

1

حالت	S_1	S_2	S_3	S_4	S_5	S_6
$\pi(i)$	رخس	رخس	آسم	آسم	آسم	آسم
$v(\pi(i))$	3	3	3	4	2	4

حالت	S_1	S_2	S_3	S_4	S_5	S_6
$\pi(i)$	رخس	رخس	آسم	آسم	آسم	آسم
$\pi(i+1)$	رخس	رخس	آسم	آسم	آسم	آسم

همان طور که در جدول مشخص می باشد Policy 1 در حالت 1 و حالت 2 + 1 تغییر خاصی نکرده است.
 در حالت 3 و 4 هر دو تصمیم می گیرند در رفتن
 نمی دانند از آن به عنوان تغییر Policy می توان برد. در نتیجه Policy همرا با Converge شده است.
 البته امکان دارد Policy باشد ولی Value ها بعد از Policy همرا شوند. اصولاً در حالت کاربردی
 Policy زودتر از Value ها همرا می شود.

از این فرین ها را Policy Libration, Policy Improvement, Policy evaluation

$$r_{k+1}^{(i)}(s) = \sum_{s'} \pi(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{(i)}(s')]$$
$$\pi_{i+1}(s) = \arg \max_a \{ \pi(s_a, s'), [R(s_a, s') + \gamma V^{\pi_i}(s')] \}$$

ω	S_1	S_2	S_3	S_4	S_5	S_6
$N_0(S)$	0	0	0	0	0	0
π_0 $r_1^{\pi_0}(S)$	f 1	f 2	f 3	f 4	f 5	f 6
α_1 $r_1^{\alpha_1}(S)$	T 1/2	T 1/2	f 3	f 4	f 5	f 6
π_2 $r_1^{\pi_2}(S)$	T 1/2	T 1/2	f 3	f 4	f 5	f 6
α_2 $r_2^{\alpha_2}(S)$	T 1/2	T 1/2	f 3	f 4	f 5	f 6
π_3 $r_1^{\pi_3}(S)$	T 1/2	T 1/2	f 3	f 4	f 5	f 6
α_3 $r_3^{\alpha_3}(S)$	T 1/2	T 1/2	f 3	f 4	f 5	f 6
π_4 $r_1^{\pi_4}(S)$	T 1/2	T 1/2	f 3	f 4	f 5	f 6
α_4 $r_4^{\alpha_4}(S)$	T 1/2	T 1/2	f 3	f 4	f 5	f 6
π_5 $r_1^{\pi_5}(S)$	T 1/2	T 1/2	f 3	f 4	f 5	f 6
α_5 $r_5^{\alpha_5}(S)$	T 1/2	T 1/2	f 3	f 4	f 5	f 6
π_6 $r_1^{\pi_6}(S)$	T 1/2	T 1/2	f 3	f 4	f 5	f 6
α_6 $r_6^{\alpha_6}(S)$	T 1/2	T 1/2	f 3	f 4	f 5	f 6

حکایت اول: $V(S) = 0$ $u_0(S) = f$

در این بخش، به بررسی قیمت $V(S_t)$ و $V(S_1)$ می‌پردازیم.

$$V^{\pi_1}(S_1) = \frac{1}{4} [(-1 + 1 \times 3) + (-1 + 1 \times 3) + (-1 + 1 \times 3) + (-1 + 1 \times 3) + (-1 + 1 \times 3) + (-1 + 1 \times 3)]$$

در S_1 و S_2

$$V^{\pi_1}(S_1) = \frac{1}{4} \times 11 = 2.75$$

در S_1 $Q = 3$ na $Q = 1$ $Q(S, f)$

$$V^{\pi_1}(S_2) = 3$$

S_2 $\left\{ \begin{array}{l} T: Q = \frac{1}{4} (11) = 2.75 \text{ } naa \\ F: Q = 1(0 + 1 \times 3) \end{array} \right.$

$$V^{\pi_1}(S_2) = 3$$

S_2 $\left\{ \begin{array}{l} T: Q = \frac{1}{4} (11) = 2.75 \\ F: Q = 1(0 + 1 \times 3) = 3 \end{array} \right.$

در S_2 و S_3

$$V^{\pi_1}(S_3) = 1(0 + 1 \times 3)$$

S_3 $\left\{ \begin{array}{l} T: Q = 3 \\ F: Q = 3 \text{ } naa \end{array} \right.$

$$V^{\pi_1}(S_4) = 1(0 + 1 \times 3) = 3$$

S_4 $\left\{ \begin{array}{l} T: Q = \frac{1}{4} (11) = 2.75 \\ F: Q = 0 \text{ } naa \end{array} \right.$

$$V^{\pi_1}(S_4) = 1(0 + 1 \times 3) = 3$$

S_4 $\left\{ \begin{array}{l} T: Q = 3 \end{array} \right.$

$$V^{\pi_1}(S_5) = 1(0 + 1 \times 3) = 3$$

$\left\{ \begin{array}{l} F: Q(S_4, f) = 3 \text{ } naa \end{array} \right.$

② - 1 Policy: Greedy Loop State

1	2	3	...	n-2	n-1	n
R	R	R	...	R	R	loop

2- Value Iteration

$$V_0(n) = 0$$

Initial State(0)

$$V_1(n) = 10 + \frac{1}{\gamma} \alpha 0 = 10$$

$$V_2(n) = 10 + \frac{1}{\gamma} + \frac{1}{\gamma^2}$$

$$V_3(n) = 10 + \frac{1}{\gamma} \alpha 10 = 10$$

$$V_4(n) = 10 + \frac{1}{\gamma} + \frac{1}{\gamma^2} + \frac{1}{\gamma^3}$$

$$V_n(n) = 10 \alpha \left(1 + \frac{1}{\gamma} + \frac{1}{\gamma^2} + \dots + \frac{1}{\gamma^{n-1}} \right) \xrightarrow{n \rightarrow \infty} 10 \alpha 2 = 20$$

20 هجری لور

جمع می: 2 می 3 اند

3- Value Iteration: $V^*(k)$ برای $k=1, \dots, n-1, 10$

$$V^*(1) = V^*(n - (n-1)) = 1 + \frac{1}{\gamma} + \dots + \left(\frac{1}{\gamma}\right)^{n-2} + \frac{V_0}{\gamma^{n-1}} = 2$$

$$V^*(2) = V^*(n - (n-2)) = 1 + \frac{1}{\gamma} + \dots + \left(\frac{1}{\gamma}\right)^{n-3} + \frac{V_0}{\gamma^{n-2}} = 2$$

$$V^*(n-1) = V^*(n-1) = \left(\frac{1}{\gamma}\right)^0 + \frac{V_0}{\gamma} = 10$$

$$V(n-k) = 1 + \frac{1}{\gamma} + \dots + \left(\frac{1}{\gamma}\right)^{k-1} + \frac{V_0}{\gamma^k} = 2$$

10 برای State زین هجری

4- Iteration: 2 در State می تواند هجری

	1	2	3	...	n-2	n-1	n
V_0	10	0	0	...	0	0	0
V_1	1	1	1	...	1	1	10
V_2	1/5	1/5	1/5	...	1/5	4	15

در State های 2 و 3 می تواند هجری
می تواند هجری 2 و 3 می تواند هجری
می تواند هجری 2 و 3 می تواند هجری

State می تواند هجری 2 و 3 می تواند هجری