

April 19, 2019

## Policy gradient HRL

update  
April 29. $g_i$ : Subgoal @ local  $i$ , Top agent considers local  $i$  subgoal (sub agent) as 'action'. $s_i$ : State @ local  $i$  $\theta_i$ : policy parameter for local  $i$  $\pi_i(g_i | s_i) \doteq \pi_{i\theta_i}(g_i | s_i)$ : policy for local  $i$  $\pi_i(v_i, \theta_{iw}) \doteq \pi_{i\theta_{iw}}(v_i, \theta_{iw})$ : policy for local  $i$  value weight $q_{it}$ : return, Sample of  $q_{\pi_i}(s_i, g_i)$  $\pi_i(v_i, \theta_{iw})$ : Policy for local  $i$  value weightnotation ↑  $\theta_{iw}$ : policy parameter for local  $i$  value weight

$$\begin{aligned} \nabla_{\theta_i} J(\theta_i, \theta_{iw}) &= \frac{\partial}{\partial \theta_i} J(\theta_i, \theta_{iw}) = \sum_{i=1}^N \frac{\partial}{\partial \theta_i} V_{i\pi_{\theta_i}}(s_{i0}) \cdot \pi_i(v_i, \theta_{iw}) = \sum_{i=1}^N \sum_S d^{\pi_i}(s_i) \underbrace{\frac{\partial \pi_i(g_i | s_i)}{\partial \theta_i} q_{\pi_i}(s_i, g_i) \cdot \pi_i(v_i, \theta_{iw})}_{\text{item for updating } \theta_i} \\ &= \sum_{i=1}^N \sum_S d^{\pi_i}(s_i) \underbrace{\sum_g \pi_i(g_i | s_i) \frac{1}{\pi_i(g_i | s_i)} \frac{\partial \pi_i(g_i | s_i)}{\partial \theta_i} q_{\pi_i}(s_i, g_i) \cdot \pi_i(v_i, \theta_{iw})}_{\text{item for updating } \theta_i} \end{aligned}$$

$$= \mathbb{E}_{\substack{\pi \\ \text{global, system}}} \frac{\frac{\partial}{\partial \theta_i} (\ln [\pi_i(g_i | s_i)] q_{\pi_i}(s_i, g_i)) \cdot \pi_i(v_i, \theta_{iw})}{\text{item for updating } \theta_i}$$

$$\therefore \theta_i \leftarrow \theta_i + \alpha \nabla_{\theta_i} J(\theta_i, \theta_{iw})$$

$$\therefore \theta_i \leftarrow \theta_i + \alpha \{ \nabla_{\theta_i} [\ln \pi_i(g_i | s_i)] q_{it} \} \pi_i(v_i, \theta_{iw})$$

$$\nabla_{\theta_{iw}} J(\theta_i, \theta_{iw}) = \frac{\partial}{\partial \theta_{iw}} J(\theta_i, \theta_{iw}) = \sum_{i=1}^N V_{i\pi_{\theta_i}}(s_{i0}) \cdot \frac{\partial}{\partial \theta_{iw}} \pi_i(v_i, \theta_{iw}) = \sum_{i=1}^N V_{i\pi_{\theta_i}}(s_{i0}) \cdot \nabla_{\theta_{iw}} \pi_i(v_i, \theta_{iw})$$

$$= \sum_{i=1}^N \sum_S d^{\pi_i}(s_i) \underbrace{\sum_g \pi_i(g_i | s_i) \{ \pi_i(s_i, g_i) \cdot \nabla_{\theta_{iw}} \pi_i(v_i, \theta_{iw}) \}}_{\text{item for updating } \theta_{iw}}$$

$$= \mathbb{E}_{\pi} [q_{\pi_i}(s_i, g_i) \cdot \nabla_{\theta_{iw}} \pi_i(v_i, \theta_{iw})]$$

$$\therefore \theta_{iw} \leftarrow \theta_{iw} + \beta \nabla_{\theta_{iw}} J(\theta_i, \theta_{iw})$$

$$\therefore \theta_{iw} \leftarrow \theta_{iw} + \beta [q_{it} \cdot \nabla_{\theta_{iw}} \pi_i(v_i, \theta_{iw})]$$

Dimitri Bertsekas, John Tsitsiklis. Neuro-Dynamic Programming. 1996

From the first order Taylor expansion:  $f(r_{t+1}) = f(r_t) + \nabla f(r_t)^T S_t + O(\gamma_t)$

Pg9-94-96  
Gradient  
Methods

$S_t$ : descent direction

↓  
正交

$\nabla f(r_t)^T S_t < 0$  ↓ 降方向  
Transpose?

