

سوال ۱:

الف) با توجه به اینکه فضای transition ها کامل^۱ probabilistic هستند،

برای اینکه مجموع پاداش را بخواهیم بیشینه کنیم، ممکن است که اگر در ایست ۵ باشیم و

یک action مانند a انجام بدهیم، ممکن است دو برابر متفاوت بگیریم.

($P(a|s, a)$ وجود دارد). همچنین می توان به این اشاره کرد که در یک State

بیشتر $P(a|s)$ وجود دارد و رندومس وجود دارد.

$$J(\theta) = \frac{1}{|P_\theta|} \sum_{\tau} \gamma^t r(s_t, a_t)$$

اگر بخواهیم کمترین گمراهی را داشته باشیم، از رابطه بالا که بیان نسبت به θ می گیرد.

با توجه به اینکه trajectory ها از یک توزیع P_θ می آیند و reward ها در طی مسیر

خودشان به صورت مستقیم از این عبارت که بیان گرفت در بایه این عبارت را باز کرده

و پس که اگر بخواهیم θ را تابع θ محاسبه شود و قابل محاسبه شود.

برای رفع این مشکل می توان پاداش را به صورتی تعریف کنیم که تابعی از θ داشته باشیم.

برای این کار می توان پارامتر جدیدی تعریف کنیم و طبق یک تابع مانند $r(s_t, a_t, \theta)$

اکنون را انتخاب کنیم و پس می توان θ را به صورتی تعریف کرد که تابعی از θ باشد.

$$J(\theta) = \frac{1}{|P_\theta|} \sum_{\tau} \gamma^t r(s_t, a_t, \theta)$$

مسئله 2

$$D_Q(J_Q) = D_Q \left\{ E \left\{ \sum_t \gamma^t r_{(s_t, a_t)} \right\} \right\} \rightarrow R = \sum_{(z)} \sum_{t=0}^{\infty} \gamma^t r_{(s_t, a_t)} \quad (الف)$$

$$Z = (s_0, a_0, s_1, a_1, \dots)$$

$$\begin{aligned} D_Q(J_Q) &= D_Q \left\{ \int P_Q(z) R(z) dz \right\} = \int D_Q P_Q(z) R(z) dz \\ &= \int D_Q \log P_Q(z) P_Q(z) R(z) dz \\ &= E_Q \left\{ D_Q \log P_Q(z) R(z) \right\} \end{aligned}$$

$$\Rightarrow D_Q(J_Q) = E_Q \left\{ \dots \right\} \quad \text{حال در نظر بگیریم که داریم:}$$

$$P_Q(z) = P(s_0) P(s_1 | s_0, a_0) P_Q(a_0 | s_0) \dots \Rightarrow$$

$$\Rightarrow D_Q \log P_Q(z) = \sum_{t=0}^{\infty} D_Q \log \pi_Q(a_t | s_t)$$

$$D_Q(J_Q) = E_Q \left\{ \left(\sum_t D_Q \log \pi_Q(a_t | s_t) \right) \left(\sum_t \gamma^t r_{(s_t, a_t)} \right) \right\} \quad \text{پس داریم که:}$$

حال بخواهیم در نظر می گیریم:

$$E(D_Q(J_Q)) = E \left\{ \frac{1}{N} \sum_{i=1}^N \left\{ \left(\sum_{t=0}^{\infty} D_Q \log \pi_Q(a_t | s_t) \right) \left(\sum_t \gamma^t r_{(s_t, a_t)} \right) \right\} \right\} =$$

$$= \frac{1}{N} \sum_{i=1}^N \left\{ E \left\{ \left(\sum_{t=0}^{\infty} D_Q \log \pi_Q(a_t | s_t) \right) \left(\sum_t \gamma^t r_{(s_t, a_t)} \right) \right\} \right\} =$$

$$\text{s.a.m} \quad = \frac{1}{N} \sum_{i=1}^N D_Q J_Q = D_Q J_Q \rightarrow \text{unbiased}$$

