

Contraction mapping $\rightarrow \|U(v_1) - U(v_2)\| \leq q \|v_1 - v_2\|$ (مقبوض المape)

$$\begin{aligned} \|U(v_1) - U(v_2)\| &= \|R + \gamma P v_1 - (R + \gamma P v_2)\| = \|\gamma P(v_1 - v_2)\| \\ &= \gamma \|P(v_1 - v_2)\| \leq \gamma \|P\| \|v_1 - v_2\| < q \|v_1 - v_2\| \text{ where } q = \gamma \|P\| \\ &\text{و } U(v) \text{ هي نقطة انقباضات.} \end{aligned}$$

$$\begin{aligned} U(v_1(s)) - U(v_2(s)) &= \sum_{s'} P(s'|s, a) \{R(s) + \gamma v_1(s')\} - \sum_{s'} P(s'|s, a) \{R(s) + \gamma v_2(s')\} \\ &= \sum_{s'} P(s'|s, a) \gamma (v_1(s') - v_2(s')) \leq \gamma \max_{s'} \{v_1(s') - v_2(s')\} \sum_{s'} P(s'|s, a) \\ &\quad \text{Contraction Mapping } U(v) \text{ في كائين } \|v_1 - v_2\|_{\infty} \end{aligned}$$

~~دالة~~ $\rightarrow U(v^k) = v^k$ (نقطة ثابتة)

$$d(U(x), U(y)) \leq q d(x, y) \quad \forall x, y \in X$$

$$\hookrightarrow d(U(v^k), U(v^{\infty})) \leq q d(v^k, v^{\infty})$$

$$d(v^k, v^{\infty}) = d(U(v^{k-1}), U(v^{\infty})) \leq q d(v^{k-1}, v^{\infty}) = q d(U(v^{k-2}), U(v^{\infty}))$$

$$\leq q^2 d(v^{k-2}, v^{\infty}) = q^2 d(U(v^{k-3}), U(v^{\infty}))$$

$$\Rightarrow d(v^k, v^{\infty}) \leq q^k d(v^0, v^{\infty})$$

$$\hookrightarrow \lim_{k \rightarrow \infty} d(v^k, v^{\infty}) \leq \lim_{k \rightarrow \infty} q^k d(v^0, v^{\infty}) = 0$$

$$\lim_{k \rightarrow \infty} d(r^k, r^\infty) \leq \lim_{k \rightarrow \infty} \gamma^k d(r^0, r^\infty) = 0$$

$$\lim_{n \rightarrow \infty} U^n(r) = r^\infty$$

نقطه ثابت

$$\|U(r_1) - U(r_2)\|_\infty \leq \gamma \|r_1 - r_2\|_\infty$$

(-)

نرمالیزاسیون norm

$$\|x - y\|_\infty \leq \|x - z\|_\infty + \|z - y\|_\infty \quad (I)$$

$$\|r^\infty - U^k(r)\|_\infty \leq \|U^{k+1}(r) - U^k(r)\|_\infty$$

داده

$$\begin{aligned} \|U^{(k+1)}(r) - U^k(r)\|_\infty &\leq \gamma \|U^k(r) - U^{k-1}(r)\|_\infty < \gamma \epsilon \\ \|U^{(k+2)}(r) - U^{k+1}(r)\|_\infty &\leq \gamma \|U^{k+1}(r) - U^k(r)\|_\infty < \gamma^2 \epsilon \\ &\vdots \end{aligned} \Rightarrow$$

$$\|U^{k+n}(r) - U^{k+n-1}(r)\|_\infty \leq \gamma^n \epsilon$$

$$\Rightarrow \sum_{i=0}^n \|U^{k+i}(r) - U^{k+i-1}(r)\|_\infty \leq \sum_{i=0}^n \gamma^i \epsilon = \frac{\gamma^{n+1} \epsilon}{1-\gamma}$$

مجموع جیب

$$\|r^\infty - U^k(r)\|_\infty \leq \lim_{n \rightarrow \infty} \sum_{i=1}^n \|U^{k+i}(r) - U^{k+i-1}(r)\|_\infty \leq \lim_{n \rightarrow \infty} \sum_{i=1}^n \gamma^i \epsilon$$

$$\Rightarrow \|r^\infty - U^k(r)\|_\infty \leq \frac{\gamma \epsilon}{1-\gamma} < \frac{\epsilon}{1-\gamma}$$

s.a.m

Every visit Monte-Carlo:

مسئله (2) (1)

Sample Episode:

