

## تصمیم‌گیری

(الف) احتمال اینکه اولین شیر در  $n$  این پرتاب ظاهر شود،  $2^{-n}$  است، پس:

$$EMV(L) = \sum_{n=1}^{\infty} 2^{-n} \cdot 2^n = \sum_{n=1}^{\infty} 1 = \infty$$

(ب) پانچ‌های نوعی در بازه‌ی  $\$4$  تا  $\$100$  است.

(ج) فرض می‌کنیم سرمایه اولیه (پس از پرداخت  $c$  برای ورود به بازی) برابر با  $\$(k-c)$  باشد.

در این صورت:

$$U(L) = \sum_{n=1}^{\infty} 2^{-n} \cdot (a \log_2(k-c+2^n) + b)$$

برای سادگی فرض می‌کنیم  $k-c=0$  در این صورت:

$$\begin{aligned} U(L) &= \sum_{n=1}^{\infty} 2^{-n} \cdot (a \log_2(2^n) + b) \\ &= \sum_{n=1}^{\infty} 2^{-n} \cdot an + b \\ &= 2a + b \end{aligned}$$

(د) حداکثر مبلغ  $c$  از حل معادله زیر بدست می‌آید:

$$a \log_2 k + b = \sum_{n=1}^{\infty} 2^{-n} \cdot (a \log_2(k-c+2^n) + b)$$

برای حالت ساده‌ی فرض شده:  $(k=c)$

$$a \log_2 k + b = \sum_{n=1}^{\infty} 2^{-n} \cdot (a \log_2 2^n + b) \Rightarrow$$

$$a \log_2 c + b = 2a + b \Rightarrow$$

$$c = 4$$

۲) ایتان بودن لازم می‌دارد که عامل ترجیح‌های یکسانی بین صفت دنباله‌های

$$[S_0, S_1, S_2, \dots] \text{ و } [S_0, S'_1, S'_2, \dots]$$

و بین صفت دنباله‌های

$$[S_1, S_2, \dots] \text{ و } [S'_1, S'_2, \dots]$$

داشته باشد. اگر سودمندی یک دنباله ماژیم پاداش باشد، ایتان بودن به سادگی نقض می‌شود.  
براشال:

$$[4, 3, 0, 0, 0, \dots] \sim [4, 0, 0, 0, 0, \dots]$$

یا

$$[3, 0, 0, 0, \dots] > [0, 0, 0, 0, \dots]$$

هنوز می‌توانیم  $U^{\pi}(s)$  را به عنوان ماژیم پاداش مورد انتظار با اجرای  $\pi$  با شروع از  $s$  تعریف کنیم.  
با این حال ترجیح‌های عامل عجیب و غریب به نظر می‌رسد.  
براشال اگر حالت جاری  $s$  پاداش  $R_{max}$  داشته باشد، عامل بین همی‌کش‌ها بی تفاوت خواهد بود اما همین  
که کش اجرا شود و عامل دیگر در  $s$  باشد، ناگهان شروع به مراقبت در مورد آنچه در بعد اتفاق می‌افتد، خواهد کرد.

۳) یک مسئله جستجوی شاخه با موارد زیر تعریف می‌شود:

$S_0$  حالت آغازین

$S(s)$  تابع مابعد (که مجموعه‌ای از زوج‌های کش-حالت  $(a, s')$  را برمی‌گرداند)

$C(s, a, s')$  یک تابع هزینه گام

GOAL-TEST (S) تابع آزمون هدف

یک راه حل بهینه یک مسیر کمترین هزینه از  $S_0$  به هر هدفی است.

برای ساخت MDP شناظر، تعریف می‌کنیم:  $R(s, a, s') = -C(s, a, s')$  اگر اندک  $s$  یک حالت هدف

باشد که در این صورت  $R(s, a, s') = 0$ .

همچنین تعریف می‌کنیم  $T(s, a, s') = 1$  اگر  $(a, s') \in S(s)$  و غیر  $\gamma = 1$ .

راه حل بهینه برای این MDP یک سیاست (policy) است که مسیر کمترین هزینه از هر حالت تا نزدیک‌ترین  
حالت هدف به آن را می‌یابد.

الف) عامل به صورت محدودی خواهد تا جای ممکن سریع تر به حالت 3 برسد، زیرا برای هر گام زمانی که حالت های 1 و 2 صرف خواهد کرد، هزینه می پردازد.

اما تخفیفی که به حالت 3 می رسد، (کنش b) دارای احتمال پایین است، بنابراین عامل باید هزینه ای که متحمل می شود (در هنگام تلاش برای رسیدن به حالت پایانی) را می نیمم کند.

