

# Markov Decision Process

Q1.  $S = 1(1,1), 2(1,2), 3(1,3), 4(2,1), 5(2,2), 6(2,3)$

1	2	3	↗ +5
4	5	6	↘ -5

$$V_{t+1}(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_t(s')]$$

$V_0(s) = 0$  for each  $s \in S$  if  $s$  is not in final states

از آنجایی که reward رسیدن به هر خانه‌ای جز 3 و 6 صفر است، اگر همچنین <sup>value</sup> خانه‌های همسایه صفر باشند پس Value خانه کنونی نیز صفر می‌شود. برای مثال می‌دانیم  $V_1$  و  $V_2$  و  $V_4$  در زمان صفر هستند. می‌توان گفت در حالتی که در  $s=1$  باشیم و به هر جهتی برویم، در یکی از خانه‌های 1 یا 2 یا 4 خواهیم رفت که Value همه آن‌ها و reward صفر است. پس  $V_{t=1}(s=1)$  هم صفر خواهد بود. با این استدلال یک به یک خانه‌ها را بررسی می‌کنیم:

$V_1(4)$  و  $V_1(5)$  خاصیت گفته شده را دارند.

$$V_1(2) = \max \left\{ \begin{array}{l} \text{right} \uparrow 0.1 \times (5 + 0.9 \times 0) \\ \text{up \& down} \uparrow 0.1 \times (5 + 0.9 \times 0) \end{array} \right\} = 0.1 \times 5 = 0.5$$

$$V_1(5) = -0.5 \leftarrow \text{شرایط کاملاً مشابه } V_1(2) \text{ است با این تفاوت که به جای } +5 \text{ در } -5 \text{ سایه بیافتد}$$

$V_1(3)$  و  $V_1(6)$  ثابت هستند (مانند Value در  $t=0$ )

$V_2(1)$ : فقط آن  $Q_2(1)$  ای صفر است که agent به سمت left برود. مانده حالت دقتی رخ می‌دهد که به right برود (با احتمال

$$V_2(1) = 0.1 \times (5 + 0.9 \times 0.5) = 0.475$$

بیشتری در خانه 2 می‌افتد تا وقتی که به up یا down برود)

$$V_2(4) = -0.475$$

$V_2(4)$ : شرایط کاملاً مشابه  $V_2(1)$  اما قرینه آن. در نتیجه

$V_2(2)$ : دنبال مانده هستیم. واضحاً حرکت به پایین (مقدار منفی) و چپ (مقدار صفر) ما را به هدف نمی‌رساند. دو حرکت

$$V_2(2) = \max \left\{ \underbrace{0.1 \times (5 + 0.9 \times 0.5) + 0.1 \times (5 + 0.9 \times 0.5)}_{0.475}, \underbrace{0.1 \times (5 + 0.9 \times 0.5) + 0.1 \times (0.5 + 0.9 \times 0.5)}_{0.5} \right\} = 0.5$$

$$V_2(5) = -0.5$$

$V_2(5)$ : شرایط دقیقاً مشابه  $V_2(2)$  اما قرینه آن است

$V_2(3)$  و  $V_2(6)$  ثابت هستند (مانند Value در  $t=0$ )

پس جدول نهایی به این صورت است

s	(1,1)	(1,2)	(1,3)	(2,1)	(2,2)	(2,3)
$V_0$	0	0	0	0	0	-5
$V_1$	0	+0.5	0	0	-0.5	-5
$V_2$	+0.475	+0.5	0	-0.475	-0.5	-5

1,1	1,2	1,3
2,1	2,2	2,3

# Markov Decision Process

Q2.

۵,۴۷۲	۷,۴	۵
-۵,۴۷۲	-۷,۴	-۵

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_T^*(s')] \rightarrow T=2$$

$\pi_1^*$  : با حرکت به down مقدار منفی دریافت خواهیم کرد. همچنین حرکت به left مطمئناً نتیجه کمتری نباشد. ~~left~~  
 right یا up خواهد داشت. دو حالت right و up را می‌سنجیم:

$$\pi^*(1) = \underset{a}{\operatorname{argmax}} \left\{ \underbrace{(0.1 + 0.1)(0 + 0.9 \times 5,472) + 0.1(0 + 0.9 \times 7,4)}_{5,114} \text{ و } \underbrace{(0.1 - 0.1)(0.9 \times 5,472) + 0.1 \times 0.9 \times 7,4}_{5,472} \right\} = \text{right} (\rightarrow)$$

$\pi_2^*$  : شرایط کاملاً مشابه  $\pi_1^*$  اما گزینه آن است. پس با حرکت به پایین یا راست مقدار ۵,۱۱۴ و ۵,۴۷۲ را خواهد گرفت. ~~left~~ با حرکت به چپ هم مطمئناً مقداری منفی خواهد گرفت. اما حرکت به بالا منفی نباشد.

