

① از فرمول‌های زیر برای Policy Iteration و Policy Improvement در الگوریتم

policy iteration استفاده می‌کنیم.

Evaluation:

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$

Improvement:

$$\pi_{i+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

~~برای استفاده از این سیاست (π_i) ابتدا باید مقادیر همگرا بشوند لذا ابتدا هم $V(s)$ ها~~
~~محاسبه می‌کنیم و چندین تکرار انجام می‌دهیم~~

ابتدا $V(s) \approx 0$ می‌گذاریم برای هم حالت‌ها و سپس چندین تکرار انجام می‌دهیم
 تا به π_i برسیم. در ابتدا سیاست هم حالت‌ها را انجام بازی می‌گذاریم.

init:	حالت	s_1	s_2	s_3	s_4	s_5	s_6
$V_s: V(s) = 0$	$V_0(s)$	0	0	0	0	0	0
$\pi_0(s) = \text{finish}$	π_0	اتمام	اتمام	اتمام	اتمام	اتمام	اتمام
iterate:	$V_1^{\pi_0}(s)$	1	2	3	4	5	6
	π_1	تاس	تاس	اتمام	اتمام	اتمام	اتمام
	$V_2^{\pi_1}(s)$	2.5	2.5	3	4	5	6
	π_2	تاس	تاس	اتمام	اتمام	اتمام	اتمام
	$V_3^{\pi_2}(s)$	2.83	2.83	3	4	5	6
	π_3	تاس	تاس	اتمام	اتمام	اتمام	اتمام
	$V_4^{\pi_3}(s)$	2.94	2.94	3	4	5	6
$V(s_1), V(s_2)$	π_4	تاس	تاس	اتمام	اتمام	اتمام	اتمام
بعد از چندین	$V_5^{\pi_4}(s)$	2.98	2.98	3	4	5	6
مرحله به 3	π_5	تاس	تاس	اتمام	اتمام	اتمام	اتمام
همگرا می‌شوند.	$V_6^{\pi_5}(s)$	2.99	2.99	3	4	5	6

حالت	s_1	s_2	s_3	s_4	s_5	s_6 .1
π_i	تاس ریختن	تاس ریختن	اتمام بازی	اتمام بازی	اتمام بازی	اتمام بازی
$V^{\pi_i}(s)$	3	3	3	4	5	6

$$V^{\pi_i}(s_1) = \frac{1}{6}((-1+1 \times 3) + (-1+1 \times 3) + (-1+1 \times 3) + (-1+1 \times 4) + (-1+1 \times 5) + (-1+1 \times 6))$$

$$V^{\pi_i}(s_1) = \frac{1}{6} \times 18 = 3 \quad V^{\pi_i}(s_2) = s_1 \text{ مشابه } = 3$$

محاسبه می شود

$$V^{\pi_i}(s_3) = 1(0+1 \times 3) = 3 \quad V^{\pi_i}(s_4) = 1(0+1 \times 4) = 4 \quad V^{\pi_i}(s_5) = 1(0+1 \times 5) = 5$$

$$V^{\pi_i}(s_6) = 1(0+1 \times 6) = 6$$

حالت	s_1	s_2	s_3	s_4	s_5	s_6 .2
π_i	تاس ریختن	تاس ریختن	اتمام بازی	اتمام بازی	اتمام بازی	اتمام بازی
π_{i+1}	تاس ریختن	تاس ریختن	اتمام بازی یا تاس ریختن	اتمام بازی	اتمام بازی	اتمام بازی

$s_1 \xrightarrow{\text{actions}} \text{تاس: } Q=3 \rightarrow \text{argmax} = \text{تاس}$
 $\text{اتمام: } Q=1$
 $Q(s_1, \text{finish})$

$s_2 \xrightarrow{\text{actions}} \text{تاس: } Q = \frac{1}{6}(18) = 3 \rightarrow \text{argmax} = \text{تاس}$
 $\text{اتمام: } Q = 2 = 1(0+1 \times 2)$

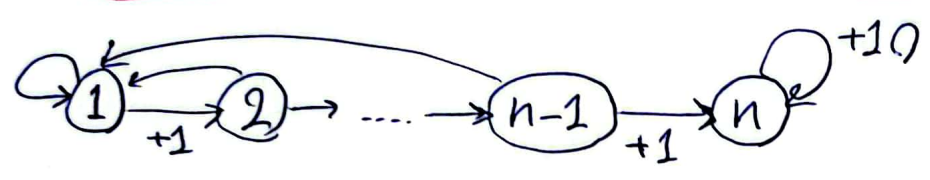
$s_3 \xrightarrow{\text{actions}} \text{تاس: } Q = \frac{1}{6}(18) = 3$
 $\text{اتمام: } Q = 1(0+1 \times 3) = 3$
 $\text{argmax} = \text{تاس / اتمام}$
 $\text{هر دو تصمیم optimal اند.}$

$s_4 \xrightarrow{\text{actions}} \text{تاس: } Q=3$
 $\text{اتمام: } Q=4 \rightarrow \text{argmax} = \text{اتمام}$

$s_5 \xrightarrow{\text{actions}} \text{تاس: } Q = \frac{1}{6}(18) = 3$
 $\text{اتمام: } Q=5 \rightarrow \text{argmax} = \text{اتمام}$

