

① از رابطه بدین معنی داریم $V^{\pi_\theta}(s) = \sum_a \pi_\theta(a|s) Q^{\pi_\theta}(s, a)$. با مشتق گرفتن از این رابطه خواهیم داشت:

$$\nabla_\theta V^{\pi_\theta}(s) = \sum_a \left[\left(\nabla_\theta \pi_\theta(a|s) \right) Q^{\pi_\theta}(s, a) + \pi_\theta(a|s) \nabla_\theta Q^{\pi_\theta}(s, a) \right] \quad \text{I}$$

از رابطه بدین برای Q می داریم:

$$\nabla_\theta Q^{\pi_\theta}(s, a) = \nabla_\theta \left[r(s, a) + \gamma \sum_{s'} P(s'|s, a) V^{\pi_\theta}(s') \right] = \gamma \sum_{s'} P(s'|s, a) \nabla_\theta V^{\pi_\theta}(s')$$

همچنین می داریم که $\nabla_\theta \pi_\theta(a|s) = \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s)$. اکنون با جایگذاری مشتق‌های $Q^{\pi_\theta}(s, a)$ و $\pi_\theta(a|s)$ در رابطه I خواهیم داشت:

$$\begin{aligned} \nabla_\theta V^{\pi_\theta}(s) &= \sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) Q^{\pi_\theta}(s, a) + \gamma \sum_a \pi_\theta(a|s) \sum_{s'} P(s'|s, a) \nabla_\theta V^{\pi_\theta}(s') \\ &= E_{a \sim \pi_\theta(a|s)} \left[\nabla_\theta \log \pi_\theta(a|s) Q^{\pi_\theta}(s, a) \right] + \gamma \sum_{s'} P_r^{\pi_\theta}(s'=s' | s=s) \nabla_\theta V^{\pi_\theta}(s') \end{aligned}$$

توجه کنید عبارت دوم با کمک $P_r^{\pi_\theta}(s'=s' | s=s) = \sum_a \pi_\theta(a|s) P(s'=s' | s, a)$ نتیجه شد و در واقع احتمال رفتن از حالت s به s' را نشان می دهد (بر اساس سیاست π_θ) . عبارت اول را نیز از این به بعد با $J(s)$ نشان می دهیم.

$$\nabla_\theta J(\theta) = \nabla_\theta V^{\pi_\theta}(s) = J(s) + \gamma \sum_{s'} P_r^{\pi_\theta}(s'=s' | s) \nabla_\theta V^{\pi_\theta}(s')$$

آنگاه $\nabla_\theta V^{\pi_\theta}(s)$ (نیز با رابطه ای که داریم) می نویسیم:

$$\nabla_\theta J(\theta) = J(s) + \gamma \sum_{s'} P_r^{\pi_\theta}(s'=s' | s) \left[J(s') + \gamma \sum_{s''} P_r^{\pi_\theta}(s''=s'' | s') \nabla_\theta V^{\pi_\theta}(s'') \right] \quad \text{II}$$

با توجه به ویژگی مارکوفی، می داریم $P_r^{\pi_\theta}(s''=s'' | s) = \sum_{s'} P_r^{\pi_\theta}(s'=s' | s) P_r^{\pi_\theta}(s''=s'' | s')$. اثبات این نکته را در زیر آورده ام:

$$P_r^{\pi_\theta}(s''=s'' | s) = \sum_{s'} P_r^{\pi_\theta}(s'=s' | s) P_r^{\pi_\theta}(s''=s'' | s')$$

$$\sum_{a, s', a'} P(s'=s' | s, a) \times \pi(a|s) \times P(s''=s'' | s', a') \times \pi(a'|s')$$

$$= \sum_{s', a'} P_r^{\pi_\theta}(s'=s' | s) \times P(s''=s'' | s', a') \times \pi(a'|s')$$

$$= \sum_{s'} P_r^{\pi_\theta}(s'=s' | s) \sum_{a'} P(s''=s'' | s', a') \pi(a'|s') = \sum_{s'} P_r^{\pi_\theta}(s'=s' | s) \times P_r^{\pi_\theta}(s''=s'' | s')$$

بنا بر این با کمک این قضیه و همچنین مجزا باز نویسی $D_{\theta} V^{\pi_{\theta}}(s_p)$ در معادله (II) به کمک رابطه‌ای که داریم، می‌توان به نتیجه زیر رسید :

$$\begin{aligned} D_{\theta} J(\theta) &= f(s) + \sum \gamma P_r^{\pi_{\theta}}(s_t = s | s_t) f(s) + \sum \gamma^2 P_r^{\pi_{\theta}}(s_p = s | s_t) f(s) + \dots \\ &= \sum_{t=0}^{\infty} \sum_s \gamma^t P_r^{\pi_{\theta}}(s_t = s | s_t) f(s) = \sum_s \sum_{t=0}^{\infty} \gamma^t P_r^{\pi_{\theta}}(s_t = s | s_t) f(s) \quad \text{(III)} \end{aligned}$$

