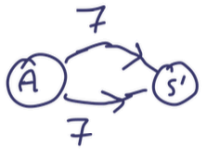


①

الف) غلط - می‌تواند آموختن هم از ریداردها قبلی در صورتی که محیط تغییر نکند استفاده کند ،
برای مثال با importance sampling از ریداردها یک پالیسی دیگری استفاده کند و خود را ترین کند.
و با از بین کلی الگوریتم‌های offline RL هست غیر به خط اند.



ب) غلط - ریداردها از اکشن مختلف می‌تواند یکبار به یکبار باشد →
که در این صورت هر دو اکشن پالیسی بهینه در این حالت اند.

ج) غلط - REINFORCE که نمونه کارهای پالیسی گزینش مسیحا "فرد پالیسی را
تغییر می‌دهد و استفاده ای از مدل مثل تدریس ندارد. پس model-free است.

د) صحیح - اثبات در note ها لکچر ۲ کدیس CS234 2019 صفحه ۱۶، ۱۷

و همچنین در Sutton & Barto موجود است
[اثبات از آنجا گرفته شده]
[در صورت سوال ذکر شده]
برای غلط ها دلیل بیاوریم به نایزدها برآید
روشن می‌شود [در]

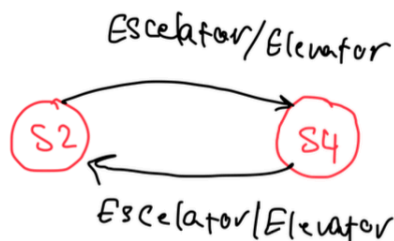
$$\pi(s) = \operatorname{argmax}_a Q(s, a)$$

(2)
(۲)

$$\pi(s_1) = \operatorname{argmax} [1, 5, -1, 5] = \text{Elevator}$$

$$\pi(s_2) = \operatorname{argmax} [+2, 1, -0, 3] = \text{Elevator}$$

$$\pi(s_3) = \operatorname{argmax} [+0, 8, +0, 9] = \text{Escalator}$$



(۳)

$$\pi(s_2) = \operatorname{argmax} [+0, 10, 8, +2, 10, 8] = \text{Escalator}$$

(۴)

$$\pi(s_4) = \operatorname{argmax} [+1, 1, 8, +1, 3, 8] = \text{Escalator}$$

(۵) روبرو حرکت S4 به S2 بسته به حالتی که بین S1 یا S3 بودیم فرق می‌کند

د به نوعی اکشن تبلی می‌شود S2 به S4 که با Escalator یا Elevator بوده رد روبرو حرکت S4 به S2 تأثیر داشته‌س MDP مفروض ماکروین نیست و در طول این تبلی information loss خواهد که در تبلی سیاست بهینه هم لزوماً برابر نیست. (که شده است)

3

(الف)

$$\mathbb{E}_{s \sim p(s), a \sim \pi_1(s, a)} R(s, a) = \sum R(s, a) p(s) p(a|s) = \sum R(s, a) p(s) \pi_1(s, a)$$

$$= \sum R(s, a) p(s) \pi_0(s, a) \frac{\pi_1(s, a)}{\pi_0(s, a)} = \sum p(s) \pi_0(s, a) \left(R(s, a) \frac{\pi_1(s, a)}{\pi_0(s, a)} \right)$$

$$\approx \mathbb{E}_{s \sim p(s), a \sim \pi_0(s, a)} R(s, a) \frac{\pi_1(s, a)}{\hat{\pi}_0(s, a)}$$

$\infty \rightarrow$ تعداد بسیار زیاد

نمونه ادرج
طبیعتاً متن اعداد بزرگی

$$\mathbb{E}_{s \sim p(s), a \sim \pi_1(s, a)} R(s, a)$$

با این روش ها به ∞ می رسد که باید آن

باید ثابت کنیم \mathbb{E} خارج نمی شود

$$\mathbb{E}_{s \sim p(s), a \sim \pi_0(s, a)} \frac{\pi_1(s, a)}{\hat{\pi}_0(s, a)} \approx \sum_{\# \text{ بسیار زیاد}} p(s) \pi_0(s, a) \frac{\pi_1(s, a)}{\pi_0(s, a)}$$

$$= \sum p(s) \pi_1(s, a) = \sum p(s) p(a|s) = \mathbb{E}_{s \sim p(s), a \sim \pi_1(s, a)} 1 = 1$$

