

(الف) درست؛ چرا ضرب تخفیف کوچک موجب کوچک شدن  $rewards$  می شود که این مورد در هر

قدم از بازی موجب رفتار بینه یابی محلی (greedy) می شود

(ب) درست؛ پاداش منفی زندگی موجب در نظر گرفتن جریمه برای انجام تعداد حرکت بیشتر می شود که بزرگی آن معادل این است که  $reward$  در جرمای مراحل بعدی در تصمیم گیری بی اثر تر شده و حرکتی که موقعیت فعلی را به پاداش می رساند انتخاب شود

(ج) نادرست؛ تاثیر این دو پارامتر به شکل متفاوتی اعمال می شود  $\gamma$  ضرب تخفیف منفی  $\gamma$  رفتار ضریبی روی پاداش خواهد داشت اما پارامتر منفی  $\gamma$  تاثیر جمعی روی پاداش خواهد داشت

(د) نادرست؛ طبق توضیحات بخش قبلی رفتار این دو پارامتر به دو شکل متفاوت اعمال می شود

2 (الف) فرض می کنیم که بازی به شکل  $episodic$  باشد که در هر مرحله مقدار پاداش دریافتی برابر  $R_t$  باشد. برای پاداش کلی داریم

$$R_{total} = R_t + R_{t+1} + R_{t+2} + \dots \xrightarrow[\text{factor}]{\text{discount}} R_t + R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

پس مقدار ارزش هر  $state$  با توجه به  $policy$ ،  $\pi$  به شکل  $E[R_t]$  تعریف می شود

$$V^{\pi}(s) = E_{\pi} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$$

طبق این رابطه مقدار پاداش تاجایی که بازی ادامه داشته باشد، روی مقدار پاداش هر مرحله به شرط آنکه  $s$  شروع کننده بازی باشد، امید ریاضی گرفته می شود و محاسبه می شود در صورت تغییر وضعیت بین استیج مقدار

$$V(s) = E_{\pi} [R_{t+1} + \gamma V(s_{t+1}) | S_t = s]$$

$$\Rightarrow V^{\pi}(s) = \sum_{\alpha \in A} \pi(s, \alpha) (\tilde{R}(s, \alpha) + \gamma \sum_{s' \in S} T(s, \alpha, s') V^{\pi}(s'))$$

که در این رابطه  $\tilde{R}(s, \alpha) = \sum_{s' \in S} T(s, \alpha, s') R(s, \alpha, s')$  و  $\pi(s, \alpha)$  احتمال انتخاب حرکت  $\alpha$

در استیت S است (چرا که به دلیل وجود حالت C، استراتژی deterministic نخواهیم داشت) و اثری با mix خواهد بود

(ج) در policy evaluation داریم:

$$V_{k+1}^{\pi}(S) = \sum_{S' \in S} T(S, \pi(S), S') [R(S, \pi(S), S') + \gamma V_k^{\pi}(S')]$$

که با توجه به انتخاب استراتژی mix باید با احتمال هر S وزن دهی شود  $(E[V_{k+1}^{\pi}(S)])$

$$V_2^{\pi}(A) = -8 + 0.5 \times V_1(B) = -7 \quad V_2^{\pi}(B) = 0.5(2 + 0.5(V_1(A))) + 0.5$$

$$V_2^{\pi}(C) = 0.5(8 + 0.5 V_1(B)) + 0.5(4 + 0.5) \rightarrow x(-2 + 0.5 V_1(C)) = 1$$

$$\rightarrow x \left( \frac{1}{4} V_1(A) + \frac{3}{4} V_1(C) \right) = 7$$

یک انتخاب اولیه

$$\pi_2(A) = ab$$

(د) برای هر State باید مقداری که Q را بیشینه می کند را بدست آوریم

$$Q_2(B, ba) = -2 + 0.5 V_2(A) = -5.5 \quad \& \quad Q_2(B, bc) = -2 + 0.5 V_2(C) = 1.5$$

$$\Rightarrow \pi_2(B) = bc$$

$$Q_2(C, ca) = 4 + 0.5 \left( \frac{1}{4} V_2(A) + \frac{3}{4} V_2(B) \right) = -5.75 \quad \& \quad Q_2(C, cb) = 8 + 0.5 V_2(B) = 8.5$$