از رابطه‌ای که در فرضیه‌ها داده شده می‌دانیم :

$$d_{s_t}^{\pi_{\theta}}(s) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t P_r^{\pi_{\theta}}(s_t = s | s_t)$$

بنابراین، رابطه بالا را در معادله (III) جایگذاری کنید :

$$D_{\theta} J(\theta) = \sum_s \frac{d_{s_t}^{\pi_{\theta}}(s)}{1-\gamma} f(s) = \frac{1}{1-\gamma} \sum_s d_{s_t}^{\pi_{\theta}}(s) E_{\alpha \sim \pi_{\theta}(\cdot | s)} [D_{\theta} \log \pi_{\theta}(\alpha | s) Q^{\pi_{\theta}}(s, \alpha)]$$

که اگر بر اساس امید ریاضی از توزیع $d_{s_t}^{\pi_{\theta}}$ باز نویسی کنیم، به همان رابطه گفته شده در

حکم سوال می‌رسیم :

$$D_{\theta} J(\theta) = \frac{1}{1-\gamma} E_{s \sim d_{s_t}^{\pi_{\theta}}} [E_{\alpha \sim \pi_{\theta}(\cdot | s)} [D_{\theta} \log \pi_{\theta}(\alpha | s) Q^{\pi_{\theta}}(s, \alpha)]]$$

(۲) قضیه Compatible Function Approximation

برای اثبات حکم کافیست برای مقادیری که در سوال قبل برای $D_{\theta} J(\theta)$ بدست آوردیم، رابطه جدیدی که در سوال جدید گفته شده، ثابت کنیم. برای این دو، معادل با حکم زیر می‌شود :

$$E_{s \sim d_{s_t}^{\pi_{\theta}}} E_{\alpha \sim \pi_{\theta}(\cdot | s)} [D_{\theta} \log \pi_{\theta}(\alpha | s) (Q^{\pi_{\theta}}(s, \alpha) - Q_{\phi}(s, \alpha))] = 0 \quad \text{(I)}$$

شرط دوم سوال می‌گوید که ϕ مقدار ϵ را کمینه می‌کند. بنابراین نتیجه می‌شود مشتق ϵ نسبت به ϕ برابر صفر است. آنگاه نشان می‌دهیم $D_{\phi} \epsilon = 0$ و رابطه (I) معادل هستند و بنابراین حکم اثبات می‌شود.

باتوجه به اینکه توزیع $d_{s_t}^{\pi_{\theta}}$ ، $\pi_{\theta}(\cdot | s)$ وابسته نیستند، پس می‌توان D_{θ} را داخل امید ریاضی برد :

$$\begin{aligned} D_{\theta} \epsilon &= D_{\phi} E_{s \sim d_{s_t}^{\pi_{\theta}}} E_{\alpha \sim \pi_{\theta}(\cdot | s)} [(Q^{\pi_{\theta}}(s, \alpha) - Q_{\phi}(s, \alpha))] \\ &= E_A E_B [D_{\phi} (Q^{\pi_{\theta}}(s, \alpha) - Q_{\phi}(s, \alpha))] \quad \text{نم‌گذاری مختصر توابع} \end{aligned}$$

$$= E_A E_B [\gamma (Q^{\pi_{\theta}}(s, \alpha) - Q_{\phi}(s, \alpha)) (-D_{\phi} Q_{\phi}(s, \alpha))]$$

$$= E_A E_B [\gamma (Q^{\pi_{\theta}}(s, \alpha) - Q_{\phi}(s, \alpha)) (-D_{\theta} \log \pi_{\theta}(\alpha | s))] \quad \text{استفاده از شرط اول}$$

الفن اگر یک ۲- از داخل امید ریاضی فاکتور بگیریم خواهیم داشت :

$$-r E_{s \sim p} \pi_\theta E_{\alpha \sim \pi(\cdot|s)} \left[\nabla_{\theta} \log \pi_\theta(\alpha|s) (Q^{\pi_\theta}(s, \alpha) - V_\pi(s, \alpha)) \right] = 0$$

