-١

(α

$$\nabla_\theta J(\theta) = \nabla_\theta \mathop{E}_{T\sim\pi_\theta}\left\{\sum_{t=0}^{\infty}\gamma^t r(s_t, a_t)\right\} = \nabla_\theta \int P(T|\theta)\sum_{t=0}^{\infty}\gamma^t r(s_t, a_t)\,dT$$

$$= \partial \int \left((\nabla_\theta P(T|\theta))\sum_{t=0}^{\infty}\gamma^t r(s_t, a_t)\right)dT$$

$$P(T|\theta) = P(s_0)\prod_{t=0}^{\infty}\pi_\theta(a_t|s_t)P(s_{t+1}|s_t, a_t) \xrightarrow{\ Ln\ }$$

$$Ln\,P(T|\theta) = Ln\,P(s_0) + \sum_{t=0}^{\infty}Ln\,P(s_{t+1}|s_t, a_t) + \sum_{t=0}^{\infty}Ln\,P(s_t|s_t)$$

$$+ \sum_{t=0}^{\infty}Ln\,\pi_\theta(a_t|s_t)$$

$$\xrightarrow{\ \partial/\partial\theta\ } \sum_{t=0}^{\infty}\nabla_\theta\,Ln\,\pi_\theta(a_t|s_t)\Bigg|$$

$$\Rightarrow \int \nabla_\theta\,Ln\,P(T|\theta)\,P(T|\theta)\left(\sum_{t=0}^{\infty}\gamma^t r(s_t, a_t)\right)dT =$$

$$\mathop{E}_{T\sim\pi_\theta}\left\{\sum_{t=0}^{\infty}\gamma^t r(s_t, a_t)\sum_{k=0}^{\infty}\nabla_\theta\,Ln\,\pi_\theta(a_k|s_k)\right\}$$

causality از این رابطه می‌توان با استفاده از داریم

$$\nabla_\theta J(\theta) = \mathop{E}_{T\sim\pi_\theta}\left\{\int\sum_{t=0}^{\infty}\gamma^t r(s_t, a_t)\sum_{k=t}^{\infty}\nabla_\theta\,Ln\,\pi_\theta(a_k|s_k)\right\}$$

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left\{ \sum_{t=0}^{\infty} \gamma^t \, \nabla_\theta \log \pi_\theta(a_t | s_t) \left( \sum_{k=0}^{\infty} \gamma^k r(s_{t+k}, a_{t+k}) \right) \right\}$$

$$Q^{\pi_\theta}(s_t, a_t) = \mathbb{E}_{\tau \sim \pi_\theta} \left\{ \sum_{k=0}^{\infty} \gamma^k r(s_{t+k}, a_{t+k}) \Big| s_t, a_t \right\}$$

$$\Rightarrow \nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left\{ \sum_{t=0}^{\infty} \gamma^t \nabla_\theta \log \pi_\theta(a_t | s_t) \, Q^{\pi_\theta}(s_t, a_t) \right\}$$

$$\mu_{s_0}^{\pi_\theta}(s) := \sum_{t=0}^{\infty} \gamma^t \, \overset{\pi_\theta}{Pr}(s_t = s | s_0)$$

↑ موو trajectory-based يكمين موو حالت اسکراذ

$$\rightsquigarrow \nabla_\theta J(\theta) = \sum_s \mu_{s_0}^{\pi_\theta}(s) \sum_a \pi_\theta(a|s) \, \nabla_\theta \log \pi_\theta(a|s) \, Q^{\pi_\theta}(s, a)$$

از آنجا که $\mu_{s_0}^{\pi_\theta}(s)$ يک تورِیع نیست و به یکِ جمع نمی باشِ با پشِایا ١

$$\sum_s \mu_{s_0}^{\pi_\theta}(s) = \frac{1}{1-\gamma} \longrightarrow \sum_{t=0}^{\infty} \gamma^t \sum_s \overset{\pi_\theta}{Pr}(s_t = s | s_0) =$$

$$\sum_{t \leq 0}^{\infty} \gamma^t = \frac{1}{1-\gamma} \longrightarrow d_{s_0}^{\pi_\theta}(s) = (1-\gamma) \, \mu_{s_0}^{\pi_\theta}(s) \checkmark$$

حال داریم ؛

$$\nabla_\theta J(\theta) = \frac{1}{1-\gamma} \sum_s d_{s_0}^{\pi_\theta}(s) \sum_a \pi_\theta(a|s) \, \nabla_\theta \log \pi_\theta(a|s) \, Q^{\pi_\theta}(s, a)$$

$$= \frac{1}{1-\gamma} \, \mathbb{E}_{s \sim d_{s_0}^{\pi_\theta}(s)} \left\{ \mathbb{E}_{a \sim \pi_\theta(\cdot | s)} \, \nabla_\theta \log \pi_\theta(a|s) \, Q^{\pi_\theta}(s, a) \right\}$$