$$= 8.5 \Rightarrow \pi_2(C) = cb$$

(ه) اگر  $\pi'$ ، greedy رفتار کند یعنی  $\pi'(S) = \arg \max_{\alpha \in A} Q(S, \alpha)$

$$Q^{\pi}(S, \pi(S)) = \max_{\alpha \in A} Q^{\pi}(S, \alpha) \geq Q^{\pi}(S, \pi'(S)) = V^{\pi'}(S) \Rightarrow \text{increase reward}$$

$$V^{\pi}(S) \leq Q^{\pi}(S, \pi'(S)) = E_{\pi} [R_{t+1} + \gamma V^{\pi}(S_{t+1}) | S_t = S] \leq E_{\pi'} [R_{t+1} + \gamma Q^{\pi}(S_{t+1}, \pi'(S_{t+1})) | S_t = S]$$

$$\rightarrow (S_{t+1}, \pi'(S_{t+1})) | S_t = S =$$

$$\leq \dots \leq E_{\pi} [R_{t+1} + \gamma R_{t+2} + \dots | S_t = S] \cdot V^{\pi}(S)$$

اگر عبارت فوق  $Q^{\pi}(S, \pi(S)) = V^{\pi}(S) \leq Q^{\pi}(S, \pi'(S)) = \max_{\alpha \in A} Q^{\pi}(S, \alpha)$  شود

(توقف)؛  $V^{\pi}(S)$  شرایط felman را ارضا می کند و  $\pi'$ ، هر دو optimal خواهند بود



$$Q_{k+1}(s, \alpha) = \sum_{s'} T(s, \alpha, s') [R(s, \alpha, s') + \gamma \max_{\beta} (Q_k(s', \beta))]$$

$$\Rightarrow Q(3, left) = \sum_{s'=\{2\}} T(3, left, s') [R(3, left, s') + 0.9 \max_{\beta=\{U, R, L\}} (Q(s', \beta))]$$

$$\Rightarrow Q(3, left) = -1 + 0.9 \max(8, 3, 6) = 6.2$$

(ب)

در استراتژی حریصانه  $\epsilon$  در هر  $state$  بین حالت های شناخته شده آنکه بیشترین پاداش را می دهد انتخاب خواهد شد اما در اینجا که شناخت از محیط وجود ندارد، ممکن است  $action$  ای باشد  $state$  معقد

- آن  $reward$  آن ناشناخته (اما بسیار زیاد تر از بقیه پاداش ها باشد) باشد و در الگوریتم  $greedy$  انتخاب نشود که باعث ناهوشمندی این الگوریتم در محیط ناشناخته خواهد بود

برای برقراری تعادل  $exploration$  باید بین  $exploration$  و  $exploitation$  تعادل برقرار کرد و گاهی با عمل غیر حریصانه باید  $exploration$  را افزایش داد تا به سیاستی (در احوی) با  $reward$  بیشینه  $global$  نزدیک شود.

چرا که استفاده از  $Q$ -value بدون نیاز به مدل محیط می شود و می توانیم با  $\pi(s) = \arg \max_{\alpha} (Q(s, \alpha))$

(ج)

- حرکت بجهت را پیدا کنیم اما اگر از  $V$ -value استفاده کنیم محاسبات پیچیده تر خواهد بود و نمی توانیم سیاست بجهت را نیز بدست آوریم

(د)

با این کار همه  $action$  ها احتمال انتخاب می گیرند که باعث می شود که حرکتی احتمالی انتخاب داشته باشد که معادل برآورد کردن  $exploration$  است که در حالت  $greedy$  آن را نداریم.