\Rightarrow Q.E.D

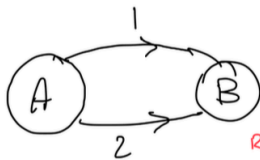
$E(C) = C$
ثابت است

(ع) اگرندۀ سفایک دیا جائے باقیم (S_1, a_1)

$$\frac{E_{S \sim p(S), a \sim \pi_0(S, a)} \frac{\pi_1(S, a)}{\pi_0(S, a)} R(S, a)}{E_{S \sim p(S), a \sim \pi_0(S, a)} \frac{\pi_1(S, a)}{\pi_0(S, a)}} = \frac{\frac{\pi_1(S_1, a_1)}{\pi_0(S_1, a_1)} R(S_1, a_1)}{\frac{\pi_1(S_1, a_1)}{\pi_0(S_1, a_1)}} = R(S_1, a_1)$$

$E(R(S_1, a_1)) = R(S_1, a_1)$

✓ براستیت $R(S, a)$ برقی π_0 سفای دار R بہت آسہ π_0 ،
ی کھد کہ لزوماً برا π_1 نیست۔ مثلاً



$R(B) = 0$ ✓
Stochastic
صرفاً A ہائے
 $R(A, 1) = 4$
 $R(A, 2) = 0$

$$\pi_0(1|A) = 0.8$$

$$\pi_0(2|A) = 0.2$$

$$\pi_1(1|A) = 0.5$$

$$\pi_1(2|A) = 0.5$$

نرخ گذشتن π_0

$(A, 1, 4)$

$$\Rightarrow R(S, a) = 0.5 \times 0 + 0.5 \times 4 = 2$$

π_1

الکھد استیصیر 4 نہ 2
 $2 \neq 4$

$\pi(s) = \text{مبلغ}$

(۶)

$$V = R + \gamma P V \Rightarrow V = (I - \gamma P)^{-1} R$$

Transition
Matrix
احتمالات
گذار مجدد
نشان داده شده است

$$\begin{bmatrix} V(M) \\ V(R) \\ V(D) \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} + \gamma \begin{bmatrix} 0 & 0.5 & 0.5 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} V(M) \\ V(R) \\ V(D) \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} V(M) \\ V(R) \\ V(D) \end{bmatrix} = \begin{bmatrix} 1 - 0.5\gamma & -0.5\gamma \\ 0 & 1 - \gamma \\ 0 & 0 & 1 - \gamma \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix}$$

ایجاد با (ماتریس) فراموش کردن یا صاف کردن (ماتریس) (ماتریس)

$$\text{ماتریس} \begin{bmatrix} 1 & -0.5\gamma & -0.5\gamma \\ 0 & 1 - \gamma & 0 \\ 0 & 0 & 1 - \gamma \end{bmatrix} \begin{bmatrix} 1 & \frac{\gamma}{2-2\gamma} & \frac{\gamma}{2-2\gamma} \\ 0 & \frac{1}{1-\gamma} & 0 \\ 0 & 0 & \frac{1}{1-\gamma} \end{bmatrix} = I$$

$$\Rightarrow \begin{bmatrix} V(M) \\ V(R) \\ V(D) \end{bmatrix} = \begin{bmatrix} 1 & \frac{\gamma}{2-2\gamma} & \frac{\gamma}{2-2\gamma} \\ 0 & \frac{1}{1-\gamma} & 0 \\ 0 & 0 & \frac{1}{1-\gamma} \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 + \frac{\gamma}{1-\gamma} - \frac{\gamma}{2-2\gamma} \\ \frac{2}{1-\gamma} \\ \frac{1}{1-\gamma} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{2-\gamma}{2-2\gamma} \\ \frac{2}{1-\gamma} \\ \frac{1}{1-\gamma} \end{bmatrix}$$

$$\left[\begin{array}{ccc|ccc} 1 & -0.5\gamma & -0.5\gamma & 1 & 0 & 0 \\ 0 & 1-\gamma & 0 & 0 & 1 & 0 \\ 0 & 0 & \gamma & 0 & 0 & 1 \end{array} \right]$$

$$\left[\begin{array}{ccc|ccc} 1 & 0 & -0.5\gamma & 1 & -\frac{0.5\gamma}{1-\gamma} & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & \gamma & 0 & 0 & 1 \end{array} \right]$$

$$\left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 1 & \frac{-0.5\gamma}{1-\gamma} & \frac{-0.5\gamma}{1-\gamma} \\ 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & \gamma & 0 & 0 & 1 \end{array} \right]$$

