

۱۱۵۱۰۵۰۵۰

مددا هاجایی

سوال شماره ۱: Markov Decision Process

در این مدل یک جیل مارکوف زمانی است که با توجه به وضعیت فعلی و عملی که انجام می‌دهیم، به یک وضعیت جدید می‌رویم. (۱) مقدار انتظارت که با توجه به این دانسته که تمام وضعیت‌ها و عمل‌ها را می‌توانیم به صورت یک مدل درآوریم.

حل تمام Transition را با این گونه می‌کنیم که Action باید به افعال $\frac{1}{10}$ به سمت راست می‌خواهیم به در برور و با افعال $\frac{1}{10}$ می‌توانیم به سمت چپ می‌رویم یا می‌توانیم به سمت راست می‌رویم. به این یک مدل Non deterministic داریم. حال اگر بخواهیم به در برور Action را به همان حالت بماند.

فصل ۱: یابی نتایج در حالت Value Iteration را می‌توانیم به این روش $\frac{9}{10}$ می‌توانیم به مقادیر مشخصی $\frac{1}{10}$ و $\frac{2}{10}$ تغییر می‌دهیم و می‌توانیم به مقادیر را به این گونه قرار می‌دهیم:

$$V_{i+1}(s) = \max \left(\sum_{s'} T(s, a, s') R(s, a, s') + \gamma V_i(s') \right)$$

Value Iteration یا به همین شکل از Value Iteration استفاده می‌کنیم.

این در

$$V_{i+1}(s) = \max_a \left(\sum_{s'} T(s, a, s') (R(s, a, s') + \gamma V_i(s')) \right)$$

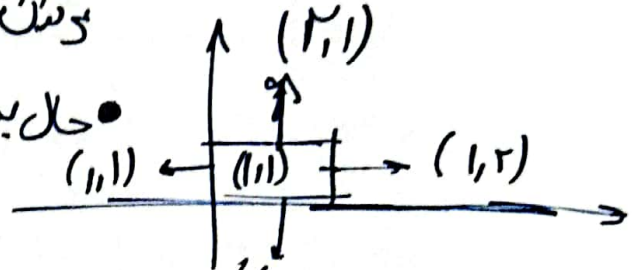
توجه: مقدار γ را می‌توانیم به این گونه قرار دهیم که مقدار انتظارت که با توجه به این دانسته که تمام وضعیت‌ها و عمل‌ها را می‌توانیم به صورت یک مدل درآوریم.

| | |
|--------|--------|
| (۲, ۲) | (۲, ۱) |
| (۱, ۱) | (۱, ۲) |

• حال بی مقدار از است معلوم.

$$V(1,2) = \max_a \left(\sum_{s'} T(s,a|s') (R(s,a|s') + \gamma V(s')) \right)$$

این ها که از است (1,1)
 و (2,1) و (1,2) و (1,1) جواب
 حالت بی مقدار از
 و (1,1) و (1,2) و (2,1) و (1,1) جواب
 و (1,1) و (1,2) و (2,1) و (1,1) جواب



• در صورتی که در هر دو حالت
 و (1,1) و (1,2) و (2,1) و (1,1) جواب
 و (1,1) و (1,2) و (2,1) و (1,1) جواب

•

① up ↑

$$\sum_{s'} T((1,1), up, s') (R(s,a|s') + \gamma V(s'))$$

• مقدار را می گیریم
 جای گذشتیم به اینجا
 Reward = 0

$$= \left(\frac{1}{10} \right) (0 + \frac{9}{10}) + \left(\frac{1}{10} \right) (0 + \frac{9}{10}) + \left(\frac{1}{10} \right) (0 + \frac{9}{10}) = 0$$

② down ↓

$$\sum_{s'} T((1,1), down, s') (R(s,a|s') + \gamma V(s'))$$

• در این حالت هم چون مانند قبل به ازای تمام حالت ها که در آن قرار
 و (1,1) و (1,2) و (2,1) و (1,1) جواب
 و (1,1) و (1,2) و (2,1) و (1,1) جواب
 و (1,1) و (1,2) و (2,1) و (1,1) جواب

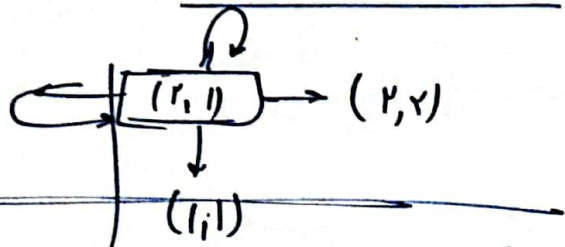
$$\rightarrow \sum_{s'} T((1,1), down, s') (R(s,a|s') + \gamma V(s'))$$

