باشد.

DQN ترکیبی از  $Q$ -learning و شبکه عصبی عمیق است هدف  $Q$ -learning یافتن تابع  $Q$  است که مدائن مقداردهی مورد انتظار را برای هر حالت و اقدام ممکن محاسبه و کند. این تابع با استفاده از معادله بلان به روز می‌شود. در DQN از شبکه‌های عصبی عمیق برای تقریب این تابع  $Q$  استفاده می‌شود.

کاربرد DQN: بازی دیدنی، رباتیک، مدیریت منابع: در بازی مانند هرچس ترانزیک شبکه‌های کلیدی یا تخصیص منابع در سیستم‌های بزرگ، DQN می‌تواند کمک کند.

مزایا: توانایی یادگیری از تجربیات خام و داده‌های پیچیده.  
قلیت استفاده در محیط‌های متنوع و پویا.

معایب: نیاز به حافظه زیاد برای ذخیره Replay Buffer  
چالش یادگیری و همگرایی به دلیل وابستگی زیادی به راه‌ها.