Define $G_{i,t} = r_{i,t} + \gamma r_{i,t+1} + \dots + \gamma^{T_i-t} r_{i,T_i}$

For each step t until T_i :

state s is the state visited at step t

$$N(s) = N(s) + 1$$

$$G(s) = G(s) + G_{i,t}$$

$$V^{\pi}(s) = G(s) / N(s)$$

در ابتدا $R = -10$

و در انت $R = 10$

در بقیه ایست ها هیچ پواردی ندارد.

$$V_{(1)} = -10 \quad V_{(2)} = -10 \quad V_{(3)} = -10 \quad V_{(8)} = -10$$

$$V_{(7)} = -10 \quad V_{(12)} = -10 \quad V_{(16)} = -10 \quad V_{(20)} = -10$$

$$\cancel{V_{(21)} = -10} \quad V_{(21)} = \frac{-10+0}{2} = -5 \quad V_{(22)} = \frac{0+10}{2} = 5$$

$$V_{(17)} = 10 \quad V_{(18)} = 10 \quad V_{(23)} = 10$$

و Value مربوط به بقیه ایست ها برابر با 0 باقی می ماند.

* برای هر دو قسمت سوال فرض گرفته شده است که $\gamma = 1$ باشد.

First-Visit Monte Carlo:

(2)

Sample Episode

$$G_{i:t} = r_{i,t} + \gamma r_{i,t+1} + \dots + \gamma^{T_i-t} r_{i,T_i}$$

For each time step t , until T_i :If this is the first visit of state s in this episode

$$N(s) = N(s) + 1$$

$$G(s) = G(s) + G_{i:t}$$

$$V(s) = G(s) / N(s)$$

$$V_{(1)} = -10 \quad V_{(2)} = -10 \quad V_{(3)} = -10 \quad V_{(8)} = -10$$

$$V_{(7)} = -10 \quad V_{(12)} = -10 \quad V_{(16)} = -10 \quad V_{(20)} = -10$$

$$V_{(21)} = -10 \quad V_{(22)} = 0 \quad V_{(27)} = 10 \quad V_{(15)} = 10$$

$$V_{(23)} = 10$$

into 0 to 10 value

سوال 3 -

$$V^{\pi}(s) = E^{\pi} \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \mid S_0 = s \right\} = \sum_{t=0}^{\infty} \gamma^t E \{ r_t \mid S_0 = s \} \quad (1)$$

$$= \sum_{a \in A} P_{ts} \pi_{(s,a)} \left\{ R_{(s,a)} + \sum_{t=1}^{\infty} \gamma^t E \{ r_t \mid S_0 = s \} \right\}$$

$$= \sum_{a \in A} \pi_{(s,a)} \left\{ R_{(s,a)} + \sum_{s' \in S} P_{(s,a,s')} \sum_{a' \in A} \pi_{(s',a')} \left\{ \gamma R_{(s',a')} + \sum_{t=2}^{\infty} \gamma^t E \{ r_t \mid S_0 = s \} \right\} \right\}$$

$$= \left(\sum_{a \in A} \pi_{(s,a)} \left\{ R_{(s,a)} + \sum_{s' \in S} P_{(s,a,s')} \sum_{a' \in A} \pi_{(s',a')} \right\} \right)$$

$$= \sum_{a \in A} \pi_{(s,a)} R_{(s,a)} + \sum_{a \in A} \pi_{(s,a)} \sum_{s' \in S} P_{(s,a,s')} \sum_{a' \in A} \gamma \pi_{(s',a')} R_{(s',a')}$$

$$+ \sum_{a \in A} \pi_{(s,a)} \sum_{s' \in S} P_{(s,a,s')} \sum_{a' \in A} \pi_{(s',a')} \sum_{s'' \in S} P_{(s',a',s'')} \sum_{a'' \in A} \gamma^2 \pi_{(s'',a'')} R_{(s'',a'')} + \dots$$

با توجه به اینکه M_0, M هر دو R, P یکسانی دارند، می توان

$$V^{\pi}(s) = V_0^{\pi}(s) \quad \text{گفت}$$

(۱) این عبارت در ت است.