A technical condition on the direction  $S_t$ :  $\forall t$ ,

$$\rightarrow C_1 \|\nabla f(r_t)\|^2 \leq -\nabla f(r_t)^T S_t, \Rightarrow \|S_t\| \leq C_2 \|\nabla f(r_t)\|, \quad (3.36)$$

↑  
positive scalar

slope of  $f$  at  $r_t$  along the direction  $S_t$

↑ guarantees that the vectors  $S_t$  and  $\nabla f(r_t)$  will not become asymptotically orthogonal.  
即保证  $\nabla f(r_t)^T S_t$  不为 0.

### ① proposition 3.4 [constant stepsize]

proof:  $g(\xi) = f(r + \xi z) \quad \frac{d}{d\xi} g(\xi) = z^T \nabla f(r + \xi z) \quad (\text{chain rule})$

↓  
scalar parameter

$$\begin{aligned} f(r+z) - f(r) &= f(r+1z) - f(r+0z) = g(1) - g(0) = \int_0^1 \frac{d}{d\xi} g(\xi) d\xi = \int_0^1 z^T \nabla f(r+\xi z) d\xi = \int_0^1 [z^T \nabla f(r) + z^T \nabla f(r+\xi z) - z^T \nabla f(r)] d\xi \\ &\leq \int_0^1 z^T \nabla f(r) d\xi + \left| \int_0^1 z^T (\nabla f(r+\xi z) - \nabla f(r)) d\xi \right| \leq z^T \nabla f(r) \int_0^1 d\xi + \int_0^1 \|z\| \cdot \|\nabla f(r+\xi z) - \nabla f(r)\| d\xi \\ &\leq z^T \nabla f(r) + \|z\| \int_0^1 L \xi \|z\| d\xi \quad \xrightarrow{\nabla f \text{ 满足 Lipschitz 连续条件.}} \\ &= z^T \nabla f(r) + \|z\| \cdot \|z\| \cdot L \int_0^1 \xi d\xi = z^T \nabla f(r) + \frac{1}{2} L \|z\|^2 \end{aligned} \quad (3.37)$$

$$\Rightarrow f(r_t + \gamma S_t) - f(r_t) \leq \nabla f(r_t)^T \gamma S_t + \frac{1}{2} L \|\gamma S_t\|^2 = \underbrace{\gamma \nabla f(r_t)^T S_t}_{\downarrow} + \underbrace{\frac{1}{2} \gamma^2 L \|S_t\|^2}_{\downarrow} \quad (3.39)$$

$r_t$ : sequence generated by gradient method  
 $r_{t+1} = r_t + \gamma S_t$

$$\begin{aligned} \Rightarrow f(r_t) - f(r_t + \gamma S_t) &\geq \gamma C_1 \|\nabla f(r_t)\|^2 - [C_2 \|\nabla f(r_t)\|]^2 \frac{1}{2} \gamma^2 L = C_1 \gamma \|\nabla f(r_t)\|^2 - \frac{1}{2} \gamma^2 L C_2^2 \|\nabla f(r_t)\|^2 \\ &= \frac{1}{2} \gamma L C_2^2 \left( \frac{2C_1}{LC_2^2} - \gamma \right) \|\nabla f(r_t)\|^2 \end{aligned} \quad (3.40)$$

monotonically nonincreasing.

$$\Downarrow 0 < \gamma < \frac{2C_1}{LC_2^2} \quad (3.38)$$

$\downarrow$   
 $f(r_t) \rightarrow -\infty$  or finite value, in ②,  $f(r_t) - f(r_{t+1}) \rightarrow 0$ , so, 3.40 implies  $\nabla f(r_t) \rightarrow 0$ .

② proposition 3.5 | Diminishing Step size

$$(3.39) \quad \begin{cases} f(r+z) - f(r) \leq z^T \nabla f(r) + \frac{1}{2} L \|z\|^2 \\ z = \gamma_t s_t \end{cases}$$

$$\Rightarrow f(r_{t+1}) \leq$$

$$(3.36) \quad \begin{cases} c_1 \|\nabla f(r_t)\|^2 \leq -\nabla f(r_t)^T s_t \\ \|s_t\| \leq c_2 \|\nabla f(r_t)\| \end{cases}$$

$$\begin{aligned} & \leq f(r_t) + (\gamma_t s_t)^T \nabla f(r_t) + \frac{1}{2} L \|s_t\|^2 \\ &= f(r_t) + \gamma_t \left( \frac{1}{2} \gamma_t L \|s_t\|^2 - \|\nabla f(r_t)^T s_t\| \right) \\ &\leq f(r_t) - \gamma_t \left( c_1 - \frac{1}{2} \gamma_t^2 L \right) \|\nabla f(r_t)\|^2 \\ &\rightarrow 0 \end{aligned}$$

