



$$V_0^{\text{stay}}(1) = 0$$

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

$$\Rightarrow V_1^{\text{stay}}(1) = \sum_{s'} T(1, \text{stay}, s') [R(1, \text{stay}, s') + \gamma V_0(s')]$$

$$= \frac{1}{2} T(1, \text{stay}, 1) [R(1, \text{stay}, 1) + \frac{1}{2} V_0(1)] = \frac{1}{2} [1 + 0] = \frac{1}{2}$$

$$V_2^{\text{stay}}(1) = \frac{1}{2} [1 + \frac{1}{2} + 1] = \frac{3}{4}$$

$$V_3^{\text{stay}}(1) = \frac{1}{2} [1 + \frac{1}{2} (1 + \frac{1}{2})] = \frac{7}{8}$$

$$\therefore V^{\text{stay}}(1) = 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots = \sum_{n=0}^{\infty} \frac{1}{2^n} = 2$$

نیم از آنهایی که r_i تمام state ها برابر 1 است V^* به استثناء

ما با هم برابر نخواهد بود زیرا وضعیت مشابهی دارند. $V^*(1) = V^*(2) = \dots = V^*(N)$

ما در هر استیج با هم داریم دو انتخاب داریم که برای ما r_{21}

10.

به همراه دارد و West یا East که برای ما r_{20} دارد و ما را به استیج می برد

که دقیقاً از نقل V مشابه state فعلی است. بنابراین در هر حرکت بهتر و

15.

است star کنیم و سود r_{21} را به r_{20} ترمیم دهیم.

$$\begin{aligned}
 V^*(1) &= \max_a Q(s, a, s') = \max \left(\underbrace{1 \times (1 + \frac{1}{2} V^*(1))}_{\text{best action} \rightarrow \text{stay}}, 1 \times (0 + \frac{1}{2} V^*(2)), 1 \times (0 + \frac{1}{2} V^*(0)) \right) \\
 &= 1 + \frac{V^*(1)}{2}
 \end{aligned}$$

پس می توانیم بگوییم سیاست بهینه برای (مسئله) V^* در هر State

stay کردن در آن State است. بنابراین در این مسئله

$$V^*(s) = V^{stay}(s) = 2 \Rightarrow V^*(1) = 2$$

محاسبه شده در قسمت

(ب) سؤال

* فرض این است که سوال MDP است و r_i ها را می دانیم

↓ (ت) سیاست بهتر این است که ابتدا به State ای که بیشترین r_i را دارد بروی
سپید برای همیشه در آن بماند Stay کنی.

5. اثبات و در هر State ای که باشیم می توانیم به راست یا چپ حرکت کنیم:

فرض کنیم در Start S_i هستیم و استیج S_p بیشترین مقدار r را برای

Stay دارد. چطور فرض کنید در K اکشن می توانیم از S_i به S_p برویم.

10. برای شروع از S_i دو سیاست تعریف می کنیم و

1- از S_i با K اکشن به S_p می رویم و در آنجا Stay می کنیم

15. 2- هر سیاست دیگری جز سیاست ①

برای محاسبه Utility ~~است~~ پس از M مرحله داریم

~~در مرحله~~

① $Utility = (m-k) R(S_p, Stay, S_p)$

② ~~$Utility = R(S_p, Stay, S_p)$~~

20. ② در بهترین حالت در ~~هر مرحله~~ یک r کسب کرده ایم که $r = R(S_p, Stay, S_p)$ است.

① است زیرا اگر $r = R(S_p, Stay, S_p)$ باشد و این سیاست همان سیاست ①

25. $Utility = m R(S_p, Stay, S_p)$ و

$R(S_p, Stay, S_p) \geq R(S_p, Stay, S_p)$ خواهد بود.

در حالتی است که $m \rightarrow \infty$ میل کند

$$\text{utility ①} = \lim_{m \rightarrow \infty} \frac{m R(S_F, \text{stay}, S_F) - K(R(S_F, \text{stay}, S_F))}{m R(S_N, \text{stay}, S_N)}$$

$$= \frac{R(S_F, \text{stay}, S_F)}{R(S_N, \text{stay}, S_N)} > 1 \Rightarrow \text{utility ①} > \text{utility ③}$$

10

(1)

15

(S, a, r, S')	$Q(1, \text{stay})$	$Q(1, \text{East})$	$Q(2, \text{West})$	$Q(2, \text{stay})$
Initial	0	0	0	0
$(1, \text{stay}, 4, 1)$	2	0	0	0
$(1, \text{East}, 0, 2)$	2	0	0	0
$(2, \text{stay}, 6, 2)$	2	0	0	3
$(2, \text{West}, 0, 1)$	2	0	0.5	3
$(1, \text{stay}, 9, 1)$	3.5	0	0.5	3