این پیشنهاد می کند که عامل باید به طور معین در حالت 1 کنش a را آزمایش کند. در حالت 2 می تواند بهتر باشد که کنش a را برآورد و در حالت 1 آزمایش کند (که گمان بهتری برای انتظار پذیرش به حالت 3 است) به جای اینکه مستقیماً نخواهد وارد حالت 3 شود. این تصمیم در حالت 2 شامل یک بده ستان عددی است.

ب) به کارگیری الگوریتم تکرار سیاست در گام های تناوب؛ تعیین ارزش و به هنگام سازی سیاست

Initialization:  $U \leftarrow \langle -1, -2, 0 \rangle$ ,  $\pi \leftarrow \langle b, b \rangle$ .

Value determination:

$$\begin{cases} u_1 = -1 + 0.1 u_3 + 0.9 u_1 \\ u_2 = -2 + 0.1 u_3 + 0.9 u_2 \\ u_3 = 0 \end{cases} \Rightarrow \begin{aligned} u_1 &= -10 \\ u_2 &= -20 \end{aligned}$$

Policy update:

In state 1,  $\sum_j T(1, a, j) u_j = 0.8 \times (-20) + 0.2 \times (-10) = -18$

$\sum_j T(1, b, j) u_j = 0.1 \times 0 + 0.9 \times (-10) = -9$

پس در حالت 1 هنوز کنش a ترجیح داده می شود.

In state 2,  $\sum_j T(2, a, j) u_j = 0.8 \times (-10) + 0.2 \times (-20) = -12$

$\sum_j T(2, b, j) u_j = 0.1 \times (0) + 0.9 \times (-20) = -18$

پس در حالت 2 کنش a ترجیح داده می شود.

قراری دهیم: unchanged?  $\leftarrow$  false و ادامه دهیم:

Value determination:

$$\begin{cases} u_1 = -1 + 0.1 u_3 + 0.9 u_1 \\ u_2 = -2 + 0.8 u_1 + 0.2 u_2 \\ u_3 = 0 \end{cases} \Rightarrow$$

$$u_1 = -10, u_2 = -15$$

Policy update:

In state 1,  $\sum_j T(1, a, j) u_j = 0.8 \times (-15) + 0.2 \times (-10) = -14$

$$\sum_j T(1, b, j) u_j = 0.1 \times 0 + 0.9 \times (-10) = -9$$

پس کنش  $b$  در حالت 1 هنوز ترجیح داده می شود.

In state 2,  $\sum_j T(2, a, j) u_j = 0.8 \times (-10) + 0.2 \times (-15) = -11$

$$\sum_j T(2, b, j) u_j = 0.1 \times 0 + 0.9 \times (-15) = -13.5$$

پس کنش  $a$  هنوز در حالت 1 ترجیح داده می شود.

unchanged? بدون تغییر و true می ماند و خاتمه می یابد.

\* توجه کنید که سیاست حاصل باشد مطابق دارد: وقتی در حالت 2 هستیم سعی می کنیم به حالت 1 برویم و وقتی در حالت 1 هستیم سعی می کنیم وارد حالت 3 شویم.

ج) سیاست آغازین با کنش  $a$  در هر دو حالت،  $\pi = \langle a, a \rangle$ ، یک سئدی غیر قابل حل منبری شود.  
سئدی تعیین ارزش آغازین فرم

$$\begin{cases} u_1 = -1 + 0.2u_1 + 0.8u_2 \\ u_2 = -2 + 0.8u_1 + 0.2u_2 \\ u_3 = 0 \end{cases}$$

رادر دارد که دو معادله اول آن ناسازگار هستند.

اگر سعی کنیم آنها را به طور عددی حل کنیم، مقادیری می یابیم که به  $-\infty$  میل می کنند.

استفاده از ضرب تخفیف  $\gamma$  با محدود کردن جویمه ها می تواند عامل در حرکت از حالات متغی می شود، (هزینه های تخفیف

یافته می شود انتظار) منجر به راه حل های خوش تعریف می شود. با این وجود انتخاب ضرب تخفیف بر سیاست

حاصل اثر می گذارد. برای  $\gamma$  کوچک، هزینه های تخفیف شده در آینده دور، نقش قابل چشم پوشی در محاسبه

ارزش بازی می کند، زیرا "هم نزدیک صفر است."

به عنوان یک نتیجه عامل می تولید کنش  $a$  در حالت 2 انتخاب کند، زیرا هزینه های تخفیف یافته ای کوتاه مدت باقی مانده

در حالت های غیر پایانی (حالت های 1 و 2) از هزینه های تخفیف یافته ای بلند مدت کنش  $a$  که به طور مکرر شکست

می خورد و عامل را در حالت 2 باقی می گذارد، وزن بیشتری پیدا می کند.