Value مثبتی می‌شود  $\uparrow$  (up)  $\rightarrow 0.1 \times 5,472 \times 0.9 + 0.1 \times 0.9(-5,472) + 0.1 \times 0.9 \times (-7,4) = 2,174 > 0$   
 $\pi_2^*$  : واضحاً حرکت به پایین حاصل منفی دارد، سه جهت دیگر را بررسی می‌کنیم

$$\begin{aligned} \text{up} &: 0.1 \times 0.9 \times 7,4 + 0.1 \times 0.9 \times 5,472 + 0.1 \times 0.9 \times 5 \\ \text{right} &: 0.1 \times 0.9 \times 5 + 0 \\ \text{left} &: 0.1 \times 0.9 \times 5,472 + 0 \end{aligned} \rightarrow \underset{a}{\operatorname{argmax}} : \text{up} (\uparrow)$$

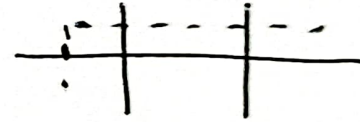
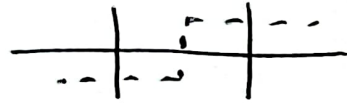
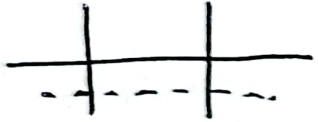
$\pi_5^*$  : حرکت به left، right و down منجر به مقادیر منفی خواهند شد و تنها حرکت به up حاصل مثبت دارد. پس با  $\underset{a}{\operatorname{argmax}}$  گرفتن،  $\uparrow$  (up) انتخاب می‌شود.

برای خانه‌های ۳ و ۶ نیاز به محاسبه  $\pi^*$  نیست. در آخر نتیجه بدین شکل می‌شود:

$\rightarrow$	$\uparrow$	۵+
$\uparrow$	$\uparrow$	۵-

Monte Carlo Simulation روشی است که با شبیه سازی بازی روی state space تخمینی از transition Function به دست می آید. هر چه تعداد simulation ها بیشتر باشد، نتیجه به دست آمده به نتیجه واقعی نزدیکتر خواهد بود اما هیچوقت نمی توان با اطمینان کامل گفت که به مقدار دقیق transition ها رسیدیم.

$(1,1) \rightarrow (1,2) \rightarrow (1,3)$        $(1,1) \rightarrow (1,2) \rightarrow (2,2) \rightarrow (2,3)$        $(1,1) \rightarrow (2,1) \rightarrow (2,2) \rightarrow (2,3)$



با مشاهده سه episode بالا، به تخمین مقادیر  $T$  می پردازیم

تخمین برای خانه  $(1,1)$

$$((1,1), \text{right}, (1,2)) = \frac{2}{3} \quad ((1,1), \text{right}, (2,1)) = \frac{1}{3}$$

تخمین برای خانه  $(2,2)$

$$((2,2), \text{right}, (2,3)) = \frac{2}{2} = 1$$

از آنجایی که Fined Policy داریم و در هر simulation ها به راست می رویم، مقدار تابع  $T$  را فقط به ازای <sup>این</sup> action می توان تخمین زد



# Markov Decision Process

Q4. learning rate = 0.1 =  $\alpha$

updating  $V(s)$  &  $V^\pi(s) \leftarrow V^\pi(s) + \alpha [\text{new sample} - V^\pi(s)] = (1-\alpha)V^\pi(s) + \alpha(\text{new sample})$  : TD-learning

$$\bar{r}_n = (1-\alpha)\bar{r}_{n-1} + \alpha r_n \quad \leftarrow \text{مؤثر کثیر تفریق ما تا الان  $\bar{r}_n$  بوده و مشاهده جدید  $r_{n+1}$  است}$$

$$\bar{r}_n = (1-\alpha) \left[ (1-\alpha)\bar{r}_{n-2} + \alpha r_{n-1} \right] + \alpha r_n = (1-\alpha)^2 \left[ (1-\alpha)\bar{r}_{n-3} + \alpha r_{n-2} \right] + (1-\alpha)\alpha r_{n-1} + \alpha r_n$$

$$= (1-\alpha)^3 \left[ (1-\alpha)\bar{r}_{n-4} + \alpha r_{n-3} \right] + \left[ (1-\alpha)^2 r_{n-2} + (1-\alpha)r_{n-1} + r_n \right] \alpha \rightarrow \dots$$

$$\bar{r}_n = \left[ \sum_{i=0}^{n-1} (1-\alpha)^i r_{n-i} \right] \alpha$$

با گذر زمان می توان  $\alpha$  را کم کرد زیرا تفریق دقیق تر شده است.