= 0

②

- ③ Right } برای این در حالت هم احتمالاً حالت
 ④ Left } عملی نیست برابر با افزودن هزینه
 در هر گام

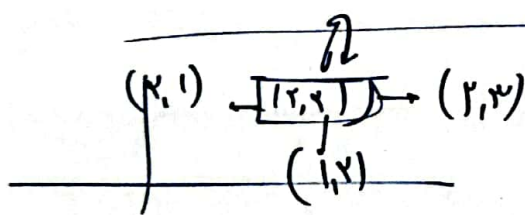
$$\Rightarrow V(1,1) = \max_a (0) = 0 \rightarrow \boxed{V(1,1) = 0}$$



• حل $(2,1)$ را بررسی می‌کنیم.

در این حالت هم ما فقط می‌توانیم به استی می
 $(2,1)$ و $(1,1)$ که خود را می‌بیند
 و در این حالت هر تار یا تارهای برابر با افزودن و هم چنین مقدارهای مختلف را برابر
 با افزودن است.

$$\Rightarrow V(2,1) = 0$$



• حل $(2,2)$ را بررسی می‌کنیم. در این حالت مقدار
 $(2,3)$ و $(2,1)$ یکسان است که برابر با 0 است.

$$\max_a \sum_{s'} T((1,2), a, s') (R((1,2), a, s') + \gamma V(s'))$$

برای حالت a بررسی می‌کنیم.

$$\textcircled{1} a \rightarrow \text{up}$$

$$\sum_{s'} T = \frac{1}{10} (0 + \frac{9}{10} (0)) + \frac{1}{10} (0 + \frac{9}{10} (0)) + \frac{1}{10} (0 + \frac{9}{10} (0)) = 0,9$$

$$\sum_{s'} = \frac{1}{10} (0 + \frac{9}{10} (0)) + \frac{1}{10} (0 + 0 (\frac{9}{10})) + \frac{1}{10} (0 + 0 (\frac{9}{10})) = 0,9$$

$$\sum_{s'} = \frac{1}{10} (0 + \frac{9}{10} (0)) + \frac{1}{10} (0 + 0 (\frac{9}{10})) + \frac{1}{10} (0 + 0 (\frac{9}{10})) = 0,9$$

$$\sum_{s'} = 0 \rightarrow \max = \boxed{V(1,1)}$$

• این مقادیر را با این نفع در یک فرم قرار دهیم مقدار Reward در رشتن به است حل پایانی باید ۰ و ۰ باشد و در صورت مقادیر ۰ و ۰ در است های آن مقادیر باید برابر ۰ و ۰ باشد و در هر کمر در صورت معادله بیان شود.

حل بدلی $\sum_{s'} T(s,a,s') (R(s,a,s') + \gamma V(s'))$ ۱,۲ ۱,۲ ۱,۲

① $\sum_{s'} (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) = 0.95$ ۰,۹۵

~~②~~ $\sum_{s'} (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) = 0.95$ ۰,۹۵

$\rightarrow \sum_{s'} (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) = 0.95$ ۰,۹۵

$\leftarrow \sum (0 + 0 + 0) = 0$

$\Rightarrow \text{max} = 0$

~~$\rightarrow \text{max} = 0$~~

$\rightarrow \text{max} (-0.95, 0.95, 0, -0.44) = 0$

سپ مقادیر این نفع را در Reward نفع لژیو در Iteration

| | | |
|---|-----|----|
| ۰ | ۰,۹ | +۵ |
| ۰ | ۰ | -۵ |

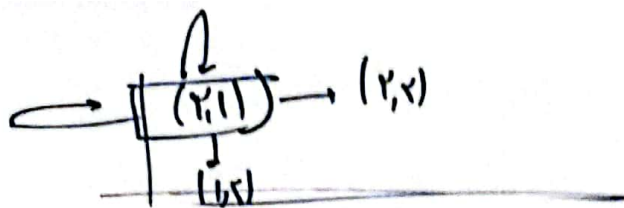
• حل در دو پایایی نیز:

ابتدا است $\sum_{s'} T(s,a,s') (R(s,a,s') + \gamma V(s'))$ ۱,۱ ۱,۱ ۱,۱

① $\sum_{s'} (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) = 0.95$

② $\rightarrow \sum_{s'} (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) = 0.95$ ۰,۹۵

③ $\rightarrow \sum_{s'} (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) + (\frac{1}{10} (0 + 0.9 \frac{9}{10})) = 0.95$ ۰,۹۵



