



$$P_i = 1, r_i = 1, \gamma = 0.5$$

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

$$\Rightarrow V^{\text{stay}}(1) = 1 [1 + 0.5 V^{\text{stay}}(1)] \rightarrow V^{\text{stay}}(1) = 1 + 0.5 V^{\text{stay}}(1)$$

$$\Rightarrow 0.5 V^{\text{stay}}(1) = 1 \Rightarrow V^{\text{stay}}(1) = 2$$

این فرآیند به صورت زیر در نظر گرفته می شود: Value iteration
این فرآیند به صورت زیر در نظر گرفته می شود: Value iteration

$$V_0(s) \rightarrow$$

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

$$V_1(1) = \max \left(1 \times \left(1 + 0.5 V_0(1) \right) + 1 \times \left(0 + 0.5 V_0(2) \right) \right)$$

$$V_1(1) = 1, \dots, V_1(2) = 1$$

$$V_2(1) = \max \left(\underbrace{(1 + 0.5 \times 1)}_{1.5}, 1 \times (0 + 0.5 \times 1) \right)$$

$$V_2(1) = V_2(2) = \dots = V_2(n) = 1.5$$

$$V_3(1) = \max \left(\underbrace{(1 + 0.5 \times 1.5)}_{1.75}, 1 \times (0 + 1.5 \times 0.5) \right)$$

الاستیسیته پ ارتدیت به نقلیله جلد به 2 جلدی شود سیدر

$$V^*(1) = 1 + \left(\left(\left(1 \times \frac{1}{2} \right) + 1 \right) \times \frac{1}{2} \right) + 1 \times \frac{1}{2} \dots$$

$$V^*(1) = 2$$

به نقلیله در نهایت لایه 1 شود

Policy iteration
سیاست لایه به لایه
Stay برابر π_0

city1 city2 ... city n
 $\pi_0 \rightarrow$ Stay Stay ... Stay

$$\rightarrow V^{\pi_0}(1) = \sum T(s, \pi_0, s') [R(s, \pi_0, s') + \gamma V^{\pi_0}(s')]$$

$$\rightarrow V^{\pi_0}(1) = 1 \times (1 + 0.5 V^{\pi_0}(1))$$

$$V^{\pi_0}(2) = 1 \times (1 + 0.5 V^{\pi_0}(2))$$

⋮

$$V^{\pi_0}(n) = 1 \times (1 + 0.5 V^{\pi_0}(n))$$

$$V^{\pi_0}(1) = V^{\pi_0}(2) = V^{\pi_0}(3) = \dots = V^{\pi_0}(n) = 2$$

Policy iteration
برای لایه

$$\pi_1(1) = \max_a \left(\overset{1.25}{\underset{\downarrow \text{Stay}}{1(1 + 0.5 \times 2)}}, \overset{0.25}{\underset{\downarrow \text{eat}}{1(1 + 0.5 \times 2)}} \right) \rightarrow \pi_1(1) \rightarrow \text{Stay}$$

برای بقیه مقادیر هم $\pi_1(k)$ - لایه $k \leq n$ برابر Stay سوره به سیاست لایه به لایه State

$$V^*(1) = 2$$

سیاست در صفر

$$\pi^*(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')] \quad (\gamma = 0.5)$$

$$\rightarrow \pi^*(city_i) = \max \left(1 \times (r_i + 1 \times V^*(city_i)), P_i (0 + 1 \times V^*(city_{i+1})) + (1 - P_i) (0 + 1 \times V^*(city_i)) \right)$$

$$, P_i (0 + 1 \times V^*(city_{i+1})) + (1 - P_i) (0 + 1 \times V^*(city_i)) = \max \left(r_i + V^*(city_i), V^*(city_i) + P_i (V^*(city_{i+1}) - V^*(city_i)) + P_i (V^*(city_{i-1}) - V^*(city_i)) \right)$$

$$= \max \left(r_i + V^*(city_i), V^*(city_i) + P_i (V^*(city_{i+1}) - V^*(city_i)), V^*(city_i) + P_i (V^*(city_{i-1}) - V^*(city_i)) \right)$$

$$V^*(city_i) + P_i (V^*(city_{i-1}) - V^*(city_i))$$

هنا طبقاً إلى هذه الصيغة $V^*(city_i)$ لا يتغير حسب P_i لأن P_i دائماً يساوي 0.5

$$\left\{ r_i, P_i (V^*(city_{i+1}) - V^*(city_i)), P_i (V^*(city_{i-1}) - V^*(city_i)) \right\}$$