در حین اثبات فرض کردیم که

(1) trajectory ها کاملاً از هم مجزا هستند

(2) در تمامی فضا π_t هایی که داریم، مشتق پذیر هستند.

$$P_{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t P(s_t = s)$$

(←)

$$B = \mathbb{E} \left\{ \left(\sum_{t=0}^{\infty} \gamma^t \log \pi_{\theta}(a_t | s_t) \right) \left(\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right) \right\}$$

$$= \sum_{t=0}^{\infty} \mathbb{E} \left\{ \gamma^t \log \pi_{\theta}(a_t | s_t) \sum_{t'=0}^{\infty} \gamma^{t'} r(s_{t'}, a_{t'}) \right\}$$

$$= \sum_{t=0}^{\infty} \left\{ \mathbb{E} \left\{ \gamma^t \log \pi_{\theta}(a_t | s_t) \sum_{t'=0}^{t-1} \gamma^{t'} r(s_{t'}, a_{t'}) \right\} + \mathbb{E} \left\{ \gamma^t \log \pi_{\theta}(a_t | s_t) \sum_{t'=t}^{\infty} \gamma^{t'} r(s_{t'}, a_{t'}) \right\} \right\}$$

با توجه به Causality می دانیم که تصمیمی که در مرحله t می گیریم، تأثیری

در $t-1, t-2, \dots$ ندارد. پس می توان جمله اول را

به این صورت نوشت:

$$E \left\{ D \log \pi_{\theta}(a_t | s_t) \sum_{t'=0}^{t-1} \gamma^{t'} r(s_{t'}, a_{t'}) \right\},$$

$$= E \left\{ D \log \pi_{\theta}(a_t | s_t) \right\} E \left\{ \sum_{t'=0}^{t-1} \gamma^{t'} r(s_{t'}, a_{t'}) \right\}$$

و همین داریم که

$$E \left\{ D \log \pi_{\theta}(a_t | s_t) \right\} = \int P_{\theta} \log \pi_{\theta}(a_t | s_t) d\pi_{\theta}$$

$$= \int \pi_{\theta} D \log \pi_{\theta} d\pi_{\theta} = \int \pi_{\theta} \frac{D \pi_{\theta}}{\pi_{\theta}} d\pi_{\theta} = D \left\{ \int \pi_{\theta} d\pi_{\theta} \right\} = D \{1\} = 0$$

پس می توان B را به این صورت نوشت:

$$\sum_{t=1}^{\infty} \left\{ E \left\{ D \log \pi_{\theta}(a_t | s_t) \sum_{t'=t}^{\infty} \gamma^{t'} r(s_{t'}, a_{t'}) \right\} \right\}$$

$$= \sum_{t=1}^{\infty} \gamma^t E \left\{ D \log \pi_{\theta}(a_t | s_t) \sum_{t'=t}^{\infty} \gamma^{t'-t} r(s_{t'}, a_{t'}) \right\},$$

$$= \sum_{t=1}^{\infty} \gamma^t E \left\{ D \log \pi_{\theta}(a_t | s_t) \hat{Q}(s_t, a_t) \right\}$$

توجه به s_t و a_t داشته باشیم

$$= \sum_{t=1}^{\infty} \gamma^t E_{s_t, a_t} \left\{ D \log \pi_{\theta}(a_t | s_t) \hat{Q}(s_t, a_t) \right\}$$

$$= \sum_{t=1}^{\infty} \gamma^t \int \int D \log \pi_{\theta}(a_t | s_t) \hat{Q}(s_t, a_t) P(s_t = s) \pi_{\theta}(s, a) da ds$$

$$= \sum_{t=1}^{\infty} \int_s P(s_t = s) \int_a D \log \pi_{\theta}(a_t | s_t) \hat{Q}(s_t, a_t) \pi_{\theta}(s, a) da ds$$

$$\int_S \sum_{t=1}^{\infty} \gamma^t P(s_t=s) E_{\pi} \{ \nabla_a \log \pi_{\theta}(a_t/s_t) \hat{Q}(s_t, a_t) \} ds =$$

$$= \int_S P_{\pi}(s) E_{\pi} \{ \nabla_a \log \pi_{\theta}(a_t/s_t) \hat{Q}(s_t, a_t) \} ds =$$

$$= E_{\pi, P_{\pi}} \{ \nabla_a \log \pi_{\theta}(a_t/s_t) \hat{Q}(s_t, a_t) \}$$

