

راه حل تکلیف شماره ی ۴  
فصل ۲۰ و ۲۱ یادگیری (یادگیری مدل ها احتمالاتی / یادگیری تقویتی)

۱) روی گرد بینی معرف کردن هر دو دارو است. روی گرد ماکزیم درست نایی معرف داروی anti-B است. در این مورد که نوشته از B وجود دارد، روی گرد بینی هنوز معرف هر دو دارو را پیشنهادی دهد در حالی که روی گرد ماکزیم درست نایی، معرف داروی anti-A است، چرا که داروی 40٪ شانس درست بودن دارد در برابر شانس 30٪ بابر هر یک از موارد B. البته این یک کار یکا تور است و شما ممکن است بپرسید یا شانس بزرگ تحت فشار باشید، حتی یک طرفدار ماکزیم درست نایی سرسخت که ممکن است مشابه این تجویز کند!

(۲) داریم:

$$L = -m(\log \sigma + \log \sqrt{2\pi}) - \sum_j \frac{(y_j - (\theta_1 x_j + \theta_2))^2}{2\sigma^2}$$

باشتق گیری نسبت به پارامترها:

$$\frac{\partial L}{\partial \theta_1} = - \sum_j \frac{x_j (y_j - (\theta_1 x_j + \theta_2))}{\sigma^2} = 0$$

$$\frac{\partial L}{\partial \theta_2} = - \sum_j \frac{(y_j - (\theta_1 x_j + \theta_2))}{\sigma^2} = 0$$

$$\frac{\partial L}{\partial \sigma} = - \frac{m}{\sigma} + \sum_j \frac{(y_j - (\theta_1 x_j + \theta_2))^2}{\sigma^3} = 0$$

که راه حل دستگاه فوق می شود:

$$\theta_1 = \frac{m(\sum_j x_j y_j) - (\sum_j y_j)(\sum_j x_j)}{m(\sum_j x_j^2) - (\sum_j x_j)^2}$$

$$\theta_2 = \frac{1}{m} \sum_j (y_j - \theta_1 x_j)$$

$$\sigma^2 = \frac{1}{m} \sum_j (y_j - (\theta_1 x_j + \theta_2))^2$$

(۳) دلف) بازنگرال گیری بر روی بازه‌ی  $[0, 1]$ ، ثابت نرمال سازی برای توزیع  $\text{beta}[a, b]$  به صورت  
 $\alpha = \Gamma(a+b) / \Gamma(a)\Gamma(b)$

دارد می شود که در آن  $\Gamma(x)$  تابع گاما است.

برابر اعداد صحیح  $a$  و  $b$  با استقرای روی  $a$  داریم: فرض می کنیم  $\alpha(a, b)$  ثابت نرمال سازی باشد.  
 برای حالت پایه داریم:

$$\alpha(1, b) = 1 / \int_0^1 \theta^0 (1-\theta)^{b-1} d\theta = -1 / \left[ \frac{1}{b} (1-\theta)^b \right]_0^1 = b$$

$$\frac{\Gamma(1+b)}{\Gamma(1)\Gamma(b)} = \frac{b \cdot \Gamma(b)}{1 \cdot \Gamma(b)} = b.$$

برابر گام استقرای فرض می کنیم  $a$  برابر  $b$  داریم:

$$\alpha(a-1, b+1) = \frac{\Gamma(a+b)}{\Gamma(a-1)\Gamma(b+1)} = \frac{a-1}{b} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$$

حال با بازنگرال گیری مجدد جزء مقدار  $\alpha(a, b)$  را ارزیابی می کنیم:

$$\begin{aligned} 1/\alpha(a, b) &= \int_0^1 \theta^{a-1} (1-\theta)^{b-1} d\theta \\ &= \left[ \theta^{a-1} \cdot \frac{1}{b} (1-\theta)^b \right]_0^1 + \frac{a-1}{b} \int_0^1 \theta^{a-2} (1-\theta)^b d\theta \\ &= 0 + \frac{a-1}{b} \frac{1}{\alpha(a-1, b+1)} \end{aligned}$$