$s_6 \xrightarrow{\text{actions}} \text{تاس: } Q(s_6, \text{dice}) = 3$
 $\text{اتمام: } Q(s_6, \text{finish}) = 6 \rightarrow \text{argmax} = \text{finish}$
 اتمام

بله! مقادیر به دست آمده بهینه هستند و policy همگرا (Converge) شده است.
 و دیگر نیازی به بروزرسانی ندارد. چرا که policy از حالت π_i به π_{i+1} تغییر خاصی نداشته است. چرا که در حالت 3 هر دو تصمیم بهینه optimal هستند و نمی توان از آن به عنوان تغییر سیاست یاد کرد و همچنین s_1, s_2, s_4, s_5, s_6 هیچ تغییری نگرفته اند.
 اما باید توجه کنیم که ممکن است value ها دیرتر از policy ها همگرا شوند و وقتی policy ها همگرا هستند value ها همگرا ~~نیستند~~ نباشند. ~~معمولاً~~
 The policy often converges long before the values.



(۲)

حالت	1	2	3	...	n-2	n-1	n
سیاست (π)	حرکت آزادانه	حرکت آزادانه	حرکت آزادانه	حرکت آزادانه	حرکت آزادانه	حرکت آزادانه	loop (خارجی شمریم)

$$V_{k+1}^{\pi}(n) = \sum_{s'} T(s, \pi_i(s), s') \left[\underbrace{R(s, \pi_i(s), s')}_{+10} + \underbrace{\frac{1}{2} V_k^{\pi}(n)}_{\gamma} \right] = Q(n, loop)$$

$$V_{k+1}^{\pi}(n) = 1 \times \left[10 + \frac{1}{2} V_k^{\pi}(n) \right]$$

$$V_0^{\pi}(n) = 0 \quad (init)$$

$$V_1^{\pi}(n) = 1 \times \left[10 + \frac{1}{2} 0 \right] = 10$$

$$V_2^{\pi}(n) = 1 \times \left[10 + \frac{1}{2} 10 \right] = 15$$

$$V_3^{\pi}(n) = 1 \times \left[10 + \frac{1}{2} 15 \right] = 17.5$$

$$\vdots \quad \left(\frac{1}{2} \right) 10 + \left(\frac{1}{2} \right)^2 10$$

$$V_i^{\pi}(n) = 1 \times \left[10 + \frac{10}{2} + \frac{10}{2^2} + \dots + \frac{10}{2^{n-1}} \right] = 10 \times \frac{1 - (\frac{1}{2})^n}{1 - \frac{1}{2}}$$

$$V_i^{\pi}(n) = 20$$

$$V_i^{\pi}(n) = 10 \times \frac{1}{\frac{1}{2}}$$

$n \rightarrow \infty$

$$V^*(n) = 20$$

پس $V^*(n)$ به 20 همگرا می شود.

$$V^*(k) = 1 + \left(\frac{1}{p}\right)^1 + \left(\frac{1}{p}\right)^2 + \dots + \left(\frac{1}{p}\right)^i + 2 \cdot \left(\frac{1}{p}\right)^{i+1} = \sum_{i=0}^{n-1-k} \left(\frac{1}{p}\right)^i + 2 \cdot \left(\frac{1}{p}\right)^{n-k}$$

$$1 + \frac{1}{p} + \left(\frac{1}{p}\right)^2 + \dots + \left(\frac{1}{p}\right)^{(n-1)-k}$$

$$V^*(k) = 2 \left(1 - \left(\frac{1}{p}\right)^{n-k}\right) + 2 \cdot \left(\frac{1}{p}\right)^{n-k}$$

$$= \frac{1 - \left(\frac{1}{p}\right)^{n-k}}{1 - \frac{1}{p}} + 2 \cdot \left(\frac{1}{p}\right)^{n-k}$$

به صفر همگرا می شود
n-k بزرگ شود
n-k بزرگ شود به دو همگرا می شود

کوچک k
بزرگ n

$$V^*(k) = 2(1 - 0) + 2 \cdot 0 = 2 \xrightarrow{\text{مثلاً}} V^*(1) = 2 \rightarrow V^*(1) = 2$$

هر چه k
بزرگتر شود
 $V^*(k)$ بیشتر می شود.

$$V^*(n-1) = 2\left(1 - \frac{1}{p}\right) + 2 \cdot \left(\frac{1}{p}\right) = 11 \rightarrow V^*(n-1) = 11$$

state
هر چه k بزرگتر باشد ارزش بیشتر دارد.

~~$V^*(k+1) = 2 \left(1 - \left(\frac{1}{p}\right)^{k+1}\right) + 2 \cdot \left(\frac{1}{p}\right)^{k+1}$~~
 ~~$V^*(n-1) = 2 \left(1 - \left(\frac{1}{p}\right)^{n-1}\right) + 2 \cdot \left(\frac{1}{p}\right)^{n-1}$~~

حالت	1	2	3	4	...	n-3	n-2	n-1	n
$V_0^*(s)$	0	0	0	0	0	0	0	0	0
$V_1^*(s)$	1	1	1	1	1	1	1	1	10
$V_2^*(s)$	1.5	1.5	1.5	1.5	1.5	1.5	1.5	6	15

در iteration های اول دوم تمام حالت ها نمی توانند صفر باشند چرا که
باصطحت آزادانه پاداش می گیرند و همچنین حالت n نیز پاداش می گیرد.
حالت های 1 تا n-1

+10