(۵)  $a$  کلیدگار محاسباتی  $\max$  و  $\Sigma$  در جای مناسب خود است.

برابر  $R(s, a)$  داریم :

$$U(s) = \max_a [R(s, a) + \gamma \sum_{s'} T(s, a, s') U(s')]$$

برابر  $R(s, a, s')$  داریم :

$$U(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma U(s')]$$

(b) برابر این کار راه حل های مختلفی وجود دارد. یک راه ایجاد یک "pre-state"  $pre(s, a, s')$  برابر هر  $s, a, s'$  است به این صورت که اجرای  $a$  در  $s$  به  $s'$  منجر می شود بلکه منجر به  $pre(s, a, s')$  می شود. در این حالت این واقعیت که شده است که عامل از حالت  $s$  آمده است و انجام  $a$  آن را به اینجا رسانده است. از این pre-state تنها یک کنش  $b$  است که همیشه به  $s'$  منجر می شود.

فرض می کنیم MDP جدید، گذار  $T'$ ، پاداش  $R'$  و نرخ تخفیف جدید  $\gamma'$  را داشته باشد، در این صورت

$$T'(s, a, pre(s, a, s')) = T(s, a, s')$$

$$T'(pre(s, a, s'), b, s') = 1$$

$$R'(s, a) = 0$$

$$R'(pre(s, a, s'), b) = \gamma^{-\frac{1}{2}} R(s, a, s')$$

$$\gamma' = \gamma^{\frac{1}{2}}$$

(c) یا ایده ی بخش  $b$ ، می توانیم حالت های  $post(s, a)$  را برای هر  $s$  و  $a$  ایجاد کنیم به گونه ای که :

$$T'(s, a, post(s, a, s')) = 1$$

$$T'(post(s, a, s'), b, s') = T(s, a, s')$$

$$R'(s) = 0$$

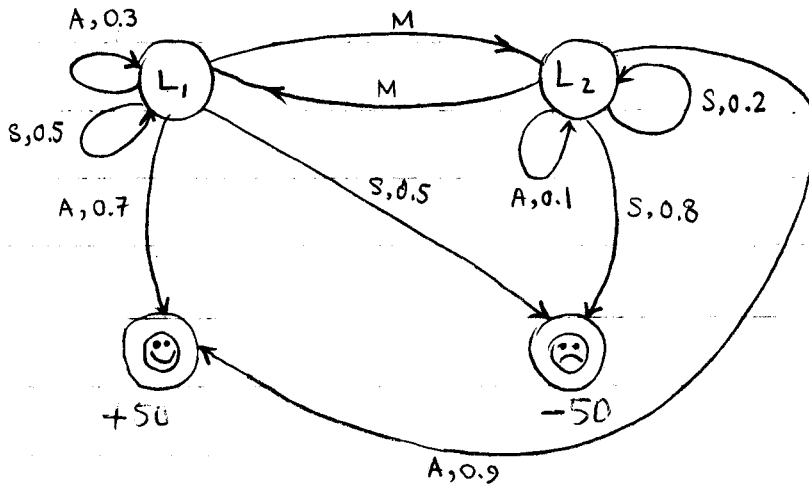
$$R'(post(s, a, s')) = \gamma^{-\frac{1}{2}} R(s, a)$$

$$\gamma' = \gamma^{\frac{1}{2}}$$

Markov Decision Process

(۶)

نمودار زیر یک مدل MDP از یک fierce battle را به تصویر کشیده است. (دعوی ضعیف)



در این دعوا

می‌توان بین مکان‌های  $L_1$  و  $L_2$  جابجا شد. یکی از این مکان‌ها به حرف نزدیک‌تر است. اگر از نزدیک‌ترین حالت عمل کنید،

- شانس بیشتری برای موفقیت دارید (90٪) (در مقابل 70٪ برای مکان دورتر از حرف)

- به هر حال ممکن است دیده شوید (با شانس 80٪) و کشته شوید. (در مقابل شانس 50٪ برای مکان

دورتر).

شما می‌توانید دیده شوید اگر در همان مکان خود باقی بمانید.

شما لازم دارید که یک نقشه action برای این موقعیت داشته باشید.

بیمکان‌ها، عمل‌های ممکن را نشان می‌دهند: M : Move : قطعی

A : Attack : تصادفی

S : stay : تصادفی

بر روی مکان‌ها،  $(a,p)$  نوشته شده است. (عمل و احتمال آن)

همه پاداش‌ها منفی هستند به جز برای حالات پایانی که در آن موقعیت شما با پاداش +50 نشان داده می‌شود،

در حال که موفقیت حرف، با پاداش -50 برای شما نشان داده می‌شود.