ب) داریم،

$$\nabla_\phi(\varepsilon) = \nabla_\phi \mathop{E}_{S \sim d_{s_o}^{\pi_\theta}} \mathop{E}_{a \sim \pi_\theta(\cdot|s)} \left[ \left( Q^{\pi_\theta}(s,a) - Q_\phi(s,a) \right)^\gamma \right] =$$

$$-\gamma \mathop{E}_{S \sim d_{s_o}^{\pi_\theta}} \mathop{E}_{a \sim \pi_\theta(\cdot|s)} \left\{ \left( Q^{\pi_\theta}(s,a) - \nabla_\phi Q_\phi(s,a) \right) \right\} = 0$$

$$Q_\phi(s,a)$$

$$\Rightarrow \mathop{E}_{S \sim d_{s_o}^{\pi_\theta}} \mathop{E}_{a \sim \pi_\theta(\cdot|s)} \left[ Q^{\pi_\theta}(s,a) \nabla_\phi Q_\phi(s,a) \right] =$$

$$\mathop{E}_{S \sim d_{s_o}^{\pi_\theta}} \mathop{E}_{a \sim \pi_\theta(\cdot|s)} \left[ Q_\phi(s,a) \nabla_\phi Q_\phi(s,a) \right] =$$

$$\mathop{E}_{S \sim d_{s_o}^{\pi_\theta}} \mathop{E}_{a \sim \pi_\theta(\cdot|s)} \left\{ Q_\phi(s,a) \nabla_\theta \log \pi_\theta(a|s) \right\}$$

اینم! $\frac{1}{1-\gamma}$ هم این اسکاله میشه حذف کواله لثبات - المسئله.

٢) الف) داریم،

$$A_\pi(s_t, a_t) = Q_\pi(s_t, a_t) - V_\pi(s_t).$$

طبق تعریف داریم / باید نشون بدیم،

$$Q_\pi(s_t, a_t) = r(s_t, a_t) + \gamma \mathop{E}_{s_{t+1} \sim P(\cdot|s_t, a_t)} \left[ V_\pi(s_{t+1}) \right]$$

$\Rightarrow \quad \underbrace{\cancel{Q_\pi(s_t,a_t)}}_{A_\pi(s_t,a_t)} = r(s_t, a_t) + \gamma \underset{s_{t+1}}{\mathbb{E}} V_\pi(s_{t+1}) - V_\pi(s_t)$

$\longrightarrow \quad r(s_t, a_t) = A_\pi(s_t, a_t) + V_\pi(s_t) - \gamma \underset{s_{t+1}}{\mathbb{E}} V_\pi(s_{t+1})$

expected return of $\pi'$

$\underset{\mathbb{E}\{R(\pi')\}}{\Big(} = \underset{\pi'}{\mathbb{E}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] =$

$\underset{\pi'}{\mathbb{E}} \left[ \sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right] + \underset{\pi'}{\mathbb{E}} \left[ \sum \gamma^t V_\pi(s_t) \right] -$

$\underset{\pi'}{\mathbb{E}} \left[ \sum_{t=0}^{\infty} \gamma^{t+1} V_\pi(s_{t+1}) \right]$

$\underset{s_0 \sim P_0}{E} \left[ V_\pi(s_0) \right]$

$= \underset{\pi'}{E} \left[ \sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right] + \underset{s_0 \sim P_0}{\mathbb{E}} \left[ V_\pi(s_0) \right] \Rightarrow$

$\eta(\pi') = \eta(\pi) + \underset{\pi'}{E} \left[ \sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right]$

$\Rightarrow \eta(\pi') = \eta(\pi) + \underset{\substack{s_0, a_0, \\ \cdots \sim \pi'}}{E} \left\{ \sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right\} \checkmark$

$$\mathbb{E}_{\pi'}\left[\sum_{t=0}^{\infty}\gamma^t A_\pi(s_t,a_t)\right] = \sum_{t=0}^{\infty}\sum_{s_t}\gamma^t P(s_t)\sum_{a_t}\pi'(a|s)A_\pi(s_t,a_t)$$

$$= \sum_{s} P_{\pi'}(s_t)\sum_{a}\pi'(a|s)A_\pi(s,a)$$

$$\mathbb{E}_{a\sim\pi(\cdot|s)}\left[A_\pi(s,a)\right] \overset{?}{=} 0 \quad , \quad V_\pi(s) = \mathbb{E}_{a\sim\pi(\cdot|s)}\left[Q_\pi(s,a)\right] \quad , \quad داریم$$