$$v^{\pi}(s) = E^{\pi} \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \mid S_0 = s \right\} \Rightarrow \pi^* = \operatorname{argmax}_{\pi} \{ v^{\pi}(s) \}$$

$$r_t \rightarrow cr_t \Rightarrow v'^{\pi}(s) = E^{\pi} \left\{ \sum_{t=0}^{\infty} \gamma^t cr_t \mid S_0 = s \right\} \\ = c E^{\pi} \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \mid S_0 = s \right\} = c v^{\pi}(s)$$

$$\pi'^* = \operatorname{argmax}_{\pi} \{ c v^{\pi}(s) \} = \operatorname{argmax}_{\pi} \{ v^{\pi}(s) \} = \pi^*$$

(۲) ثابت

$$v^{\pi}(s) = E^{\pi} \left\{ \sum_{t=0}^T \gamma^t r_t \mid S_0 = s \right\} \Rightarrow \pi^* = \operatorname{argmax}_{\pi} \{ v^{\pi}(s) \}$$

$$r_t \rightarrow cr_t + c \Rightarrow v'^{\pi}(s) = E^{\pi} \left\{ \sum_{t=0}^T \gamma^t (r_t + c) \mid S_0 = s \right\}$$

$$= E^{\pi} \left\{ \sum_{t=0}^T \gamma^t r_t \mid S_0 = s \right\} + \sum_{t=0}^T \gamma^t c$$

با توجه به عبارت به دست آمده $v^{\pi}(s)$ مقدار $\sum_{t=0}^T \gamma^t c$ اضافه می شود.

که با توجه به اینکه طول هر trajectory به اندازه T است در این مقادیر با هم

تفاوت دارد، باعث می شود که argmax مقدار دیگری را به عنوان π

ببیند.

$$v^{\pi}(s) = E^{\pi} \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \mid S_0 = s \right\} \Rightarrow \pi^* = \underset{\pi}{\operatorname{argmax}} \{ v^{\pi} \} \quad (1)$$

$$r_t \rightarrow r_t + c$$

درست است.

$$\begin{aligned} v'^{\pi}(s) &= E^{\pi} \left\{ \sum_{t=0}^{\infty} \gamma^t \{ r_t + c \} \mid S_0 = s \right\} \\ &= E^{\pi} \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \mid S_0 = s \right\} + \sum_{t=0}^{\infty} \gamma^t c = v^{\pi}(s) + \frac{c}{1-\gamma} \\ &\Rightarrow \pi'^* = \underset{\pi}{\operatorname{argmax}} \{ v'^{\pi} \} = \underset{\pi}{\operatorname{argmax}} \left\{ v^{\pi} + \frac{c}{1-\gamma} \right\} \\ &= \underset{\pi}{\operatorname{argmax}} \{ v^{\pi} \} = \pi^* \Rightarrow \pi'^* = \pi^* \end{aligned}$$

(2) این عبارت درست است.

$$a = \underset{\pi}{\operatorname{argmax}} \{ v^{\pi}(s) \}, \pi^* = \underset{\pi}{\operatorname{argmax}} \{ v^{\pi} \} \quad \text{در نظر بگیرید که}$$

$$\cdot \text{مثلاً } v^{\pi}(s) = E^{\pi} \left\{ \sum_{t=0}^{\infty} \gamma^t r_t \mid S_0 = s \right\} \quad \text{حل}$$

~~این action~~ ~~در MDP~~ ~~است~~

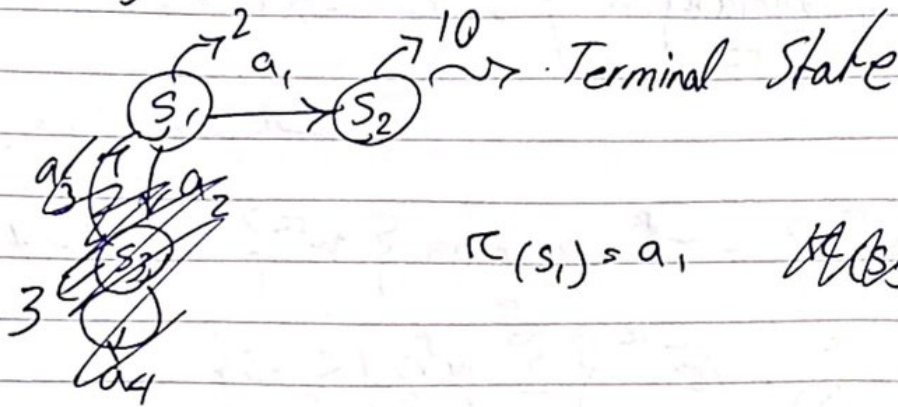
$$\begin{aligned} v^{\pi}(s) &= \sum_{a \in A} \pi(s, a) R(s, a) + \sum_{a \in A} \pi(s, a) \sum_{s' \in S} P(s, a, s') \sum_{a' \in A} \gamma \pi(s', a') R(s', a') + \\ &+ \sum_{a \in A} \pi(s, a) \sum_{s' \in S} P(s, a, s') \sum_{a' \in A} \gamma \pi(s', a') \sum_{s'' \in S} P(s', a', s'') \sum_{a'' \in A} \gamma^2 \pi(s'', a'') R(s'', a'') + \\ &\dots \end{aligned}$$

سوال 4-
الف) نامت ات.