با بکارگیری فاکتور تخفیف  $\gamma = 0.9$ ، سیاست بهینه (optimal policy) یا همان نقشه‌نگاشی‌ها را

محاسبه کنید.

راه حل :

نماینده action-value ها برای همه حالات و کنش ها لازم است.

$$V^*(Far) = \max \{ Q^*(Far, M), Q^*(Far, S), Q^*(Far, A) \}$$

$$V^*(Close) = \max \{ Q^*(Close, M), Q^*(Close, S), Q^*(Close, A) \}$$

$$V^*(\text{😊}) = +50$$

$$V^*(\text{😞}) = -50$$

در value-iteration با برآوردهای آغازین زیر شروع می‌کنیم :

$$\forall s \forall a \quad Q_0(s, a) = 0$$

پس همه action-value ها را با توجه به قاعده زیر به‌کم می‌کنیم :

$$Q_{n+1}(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) V_n(s')$$

$$V_n(s') = \max_{a'} Q_n(s', a') \quad \text{که در آن}$$

بنابراین در تکمیل اول الگوریتم به دست می‌آوریم :

$$Q_1(Far, A) = 0 + 0.9 [0.7(+50) + 0.3(0)] = 31.5$$

$$Q_1(Far, S) = 0 + 0.9 [0.5(-50) + 0.5(0)] = -22.5$$

$$Q_1(Close, A) = 0 + 0.9 [0.9(+50) + 0.1(0)] = 40.5$$

$$Q_1(Close, S) = 0 + 0.9 [0.8(-50) + 0.2(0)] = -36$$

ادزش عمل‌های M همان صفر باقی می‌ماند :

$$Q_1(Far, M) = 0 + 0.9 (1 \times 0) = 0$$

$$Q_1(Close, M) = 0 + 0.9 (1 \times 0) = 0$$

پس از این‌ها مقدار این دو حالت خواهد بود :

$$V_1(Far) = \max_a Q(Far, a) = 31.5$$

$$V_1(Close) = \max_a Q(Close, a) = 40.5$$

که متناظر با کنش همه (A) در هر دو حالت می‌باشد.

در تکرار بعدی خواهیم داشت :

$$Q_2(\text{Far}, A) = 0 + 0.9 [0.7(+50) + 0.3(31.5)] = 40.5$$

$$Q_2(\text{Far}, S) = 0 + 0.9 [0.5(-50) + 0.5(31.5)] = -8.325$$

$$Q_2(\text{Far}, M) = 0 + 0.9 V_1(\text{Close}) = 0.9(40.5) = 36.45$$

$$Q_2(\text{Close}, A) = 0 + 0.9 [0.9(+50) + 0.1(40.5)] = 44.145$$

$$Q_2(\text{Close}, S) = 0 + 0.9 [0.8(-50) + 0.2(40.5)] = -28.71$$

$$Q_2(\text{Close}, M) = 0 + 0.9 V_1(\text{Far}) = 0.9(31.5) = 28.35$$

د

$$V_2(\text{Far}) = 40$$

$$V_2(\text{Close}) = 44.145$$

که تناظر با فعل عمل در هر دو حالت است.

این فرایند می تواند ادامه یابد تا مقادیر بین تکرارهای متوالی تغییر چندانی نداشته باشند.

بر اساس آنچه تا این نقطه دیده ایم، بهترین نقطه کنش عمل کردن در هر زمان ها بنظر می رسد.

آیا می توان بدون بسط بازی کاربستری چیز بهتری گفت ؟

بدیهی است که عمل "Stay" در هر دو حالت زیر بهینه است.

در حالت Close، بدیهی است که بهترین کار برای انجام به طور مداوم عمل کردن "Attack" است.

(که برای آن هزینه ای نمی پردازیم)

در واقع می توانیم ارزش ها را به طور تحلیلی در حد مناسب کنیم (تغییرات در هنگام بازی از یک تکرار به تکرار دیگر)

$$Q^*(\text{Close}, A) = \gamma \times 0.9 \times 50 + \gamma^2 (1 - 0.9)(0.9) \times 50$$

$$+ \gamma^3 (1 - 0.9)^2 (0.9) \times 50$$

+ ...

$$= \gamma \times 0.9 \times 50 (1 + \gamma (1 - 0.9) + (\gamma (1 - 0.9))^2 + \dots)$$

$$= \frac{0.9 \gamma}{1 - \gamma (1 - 0.9)} \times 50 \approx 39.1304$$



حال برای حالت Far، سؤال بین Attack یا Move به مکان نزدیک تر Close است.  
 ارزش‌ها را برای هر دوی این کنش‌ها محاسبه کنید (به همان روش قبلی):

$$Q^*(Far, M) = 0.9 V^*(Close) \approx 44.5055$$

$$Q^*(Far, A) = \frac{0.9 \times 0.7}{1 - 0.9 \times 0.3} \times 50 = 43.1507$$

بنابراین رفتن به Close بهتر است.