$(2,1) = \text{load}$

$$\downarrow \textcircled{1} \sum \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 1(1)) = \underline{\underline{0.418}}$$

$$\uparrow \textcircled{2} \sum \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 4(1)) = 0.418$$

$$\rightarrow \textcircled{3} \sum \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) = \underline{\underline{0.147}}$$

$$\leftarrow \textcircled{4} \sum \frac{1}{1} (0 + 0(\frac{9}{1})) + 0 + 0 = 0$$

$$\rightarrow \text{max} = \underline{\underline{0.147}} \rightarrow \underline{\underline{V(2,1) = 0.147}}$$

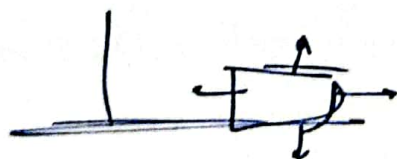
$$\downarrow \textcircled{1} \sum \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) = \underline{\underline{0.147}}$$

$$\uparrow \textcircled{2} \sum \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) = 0.147$$

$$\rightarrow \textcircled{3} \sum \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) = \underline{\underline{0.147}}$$

$$\leftarrow \textcircled{4} \sum \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) = \underline{\underline{0.147}}$$

$$\rightarrow \text{max} = \underline{\underline{0.147}} \rightarrow \underline{\underline{V(2,1) = 0.147}}$$



$(1,2) = \text{load}$

$$\downarrow \textcircled{1} \sum \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) = -0.90$$

$$\uparrow \textcircled{2} \sum \frac{1}{1} (0 + 4(1)) + \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) = \underline{\underline{0.147}}$$

$$\rightarrow \textcircled{3} \sum \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) = -0.90$$

$$\leftarrow \textcircled{4} \sum \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) + \frac{1}{1} (0 + 0(\frac{9}{1})) = 0.147$$

$$\Rightarrow \text{max} = \underline{\underline{0.147}} \rightarrow \underline{\underline{V(1,2) = 0.147}}$$

(5)

سپتقایریهای بین کشورهاست: در دردم از Value Iteration

| | | |
|-------|-------|-----|
| ۵,۴۷۲ | ۸,۲۸۴ | + ۵ |
| ۰ | ۴,۵۲۲ | - ۵ |

که البته باید فرمود که
ما مقدار Demand را بر
رفتاریهای پایانی داریم
و هر چه مقدار
پایانی ثابت بداند:

| S | (۱) | (۱,۲) | (۱,۳) | (۲,۱) | (۲,۲) | (۲,۳) |
|-------|-----|-------|-------|-------|-------|-------|
| V_0 | ۰ | ۰ | -۵ | ۰ | ۰ | ۵ |
| U_1 | ۰ | ۰ | -۵ | ۰ | ۷,۱ | ۵ |
| V_1 | ۰ | ۴,۵۲۲ | -۵ | ۵,۴۷۲ | ۸,۲۸۴ | ۵ |

حال مقدار Policy برای مدت ۱۰

$\pi^*(s) =$
Policy این نوعی است که
چون که برابر با مقدار Q ها اندام بر مقدار

• $R_{i+1}(s) = \arg \max_a \sum_s T(s,a,s') [R(s,a,s') + \gamma V^{\pi}(s')]$
(در این باره است که می دانیم)

$R_{i+1}(1,1) = \arg \max:$

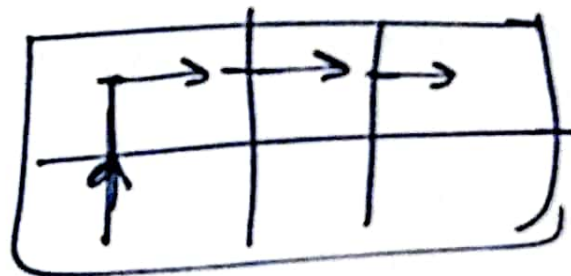
$\rightarrow \frac{1}{10} (0 + \frac{9}{10} (5,472)) + \frac{1}{10} (0 + \frac{9}{10} (0)) + \frac{1}{10} (0 + \frac{9}{10} (4,522)) = 4,34$
 $\rightarrow \frac{1}{10} (0 + \frac{9}{10} (4,522)) + \frac{1}{10} (0 + \frac{9}{10} (5,472)) + \frac{1}{10} (0 + \frac{9}{10} (0)) = 3,74$
 $\rightarrow \frac{1}{10} (0 + \frac{9}{10} (0)) + \frac{1}{10} (0 + \frac{9}{10} (0)) + \frac{1}{10} (0 + \frac{9}{10} (5,472)) = 0,492$
 ④ $\rightarrow \frac{1}{10} (0 + \frac{9}{10} (0)) + \frac{1}{10} (0 + \frac{9}{10} (0)) + \frac{1}{10} (0 + \frac{9}{10} (4,522)) = 0,4092$