صفرالست ، خواهیم داشت :

$$\bar{A}(s) = \mathbb{E}_{a'\sim\pi'(\cdot|s)}\left[A_\pi(s,a')\right] =$$

$$\mathbb{E}_{(a,a')\sim(\pi,\pi')|s}\left[A_\pi(s,a')-A_\pi(s,a)\right] =$$

$$P(a\ne a'|s)\,\mathbb{E}_{(a,a')\sim(\pi,\pi')|s;\,a\ne a'}\left[A_\pi(s,a')-A_\pi(s,a)\right]$$

$$\le \alpha\,\mathbb{E}_{(a,a')\sim(\pi,\pi')|s,\,a=a'}\left[A_\pi(s,a')-A_\pi(s,a)\right]$$

$$\Rightarrow |\bar{A}(s)| \le \alpha\left(\max_{s,a'}\left|A_\pi(s,a')\right| + \max_{s,a}\left|A_\pi(s,a)\right|\right)$$

$$= 2\alpha\,\max_{s,a}\left|A_\pi(s,a)\right|$$

d) با احتمال $(1-\alpha)^t$ تمام action ها / policy دوم یکسان خواهد بود

← با احتمال $1-(1-\alpha)^t$ متفاوت خواهد بود.

$$\left| \mathbb{E}_{S_t \sim \pi'}\left[\bar{A}(S_t)\right] - \mathbb{E}_{S_t \sim \pi}\left[\bar{A}(S_t)\right] \right| = \left(1-(1-\alpha)^t\right)\left| \mathbb{E}_{\substack{S_t \sim \pi' \mid \\ a \neq a'}}\left[\bar{A}(S_t)\right] - \mathbb{E}_{\substack{S_t \sim \pi \mid a \neq a'}}\left[\bar{A}(S_t)\right] \right|$$

$$\leq \left(1-(1-\alpha)^t\right)\left( \left| \mathbb{E}_{\substack{S_t \sim \pi' \mid a \neq a'}}\left[\bar{A}(S_t)\right]\right| + \left|\mathbb{E}_{\substack{S_t \sim \pi \mid a \neq a'}}\left[\bar{A}(S_t)\right]\right| \right)$$

$$\leq \left(1-(1-\alpha)^t\right)\left( 2\alpha \max_{s,a}\left|A_\pi(s,a)\right| + 2\alpha \max_{s,a}\left|A_\pi(s,a)\right| \right)$$

$$= 4\alpha\left(1-(1-\alpha)^t\right)\max\left|A_\pi(s,a)\right|$$

e)

$$\left| \eta(\pi') - L_\pi(\pi') \right| = \left| \sum_{t=0}^{\infty} \gamma^t \left( \mathbb{E}_{S_t \sim \pi'}\left[\bar{A}(S_t)\right] - \mathbb{E}_{S_t \sim \pi}\left[\bar{A}(S_t)\right]\right) \right|$$

$$\leq \sum_{t=0}^{\infty} \gamma^t \left| \mathbb{E}_{S_t \sim \pi'}\left[\bar{A}(S_t)\right] - \mathbb{E}_{S_t \sim \pi}\left[\bar{A}(S_t)\right]\right| \leq$$

$$\sum_{t=0}^{\infty} \gamma^t \cdot 4\alpha\left(1-(1-\alpha)^t\right)\varepsilon = 4\alpha\varepsilon\left( \sum_{t=0}^{\infty}\gamma^t - \sum_{t=0}^{\infty}\gamma^t(1-\alpha)^t \right)$$

$$= 4\alpha\varepsilon\left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha)} \right)$$

$$\leq 4\alpha\varepsilon\left( \frac{\alpha\gamma}{(1-\gamma)^2} \right) = \frac{4\alpha^2\varepsilon\gamma}{(1-\gamma)^2}$$

$$\left| \eta(\pi') - L_\pi(\pi') \right| \leq \frac{4\alpha^2 \varepsilon \gamma}{(1-\gamma)^2} \implies$$

$$\eta(\pi') - L_\pi(\pi') \geq -\frac{4\alpha^2 \varepsilon \gamma}{(1-\gamma)^2} \implies$$

$$\eta(\pi') \geq L_\pi(\pi') - \frac{4\alpha^2 \varepsilon \gamma}{(1-\gamma)^2} \quad \text{(بتا(π))} \quad \text{(بتا)}$$

$$= L_\pi(\pi') - \frac{4 \varepsilon \gamma}{(1-\gamma)^2} D_{TV}^{max}(\pi, \pi')^2 \geq$$

$$L_\pi(\pi') - \frac{4 \varepsilon \gamma}{(1-\gamma)^2} D_{KL}^{max}(\pi, \pi')$$

$$\boxed{\begin{array}{c} D_{TV}\left( \pi(\cdot|s) \parallel \pi'(\cdot|s) \right)^2 \leq D_{KL}\left( \pi(\cdot|s) \parallel \pi'(\cdot|s) \right) \\[3mm] \implies D_{TV}^{max}(\pi, \pi') \leq D_{KL}^{max}(\pi, \pi') \end{array}}$$

به خاطر اینکه نامساوی بالا برقرار است.