

$s$	(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	$\frac{1}{8} \times (0.9 \times 5) = 3.6$	-5
$V_2$	0	$\frac{1}{8} \times (0.9 \times 3.6) = 2.59$	5	$\frac{1}{8} \times (0.9 \times 3.6) = 2.59$	$\frac{1}{8} \times (0.9 \times 5) = 3.6$	-5

$+ 0.1 \times (0.9 \times -5) = 2.14$ 
 $+ 0.1 \times (0.9 \times 3.6) = 3.92$

با توجه به حدود، ارزش، هاء، بالا محاسبه کنند (باء، خانههای)، که ده با حند Action با مقدار

<del>2,59</del>	<del>3,6</del>	3,92
s	2,14	-5

طبق ضریل دارم و مسأله مسأله

$$\begin{aligned} V_1((2,2)) &= 0,8(0 + 0,9(5)) + 0,1(0 + 0,9(0)) + 0,1(0 + 0,9(0)) \\ &= 3,6 \end{aligned}$$

$$\begin{aligned} V_2((2,1)) &= 0,8(0 + 0,9(3,6)) + 0,1(0 + 0,9(0)) + 0,1(0 + 0,9(0)) \\ &= 2,59 \end{aligned}$$

$$\begin{aligned} V_2((1,2)) &= 0,8(0 + 0,9(3,6)) + 0,1(0 + 0,9(-5)) + 0,1(0 + 0,9(0)) \\ &= 2,14 \end{aligned}$$

بهمت یا بود باز همینجا نباشد

$$\begin{aligned} V_2((2,2)) &= 0,8(0 + 0,9(5)) + 0,1(0 + 0,9(3,6)) + 0,1(0 + 0,9(0)) \\ &= 3,92 \end{aligned}$$

2

S	(1,1)	(1,2)	(1,3)	(2,1)	(2,2)	(2,3)
$\pi^*(S)$	$UP \uparrow$	$UP \uparrow$	-	$Right \rightarrow$	$Right \rightarrow$	-

<del>2,59</del> $\rightarrow$	<del>3,92</del> $\rightarrow$	+5
0 S $\uparrow$	$\uparrow$ 2,14	-5

$2,59 > 2,14 \xrightarrow{(1,1)} \text{So move up}$   
 $3,92 > -5 \xrightarrow{(1,2)} \text{So move up}$

$3,92 > 9 \xrightarrow{(2,1)} \text{So move right}$   
 $5 > 2,14 \xrightarrow{(2,2)} \text{So move right}$

: monte carlo  $\tilde{\pi}_{\text{target}}$

<del>2,59</del>	<del>3,92</del>	+5
S	<del>2,14</del>	-5

$\approx \text{avg! } 3$

Summing all the rewards coming after the first visit to  $(S_t)$

$$(1,1): \text{ For episode } E_1: G = 2,14 + (-5) = -2,86$$

$$\sim \sim \quad E_2: G = 2,14 + 3,92 + 5 = 11,06$$

$$\sim \sim \quad E_3: G = 2,59 + 3,92 + 5 = 11,51$$

$$V((1,1)) = \frac{-2,86 + 11,06 + 11,51}{3} = 6,57$$


---

$$(2,2): \text{ For episode } E_1: G = \text{null}$$

$$\sim \sim \quad E_2: G = 5$$

$$\sim \sim \quad E_3: G = 5$$

$$V((2,2)) = \frac{5+5}{3} = 3,33$$

در این قسمت معرفه کردیم چهار گزینه exploration policy

• نرم خواهد بود که استخراج ایجاد کند action

: پایه Valve میگیرد

-1	-1	5
2	-1	-5

$$\begin{array}{c|c} \leq & V(s) \\ \hline (1,1) & 0 \\ (1,2) & 0 \\ \vdots & i \\ (2,3) & 0 \end{array}$$

$$V(s_1)_{obs} = R(s_2) + V(s_2)$$

$$TD_{Error} = V(s_1)_{obs} - V(s_1)^{exp}$$

$$V(s_1) = V(s_1) + \alpha \times TD_{Error}$$

حل > مسأله نمر از روشی که می خواهد

$$V((1,1)) = -0,1$$

$$V((1,2)) = -0,1$$

$$\begin{aligned} V((1,3)) &= -5 + 0,1(45+5) \\ &= -4,05 \end{aligned}$$

$s$	$V(s)$
$(1,1)$	-0,1
$(1,2)$	-0,1
$(1,3)$	-4,05
$(2,1)$	0
$(2,2)$	0
$(2,3)$	0

نایابی در مس:



$$v((1,1)) = -0,3$$
$$v((2,1)) = -0,1$$

$$v((2,2)) = -0,1$$

(1,1)	-0,3
(1,2)	-0,1
(1,3)	-4,05
(2,1)	-0,1
(2,2)	-0,1
(2,3)	0

مکانیزم