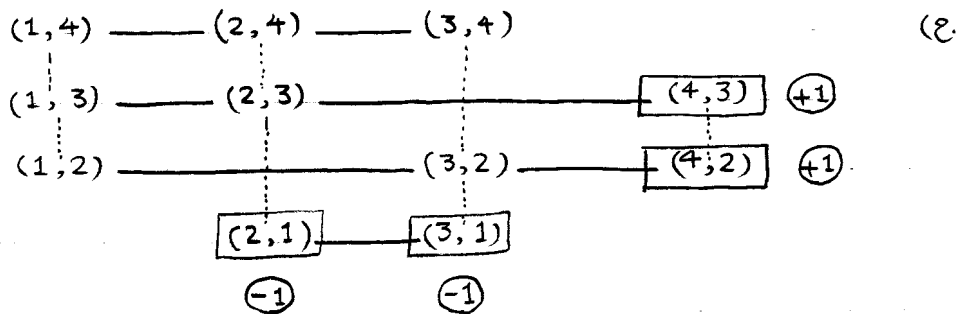
سیاست بهینه برای این صورت مساله (۱) رفتن به Close (تندیک تر شدن) و عمل کردن از آنجاست.  
 می‌توان تصور کرد که یک سیاست متفاوت حس بهتری برای این مساله ایجاد می‌کند؟

(۷)

$$U_A(s) = R(s) + \max_{\alpha} \sum_{s'} P(s'|a, s) U_B(s') \quad : U_A \text{ برابر}$$

$$U_B(s) = R(s) + \min_{\alpha} \sum_{s'} P(s'|a, s) U_A(s') \quad : U_B \text{ برابر}$$

ب) برای انجام تکرار ارزش، به طور ساده هر یک از معادلات فوق را به معادله‌ی جنگام سازی بهین تبدیل می‌کنیم و آنها را به طور متناوب به کار می‌بریم و هر یک را برای همی حالت‌ها به طور همزمان به کار می‌بریم. این فرآیند زمانی خاتمه می‌یابد که بردار سود مندی برای یک بازیکن، مثلاً بردار سود مندی قبلی برای همان بازیکن شود (یعنی دوگام قبل تر). [ توجه کنید که معمولاً  $U_B$  و  $U_A$  در ترتیب متعادل یکسان نمی‌باشند.]



د) اجرای الگوریتم تکرار ارزش در زیر نشان داده شده است. ارزش حالت پایان برای هر بازیکن با دایره دور آن نشان داده شده است. سایر ارزش‌ها با 0 مقدار دهی اولیه شده‌اند. این الگوریتم به صورت زیر ادامه می‌یابد:

	(1, 4)	(2, 4)	(3, 4)	(1, 3)	(2, 3)	(4, 3)	(1, 2)	(3, 2)	(4, 2)	(2, 1)	(3, 1)
$U_A$	0	0	0	0	0	(+1)	0	0	(+1)	(-1)	(-1)
$U_B$	0	0	0	0	-1	(+1)	0	-1	(+1)	(-1)	(-1)
$U_A$	0	0	0	-1	+1	(+1)	-1	+1	(+1)	(-1)	(-1)
$U_B$	-1	+1	+1	-1	-1	(+1)	-1	-1	(+1)	(-1)	(-1)
$U_A$	+1	+1	+1	-1	+1	(+1)	-1	+1	(+1)	(-1)	(-1)
$U_B$	-1	+1	+1	-1	-1	(+1)	-1	-1	(+1)	(-1)	(-1)

سیاست بهینه برای هر بازیکن به صورت زیر است:

	(1, 4)	(2, 4)	(3, 4)	(1, 3)	(2, 3)	(4, 3)	(1, 2)	(3, 2)	(4, 2)	(2, 1)	(3, 1)
$\pi_A^*$	(2, 4)	(3, 4)	(2, 4)	(2, 3)	(4, 3)		(3, 2)	(4, 2)			
$\pi_B^*$	(1, 3)	(2, 3)	(3, 2)	(1, 2)	(2, 1)		(1, 3)	(3, 1)			