

① الف) حالت : موقعیت کور  $(X, Y)$  در صفحه

اکشن : چه حش می‌دهد که باعث تغییر حالت صفحه می‌شود. (زاویه، دور و پنج صفحه)

پاداش : مقدار یک پاداش به ازای هر صفحه به خط بر

تمام کور از محدوده  $(\sqrt{x^2 + y^2})$

و همچنین یک پاداش زیاد به ازای هر کور که در

محدوده حالت تعادل  $(1)$  و یک پاداش متنوع

به ازای افتادن کور از صفحه

ب) حالت : مقدار سه مایه (مبارک)

اکشن : خدیه یافته شدن / یا (در وقت دیگر)

پاداش : پاداش مثبت به ازای اتمه شدن مایه و پاداش متنوع  
به ازای کاهش آن  $(\Delta X)$

② ب)  $right + 1$   $Y$   $right + 4$

$X, right, +1, X, \dots, terminal$

$X, left, 0, X, left, 0, \dots$  الف)

$$G_0 = R_1 + \gamma R_2 + \gamma^2 R_3 = 1 + \frac{1}{2} \times 1 + \frac{1}{2}^2 \times 4 = 1.5$$



Subject :

Year :

Month :

Date :

$$V_{\pi}(s) = \sum_a P(a|s) \sum_{s',1} P(s',1|s,a) [1 + \gamma V_{\pi}(s')] \quad \rightarrow$$

$$V_{\pi_1}(Y) = 4 + 0.5 \times V_{\pi_1}(\text{terminal}) = 4$$

$$V_{\pi_2}(X) = \frac{3}{4} (1 + 0.5 \times V_{\pi_2}(X)) \quad (\leftarrow)$$

$$+ \frac{1}{4} (-1 + 0.5 \times V_{\pi_2}(Y))$$

$$V_{\pi_2}(Y) = 4 + 0.5 \times V_{\pi_2}(\text{terminal}) = 4$$

$$\Rightarrow V_{\pi_2}(X) = \frac{3}{4} + \frac{3}{8} V_{\pi_2}(X) + \frac{1}{4} (1)$$

$$\Rightarrow \frac{5}{8} V_{\pi_2}(X) = 1 \Rightarrow V_{\pi_2}(X) = \frac{8}{5} = 1.6$$

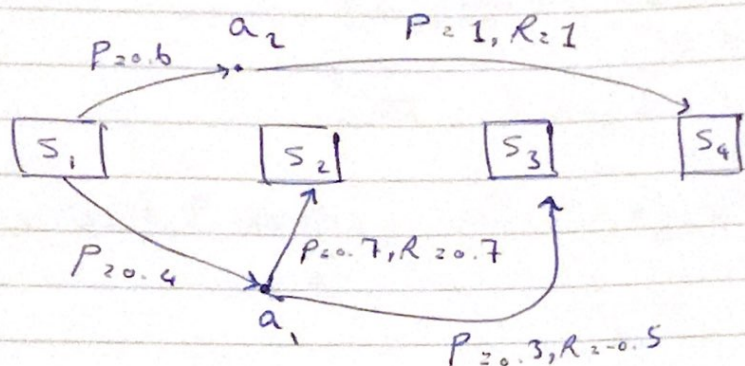
$$G_0 = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots \quad (\text{r})$$

$$= 5 + 10 (\gamma + \gamma^2 + \gamma^3 + \dots) = 5 + 10 \frac{\gamma}{1-\gamma}$$

$$= 5 + 10 \frac{0.9}{0.1} = 95$$

$$G_0 = \underbrace{R_1}_5 + \gamma \overset{0.9}{G_1} = 95 \Rightarrow 90 = 0.9 G_1 \Rightarrow G_1 = 100$$

(15)



$$V_{\pi}(s) = \sum_a \alpha(a|s) \sum_{s',r} P(s',r|s,a) \left[ r + \gamma V_{\pi}(s') \right]$$

$$V_{\pi}(s_1) = 0.4 \times \left[ 0.7 (0.5 + \gamma V_{\pi}(s_2)) + 0.3 (-0.5 + \gamma V_{\pi}(s_3)) \right] + 0.6 (1 + \gamma V_{\pi}(s_4))$$

$$q_{\pi}(s_1, a_1) = 0.7 (0.5 + \gamma V_{\pi}(s_2)) + 0.3 (-0.5 + \gamma V_{\pi}(s_3))$$

$$q_{\pi}(s_1, a_2) = 1 (1 + \gamma V_{\pi}(s_4))$$

$$\Rightarrow V_{\pi}(s_1) = 0.6 q_{\pi}(s_1, a_2) + 0.4 q_{\pi}(s_1, a_1)$$

$$q'_{\pi}(s,a) = q_{\pi}(s,a) + P(s) \rightarrow P(s) = s, a \quad f: s \rightarrow R \quad (16)$$

تابع ادالاف ، تابع ارتباط الکت و تقیید عمل بر روی سته و خطای آن تا تیره ندارد

حال با استفاده از رابطه زیر می توان به سینه نه به طریقی که به خطای خود



$$\pi(s) = \arg \max_a q_\pi(s, a)$$

$$\arg \max_a q_\pi(s, a) = \pi$$

$$\arg \max_a q_\pi(s, a) + f(s) = \arg \max_a q_\pi(s, a) = \pi$$

$$\Rightarrow \pi' = \pi$$

$$V_\pi(x) = \frac{1}{3} \left[ \begin{array}{l} \text{a} \quad (4) \\ 0.25(-1.6 + 0.9 \times 8) + 0.75(0.2 + 0.9 \times 2) \\ \text{b} \quad + (-2.7 + 0.9 \times 3) + 0.2(-5 + 0.9 \times 0) \\ \text{c} \quad + 0.8(1.6 + 0.8 \times V_\pi(x)) \end{array} \right]$$

$$\Rightarrow V_\pi(x) = \dots$$

$$V_\pi^*(x) = \max \left[ \frac{1}{3} a, \frac{1}{3} b, \frac{1}{3} c \right] = \max [0.96, 0, 0.85]$$

$$= 0.96 \Rightarrow L$$

$$V_\pi(s_{13}) = \frac{1}{4} \left( \begin{array}{l} 0 + 0.92 \times 2.3 + 0.92 \times 0.4 + 0.92 \times 0.7 \\ + 0.92 \times (-0.4) \end{array} \right) \quad (V)$$

$$\approx 0.7$$