که همان رابطه (I) را نتیجه می دهد و حکم ثابت شد

بخش 2.1

$$\eta(\pi) = E_{s, \alpha, s_1, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right] = E_{s, \alpha, s_1, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right] = E_{s, \alpha, s_1, \dots} [V_\pi(s)]$$

روابط را به نحو زیر ادامه می دهیم :

$$\begin{aligned} \eta(\pi) &= E_{s, \alpha, s_1, \dots} [V_\pi(s)] = E_{\gamma \sim \pi} [V_\pi(s)] \\ &= E_{\gamma \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t V_\pi(s_t) - \sum_{t=1}^{\infty} \gamma^t V_\pi(s_t) \right] \\ &= E_{\gamma \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t (r V_\pi(s_1) - V_\pi(s_t)) \right] \quad \textcircled{I} \end{aligned}$$

$$\eta(\pi') = E_{\gamma \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right] \quad \textcircled{II}$$

همچنین می دانیم که :

$$\textcircled{I}, \textcircled{II} \Rightarrow \eta(\pi') - \eta(\pi) = E_{\gamma \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t) + r V_\pi(s_{t+1}) - V_\pi(s_t)) \right] \quad \textcircled{III}$$

توجه کنید بر اساس تعاریف به سادگی می توان دید
و بنا براین :

$$A_\pi(s, \alpha) = Q_\pi(s, \alpha) - V_\pi(s) = r(s) + r V_\pi(s_{t+1}) - V_\pi(s_t) \quad \textcircled{IV}$$

$$\textcircled{III}, \textcircled{IV} \Rightarrow \eta(\pi') = \eta(\pi) + E_{\gamma \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, \alpha_t) \right]$$

$$\eta(\pi') - \eta(\pi) = E_{\gamma \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, \alpha_t) \right] \quad \textcircled{b}$$

(b) از نتیجه بخش قبل استفاده می کنیم :

$$= \sum_{t=0}^{\infty} E_{s_t \sim p_{\pi}(s_t)} \left[E_{\alpha_t \sim \pi'(\cdot|s_t)} [\gamma^t A_\pi(s_t, \alpha_t)] \right]$$

$$= \sum_{t=0}^{\infty} \sum_s p(s_t=s) \sum_{\alpha} \pi'(\alpha|s) \gamma^t A_\pi(s, \alpha)$$

$$= \sum_s \left(\sum_{t=0}^{\infty} \gamma^t p(s_t=s) \right) \sum_{\alpha} \pi'(\alpha|s) A_\pi(s, \alpha)$$

$$\Rightarrow \eta(\pi') = \eta(\pi) + \sum_s p_{\pi'}(s) \sum_{\alpha} \pi'(\alpha|s) A_\pi(s, \alpha)$$

حکم ثابت شد .

Subject :

Year . Month . Date . ()

(c) از تعریف، می دانیم :

$$\bar{A}(s) = \sum_{\alpha'} \pi'(\alpha' | s) A_{\pi}(s, \alpha') = \sum_{\alpha'} \sum_{\alpha} P(\alpha, \alpha' | s) \bar{A}_{\pi}(s, \alpha')$$

تساوی درم از اینجا آمد که π' ، marginal dist. روی α' از توزیع $\alpha, \alpha' | s$ است.
همچنین توجه کنید که $E_{\alpha \sim \pi(s)} [A_{\pi}(s, \alpha)] = 0$ زیرا :

$$\begin{aligned} E_{\alpha \sim \pi(s)} [A_{\pi}(s, \alpha)] &= E_{\alpha \sim \pi(s)} [Q_{\pi}(s, \alpha)] - E_{\alpha \sim \pi(s)} [V_{\pi}(s)] \\ &= V_{\pi}(s) - V_{\pi}(s) = 0 \end{aligned}$$

پس می توان نوشت :

$$0 = \sum_{\alpha} \pi(\alpha | s) A_{\pi}(s, \alpha) = \sum_{\alpha} \sum_{\alpha'} P(\alpha, \alpha' | s) A_{\pi}(s, \alpha)$$

با توجه به 0 بودن این عبارت، می توان $\bar{A}(s)$ را بدینوسیله کرد :

$$|\bar{A}(s)| = \left| \sum_{\alpha', \alpha} P(\alpha, \alpha' | s) A_{\pi}(s, \alpha') - \sum_{\alpha, \alpha'} P(\alpha, \alpha' | s) A_{\pi}(s, \alpha) \right|$$

$$= \left| \sum_{\alpha, \alpha'} P(\alpha, \alpha' | s) (A_{\pi}(s, \alpha') - A_{\pi}(s, \alpha)) \right|$$