→ فرض کنید در تک گامی

$$h = \nabla_a \log \pi_{\theta}(a, s) (\hat{Q}(s, a) - b(s))$$

$$\text{Var}(h) = E_{P_{\pi, \pi}}(h^T h) - E(h)^T E(h)$$

$$= E_{P_{\pi, \pi}} \left(\left(\nabla_a \log \pi_{\theta}(a, s) \right)^T \nabla_a \log \pi_{\theta}(a, s) \right) b(s)^2$$

$$- 2 E_{P_{\pi, \pi}} \left(\left(\nabla_a \log \pi_{\theta}(a, s) \right)^T \left(\nabla_a \log \pi_{\theta}(a, s) \right) \hat{Q}(s, a) \right) b(s)$$

$$+ E_{P_{\pi, \pi}} \left(\left(\nabla_a \log \pi_{\theta}(a, s) \right)^T \left(\nabla_a \log \pi_{\theta}(a, s) \right) \hat{Q}(s, a)^2 \right)$$

...

با توجه که تخمیندها ناآرپ است تقریباً جلاست به $b(s)$
وابسته نیست و در مشتق گرفتن تأثیری ندارند.

$$\frac{\partial(\text{Var}(h))}{\partial(b(s))} = 0 \Rightarrow 2b(s) \mathbb{E}_{\pi_Q} \left\{ (V_Q \log \pi_Q(a, s))^T V_Q \log \pi_Q(a, s) \right\} - 2 \mathbb{E}_{\pi_Q} \left\{ (V_Q \log \pi_Q(a, s))^T V_Q \log \pi_Q(a, s) Q'(s, a) \right\} = 0$$

$$\Rightarrow b(s) = \frac{\mathbb{E}_{\pi_Q} \left\{ (V_Q \log \pi_Q(a, s))^T V_Q \log \pi_Q(a, s) Q'(s, a) \right\}}{\mathbb{E}_{\pi_Q} \left\{ (V_Q \log \pi_Q(a, s))^T V_Q \log \pi_Q(a, s) \right\}}$$

در نظر بگیریم؛ داشتن π_Q دو عبارت است: نیاز داریم از هم مستقل باشند.
 $Q'(s, a)$ و $(V_Q \log \pi_Q(a, s))^T V_Q \log \pi_Q(a, s)$ از هم مستقل هستند.

$$b(s) = \frac{\mathbb{E}_{\pi_Q} \left\{ (V_Q \log \pi_Q(a, s))^T V_Q \log \pi_Q(a, s) \right\} \mathbb{E}_{\pi_Q} \left\{ Q'(s, a) \right\}}{\mathbb{E}_{\pi_Q} \left\{ (V_Q \log \pi_Q(a, s))^T V_Q \log \pi_Q(a, s) \right\}} \Rightarrow$$

$$\Rightarrow b(s) = \mathbb{E}_{\pi_Q} \left\{ Q'(s, a) \right\} = V^{\pi}(s)$$

سوال ۴

الف) فرض کنید μ و ν دو توزیع احتمالی به روی فضای X باشند. در این صورت داریم که:

$$\|\mu - \nu\|_{TV} = \inf_{\gamma} \{P(X \neq Y) \mid \gamma \text{ is a coupling of } \mu, \nu\}$$