در این معادله V^{π} ، π به این صورت محاسبه می‌شود که از چه انت‌هایی می‌توان به انت s رسید، V^{π} ، π به این اساس محاسبه می‌شود.

معادله بلین اصل مکمل معادله‌ای است که در پایین داریم:

$$V^{\pi}(s) = \sum_{s'} \sum_a \pi(s, a) P(s', a, s) \{ R(s, a) + \gamma V^{\pi}(s') \}$$



$$\pi(s_1) = a_1$$

$$\gamma = 1$$

معادله بلین عادی: $V^{\pi}(s_2) = 10$ $V^{\pi}(s_1) = 12$

معادله داده شده: $V^{\pi}(s_1) = 2$ $V^{\pi}(s_2) = -8$

که به این است اگر مقدار $Value$ نمی‌تواند منفی باشد.

همچنین به این مسئله می‌توان اشاره کرد که اگر γ نزدیک به 1 باشد، این

معادله بلین، مقدار زیرین برای $Value$ به دست می‌آورد.

$$\begin{aligned}
 r^{\pi}(s) &= \mathbb{E}\{G_t | S_t = s, \pi\} + \mathbb{E}\left\{\sum_{k=0}^{\infty} \gamma^k R_{t+k} | S_t = s, \pi\right\} \\
 &= \sum_{k=0}^{\infty} \gamma^k \mathbb{E}\{R_{t+k} | S_t = s, \pi\} = \sum_{a \in A} \pi(s, a) \{R(s, a) + \sum_{k=1}^{\infty} \gamma^k \mathbb{E}\{R_{t+k} | S_t = s, \pi\}\} \\
 &= \sum_{a \in A} \pi(s, a) \{R(s, a) + \sum_{s' \in S} P(s, a, s') \sum_{a' \in A} \pi(s', a') \{ \gamma R(s', a') + \sum_{k=2}^{\infty} \gamma^k \mathbb{E}\{R_{t+k} | S_t = s, \pi\} \} \} \\
 &= \sum_{a \in A} \pi(s, a) R(s, a) + \sum_a \pi(s, a) \sum_{s'} P(s, a, s') \sum_{a'} \gamma \pi(s', a') R(s', a') \\
 &\quad + \sum_a \pi(s, a) \sum_{s'} P(s, a, s') \sum_{a'} \pi(s', a') \sum_{s'' \in S} P(s', a', s'') \sum_{a'' \in A} \gamma^2 \pi(s'', a'') R(s'', a'') + \dots
 \end{aligned}$$

$$\begin{aligned}
 r^{\pi}(s) &= \mathbb{E}\{G_0 | S_0 = s, \pi\} + \mathbb{E}\left\{\sum_{t=0}^{\infty} \gamma^t R_t | S_0 = s, \pi\right\} \\
 &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}\{R_t | S_0 = s, \pi\} = \sum_a \pi(s, a) \{R(s, a) + \sum_{t=1}^{\infty} \gamma^t \mathbb{E}\{R_t | S_0 = s, \pi\}\} \\
 &= \sum_a \pi(s, a) \{R(s, a) + \sum_{s'} P(s, a, s') \sum_{a'} \pi(s', a') \{ \gamma R(s', a') + \sum_{t=2}^{\infty} \gamma^t \mathbb{E}\{R_t | S_0 = s, \pi\} \} \} \\
 &= \sum_a \pi(s, a) R(s, a) + \sum_a \pi(s, a) \sum_{s'} P(s, a, s') \sum_{a'} \gamma \pi(s', a') R(s', a') \\
 &\quad + \sum_a \pi(s, a) \sum_{s'} P(s, a, s') \sum_{a'} \pi(s', a') \sum_{s'' \in S} P(s', a', s'') \sum_{a'' \in A} \gamma^2 \pi(s'', a'') R(s'', a'') + \dots
 \end{aligned}$$