نامساوی مثلث

$$\leq \sum_{\alpha, \alpha'} P(\alpha, \alpha' | s) \frac{1}{2} |A_{\pi}(s, \alpha') - A_{\pi}(s, \alpha)|$$

حالا توجه کنید اگر $\alpha = \alpha'$ ، آنگاه $A_{\pi}(s, \alpha') - A_{\pi}(s, \alpha) = 0$

در غیر این صورت نیز می توان کران بالایی را داد :

$$|A_{\pi}(s, \alpha') - A_{\pi}(s, \alpha)| \leq |A_{\pi}(s, \alpha')| + |-A_{\pi}(s, \alpha)| \leq 2 \max_{s, \alpha} |A_{\pi}(s, \alpha)|$$

در نتیجه می توان نوشت :

$$|\bar{A}(s)| \leq \sum_{\alpha \neq \alpha'} P(\alpha, \alpha' | s) \times 2 \max_{s, \alpha} |A_{\pi}(s, \alpha)| \Rightarrow |\bar{A}(s)| \leq 2 \max_{s, \alpha} |A_{\pi}(s, \alpha)|$$

$$\text{می دانیم: } \sum_{\alpha \neq \alpha'} P(\alpha, \alpha' | s) = P(\alpha \neq \alpha' | s) \leq \alpha$$

حکم ثابت شد

$$|E_{s_t \sim \pi} [\bar{A}(s_t)] - E_{s_t \sim \pi} [\bar{A}(s_t)]| = \left| \sum_s P_{\pi}(s_t = s) \bar{A}(s) - \sum_s P_{\pi}(s_t = s) \bar{A}(s) \right| \quad (d)$$

نامساوی مثلث

$$\leq \sum_s |P_{\pi'}(s_t = s) - P_{\pi}(s_t = s)| |\bar{A}(s)|$$

لم قبلی

$$\leq \sum_s |P_{\pi'}(s_t = s) - P_{\pi}(s_t = s)| \times 2 \max_{s, \alpha} |A_{\pi}(s, \alpha)|$$

بنابراین کیفیت ثابت کنیم $\sum_s |P_{\pi'}(s_t = s) - P_{\pi}(s_t = s)| \leq 2(1 - (1 - \alpha)^t)$ تا حکم نهایی

اثبات شود این موضوع نیز در کلاس ثابت شده بود و اثبات آن را در ادامه تکراری کنیم :

(تقریباً)

فرض کنید S'_n امتیاز در زمان n بر اساس پالیسی π' ، S_n امتیاز در زمان n بر اساس پالیسی π باشد، با این فرض که هر دو از حالت اولیه یکسان شروع شوند: $S'_1 = S_1 = p(s)$

در هر لحظه، اکشن‌های (a_n, a'_n) را از توزیع $P(a_n, a'_n | S_n)$ انتخاب می‌کنیم (اگر $S_n = S'_n$ بود). اگر $a_n = a'_n$ باشد، آنگاه $S_{n+1} = S'_{n+1}$ خواهد بود (چون از transition و random seed های یکسان استفاده می‌کنیم) که احتمال این اتفاق حداقل $1-\alpha$ است. اگر اینگونه نباشد، S_{n+1} و S'_{n+1} لزومی ندارد برابر شوند، هر چند باز هم ممکن است.

آنگاه احتمال اینکه در لحظه n ، $S_n \neq S'_n$ باشد به نحو زیر است:

$$\left. \begin{aligned} P(S'_n \neq S_n) &= P(S'_n \neq S_n | S_{n-1} = S'_{n-1}) \times P(S_{n-1} = S'_{n-1}) + P(S'_n \neq S_n | S_{n-1} \neq S'_{n-1}) P(S_{n-1} \neq S'_{n-1}) \\ P(S'_n \neq S_n | S_{n-1} = S'_{n-1}) &\leq P(a_n \neq a'_n | S_{n-1} = S'_{n-1}) \leq \alpha \\ P(S'_n \neq S_n | S_{n-1} \neq S'_{n-1}) &\leq 1 \end{aligned} \right\}$$

$$\Rightarrow P(S'_n \neq S_n) \leq \alpha \times (1 - P(S_{n-1} \neq S'_{n-1})) + P(S_{n-1} \neq S'_{n-1}) = \alpha + (1-\alpha)P(S_{n-1} \neq S'_{n-1})$$

$$\Rightarrow P(S'_n \neq S_n) \leq \alpha + \alpha(1-\alpha) + \alpha(1-\alpha)^2 + \dots + \alpha(1-\alpha)^{n-1} = \alpha \frac{1-(1-\alpha)^n}{1-(1-\alpha)} = 1 - (1-\alpha)^n \quad (*)$$

$$\sum_s |P_{\pi}(S_t = s) - P_{\pi'}(S_t = s)| = 2D_{TV}(P_{\pi}(S_t = \cdot) \| P_{\pi'}(S_t = \cdot))$$

همچنین می‌دانیم:

$$2P(S'_t \neq S_t) \leq 2(1 - (1-\alpha)^n) \quad (*)$$

حکم مفادگی که می‌خواهیم ثابت شد و اثبات کامل شد.