اثبات قضیه ۱:

$$\forall A \subseteq X: \mu(A) - \nu(A) = P(X \in A) - P(Y \in A) \leq P(X \in A, Y \notin A)$$

$$\leq P(X \neq Y) \Rightarrow$$

$$\Rightarrow \forall A \subseteq X: \mu(A) - \nu(A) \leq \inf_{\gamma} \{P(X \neq Y) \mid (X, Y) \text{ is a coupling of } \mu, \nu\}$$

$$\Rightarrow \|\mu - \nu\|_{TV} = \max_{A \subseteq X} |\mu(A) - \nu(A)| \leq \inf_{\gamma} \{P(X \neq Y) \mid (X, Y) \text{ is a coupling of } \mu, \nu\}$$

برای اثبات قضیه تنها کافی است که μ, ν معرفی کنیم که در آن

تساوی برقرار باشد. برای این کار به صورت زیر عمل می‌کنیم.

$$P \leq \sum \min \{ \mu_{(n)}, \nu_{(n)} \} = \sum_{n: \mu_{(n)} \geq \nu_{(n)}} \nu_{(n)} + \sum_{n: \mu_{(n)} < \nu_{(n)}} \mu_{(n)}$$

s.a.m

$$\Rightarrow P = \sum_{n: \mu_{(n)} \geq r_{(n)}} \mu_{(n)} + (1-\alpha) \sum_{n: \mu_{(n)} < r_{(n)}} \mu_{(n)} + \sum_{n: \mu_{(n)} < r_{(n)}} r_{(n)} =$$

$$= 1 + \sum_{n: \mu_{(n)} < r_{(n)}} \{r_{(n)} - \mu_{(n)}\}$$

$$|\mu - r|_{TV} = \sum_{n: \mu_{(n)} < r_{(n)}} (\mu_{(n)} - r_{(n)}) \quad \text{قضیه دوم:}$$

اثبات قضیه دوم: در نظر بگیرید یک مجموعه داده B داریم به صورتی که:

$$B = \{n \mid \mu_{(n)} \geq r_{(n)}\}$$

$$\forall A \subseteq X : \mu(A) - r(A) \leq \mu(A \cap B) - r(A \cap B) \leq \mu(B) - r(B)$$

$\forall n$ (اثبات این دو نامساوی)

$$\textcircled{1}, \forall n \in A \cap \bar{B} \Rightarrow \mu_{(n)} - r_{(n)} \leq 0$$

$$\hookrightarrow \underbrace{(\mu(A) - \mu(A \cap B))}_{\mu(A \cap \bar{B})} - \underbrace{(r(A) - r(A \cap B))}_{r(A \cap \bar{B})} \leq 0$$

$$\textcircled{2} \hookrightarrow \mu(B) (\mu(B) - \mu(A \cap B)) - (r(B) - r(B \cap B)) \geq 0$$

حال اگر دو نامساوی به دست آمده $A=B$ قرار دهیم، داریم:

$$\left. \begin{array}{l} \mu(A) - r(A) = \mu(B) - r(B) \\ \mu(B) - r(B) \geq 0 \end{array} \right\} \Rightarrow |\mu(A) - r(A)| = \mu(B) - r(B)$$

با توجه به اثبات‌های گذشته

$$\|\mu - \nu\|_{TV} = \mu(B) - \nu(B) = \sum_{n \in B} (\mu_{(n)} - \nu_{(n)})$$

$$P = 1 - \|\mu - \nu\|_{TV}$$

حال دو مقدار X و Y را به صورت زیر می‌سازیم:

① با احتمال P Z را با احتمال $\frac{\min\{\mu_{(n)}, \nu_{(n)}\}}{P}$ می‌سازیم و مقادیر دهیم

$$X = Y = Z \quad \rightarrow \quad h_3(Z)$$

② با احتمال $1-P$ X و Y را از بین احتمال‌های زیر می‌سازیم:

$$h_1(X) = \begin{cases} \frac{\mu_{(n)} - \nu_{(n)}}{1-P} & \mu_{(n)} > \nu_{(n)} \\ 0 & \text{o.w.} \end{cases}$$

$$h_2(Y) = \begin{cases} \frac{\nu_{(y)} - \mu_{(y)}}{1-P} & \nu_{(y)} > \mu_{(y)} \\ 0 & \text{o.w.} \end{cases}$$

در نتیجه داریم که:

$$X \sim P h_3(n) + (1-P) h_1 = \mu_{(n)}$$

$$Y \sim P h_3(y) + (1-P) h_2 = \nu_{(n)}$$

در نتیجه (X, Y) یک Coupling است.