$$S=1: \pi(1, up) = \frac{e^4}{e^3+e^4} \text{ و } \pi(1, right) = \frac{e^3}{e^3+e^4}$$

$$S=2: \pi(2, right) = \frac{e^8}{e^8+e^3+e^6} \text{ و } \pi(2, left) = \frac{e^3}{e^3+e^8+e^6} \text{ و } \pi(2, up) = \frac{e^6}{e^3+e^8+e^6}$$

$$S=3: \pi(3, left) = \frac{e^7}{e^7+e^9} \text{ و } \pi(3, up) = \frac{e^9}{e^7+e^9}$$

$$S=4: \pi(4, right) = \frac{e^5}{e^5+e^2} \text{ و } \pi(4, down) = \frac{e^2}{e^5+e^2}$$

$$Q(S=5) \pi(S, right) = \frac{e^8}{e^8 + e^5 + e^6} \quad \& \quad \pi(S, left) = \frac{e^5}{e^8 + e^5 + e^6} \quad \& \quad \pi(S, down) = \frac{e^6}{e^8 + e^5 + e^6}$$

$$Q(2, up) \leftarrow Q(2, up) + 0.2 [-1 + 0.8 Q(5, right) - Q(2, up)] = 4.8$$

$$Q(5, right) \leftarrow Q(5, right) + 0.2 [10 + 0.8 Q(6) - Q(5, right)] = 8.4$$

Q

(24) (II) ابتدا مقدار  $Q$  را 0 در نظر می گیریم و با هر غمزه آن را بروز می کنیم:  $\alpha = 0.1$  &  $\gamma = 0.9$  فرض

$$Q(A, 1) \leftarrow Q(A, 1) + 0.1 [-3 + 0.9 \max_{\beta} Q(B, \beta) - Q(A, 1)] = 0.3$$

$$Q(B, 1) \leftarrow Q(B, 1) + 0.1 [4 + 0.9 (\max_{\beta} Q(A, \beta)) - Q(B, 1)] = 0.427$$

$$Q(A, 2) \leftarrow Q(A, 2) + 0.1 [-4 + 0.9 (\max_{\beta} Q(A, \beta)) - Q(A, 2)] = -0.373$$

$$Q(A, 1) \leftarrow Q(A, 1) + 0.1 [-3 + 0.9 (\max_{\beta} Q(B, \beta) - Q(A, 1))] = 0.00843$$

$$Q(A, 2) \leftarrow Q(A, 2) + 0.1 [0.9 + 0.9 (\max_{\beta} Q(A, \beta) - Q(A, 2))] = -0.23$$

(ب) این سیاست می تواند آن باشد که  $S=A$ ،  $a=2$  انتخاب شود چرا که  $reward = \frac{4+1}{2} = -1.5$  اما در حالت  $a=1$ ،  $reward = -3$  و در  $S=B$  نیز  $a=1$  انتخاب کنیم تا از انتخاب رندم  $R$  بیشتری دریافت کنیم

(ج) از آنجایی که  $Q$ -learning، ~~on-policy~~ <sup>off-policy</sup> است، انتظار داریم با استفاده از حرکت به مقدار یکسان  $value$  برسیم چرا که  $update$  مستقل  $policy$  است. سیاست  $random$  حرکت در  $exploitation$  نخواهد شد و محیط را  $explore$  خواهد کرد بدون توجه به مقدار  $reward$ . در مقابل  $policy$  سیاست  $explore$  (حرکت حیرانانه)،  $exploit$  می کند و  $state$  های ناشناخته (بدون تجربه) را  $explore$  می کند و  $reward$  بیشتری را تجربه کنیم که ممکن است بایادگیری آن ها

5 اگر بخواهیم یک الگوریتم برای کاهش این نرخ استفاده کنیم، می‌توانیم با ضرب کردن در انجام

action،  $\epsilon$  را به شکل  $\exp(-\frac{n}{10})$  کاهش دهیم  $\frac{\epsilon}{n}$  اما این نسبت به تغییر

استراتژی حریف امن نخواهد بود برای اینکه از این مورد نیز امان باشیم می‌توانیم در هر مرحله  $\alpha$  را  
باضرب بسیار کوچکی از تفاوت  $estimated R$  با پاداش دریافتی از بازی افزایش دهیم تا در صورتی  
که سیاست حریف تغییر کرد (که معادل عدم همخوانی پاداش دریافتی با پاداش تخمینی است)، مقدار  $\alpha$  افزایش  
یابد