$$\left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 1 & \frac{-0.5\gamma}{1-\gamma} & \frac{-0.5\gamma}{1-\gamma} \\ 0 & 1 & 0 & 0 & \frac{1}{1-\gamma} & 0 \\ 0 & 0 & \gamma & 0 & 0 & 1 \end{array} \right]$$

$$\left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 1 & \frac{-0.5\gamma}{1-\gamma} & \frac{-0.5\gamma}{1-\gamma} \\ 0 & 1 & 0 & 0 & \frac{1}{1-\gamma} & 0 \\ 0 & 0 & 1 & 0 & 0 & \frac{1}{1-\gamma} \end{array} \right]$$

نظام التحكم في السياسة. كيف نحل policy improvement، الاجابة كالتالي.

$\gamma = 0.9$

$$\hat{\pi}(s) = \operatorname{argmax}_{a \in A} [R(s) + \gamma \sum_{s' \in S} P(s'|sa) V^{\pi}(s')]$$

$$V(M) = \frac{2 - 0.9}{2 - 2 \times 0.9} = \frac{1.1}{0.2} = 5.5$$

$$V(R) = \frac{2}{1 - 0.9} = 20$$

$$V(D) = \frac{-1}{1 - 0.9} = -10$$

$$\hat{\pi}(\text{Mountain}) = \operatorname{argmax}_{P_{i,w}} [1 + \gamma [0.5 V(R) + 0.5 V(D)], 1 + \gamma [0.1 V(M) + 0.2 V(D) + 0.7 V(R)]]$$

$$= \operatorname{argmax}_{P_{i,w}} [1 + \frac{0.9}{2} (20 - 10), 1 + 0.9 (.55 + (-2) + 14)]$$

$$= \operatorname{argmax}_{P_{i,w}} [5.5, 12.25] = \text{War}$$

$$\hat{\pi}(\text{Riverside}) = \operatorname{argmax}_{P_{i,w}} [2 + \gamma (1 \cdot V(R)), 2 + \gamma (0.2 V(R) + 0.8 V(M))]$$

$$= \operatorname{argmax}_{P_{i,w}} [2 + 0.9 (20), 2 + 0.9 (0.2 \times 20 + 0.8 \times 5.5)]$$

$$= \operatorname{argmax}_{P_{i,w}} [20, 9.56] = \text{Peace}$$

$$\hat{\pi}(\text{Desert}) = \operatorname{argmax}_{P_{i,w}} [-1 + \gamma (1 \cdot V(D)), -1 + \gamma (1 \cdot V(M))] =$$

$$\operatorname{argmax}_{P_{i,w}} [-1 + 0.9 \times (-10), -1 + 0.9 (5.5)] = \operatorname{argmax}_{P_{i,w}} [-10, 3.25] = \text{War}$$

(C)

M-P	M-W	R-P	R-W	D-P	D-W
0	0	0	0	0	0
0	0	0	0	0	-1
0	0	0	0	0	1
0	0	0.5	0	0	1
0.75	0	0.5	0	0	1

$$1: Q(D, W) = Q(D, W) + \alpha (R + \lambda \max_a Q(D, a) - Q(D, W))$$

$$Q(D, W) = 0 + 0.5 (-2 + 1 \max(0, 0)) = 0.5 \times -2 = -1$$

$$2: Q(D, W) = Q(D, W) + \alpha (R + \lambda \max_a Q(R, a) - Q(D, W))$$

$$Q(D, W) = -1 + 0.5 (3 + 1 \max(0, 0) - (-1)) =$$

$$-1 + 0.5 (3 + 1) = -1 + 2 = 1$$

$$3: Q(R, P) = Q(R, P) + \alpha (R + \lambda \max_a Q(M, a) - Q(R, P))$$

$$Q(R, P) = 0 + 0.5 (1 + 1 \max(0, 0) - 0) = 0.5$$

$$4: Q(M, P) = Q(M, P) + \alpha (R + \lambda \max_a Q(R, a) - Q(M, P))$$

$$Q(M, P) = 0 + 0.5 (1 + 1 \max(0.5, 0) - 0) =$$

$$0.5 (1 + 0.5) = 0.5 (1.5) = 0.75$$

Mountain-peace	Riverside-peace	Desert-war
0	0	0
0	0	-1
0	0	1
0	0.5	1
0.75	0.5	1