پس این X و Y که دقیقاً نام‌های داخل قضیه ① را به‌کار می‌بریم، از این قضیه ① ثابت می‌شود.

ما به این نتایج در مورد می پردازیم. در نظر بگیریم دو توزیع مانند
 $\pi_Q(\cdot | s_t)$ و $\pi_{Q'}(\cdot | s_t)$ داریم. در نتیجه طبق قضایای ثابت کردیم
 $X \sim \pi_Q(\cdot | s_t), Y \sim \pi_{Q'}(\cdot | s_t)$ و X و Y ای کار به صورتی که

$$\Rightarrow P(X \neq Y) = \underbrace{\|\pi_{Q'}(\cdot | s_t) - \pi_Q(\cdot | s_t)\|_{TV}}_{\leq \epsilon}$$

$$P_{Q'}(s_t) = (1 - \epsilon)^t P_Q(s_t) + (1 - (1 - \epsilon)^t) P_{\text{mistake}}(s_t) \rightarrow$$

$$\Rightarrow |P_{Q'}(s_t) - P_Q(s_t)| = (1 - (1 - \epsilon)^t) |P_{\text{mistake}}(s_t) - P_Q(s_t)|$$

$$\leq 2(1 - (1 - \epsilon)^t)$$

$$\epsilon \leq \epsilon \Rightarrow (1 - (1 - \epsilon)^t) \leq (1 - (1 - \epsilon)^t)$$

همین طبق بحث میگویم، داریم که $\epsilon t \gg 1 - (1 - \epsilon)^t$
 در نتیجه در نهایت داریم که

$$|P_{Q'}(s_t) - P_Q(s_t)| \leq 2\epsilon t$$

ب) قضیه ۳: فرض کنید تابع به صورت $h(s_t)$ باشد

$$E_{P_{\theta'}(s_t)} \{h(s_t)\} \geq E_{P_{\theta}(s_t)} \{h(s_t)\} - 2\epsilon t \max_{s_t} |h(s_t)|$$

اثبات قضیه ۳:

$$E_{P_{\theta'}(s_t)} \{h(s_t)\} = \sum_{s_t} P_{\theta'}(s_t) h(s_t) =$$

$$= \sum_{s_t} P_{\theta}(s_t) h(s_t) + \sum_{s_t} (P_{\theta'}(s_t) - P_{\theta}(s_t)) h(s_t)$$

$$\geq E_{P_{\theta}(s_t)} \{h(s_t)\} - |P_{\theta'}(s_t) - P_{\theta}(s_t)| \max_{s_t} |h(s_t)|$$

$$\geq E_{P_{\theta}(s_t)} \{h(s_t)\} - 2\epsilon t \max_{s_t} |h(s_t)|$$

حال به اثبات صورت سؤال می پردازیم:

۱) حالت اول (معروف):

$$h(s_t) = E_{a_t \sim \pi_{\theta'}(\cdot|s_t)} \left\{ \frac{\pi_{\theta'}(a_t|s_t)}{\pi_{\theta}(a_t|s_t)} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \right\}$$

$$\sum_t E_{\theta'} \{h(s_t)\} \geq \sum_t E_{\theta}(h(s_t)) - \sum_{t=0}^T 2\epsilon t \max_{s_t} |h(s_t)|$$

مرتبه اختلاف به این است:

$$\sum_{t=0}^T 2\epsilon t \max_{s_t} |f(s_t)| = \epsilon T(T+1) \times |f(s_t)|$$

با توجه به اینکه $f(s_t)$ تقریباً $O(r_{\max})$ است، در نتیجه مرتبه اختلاف به این است با $O(\epsilon T^2 r_{\min})$ maximum reward