$\rightarrow \text{max}_s \sum_i \pi_i \rightarrow \text{arg max}_s = \uparrow \rightarrow \underline{\underline{u_p}}$
 (پس از آنکه u_p انتخاب شود)

$$R_{i+1}(s) = \text{arg max}_a \sum_{s'} T(s, a, s') (R_i(s, a, s') + \gamma V^p(s'))$$

که گاهی بقیه است - هام ی ب به سبب این داریم

| | | |
|-------------------------------|-------------------------------|--|
| $\underline{u_R} \rightarrow$ | $\underline{u_R} \rightarrow$ | |
| $\uparrow R_1$ | $\uparrow u$ | |



و می آید به نتیجه خود

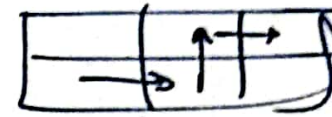
$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots$$



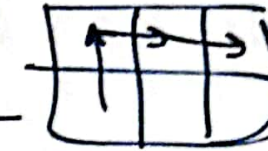
→ Rewards

Path ① $\rightarrow (1,1) \rightarrow (1,2) \rightarrow (1,3)$

Path ② $\rightarrow (1,1) \xrightarrow{0} (1,2) \xrightarrow{-5} (2,2) \rightarrow (2,3)$



Path ③ $\rightarrow (1,1) \xrightarrow{0} (2,1) \xrightarrow{5} (2,2) \rightarrow (2,3)$



~~Monte Carlo~~ • Monte Carlo Stimulation:

Path ① $(1,1) \rightarrow -5 \quad (1,2) \rightarrow -5$

Path ② $(1,1) \rightarrow 5 \quad (1,2) \rightarrow 5 \quad (2,2) \rightarrow 5$

Path ③ $(1,1) \rightarrow 5 \quad (2,1) \rightarrow 5 \quad (2,2) \rightarrow 5$

• مقدار میانگین:

$$(1,1) \rightarrow \frac{-5 + 5 + 5}{3} = \underline{1.67}$$

$$(2,2) \rightarrow \frac{5 + 5}{2} = \underline{5}$$

$$\left. \begin{array}{l} V(1,1) = 1.67 \\ V(2,2) = 5 \end{array} \right\}$$

۴) Agar به اساس معائنہ TD یا دیگر بانسٹا از نرخ یا دیگر $\left(\frac{1}{1.0}\right)$

و بازنویس مقایسه ای هند (کریه جنبه فقهی الهامی) است. بواز (۱۲)

• 5 دلائل TD (amny) اس میں مل سکتے ہیں :

$$S_{\text{simple}} = R(s, q(s, s')) + \gamma V^R(s)$$

- $V^R(s)$ = $(1 - \alpha) V^D(s) + \frac{\alpha \text{ Sample}}{h}$

• بنویسیم تا هم میانه بین وزن های جدید و قدیم
 α Learning Rate

$$V^{\pi}(s) = (1-\alpha)V^R(s) + \alpha \left(R(s, \pi(s), s') + \gamma V^R(s') \right)$$

$$\tau \Rightarrow (1,1)$$

$$\rightarrow V(1,1), (1 - \frac{1}{1.}) (0) + \frac{1}{1.} (0 + \frac{9}{1.} (0)) = 0$$

$$V(1, r) = \left(\frac{9}{10}\right) V(1, r) + \frac{1}{10} (0 + \frac{9}{10} (0)) = 0$$

$$V_{(1,1)} = 0.1 \cdot \frac{1}{1.1} \left(\frac{9}{1.1} + \frac{1.5}{1.1} \right) = -0.049$$

$$V(1,1) = -0,45 + 0,1(0,9 \times 0 + 0,45) = -0,405$$

$$v(r) = 0 + 0,1 \left(0 + \frac{9}{1} \right) \times 10^6, \frac{9 \times 10^6}{1}$$

| S | $(1,1)$ | $(1,x)$ | $(1,y)$ | $(y,1)$ | (x,x) | (y,y) |
|-------|---------|---------|---------|---------|---------|---------|
| V_0 | 0 | 0 | -0 | 0 | 0 | 0 |
| V_1 | | -0,40 | -0 | 0 | 0 | 0 |
| V_2 | -0,040 | -0,400 | -0 | 0 | 0,400 | 0 |