پس داریم:

$$\alpha(a, b) = \frac{b}{a-1} \alpha(a-1, b+1) = \frac{b}{a-1} \frac{a-1}{b} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$$

(-) میانگین بازنگرال زیر داده می شود:

$$\begin{aligned} \mu(a, b) &= \alpha(a, b) \int_0^1 \theta \cdot \theta^{a-1} (1-\theta)^{b-1} d\theta = \alpha(a, b) \int_0^1 \theta^a (1-\theta)^{b-1} d\theta \\ &= \alpha(a, b) / \alpha(a+1, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{a\Gamma(a)\Gamma(b)}{(a+b)\Gamma(a+b+1)} = \frac{a}{a+b} \end{aligned}$$

$$d \text{beta}[a, b](\theta) / d\theta = 0$$

(ج) یافتن مقدار بیشترین:

$$\frac{d}{d\theta} (\alpha(a, b) \theta^{a-1} (1-\theta)^{b-1}) =$$

$$\alpha(a, b) [(a-1)\theta^{a-2} (1-\theta)^{b-1} - (b-1)\theta^{a-1} (1-\theta)^{b-2}] = 0 \Rightarrow$$

$$(a-1)(1-\theta) = (b-1)\theta \Rightarrow \theta = \frac{a-1}{a+b-2}$$

(۴)

فرآیند یادگیری تفاضل زمانی مقادیر Q (Q-learning) به

یادگیری تفاضل زمانی ارزش حالات V چیست؟

اگر یادگیری تفاضل زمانی بر روی ارزش های حالات استفاده شود،

استراج policy از روی ارزش های یادگرفته شده، دشوار خواهد بود.

(نیاز به مدل گذر T داریم تا بتوان سیاست را از روی ارزش ها محاسبه کنیم.)

برای یادگیری تفاضل زمانی مقادیر Q، یک سیاست می تواند مستقیماً با معادله زیر استراج شود:

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

(۵)

در برخی مسائل RL، پاداش‌ها (rewards) برای اهداف، ثبت و دیگر موارد صفر یا منفی هستند.

آیا علامت پاداش‌ها مهم است یا تنها فاصله میان آنها اهمیت دارد؟

با استفاده از تعریف دریافتی تخفیف یافته  $R_t$  (discounted reward)، ثابت کنید که افزودن یک مقدار

ثابت  $C$  به همه پاداش‌های ابتدایی، ثابت  $K$  را به ارزش همه‌ی حالات می‌افزاید و بنابراین

بر ارزش‌های نسی همه‌ی حالات تحت همه‌ی سیاست‌ها تأثیری نمی‌گذارد.

\* مقدار  $K$  را بر حسب  $C$  و  $\gamma$  محاسبه کنید.

راه حل:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$$

$$= \sum_{i=0}^{\infty} \gamma^i r_{t+i+1}$$

با اضافه کردن  $C$  به همه پاداش‌های ابتدایی، داریم:

$$r'_t = r_t + C$$

و از آنجا

$$R'_t = \sum_{i=0}^{\infty} \gamma^i r'_{t+i+1} = \sum_{i=0}^{\infty} \gamma^i (r_{t+i+1} + C) = \underbrace{\sum_{i=0}^{\infty} \gamma^i r_{t+i+1}}_{R_t} + \underbrace{\sum_{i=0}^{\infty} \gamma^i C}_K$$

در نتیجه:

$$R'_t = R_t + K$$

که در آن

$$K = C \sum_{i=0}^{\infty} \gamma^i = C \frac{1}{1-\gamma}$$

لذا تنها فاصله میان پاداش‌ها اهمیت دارند، نه مقادیر مطلق آنها.