(e) از تعاریف موجود برای γ و L_{π} می‌توان نوشت:

$$|J(\pi) - L_{\pi}(\pi)| = \left| E_{\gamma \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \bar{A}(S_t) \right] - E_{\gamma \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \bar{A}(S_t) \right] \right|$$

$$= \left| \sum_{t=0}^{\infty} \gamma^t (E_{\gamma \sim \pi} [\bar{A}(S_t)] - E_{\gamma \sim \pi} [\bar{A}(S_t)]) \right|$$

خطی بودن امید ریاضی:

$$\leq \sum_{t=0}^{\infty} \gamma^t |E_{\gamma \sim \pi} [\bar{A}(S_t)] - E_{\gamma \sim \pi} [\bar{A}(S_t)]|$$

نامساوی مثلث:

$$\leq \sum_{t=0}^{\infty} \gamma^t (\epsilon \alpha (1 - (1-\alpha)^t) \epsilon) = \epsilon \alpha \epsilon \sum_{t=0}^{\infty} \gamma^t (1 - (1-\alpha)^t)$$

لم سوال قبلی:

آنگاه کافیت ثابت کنیم $\sum_{t=0}^{\infty} \gamma^t (1 - (1-\alpha)^t) \leq \frac{\alpha \gamma}{(1-\gamma)^2}$ ماحکم ثابت شود:

$$\sum_{t=0}^{\infty} \gamma^t (1 - (1-\alpha)^t) = \sum_{t=0}^{\infty} \gamma^t - \sum_{t=0}^{\infty} (\gamma(1-\alpha))^t = \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha)}$$

$$= \frac{1-\gamma(1-\alpha) - (1-\gamma)}{(1-\gamma)(1-\gamma+\gamma\alpha)} = \frac{\gamma\alpha}{(1-\gamma)(1-\gamma+\gamma\alpha)} \leq \frac{\gamma\alpha}{(1-\gamma)^2}$$

حکم ثابت شد.

Subject :

Year . Month . Date . ()

(f) لم 25 به مای گوید اگر (π, π') یک α -coupled Policy Pair باشد (کطبق، اهنای هست) آنگاه داریم :

$$|\eta(\pi) - L_{\pi}(\pi')| \leq \frac{\alpha^2 \gamma \epsilon}{(1-\gamma)^2} \Rightarrow \eta(\pi') \geq L_{\pi}(\pi') - \frac{\alpha^2 \gamma \epsilon}{(1-\gamma)^2}$$

با توجه به اهنای، می توانیم در لم بالا قرار دهیم $\alpha = D_{TV}^{\max}(\pi, \pi')$:

$$\eta(\pi') \geq L_{\pi}(\pi') - \frac{\epsilon (D_{TV}^{\max}(\pi, \pi'))^2 \gamma}{(1-\gamma)^2} \quad (*)$$

با توجه به نامساوی $D_{TV}(\pi \parallel \pi') \leq D_{KL}(\pi \parallel \pi')$ در فرضیات داده شده، داریم :

$$D_{TV}(\pi(\cdot|s) \parallel \pi'(\cdot|s))^2 \leq D_{KL}(\pi(\cdot|s) \parallel \pi'(\cdot|s)) \quad (I)$$

$$\text{بر } b : (D_{TV}^{\max}(\pi, \pi'))^2 = \left[\max_s D_{TV}(\pi(\cdot|s) \parallel \pi'(\cdot|s)) \right]^2 = \max_s [D_{TV}(\pi(\cdot|s) \parallel \pi'(\cdot|s))]^2$$

$$\stackrel{(I)}{\leq} \max_s D_{KL}(\pi(\cdot|s) \parallel \pi'(\cdot|s)) = D_{KL}^{\max}(\pi, \pi')$$

$$\stackrel{(*)}{\Rightarrow} \eta(\pi') \geq L_{\pi}(\pi') - \frac{\epsilon \gamma}{(1-\gamma)^2} \times D_{KL}^{\max}(\pi, \pi') \quad \text{حکم اثبات شد}$$