| | +ω |
|---|---|
| $S$ | −ω |

$1 , 2 , 3$ (grid labels)

$$V_{q+1}^0(S) = \max_a \left[ \sum_{s'} T(S,a,s')\left(R(S,a,s') + \gamma V_q^{\circ}(s')\right)\right] \cdot 1 \; (\text{سوال})$$

$\gamma = 0,9$

$\forall S: \; V_0(S_{ij}) = 0$

هزینه هر حرکت = 0,1− ، مقدار حالت شروع $\int (G) = 0$

$V_1(S_{11}) =$

a=up : $0,\Lambda \times (0+0,9\times 0) + 0,1\times(0+0,9\times 0)+0,1(0+0,9\times 0)$   up / left / right

a=down : $0,\Lambda \times (0+0,9\times 0)+0,1\times(0+0,9\times 0)+0,1(0+0,9\times 0)=0$   down / left / right=0

a=left : 0

a=right : 0

$\Rightarrow \boxed{V_1(S_{11}) = 0}$

$V_1(S_{12}) =$

a=up : $0,\Lambda(0+0,9\times 0) + 0,1(0+0,9\times 0)+ 0,1(-\omega + 9(-\omega)) = -0,9\omega$   up / left / right

a=down : $0,\Lambda(0+0,9\times 0) + 0,1(0+0,9\times 0)+0,1(-\omega + 0,9\times(-\omega)) = -0,9\omega$   down / left / right

a=left : $0,\Lambda(0+0,9\times 0) + 0,1(0+0,9\times 0)+0,1(-0+0,9\times 0) = 0$   left / up / down

a=right : $0,1\Lambda(-\omega + 0,9\times(-\omega)) + 0,1(0+0,9\times 0)+0,1(0+0,9\times 0) = -V,9$   right

$\Rightarrow \boxed{V_1(S_{12}) = 0}$

$\boxed{V_1(S_{13}) = -\omega}$

$V_1(S_{11})$ چرا؟

$\boxed{V_1(S_{21}) = 0}$

$V_1(S_{22}) =$

a=up : $0,\Lambda(0+0,9\times 0) + 0,1(0+0,9\times 0)+ 0,1(+\omega + 0,\Lambda(\omega)) = +0,9\omega$   up / left / right

a=down : $0,\Lambda(0+0,9\times 0)+0,1(0+0,9\times 0)+ 0,1(\omega+0,9(\omega)) = +0,9\omega$   down / left / right

a=left : $0,\Lambda(0+0,9\times 0)+0,1(0+0,9\times 0)+0,1(0+0,9\times 0) = 0$   left / up / down

a=right : $0,\Lambda(\omega+\omega\times 0,9)+0,1(0+0\times 0,9)+0,1(0+0,9\times 0) = +V,9$   right / up / down

$\boxed{V_1(S_{22}) = +V,9} \Leftarrow$

$\boxed{V_1(S_{23}) = +\omega}$

$V_2(S_{11})$

a=up : $0,\Lambda(0+0,9\times 0)+0,1(0+0,9\times 0)+0,1(0+0,9\times 0) = 0$   up / left / right

a=down : $0,\Lambda(0+0,9\times 0)+0,1(0+0,9\times 0)+0,1(0+0,9\times 0) = 0$   down / left / right

a=left : $0,\Lambda(0+0,9\times 0)+0,1(0+0,9\times 0)+0,1(0+0,9\times 0) = 0$   left / up / down

a=right : $0,\Lambda(0+0,9\times 0)+0,1(0+0,9\times 0)+0,1(0+0,9\times 0) = 0$   right / down / up

$\boxed{V_2(S_{11}) = 0}$

$V_2(S_{12})$

a=up : $0,\Lambda(0+0,9\times V,9)+0,1(0+0,9\times 0)+0,1(-\omega+0,9(-\omega)) = +\frac{V}{?}0\, \text{YY}$   up / left / right

a=down : $0,\Lambda(0+0,9\times 0)+0,1(0+0,9\times 0)+0,1(-\omega+0,\Lambda(-\omega)) = -0,9\omega$   down / left / right

a=left : $0,\Lambda(0+0,9\times 0)+0,1(\omega+0,9\times 0)+0,1(0+0,9\times 0) = 0$   left / down / up

a=right : $0,\Lambda(-\omega+(\omega)0,9)+0,1(0+0,9\times 0)+0,1(0+0,9\times\omega) = -1,9\text{Y}$   right / down / up