انبات همگرنی Q-learning

راه حل : طرح کل انبات :

- حالت دینار اقطعی deterministic در نظر بگیرید که در آن هر  $(s, a)$  به تعداد دفعات نامتناهی مشاهده می شود.
- یک بازه کامل را به عنوان بازه ای که در خلال آن هر  $(s, a)$  مشاهده می شود در نظر بگیرید . interval
- نشان دهید که در خلال هر چنین بازه ای ، مقدار قدر مطلق خطا ( حداکثر خطا ) در جدول Q با ضرب  $\gamma$  کاهش می یابد .
- در نتیجه ، برای  $\gamma < 1$  ، پس از بی نهایت بهنگام سازی ، حداکثر خطا به سمت صفر میل می کند .

راه حل ←

رض کنید  $\hat{Q}_n$  جدولی باشد که پس از n بهنگام سازی به دست آمده است و  $e_n$  حداکثر خطا در این جدول باشد :

$$e_n = \max_s \max_a |\hat{Q}_n(s, a) - Q(s, a)|$$

حداکثر خطا را پس از (n+1) امین بهنگام سازی می یابیم :

$$\begin{aligned} e_{n+1} &= |\hat{Q}_{n+1}(s, a) - Q(s, a)| \\ &= |(r + \gamma \max_{a'} \hat{Q}_n(s', a')) - (r + \gamma \max_{a'} Q_n(s', a'))| \\ &= \gamma |\max_{a'} \hat{Q}_n(s', a') - \max_{a'} Q_n(s', a')| \\ &\leq \gamma \max_{a'} |\hat{Q}_n(s', a') - Q_n(s', a')| \\ &\leq \gamma \max_{s''} \max_{a'} |\hat{Q}_n(s'', a') - Q_n(s'', a')| \\ &= \gamma e_n \end{aligned}$$

بنابراین :

$$e_{n+1} \leq \gamma e_n , \quad e_{n+1} \geq 0 \Rightarrow$$

$$\frac{e_{n+1}}{e_n} \leq \gamma \Rightarrow \prod_{i=0}^n \frac{e_{i+1}}{e_i} \leq \prod_{i=0}^n \gamma \Rightarrow \frac{e_{n+1}}{e_0} \leq \gamma^{n+1}$$

$$\Rightarrow e_{n+1} \leq e_0 \gamma^{n+1}$$

$$\Rightarrow e_m \leq e_0 \gamma^m , \quad e_m \geq 0$$

$$0 \leq \lim_{m \rightarrow \infty} e_m \leq \lim_{m \rightarrow \infty} e_0 \gamma^m = 0 \Rightarrow$$

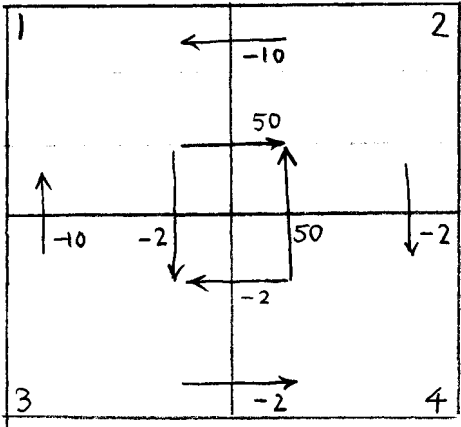
$$0 \leq \gamma < 1$$

$$\lim_{m \rightarrow \infty} e_m = 0$$

توجه می‌نمایند که هیچ فرضی بر روی دنباله‌ی action ها وجود ندارد، بنابراین Q-learning می‌تواند تابع Q را یاد بگیرد (و بنابراین سیاست بهینه را (optimal policy))، در حالی که در action های انتخاب شده به طور تصادفی آموزش می‌بیند. داده‌ی که دنباره یادگیری حاصل، هر  $(s, a)$  را به نهایت مرتبه مشاهده می‌کند.