لذلك يمكننا تجاهل P_i في الحسابات وعلينا أن نكتب $\gamma = 0.5$ فقط

$$\rightarrow \text{Sample} = R(s, a, s') + \gamma \max_{a'} (s', a')$$

نستعمل $Q(s, a)$ بدلاً من $V^*(s)$

$$Q(s, a) \leftarrow (1 - \alpha) Q(s, a) + \alpha [\text{Sample}]$$

$$\gamma = 0.5, \alpha = 0.5$$

Sample 1: (1, stay, 4, 1)

$$\text{Sample 1} = 4 + \gamma \max_{a'} Q(1, a') = 4, \quad Q(1, \text{stay}) \leftarrow (1 - 0.5) Q(1, \text{stay}) + 0.5 (4)$$

$$\rightarrow \boxed{Q(1, \text{stay}) = 2} \checkmark, \quad Q(1, \text{east}) = 0, \quad Q(2, \text{west}) = 0, \quad Q(2, \text{stay}) = 0$$

Sample 2: (1, East, 0, 2)

$$\text{Sample 2} = 0 + \gamma \max_{a'} Q(2, a') = 0, \quad Q(2, \text{east}) \leftarrow (1 - \alpha) Q(2, \text{east}) + \alpha \text{Sample 2} = 0$$

$$\rightarrow Q(1, \text{stay}) = 2, \quad Q(1, \text{east}) = 0, \quad Q(2, \text{west}) = 0, \quad Q(2, \text{stay}) = 0$$

Sample 3: (2, stay, 6, 2)

$$\rightarrow \text{Sample 3} = 6 + \gamma \max_a Q(2, a) = \underline{6}$$

$$Q(2, \text{stay}) \leftarrow (1-\alpha) Q(2, \text{stay}) + \alpha \text{Sample 3} = \underline{3}$$

$$Q(1, \text{stay}) = 2, Q(1, \text{east}) = 0, Q(2, \text{west}) = 0, Q(2, \text{stay}) = \underline{3}$$

Sample 4: (2, west, 0, 1)

$$\text{Sample 4} = R(s, a, s') + \gamma \max_a Q(s', a) = 0 + \gamma \max_a Q(1, a) = \underline{1}$$

$a = \text{stay}$

$$Q(2, \text{west}) = (1-\alpha) Q(2, \text{west}) + \alpha \text{Sample 4} = 0.5 \times 0 + \alpha \times 1 = \underline{0.5}$$

$$Q(1, \text{stay}) = 2, Q(1, \text{east}) = 0, Q(2, \text{west}) = 0.5, Q(2, \text{stay}) = \underline{3}$$

Sample 5: (1, stay, 4, 1)

$$\text{Sample 5} = 4 + \gamma \max_a Q(1, a) = \underline{5}$$

$$Q(1, \text{stay}) \leftarrow (1-\alpha) Q(1, \text{stay}) + \alpha \text{Sample 5} = 0.5 \times 2 + 0.5 \times 5 = \underline{3.5}$$

القيمة المتوقعة

| | $Q(1, \text{stay})$ | $Q(1, \text{East})$ | $Q(2, \text{west})$ | $Q(2, \text{stay})$ |
|-------------------|---------------------|---------------------|---------------------|---------------------|
| (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 |

اوله نسبت (ت) باید بقدر Max

$$\text{Max} [r_i, P_i (v^*(\text{city}(i+1)) - v^*(\text{city}(i))) \text{ و } P_i (v^*(\text{city}(i+1)) - v^*(\text{city}(i)))]$$

اینه اوله آن کسم کی لزمین به بقدر max کانه سنی به بقدر r_i , P_i بقدر

اگر r_i ها برابر باشند آن وقت $v^*(\text{city}(i)) - v^*(\text{city}(i+1))$ برابر با صفر شود دلیل اینه که

اینکه شهرها هم برابر است در صحت پ که $r_i \geq 1$ یعنی برابر هم می بود لکنه که در بیان صحت نیست

Stay انتخاب کرد اگر $v^*(\text{city}(i))$ از قیمت شهرهای مجاورش بیشتر باشد آن وقت باز

Stay انتخاب شود