$\Rightarrow \boxed{V_2(S_{12}) = +\varepsilon,\omega\text{IY}}$

$\boxed{V_2(S_{13}) = -\omega}$

$V_2(S_{21})$

a=up : $0,\Lambda(0+0,9\times 0)+0,1(0+0,9\times 0)+0,1(0+0,9\times V,9) = +0,9\Lambda\varepsilon$   up / left / right

a=down : $0,\Lambda(0+0,9\times 0)+0,1(0+0,9\times 0)+0,1(0+0,9\times V,9) = 0,9\Lambda\varepsilon$   down / left / right

a=left : $0\Lambda(0+0,9\times 0)+0,1(0+0,9\times 0)+0,1(0+0,9\times 0) = 0$   left / up / down

a=right : $0,\Lambda\times(0+0,9\times V,9)+0,1(0+0,9\times 0)+0,1(0+0,9\times 0) = 0,9\varepsilon\text{Y}$   right / up / down

$\Rightarrow \boxed{V_2(S_{21}) = 0,\varepsilon\omega\text{Y}}$

$V_2(S_{22}) =$

a=up : $0.8(0+0.9 \times V_1.9) + 0.1 \times (0+0.9 \times 0) + 0.1(\underset{right}{0+0.9 \times (0)}) = 9.442$

a=down : $0.8(0+0.9 \times 0) + 0.1 \times (0+0.9 \times 0) + 0.1(+0+0.9(0)) = 0.195$

a=left : $0.8(0+0.9 \times 0) + 0.1 \times (0+0.9 \times V_1.9) + 0.1(0+0.9 \times 0) = +0.198$

a=right : $0.8(+0+0.9 \times 0) + 0.1(0+0.9 \times V_1.9) + 0.1(0+0.9 \times 0) = +8.286$

$\Rightarrow \boxed{V_2(S_{22}) = 8.286}$

$\boxed{V_2(S_{22}) = +0}$

| 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 | 1.9 | +5 |
| $V_2$ | 0 | 4.522 | -5 | 5.472 | 8.286 | +5 |

$\pi^*(S) = \arg\max_a Q^*(S,a)$

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

$\pi^*(S_{11}) =$

a=up : $0.8(0+0.9 \times \underbrace{5.472}) + 0.1 \times (0+0.9 \times 0) + \overline{0.1 \times (0+0.9 \times 4.522)} = 4.24$

a=down : $0.8(0+0.9 \times 0) + 0.1 \times (0+0.9 \times 0) + 0.1(0+0.9 \times 4.522) = 0.499$

a=left : $0.8(0+0.9 \times 0) + 0.1 \times (0+0.9 \times 0) + 0.1(0+0.9 \times 5.472) = 0.4924$

a=right : $0.8(0+0.9 \times 4.522) + 0.1(0+0.9 \times 0) + 0.1(0+0.9 \times 5.472) = 3.7529$

$\boxed{up}$

$\pi^*(S_{12}) =$

a=up : $0.8(0+5) \times 8.286) + 0.1(0+0.9 \times 0) + 0.1(-5+0.9 \times 0)) = 5.05$

a=down : $0.8(0+0.9 \times 4.522) + 0.1(0+0.9 \times 0) + 0.1(0+0.9 \times 5)) = ----$

a=left : $0.8(0+0.9 \times 0) + 0.1(0+0.9 \times 8.286) + 0.1 \times (0+5) \times 4.522) = ---$

a=right : $0.8(-5+0.9(-5)) + 0.1(0+0.9 \times 4.522) + 0.1(0+0.9 \times 8.286) = -$

$\boxed{up}$

$\pi^*(S_{11}) \Rightarrow \boxed{right}$  $\pi^*(S_{22}) \Rightarrow \boxed{right}$