(I) از همان نتایج Simulation ها در بخش قبل استفاده می کنیم. جدول اولیه

$$\begin{array}{c|c|c} 0 & 0 & 5 \\ \hline 0 & 0 & -5 \end{array} \xrightarrow{(1,1) \rightarrow (1,2)} \begin{array}{c|c|c} 0 & 0 & 5 \\ \hline 0 & 0 & -5 \end{array} \xrightarrow{(1,2) \rightarrow (1,3)} \begin{array}{c|c|c} 0 & 0 & 5 \\ \hline 0 & 0.95 & -5 \end{array} \xrightarrow{(1,3) \rightarrow (2,3)} \begin{array}{c|c|c} 0 & 0 & 5 \\ \hline -0.1084 & -0.1084 & -5 \end{array}$$

$$\begin{array}{c|c|c} 0 & 0 & 5 \\ \hline 0 & 0 & -5 \end{array} \xrightarrow{(1,1) \rightarrow (1,2)} \begin{array}{c|c|c} 0 & 0 & 5 \\ \hline -0.1084 & -0.1084 & -5 \end{array} \xrightarrow{(1,2) \rightarrow (2,2)} \begin{array}{c|c|c} 0 & 0 & 5 \\ \hline -0.1084 & -0.1084 & -5 \end{array} \xrightarrow{(2,2) \rightarrow (2,3)} \begin{array}{c|c|c} 0 & 0 & 5 \\ \hline -0.1084 & -0.1084 & -5 \end{array}$$

$$\begin{array}{c|c|c} 0 & 0 & 5 \\ \hline -0.1084 & -0.1084 & -5 \end{array} \xrightarrow{(2,1) \rightarrow (2,2)} \begin{array}{c|c|c} 0 & 0 & 5 \\ \hline -0.1084 & -0.1084 & -5 \end{array} \xrightarrow{(2,2) \rightarrow (2,3)} \begin{array}{c|c|c} 0 & 0 & 5 \\ \hline -0.1084 & -0.1084 & -5 \end{array} \xrightarrow{(2,3) \rightarrow (3,3)} \begin{array}{c|c|c} 0 & 0 & 5 \\ \hline -0.1084 & -0.1084 & -5 \end{array}$$

$$\begin{array}{c|c|c} 0 & 0 & 5 \\ \hline -0.1084 & -0.1084 & -5 \end{array} \xrightarrow{(2,2) \rightarrow (2,3)} \begin{array}{c|c|c} 0 & 0 & 5 \\ \hline -0.1084 & -0.1084 & -5 \end{array} \xrightarrow{(2,3) \rightarrow (3,3)} \begin{array}{c|c|c} 0 & 0 & 5 \\ \hline -0.1084 & -0.1084 & -5 \end{array}$$

(II) یکبار دیگر Value ها را تخمین می زنیم (با روش value iter)

$$\begin{array}{c|c|c} 0 & 0 & 5 \\ \hline 0 & 0 & -5 \end{array} \xrightarrow{\text{First iter}} \begin{array}{c|c|c} 0 & 0.95 & 5 \\ \hline 0 & -0.95 & -5 \end{array} \xrightarrow{\text{second iter}} \begin{array}{c|c|c} 0 & 0.95 & 5 \\ \hline -0.1084 & -1.172 & -5 \end{array}$$

$$V_{t+1}(s) = (1-\alpha)V_t(s) + \sum_{s'} \alpha [R + \gamma V_t(s')]$$

AI-A6 810101465

## Deep Q-Network

Q1.

Deep Q-Network is a deep reinforcement learning algorithm that combines Q-learning algorithm with deep neural networks. DQN is a model-free algorithm.

Q-learning is used for estimating value function, and DNN is used for estimating Q Function. That's because directly calculating Q values is not efficient when working with large datasets.

- 1) init: initialize the NN with random weights.
- 2) action selection: using a greedy policy or other exploration strategies, choose an action.
- 3) perform action: execute the chosen action and receive results (next state and reward).
- 4) store result: store the result tuple in the experience replay memory.
- 5) sample a batch: randomly take a sample of transition from memory.
- 6) update Q-Network: for each sampled transition calculate the Temporal Difference error. using the error  $\delta$  update weights of NN by applying GD algorithm (use a loss function)
- 7) repeat steps 2-6 until convergence or a desired performance is reached

Applications of DQN :

- game playing
- robot control
- autonomous driving
- resource management in cloud computing
- trading strategies in finance