حالت اول نامحدود

$$\begin{aligned} \sum_{t=0}^{\infty} 2\epsilon t \max_{s_t} |f(s_t)| &\leq 2\epsilon \sum_{t=0}^{\infty} t \gamma^t r_{\max} \\ &= 2\epsilon r_{\max} \times \gamma \sum_{t=0}^{\infty} \frac{d}{d\gamma} (\gamma^t) \\ &= 2\epsilon r_{\max} \times \gamma \times \frac{d}{d\gamma} \left(\frac{1}{1-\gamma} \right) = \frac{2\epsilon \gamma r_{\max}}{(1-\gamma)^2} \end{aligned}$$

در این صورت اختلاف $O\left(\frac{2\epsilon \gamma r_{\max}}{(1-\gamma)^2}\right)$

Using Pinsker inequality (2)

$$\|\pi_{\theta'}(a_t|s_t) - \pi_{\theta}(a_t|s_t)\| \leq \sqrt{\frac{1}{2} D_{KL}(\pi_{\theta'}(\cdot|s_t) \parallel \pi_{\theta}(\cdot|s_t))} \Rightarrow$$

$$\Rightarrow D_{KL} \geq 2 \|\pi_{\theta'}(a_t|s_t) - \pi_{\theta}(a_t|s_t)\|^2 \geq 2\epsilon^2$$

$$D = \sum_t \left\{ \frac{1}{P_\theta(s_t)} \left\{ \frac{1}{\pi_\theta(a_t|s_t)} \left\{ \frac{\pi_{\theta'}(a_t|s_t)}{\pi_\theta(a_t|s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right\} \right\} \right\}$$

$$\nabla_{\theta'} D = \sum_t \left\{ \frac{1}{P_\theta(s_t)} \left\{ \frac{1}{\pi_\theta(a_t|s_t)} \left\{ \frac{\nabla_{\theta'} \pi_\theta(a_t|s_t)}{\pi_\theta(a_t|s_t)} \gamma^t A^{\pi_\theta}(s_t, a_t) \right\} \right\} \right\}$$

$$= \sum_t \left\{ \frac{1}{P_\theta(s_t)} \left\{ \frac{1}{\pi_\theta(a_t|s_t)} \left\{ \frac{\pi_\theta(a_t|s_t)}{\pi_\theta(a_t|s_t)} \nabla_{\theta'} \log \pi_\theta(a_t|s_t) \gamma^t A^{\pi_\theta}(s_t, a_t) \right\} \right\} \right\}$$

$$\Rightarrow \nabla_{\theta'} D \Big|_{\theta'=\theta} = \sum_t \frac{\nabla_{\theta'} D}{\pi_\theta}$$

$$\nabla_{\theta'} D \Big|_{\theta'=\theta} = \sum_t \left\{ \frac{1}{P_\theta(s_t)} \left\{ \frac{1}{\pi_\theta(a_t|s_t)} \left\{ \gamma^t \nabla_{\theta'} \log \pi_\theta(a_t|s_t) A^{\pi_\theta}(s_t, a_t) \right\} \right\} \right\}$$

$$D = \sum_t \left\{ \frac{1}{P_\theta(s_t)} \left\{ \frac{1}{\pi_\theta(a_t|s_t)} \left\{ \gamma^t A^{\pi_\theta}(s_t, a_t) \right\} \right\} \right\} + \nabla_{\theta'} D \Big|_{\theta'=\theta}^T (\theta' - \theta)$$

