

Gridworld:

	1	2	3
2			+5
1	S		-5

رستون بردن
(2,3)

MDP

سوال اول

سوال 1

$$V_{i+1}(S) = \max_a \left(\sum_{S'} T(S, a, S') (R(S, a, S') + \gamma V_i(S')) \right)$$

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	-7.6	+5
V_2	0	4.522	-5	5.472	8.284	+5

discount

$\gamma = 0.9$

$$V_1((1,1)) = \max_a \left\{ \begin{array}{l} \text{up: } 0.8 \times (0 + 0.9 \times 0) + 0.1 \times (0 + 0.9 \times 0) + 0.1 \times (0 + 0.9 \times 0) = 0 \\ \text{down: } 0.8 \times (0 + 0) + 0.1 \times (0 + 0) + 0.1 \times (0 + 0.9 \times 0) = 0 \\ \text{left: } 0 \rightarrow \text{بسته} \\ \text{right: } 0 \end{array} \right\} = 0$$

$V((2,1)) = 0 \rightarrow$ زیر اینتر $V((1,1))$ تمام خانه های اطراف آن
آب $(1,1) \times (1,2)$ و $(1,3)$ در سطح تبدیل بریزد

$$V((1,2)) = \max_a \left\{ \begin{array}{l} \text{up: } 0.8(0+0) + 0.1(0+0) + 0.1(-5 + 0.9(-5)) = -0.95 \\ \text{down: } 0.8(0+0) + 0.1(0+0) + 0.1(-9.5) = -0.95 \\ \text{right: } 0.8(-9.5) + 0.1(0+0) + 0.1(0+0) = -7.6 \\ \text{left: } 0 \end{array} \right\} = 0$$

از اینجا به بعد دیگر تقادیری که می بینیم
را به ترتیب می بینیم
که به ترتیب می بینیم

$$V((2,2)) = \max_a \left\{ \begin{array}{l} \text{up: } 0.1(9.5) = 0.95 \\ \text{down: } 0.1(9.5) = 0.95 \\ \text{left: } 0 \\ \text{right: } 0.8(9.5) = 7.6 \end{array} \right\} = +7.6$$

سوال ٢

$$V_2(1,1) = \max_a \left\{ \begin{array}{l} \text{up: } 0(0) = 0 \\ \text{down: } 0 \\ \text{left: } 0 \\ \text{right: } 0.8(-0.855) = 0 \end{array} \right\} = 0$$

$$V_2(2,1) = \max_a \left\{ \begin{array}{l} \text{up: } 0.1(0.9 \times 7.6) = 0.684 \\ \text{down: } 0.1(6.84) = 0.684 \\ \text{left: } 0 \\ \text{right: } 0.8(6.84) = 5.472 \end{array} \right\} = 5.472$$

$$V_2(1,2) = \max_a \left\{ \begin{array}{l} \text{up: } 0.8(6.84) + 0 + 0.1(-9.5) = 4.522 \\ \text{down: } 0.8(0) + 0 + 0.1(-9.5) = -1.634 \\ \text{left: } 0.1(0) + 0.1(6.84) = 0.684 \\ \text{right: } 0.8(-5) + 0.1(0) + 0.1(4) = -3.6 \end{array} \right\} = 4.522$$

$$V_2(2,2) = \max_a \left\{ \begin{array}{l} \text{up: } 0.8(6.84) + 0.1(9.5) = 6.422 \\ \text{down: } 0.8(0) + 0.1(9.5) = 0.95 \\ \text{left: } 0.1(6.84) + 0.1(-0.5) = 0.684 \\ \text{right: } 0.8(9.5) + 0.1(6.84) + 0.1(-0.5) = 8.284 \end{array} \right\} = 8.284$$

$$\pi^*(s) = \arg \max_a Q^*(s, a) ; Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

$$\pi^*(2,1) = \arg \max_a \left\{ \begin{array}{l} \text{up: } 0.9(4.9248) + 0.1(7.45) \\ \text{down: } 0.8(0) + 0.1(4.924) + 0.1(7.45) \\ \text{left: } 0.9(4.924) + 0.1(0) \\ \text{right: } 0.8(7.45) + 0.1(4.924) + 0.1(0) \end{array} \right\} = \text{right}$$

Subject:

Date:

0.1 + 0.8
اقتلان در خانه (1 و 2) در صورت افش بر پایین

$$\pi((1,1)) = \arg \max_a \left\{ \begin{array}{l} \text{up: } 0.8(4.924) + 0.1(0) + 0.1(4.069) \\ \text{down: } 0.9(0) + 0.1(4.069) \\ \text{left: } 0.9(0) + 0.1(4.924) \\ \text{right: } 0.8(4.069) + 0.1(0) + 0.1(4.924) \end{array} \right\} = \text{up}$$

$$\pi((2,2)) = \arg \max_a \left\{ \begin{array}{l} \text{up: } 0.8(7.45) + 0.1(4.924) + 0.1(9.5) \\ \text{down: } 0.8(4.069) + 0.1(4.924) + 0.1(9.5) \\ \text{left: } 0.8(4.924) + 0.1(7.45) + 0.1(4.069) \\ \text{right: } 0.8(9.5) + 0.1(7.45) + 0.1(4.069) \end{array} \right\} = \text{right}$$

$$\pi((1,2)) = \arg \max_a \left\{ \begin{array}{l} \text{up: } 0.8(7.378) + 0.1(9.5) + 0.1(0) \\ \text{down: } 0.8(4.069) + 0.1(0.073) + 0.1(9.5) \\ \text{left: } 0.8(0.073) + 0.1(4.069) + 0.1(7.378) \\ \text{right: } 0.8(9.5) + 0.1(4.069) + 0.1(7.378) \end{array} \right\} = \text{up}$$