(۷)

شکل زیر یک دنیای شبکه‌ای چهار حالتی را به تصویر کشیده است، که در آن حالت ۲، "طلا" را نشان می‌دهد. با استفاده از مقادیر پاداش بلافاصله (immediate rewards) نشان داده شده بر روی شکل و بکارگیری الگوریتم Q-learning، بر روی حالات به صورت پاداش متغرد حرکت کنید و جدول state-action را به‌نگام بنویسید.  $(\gamma = 0.9)$



$r(s,a)$

راه حل:

ابتدا همه درایه‌های جدول مقادیر Q را به هم‌زمان مقدار دهی می‌کنیم:

$$\forall s \forall a (s \in \{1, 2, 3, 4\} \wedge a \in \{\rightarrow, \leftarrow, \uparrow, \downarrow\} \Rightarrow Q(s, a) = 0)$$

حال فرمول زیر را تکرار می‌کنیم:

$$Q(s, a) \leftarrow r(s, a) + \gamma \max_{a'} \{Q(s', a')\}$$



دور اول:

$$Q(3, \rightarrow) = -2 + 0.9 \max \{Q(4, \leftarrow), Q(4, \uparrow)\} = -2$$

$$Q(4, \uparrow) = 50 + 0.9 \max \{Q(2, \downarrow), Q(2, \leftarrow)\} = 50$$

$$Q(2, \leftarrow) = -10 + 0.9 \max \{Q(1, \rightarrow), Q(1, \downarrow)\} = -10$$

$$Q(1, \downarrow) = -2 + 0.9 \max \{Q(3, \uparrow), Q(3, \rightarrow)\} = -2$$

$$Q(3, \rightarrow) = -2 + 0.9 \max \{Q(4, \leftarrow), \underbrace{Q(4, \uparrow)}_{50}\} = 43$$

		a			
		↑	↓	→	←
s	1	-	-2	0	-
	2	-	0	-	-10
	3	0	-	43	-
	4	50	-	-	0

دور دوم :

$$Q(4, \uparrow) = 50 + 0.9 \max \{ Q(2, \downarrow), Q(2, \leftarrow) \} = 50 + 0.9 \max \{ 0, -10 \} = 50$$

$$Q(2, \leftarrow) = -10 + 0.9 \max \{ Q(1, \rightarrow), Q(1, \downarrow) \} = -10 + 0.9 \max \{ 0, -2 \} = -10$$

$$Q(1, \downarrow) = -2 + 0.9 \max \{ Q(3, \uparrow), Q(3, \rightarrow) \} = -2 + 0.9 \max \{ 0, 43 \} = 36.7$$

$$Q(3, \rightarrow) = -2 + 0.9 \max \{ Q(4, \leftarrow), Q(4, \uparrow) \} = -2 + 0.9 \max \{ 0, 50 \} = 43$$

		$\alpha$			
		$\uparrow$	$\downarrow$	$\rightarrow$	$\leftarrow$
S	Q				
	1	-	36.7	0	-
	2	-	0	-	-10
	3	0	-	43	-
4	50	-	-	0	

دور سوم : تنها  $Q(2, \leftarrow)$  تغییر می کند :

$$Q(2, \leftarrow) = -10 + 0.9 \max \{ Q(1, \rightarrow), Q(1, \downarrow) \} = -10 + 0.9 \max \{ 0, 36.7 \} = 23.03$$

دور چهارم : تنها  $Q(3, \rightarrow)$  و  $Q(4, \uparrow)$  تغییر می کنند :

$$Q(4, \uparrow) = 50 + 0.9 \max \{ Q(2, \downarrow), Q(2, \leftarrow) \} = 50 + 0.9 \max \{ 0, 23.03 \} = 70.73$$

$$Q(3, \rightarrow) = -2 + 0.9 \max \{ Q(4, \leftarrow), Q(4, \uparrow) \} = -2 + 0.9 \max \{ 0, 70.73 \} = 61.66$$

		$\alpha$			
		$\uparrow$	$\downarrow$	$\rightarrow$	$\leftarrow$
S	Q				
	1	-	36.7	0	-
	2	-	0	-	23.03
	3	0	-	61.66	-
4	70.73	-	-	0	



RL و MDPs (۸)

یک روبات متحرک خودکار را در نظر بگیرید که می تواند FAST یا SLOW در هر گام زمان حرکت کند.

حرکت FAST به طور کلی پاداش +2 را می دهد، در حالی که

حرکت SLOW تنها پاداش +1 را می دهد.