قسم الف نوال

Marker Chains and Mining Times

(سوال 4)

$$\begin{aligned}
 \nabla_{\theta} r^{M_{\theta}}(s) &= \nabla_{\theta} Q^{M_{\theta}}(s, M_{\theta}(s)) = \nabla_{\theta} \left(r(s, M_{\theta}(s)) + \int_S \gamma P(s'|s, M_{\theta}(s)) r^{M_{\theta}}(s') ds' \right) \\
 &= \nabla_{\theta} M_{\theta}(s) \nabla_a r(s, a) \Big|_{a=M_{\theta}(s)} + \nabla_{\theta} \int_S \gamma P(s'|s, M_{\theta}(s)) r^{M_{\theta}}(s') ds' \\
 &= \nabla_{\theta} M_{\theta}(s) \nabla_a r(s, a) \Big|_{a=M_{\theta}(s)} \\
 &\quad + \int_S \gamma \left(P(s'|s, M_{\theta}(s)) \nabla_{\theta} r^{M_{\theta}}(s') + \nabla_{\theta} M_{\theta}(s) \nabla_a P(s'|s, a) \Big|_{a=M_{\theta}(s)} r^{M_{\theta}}(s') \right) ds' \\
 &= \nabla_{\theta} M_{\theta}(s) \nabla_a \left(r(s, a) + \int_S \gamma P(s'|s, a) r^{M_{\theta}}(s') ds' \right) \Big|_{a=M_{\theta}(s)} \\
 &\quad + \int_S \gamma P(s'|s, M_{\theta}(s)) \nabla_{\theta} r^{M_{\theta}}(s') ds' =
 \end{aligned}$$

$$= \nabla_{\theta} M_{\theta}(s) \nabla_a Q^{M_{\theta}}(s, a) \Big|_{a=M_{\theta}(s)} + \int_S \gamma P(s \rightarrow s', 1, M_{\theta}) \nabla_{\theta} r^{M_{\theta}}(s') ds'$$

حال به صورت بازگشتی در این رابطه جایگزینی می کنیم

$$\begin{aligned}
 \nabla_{\theta} r^{M_{\theta}}(s) &= \nabla_{\theta} M_{\theta}(s) \nabla_a Q^{M_{\theta}}(s, a) \Big|_{a=M_{\theta}(s)} \\
 &\quad + \int_S \gamma P(s \rightarrow s', 1, M_{\theta}) \nabla_{\theta} M_{\theta}(s') \nabla_a Q^{M_{\theta}}(s', a) \Big|_{a=M_{\theta}(s')} ds' + \\
 &\quad + \int_S \gamma P(s \rightarrow s', 1, M_{\theta}) \int_S \gamma P(s' \rightarrow s'', 1, M_{\theta}) \nabla_{\theta} r^{M_{\theta}}(s'') ds'' ds'
 \end{aligned}$$

s.a.m

$$\begin{aligned}
&= \nabla_Q M(s) \nabla_a Q^M(s, a) \Big|_{a=M(s)} \\
&+ \int_S \gamma P(s \rightarrow s', 1, M) \nabla_Q M(s') \nabla_a Q^M(s', a) \Big|_{a=M(s')} ds' \\
&+ \int_S \gamma^2 P(s \rightarrow s', 2, M) \nabla_Q V^M(s') ds'
\end{aligned}$$

به همین صورت اگر ادامه دهیم به این می‌رسیم که

$$\nabla_Q V^M(s) = \int_S \sum_{t=0}^{\infty} \gamma^t P(s \rightarrow s', t, M) \nabla_Q M(s') \nabla_a Q^M(s', a) \Big|_{a=M(s')} ds'$$

حال با اِصِدِ ریاض گرفتن از این عبارت می‌رسیم به

$$\begin{aligned}
\nabla_Q J(M) &= \nabla_Q \int_S P_i(s) V^M(s) ds \\
&= \int_S P_i(s) \nabla_Q V^M(s) ds
\end{aligned}$$

$$\begin{aligned}
&= \int_S \int_S \sum_{t=0}^{\infty} \gamma^t P_i(s) P(s \rightarrow s', t, M) \nabla_Q M(s') \nabla_a Q^M(s', a) \Big|_{a=M(s')} ds' ds \\
&= \int_S P^M(s) \nabla_Q M(s) \nabla_a Q^M(s, a) \Big|_{a=M(s)} ds
\end{aligned}$$