جدول نهایی:

S	(1,1)	(1,2)	(1,3)	(2,1)	(2,2)	(2,3)
$\pi^*(S)$	up	up	-	right	right	-

سوال 3) کان است میانگین مجموع reward ها برای هر خانه را بدست آوریم که می شود Value برای آن خانه:

reward:

Discount Factor $\gamma = 0.9$ I) (1,1) \rightarrow (1,2) \rightarrow (1,3)discounted reward: $1 \times r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$ II) (1,1) \rightarrow (1,2) \rightarrow (2,2) \rightarrow (2,3)III) (1,1) \rightarrow (2,1) \rightarrow (2,2) \rightarrow (2,3)

$$V^*((1,1)) = \frac{1}{3} \left[\begin{array}{l} \text{(I)} \quad (1 \times 0 + 0.9 \times (-5)) + \text{(II)} \quad (1 \times 0 + 0.9 \times 0 + (0.9)^2 \times 5) + \text{(III)} \quad (1 \times 0 + 0.9 \times 0 + (0.9)^2 \times 5) \end{array} \right] = 1.2$$

$$V^*((2,2)) = \frac{1}{2} \left[\begin{array}{l} \text{(II)} \quad (1 \times 5) + \text{(III)} \quad (1 \times 5) \end{array} \right] = 5$$

P4PCO

$$\begin{aligned} \text{sample} &= R(s_t, \pi(s_t, s')) + \gamma V^\pi(s') \\ V^\pi(s) &\leftarrow V^\pi(s) + \alpha (\text{sample} - V^\pi(s)) \end{aligned}$$

اولین ترنژیشن: $(1,1) \xrightarrow{\text{right}} (1,2)$; reward = 0

$$\text{sample} = 0 + 0.9(0) = 0 \rightarrow V^\pi((1,1)) = 0 + 0.1 = 0$$

دومین ترنژیشن: $(1,2) \xrightarrow{\text{right}} (1,3)$; reward = -5

$$\text{sample} = -5 + 0.9(-5) = -9.5 \rightarrow V^\pi((1,2)) = 0 + 0.1(-9.5 + 0) = -0.95$$

نتیجه نهایی:

s	(1,1)	(1,2)	(1,3)	(2,1)	(2,2)	(2,3)
$V^\pi(s)$	0	-0.95	-5	0	0	+5

DQN

سوال اول

DQN یکی از الگوریتم‌های پیشرفته در محیط RL است که Q-Learning را با الگوریتم‌های deep learning ترکیب می‌کند.

در واقع از یک deep neural network برای تخمین ارزش $Q(s, a)$ ها (که ذخیره می‌شود) استفاده از الگوریتم Q-Learning می‌کنیم.

استفاده می‌کند. مزیت این الگوریتم نسبت به Q-learning این است که با استفاده از آن می‌توان به حساب مسائل

با ابعاد بسیار بالا استفاده کرد (تعداد پارامترهای زیاد). حال کمی به توضیح نحوه عملکرد DQN می‌پردازیم:

Subject :
Date :

این الگوریتم سه جز اصلی دارد :

① Reply Buffer : یک بافر است که تجربیات قبل یا درون آن به صورت (S, a, r) ذخیره

می شود که در ادامه برای یافتن ارتباط های بین transition های متوالی و استیت ها استفاده می شود.

② Q Network : یک deep neural network است که استیت کنونی را به عنوان ورودی می گیرد و مقادیر Q را

برای تمام action های ممکن محاسبه می دهد. وظیفه این شبکه این است که اختلاف بین Q های پیشین و

Q های هدف را کمینه کند.

③ Target Network : یک شبکه جدا از شبکه Q است (که توفیق دارد) که برای مقادیر مستقیم از اجرا آن

دارای پارامترهای ثابت است و وظیفه آن تولید Q های هدف برای update کردن Q-network است.

حال بخواهیم به توضیحات بالا اضافه کنیم الگوریتم DQN را توضیح می دهیم. ابتدا Q-network را با وزن های رندوم و یک کمی

از آن را به عنوان target network در نظر می گیریم (مقادیر را در replay buffer هم می گیریم است). سپس شروع

به انجام transition در محیط می کنیم و هر کدام را که انجام دادیم به صورت (S, a, r) در بافر ذخیره می کنیم. سپس از تجربه

محیط به تعداد زیاد n به سراسر شبکه های Q به صورت رندوم از بافر n شروع به انتخاب batch می کنیم و برای هر

کدام از ترنسزینس ها در این batch مقدار زیر را محاسبه می کنیم :

$$y = R + \gamma \max_a Q(S', a)$$

α target

PAPCO

Subject :

Date

که مقدار Φ_{target} از target network بیست می آید را با آیدیت کردن target network weight در مرحله

بررسی match شدن با وزن های $Q\text{-network}$ L میس $Q\text{network}$ را با میسیم کردن مقدار خطا

که به صورت $\frac{1}{2}(y - Q(s, a))^2$ می توان تعریف کرد که آیدیت می کنیم و این فرایند را ادامه می دهیم

تا $Q\text{-network}$ L همگرا شود.