به هر حال، روبات باید دمای داخلی خود را به حساب آورد، که می تواند HOT یا OK باشد.

حرکت کردن در حالت SLOW بهر به دمای کمتری می شود، در حالی که

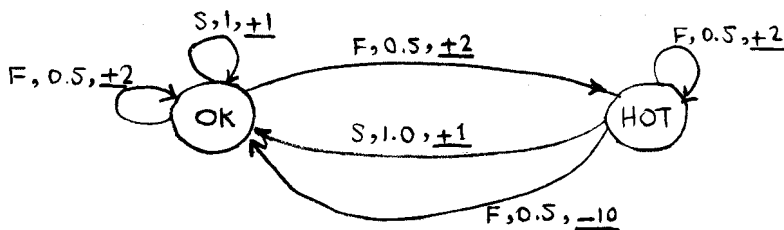
حرکت کردن در حالت FAST بهر به افزایش دما می شود.

اگر روبات HOT باشد، این خطر وجود دارد که -overheat شود که در این نقطه باید بایستد، دمایش را کم کند و

تعمیر شود. گذرهای MDP و پاداش ها به صورت زیر مشخص می شود:

S	a	S'	T(S, a, S')	R(S, a, S')	
OK	SLOW	OK	1.0	+1	SLOW S
OK	FAST	OK	0.5	+2	FAST F
OK	FAST	HOT	0.5	+2	
HOT	SLOW	OK	1.0	+1	
HOT	FAST	HOT	0.5	+2	
HOT	FAST	OK	0.5	-10	

توجه کنید که با اینکه تعمیر وقت کمی است، روبات پس از آن OK می شود (سطر آخر جدول)



(آ) دو گام از value iteration در جدول زیر را اجرا کنید. (با استفاده از discount  $\gamma = 0.8$ )

(خانه ها شور خورده را در نظر بگیرید)

S	$V_0$	$V_1$	$V_2$
OK	0	2	3.2
HOT	0	1	

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

$$V_1(OK) = \max\{1(1 + \gamma \times 0), 0.5(2 + \gamma \times 0) + 0.5(2 + \gamma \times 0)\} = \max\{1, 2\} = 2$$

$$V_1(HOT) = \max\{1(1 + \gamma \times 0), 0.5(2 + \gamma \times 0) + 0.5(-10 + \gamma \times 0)\} = \max\{1, -4\} = 1$$

$$V_2(OK) = \max\{1(1 + \gamma \times 2), 0.5(2 + \gamma \times 2) + 0.5(2 + \gamma \times 1)\} = \max\{2.6, 3.2\} = 3.2$$

ب) Q-learning را با تخفیف  $\gamma = 0.8$  اجرا کنید و از نرخ یادگیری  $\alpha = 0.5$ ، با استفاده از نمونه‌های گذر زیر کمک بگیرید. مقادیر Q ای که در یک گام تغییر نمی‌کنند را دوباره بنویسید.  
نمونه‌های زیر را با تجربیه‌های عامل زیر فرض کنید:

(OK, FAST, HOT), reward +2, calculate  $Q_1$

(HOT, FAST, OK), reward -10, calculate  $Q_2$

(OK, SLOW, OK), reward +1, calculate  $Q_3$

$$Q(s, a) \leftarrow Q(s, a) + 0.5 [R(s, a, s') + 0.8 \max_{a'} Q(s', a') - Q(s, a)]$$

$$Q_1(OK, FAST) = 0 + 0.5 [2 + 0.8 \max_{a'} Q(HOT, a') - 0] = 1$$

$$Q_2(HOT, FAST) = 0 + 0.5 [-10 + 0.8 \max_{a'} Q(OK, a') - 0] = -4.6$$

$$Q_3(OK, SLOW) = 0 + 0.5 [1 + 0.8 \max_{a'} Q(OK, a') - 0] = 0.9$$

S	a	$Q_0$	$Q_1$	$Q_2$	$Q_3$
OK	SLOW	0			0.9
OK	FAST	0	1.0		
HOT	SLOW	0			
HOT	FAST	0		-4.6	