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

.2

$$V^*(S) = \mathcal{E}\left\{\sum \gamma^i R_i\right\}$$

I) $(1,1) \underset{0}{-} (1,2) \underset{-5}{-} (1,3)$

II) $(1,1) \underset{0}{-} (1,2) \underset{0}{-} (2,2) \underset{-5}{-} (2,3)$

III) $(1,1) \underset{0}{-} (2,1) \underset{0}{-} (2,2) \underset{5}{-} (2,3)$

$$V^\pi(S_{11}) = \frac{\left(0 + (-5)\times 0.9\right) + \left(0 + 0\times 0.9 + 0.9^2(+5)\right) + \left(0 + 0\times 0.9 + 5\times(0.9)^2\right)}{3} = 1.2$$

$$V^\pi(S_{22}) = \frac{(5 + 0)}{2} = 5$$

$$Sample = R(S, \pi(S), S' + \gamma V^\pi(S'))$$

$$V^\pi_{(S)} = (1-\alpha) N^\pi_{(S)} + (\alpha)[Sample]$$

$\alpha = 0.1$

$\gamma = 0.9$

X-Random

حالة (ع

استمالي

$$V^\pi(S_{11}) = (1 - 0.1) \times 0 + 0.1 \left[0 + 0.9 \times 0\right] = 0$$

$$V^\pi(S_{13}) = (1 - 0.1) \times 0 + 0.1 \left[-5 + 0.9 \times 0\right] = -0.50$$

**DQN:** DQN یک الگوریتم قوی در زمینه یادگیری تقویتی است. این الگوریتم، اصول شبکه های عصبی را با Q-learning ترکیب میکند و به agent اجازه می‌دهد تا سیاست‌های optimal policy را یاد بگیرد.

این الگوریتم از یک رویکرد مبتنی بر شبکه‌های عصبی برای یادگیری و بهینه‌سازی توابع action value استفاده می‌کند و رویکرد را به کانال‌های اساسی شکل خلاصه میکند:

11

**State representation:** وضعیت فعلی محیط را به یک بازنمایی عددی مناسب، مانند مقادیر آن پیکسل‌های ها و ویژگی‌های دیگر تبدیل میکند.

12

**Neural network architecture:** یک شبکه عصبی طراحی میکنیم که شبکه این شبکه state را به عنوان ورودی می‌گیرد و برای هر action مقادیر متناظر action-values عرضه بدهد.

**Experience Replay (3:** تجربیات agent شامل state های، action های، reward های next state, را در حافظه replay memory buffer ذخیره میکند.

4

**Q-learning updates:** از min-batches ای از تجربیات ای که از replay memory نمونه میگیرد استفاده میکنیم تا شبکه را آپدیت کردن. در اینجا از توابع جامعی loss استفاده می‌شود که اختلافات بین مقادیر پیش‌بینی شده (predicted action value) و هدف را به حداقل می‌رساند که از معادلات Bellman بدست می‌آید.

5

**Exploration and Exploitation:** با بالانس کردن بین exploration, exploitation، تکنیک استراتژی‌های action های را برای استخراج policy استفاده می‌شود.

6

**Target network:** استفاده از یک شبکه هدف جداگانه، مشابه شبکه اصلی با فاکتور شبکه. اصل شبکه هدف را تثبیت فرآیند یادگیری بهبود می‌دهد و پایداری.

7

**Repeat steps 1 to 6:** با محیط تعامل میکنیم، تجربیات جمع آوری میکنیم؟ شبکه را آپدیت کنیم و policy را بهبود...

چند مفهوم اساسی:

**Q-learning:** از الگوریتم DQN از Q-learning استفاده می‌شود، هدف این تخمین تابع ارزش action-value (Q-function) که map state های به عملکردهای مورد انتظار پیش‌بینی شده می‌کند.

**Experience Replay:** این تکنیک شبکه مشابه را ذخیره می‌کند در replay memory buffer. ارتباطات تجربیات این memory buffer سپس با انتخاب تصادفی در طول بهبود بخشیدن شبکه برای شکستن وابستگی‌های زمانی و شبکه یادگیری تقویت برتر می‌شود.