

Subject: ()

تیرین دوازدهم هوش

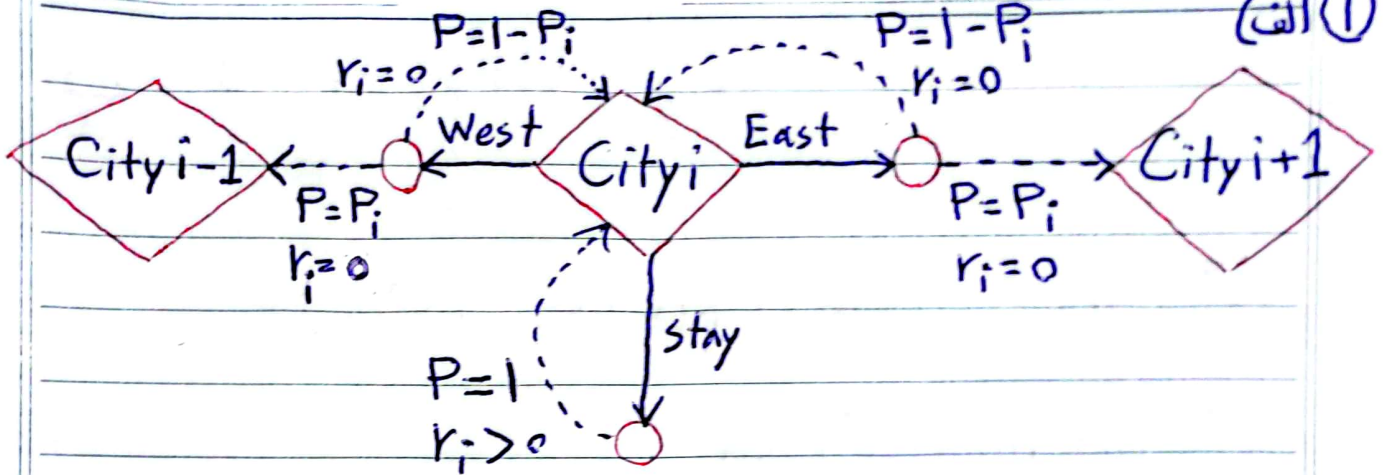
Year:

Month:

فرزان رحمانی

Day:

① الف



$$\gamma = \frac{1}{\gamma}, P_i = 1, r_i = 1 \quad \text{ب.}$$

$$V^\pi(s) = \sum_{s'} T(s, \pi(s), s') \left[R(s, \pi(s), s') + \gamma V^\pi(s') \right]$$

(P_i) فقط یک انتخاب با احتمال 1 داریم

$$V^{\text{stay}}(1) = 1 \left(1 + \frac{1}{2} V^{\text{stay}}(1) \right) \rightarrow V^{\text{stay}}(1) = 1 + \frac{1}{2} V^{\text{stay}}(1)$$

$$\frac{1}{2} V^{\text{stay}}(1) = 1 \rightarrow \boxed{V^{\text{stay}}(1) = 2}$$

با هم قابل حل است

$$V_0^{\text{stay}}(1) = 0 \quad V_1^{\text{stay}}(1) = 1 \left(1 + \frac{1}{2} \times 0 \right) = 1 \quad V_2^{\text{stay}}(1) = 1 \left(1 + \frac{1}{2} \right) = \frac{3}{2}$$

$$V_3^{\text{stay}}(1) = 1 \left(1 + \frac{3}{4} \right) = \frac{7}{4} \rightarrow V^{\text{stay}}(1) = 1 + \frac{1}{2} + \frac{1}{4} + \dots = \frac{1}{1 - \frac{1}{2}} = 2$$

TANIN

$$\gamma = \frac{1}{2}, P_i = 1, r_i = 1 \quad (\text{ب})$$

برای شرق و غرب

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

$$V^*(s) = \max_a Q^*(s, a)$$

همچنین چون برای همه آنها پارامترها یکسان هستند لذا استیت‌ها

$$V^*(1) = V^*(2) = V^*(3) = \dots = V^*(N)$$

متقارن هستند یعنی
وضعیت مشابهی دارند.

$$I) Q^*(1, \text{stay}) = 1 \times (1 + 0.5 V^*(1)) = 1 + \frac{1}{2} V^*(1)$$

$$II) Q^*(1, \text{east}) = 1 \times (0 + 0.5 V^*(2)) = \frac{1}{2} V^*(2) = \frac{1}{2} V^*(1)$$

$$III) Q^*(1, \text{west}) = 1 \times (0 + 0.5 V^*(N)) = \frac{1}{2} V^*(N) = \frac{1}{2} V^*(1)$$

$$V^*(1) = \max_a Q^*(1, a) = Q^*(1, \text{stay}) \quad \text{چون با 1 جمع می شود از بقیه بیشتر است}$$

$$I) Q^*(1, \text{stay}) = 1 + \frac{1}{2} Q^*(1, \text{stay}) \rightarrow Q^*(1, \text{stay}) = 2 = V^*(1)$$

$$II) Q^*(1, \text{east}) = \frac{1}{2} V^*(1) = \frac{2}{2} = 1 \rightarrow Q^*(1, \text{east}) = 1$$

$$III) Q^*(1, \text{west}) = \frac{1}{2} V^*(1) = \frac{2}{2} = 1 \rightarrow Q^*(1, \text{west}) = 1$$