$$\mathbb{E}\{G_t | S_t = s, \pi\} \neq \mathbb{E}\{G | S_0 = s, \pi\}$$

$$v^{\pi}(s) = \mathbb{E} \{ G_{\infty} | S_0 = s, \pi \} = \mathbb{E} \left\{ \sum_{t=0}^{\infty} R_t | S_0 = s, \pi \right\} \quad (2)$$

$$\hookrightarrow v^{\pi}(s_{L-1}) = \mathbb{E} \{ R_{L-1} \}$$

$$\begin{aligned} v^{\pi}(s_{L-2}) &= \mathbb{E} \{ R_{L-2} + R_{L-1} \} = \mathbb{E} \{ R_{L-2} \} + \mathbb{E} \{ R_{L-1} \} \\ &= \mathbb{E} \{ R_{L-2} \} + v^{\pi}(s_{L-1}) \end{aligned}$$

✓

,

$$\begin{aligned} v^{\pi}(s_n) &= \mathbb{E} \{ R_n + R_{n+1} + \dots + R_{L-1} \} = \mathbb{E} \{ R_n \} + \mathbb{E} \{ R_{n+1} + \dots + R_{L-1} \} \\ &= \mathbb{E} \{ R_n \} + v^{\pi}(s_{n+1}) \end{aligned}$$

نتیجه به ایند هر R_n یک عدد منفی است، $\mathbb{E} \{ R_n \}$ نیز منفی است در نتیجه

$$v^{\pi}(s_n) \{ v^{\pi}(s_{n+1}) \Rightarrow$$

$$\Rightarrow v^{\pi}(s_0) \{ v^{\pi}(s_1) \{ \dots \{ v^{\pi}(s_{L-1})$$

$$S_0: C \quad S_1: D \quad S_2: E$$

n-Step Returns: $G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n}$ سوال 5
الف

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t^{(n)} - V(S_t)) + \gamma^n V(S_{t+n})$$

$n=1 \Rightarrow$ Temporal-Difference:

$$G_0^{(1)} = 0 + 0,5 \rightarrow V(S_0) = V(S_0) + \alpha (0,5 - 0,5) = 0,5$$

$$G_1^{(1)} = 0 + 0,5 \rightarrow V(S_1) = V(S_1) + \alpha (0,5 - 0,5) = 0,5$$

$$G_2^{(1)} = 1 \rightarrow V(S_2) = 0,5 + \alpha (0,5 + 1) = 0,55$$

Value در هر یک به E آید می شود.

$$n=2 \rightsquigarrow G_0^{(2)} = 0 + 0 + 0,5 \rightarrow V(S_0) = 0,5$$

$$G_1^{(2)} = 0 + 0 + 1 + 0,5 \rightarrow V(S_1) = 0,55$$

$$G_2^{(2)} = 1 \rightsquigarrow V(S_2) = 0,55$$

D, E در آید می شود

برای $n \geq 2$ با توجه به اینکه sample ما تنها 3 آید دارد.

در نتیجه در هر 3 حالت هم 3 آید D, E آید می شود.

(ب) مقدار α در واقع به bias-variance Trade-off اشاره دارد. زمانی که α کم باشد نمونه‌های اولیه اهمیت بیشتری دارند و داده‌های جدید اهمیت کمتری دارند و زمانی که α زیاد باشد، نیز نمونه‌های جدید خیلی اهمیت دارند.

(ج)

(آ) با افزایش تعداد محاسبات که تکرارها و اپیزودها را بیشتر کنیم که نمونه‌های بیشتری دانسته باشیم. در غیر این صورت خطای افزایش پیدا خواهد کرد.

(ب) افزایش تعداد اپیزودها و تعداد تکرارها باعث می‌شوند که نمونه‌های بیشتری دانسته باشیم که همین امر منجر به کاهش ارور می‌شود.

$$E_t(s) = \gamma \lambda E_{t-1}(s) + \lambda (S_t = s) \quad (د)$$

$$E_t(s) = 1 + \gamma \lambda E_{t-1}(s) = 1 + \gamma \lambda + (\gamma \lambda)^2 E_{t-2}(s) + \dots$$

$$\Rightarrow E_t(s) = \sum_{n=0}^{t-1} (\gamma \lambda)^n$$

$$\lim_{t \rightarrow \infty} E_t(s) = \lim_{t \rightarrow \infty} \sum_{n=0}^{t-1} (\gamma \lambda)^n = \frac{1}{1 - \gamma \lambda} = \frac{1}{1 - 0,2} = \frac{1}{0,8} = 1,25$$