TANIN

$$V^{\text{stay}}(1) = V^*(1) = \max_a Q^*(1, a) = Q^*(1, \text{stay}) = 2$$

$$\gamma = 1, P_i > 0, r_i > 0 \quad (\text{ت})$$

$$\pi^*(s) = \arg \max_a \sum_{s'} T(s, \pi^*(s), s') [R(s, \pi^*(s), s') + \gamma V^*(s')]$$

$$\pi^*(i) = \arg \max \left\{ \underbrace{1(r_i + V^*(i))}_{\text{stay}}, \underbrace{P_i V^*(i+1) + (1-P_i) V^*(i)}_{\text{east}}, \right.$$

$$\left. \underbrace{P_i V^*(i-1) + (1-P_i) V^*(i)}_{\text{west}} \right\}$$

قرنول سیاست بهینه

$$\pi^*(i) = \arg \max \left\{ \underbrace{V^*(i) + r_i}_{\text{stay}}, \underbrace{V^*(i) + P_i (V^*(i+1) - V^*(i))}_{\text{east}}, \right.$$

$$\left. \underbrace{V^*(i) + P_i (V^*(i-1) - V^*(i))}_{\text{west}} \right\}$$

$$P_i (V^*(i-1) - V^*(i)), P_i (V^*(i+1) - V^*(i)), r_i$$

اختلاف ارزش شهر جدید
باشهر فعلی

داداش
شهر فعلی
می باشد

بسته به مقادیر بالا سیاست بهینه می تواند متفاوت باشد مثلاً اگر r_i هم

شهرها برابر باشد آنگاه $\pi^*(i) = \text{stay}$ می شود چرا که

ارزش شهرها برابر است و بماندن در یک شهر سود مشابهی می گیریم

$$V^*(i+1) - V^*(i) = 0$$

اما اگر r_i ها تفاوت داشته باشند سیاست بهینه این می شود که

TANIN

ابتداء به شهری برویم که بیشترین r_i را دارد سپس آنها را به این

و تجارت کنیم. $\max\{r_i\}$

چون r_i و P_i دقیق داده نشده اند نمی توان سیاست ثابتی

را بیان کرد و برای تعیین سیاست باید از قریب به دست آمده و

مقدارهای r_i و P_i هر شهر استفاده کنیم. مقدارهای $V^*(s)$

نیز وابسته به r_i و P_i هستند.

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

$$\pi^*(i) = \arg \max \left\{ \underbrace{V^*(i) + r_i}_{\text{stay}}, \underbrace{P_i (V^*(i+1) - V^*(i)) + V^*(i)}_{\text{east}}, \underbrace{V^*(i) + P_i (V^*(i-1) - V^*(i))}_{\text{west}} \right\}$$

Subject: ()

Year:

Month:

Day:

 $\gamma = 0.5$ $\alpha = 0.5$ - learning rate

(s, a, r, s')	$Q(1, stay)$	$Q(1, east)$	$Q(2, west)$	$Q(2, stay)$
initial	0	0	0	0
$(1, stay, 4, 1)$	2	0	0	0
$(1, east, 0, 2)$	2	0	0	0
$(2, stay, 6, 2)$	2	0	0	3
$(2, west, 0, 1)$	2	0	0.5	3
$(1, stay, 4, 1)$	3.5	0	0.5	3

$$sample = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$

از فرمول روبه رو
استفاده میکنیم

$$Q(s, a) \leftarrow (1 - \alpha) Q(s, a) + \alpha [sample]$$

$$(1, stay, 4, 1): sample = 4 + \frac{1}{2} \times 0 = 4$$

$$Q(2, stay) = \frac{1}{2} \times 0 + \frac{1}{2} \times 4 = 2$$

$$(1, east, 0, 2): sample = 0 + \frac{1}{2} \max_{a'} Q(2, a') = 0 + 0 = 0$$

$$Q(1, east) = \frac{1}{2} \times 0 + \frac{1}{2} \times 0 = 0$$

$$(2, stay, 6, 2): sample = 6 + \frac{1}{2} \times 0 = 6$$

$$Q(2, stay) = \frac{1}{2} \times 0 + \frac{1}{2} \times 6 = 3$$

$$(2, west, 0, 1): sample = 0 + \max_{a'} (1, a') = 0 + \frac{1}{2} \times 2 = 1$$

$$Q(2, west) = \frac{1}{2} \times 0 + \frac{1}{2} \times 1 = 0.5$$

$$(1, stay, 4, 1): sample = 4 + \frac{1}{2} \times 2 = 5$$

$$Q(1, stay) = \frac{1}{2} \times 2 + \frac{1}{2} \times 5 = 3.5$$

TANIN