$\nabla f(r_t)^T s_t < 0$  descent direction

$$\Rightarrow f(r_{t+1}) \leq f(r_t) - \gamma_t c_1 \|\nabla f(r_t)\|^2 \quad (3.42)$$

$$t \geq t' \quad \sum_{t=t'}^{\infty} \gamma_t \|\nabla f(r_t)\|^2 \leq f(r_{t'}) - \lim_{t \rightarrow \infty} f(r_t) < \infty$$

$$\therefore \sum_{t=t'}^{\infty} \gamma_t = \infty \quad \therefore \|\nabla f(r_t)\|^2 > \epsilon \forall t > t', \text{ should be incorrect}$$

$$\Rightarrow \lim_{i \rightarrow \infty} \|\nabla f(r_i)\| = 0$$

反证法：假設  $\lim_{i \rightarrow \infty} \|\nabla f(r_i)\| > 0$ ,  $\exists \epsilon > 0$ . let  $\|\nabla f(r_i)\| < \frac{1}{2}\epsilon$ ,  $\|\nabla f(r_{i+1})\| > \epsilon$ ,  $i < i < i+1$

implying that  $\nabla f(\bar{r}) = 0$   
(Q.E.D.).

請用  
證明

$$\begin{cases} \|\nabla f(r_t)\| < \frac{1}{2}\epsilon, \quad \|\nabla f(r_{i+1})\| > \epsilon, \\ \frac{1}{2}\epsilon \leq \|\nabla f(r_i)\| \leq \epsilon, \quad t < i < i+1 \end{cases}$$

$$\therefore \|\nabla f(r_{t+1})\| - \|\nabla f(r_t)\| \leq \|\nabla f(r_{t+1}) - \nabla f(r_t)\| \leq L \|r_{t+1} - r_t\| = \gamma_t L \|s_t\| \leq \gamma_t L c_2 \|\nabla f(r_t)\|,$$

滿足 Lipschitz 連續條件

subset of integers  
 $T, t \in T$

$$\therefore \|s_t\| \leq c_2 \|\nabla f(r_t)\|$$

$$\|\nabla f(r) - \nabla f(\bar{r})\| \leq L \|r - \bar{r}\|; \text{ Lipschitz condition (3.4)}$$

$$\text{if } i = t, t+1, \dots, i-1 \\ \|\nabla f(r_i)\| \geq \frac{\epsilon}{4}$$

$$\Rightarrow f(r_{i+1}) \leq f(r_t) - C \left( \frac{\epsilon}{4} \right) \sum_{j=t}^{i-1} \gamma_j$$

$$\begin{aligned} & \xrightarrow{\text{兩端}} \frac{1}{2Lc_2} \leq \sum_{j=t}^{i-1} \gamma_j \quad (3.43) \\ & \text{因為 } \text{① } \nabla f(r_i) \text{ 有極值} \\ & \Rightarrow t \rightarrow \infty, t \in T, = 0 \end{aligned}$$

$$\therefore \frac{\epsilon}{2} \leq \|\nabla f(r_{i+1})\| - \|\nabla f(r_t)\| \leq \|\nabla f(r_{i+1}) - \nabla f(r_t)\|$$

$$\leq L \|r_{i+1} - r_t\| \leq L \sum_{j=t}^{i-1} \gamma_j \|s_j\|$$

$$= L c_2 \sum_{j=t}^{i-1} \gamma_j \|\nabla f(r_j)\| \leq L c_2 \epsilon \sum_{j=t}^{i-1} \gamma_j$$

April 27, 2019

$\hat{\pi}_{\text{local}}(s_{it}, g_{it})$ : local  $i$  is some unbiased estimator of  $q_{\pi_i}(s_{it}, a_{it})$ .

[ $q_{\pi_i}$ ] let  $f_{x_i}: S \times A \rightarrow \mathbb{R}$  be our approximation to  $q_{\pi_i}$ , with parameter  $x_i$ . It is natural to learn  $f_{x_i}$  by following  $\pi_i(s_i | s_i)$  and updating  $x_i$  by a rule such as

$$\begin{aligned} \Delta x_{it} &\propto \frac{\partial}{\partial x_i} [\hat{q}_{\pi_i}(s_{it}, g_{it}) - f_{x_i}(s_{it}, g_{it})]^2 \\ &\propto [\hat{q}_{\pi_i}(s_{it}, g_{it}) - f_{x_i}(s_{it}, g_{it})] \frac{\partial f_{x_i}(s_{it}, g_{it})}{\partial x_i} \end{aligned}$$

When such a process has converged to a local optimum, then

$$\sum_{i=1}^N \sum_{s} d^{x_i}(s_i) \sum_g \pi_i(g_i | s_i) [q_{\pi_i}(s_i, g_i) - f_{x_i}(s_i, g_i)] \frac{\partial f_{x_i}(s_{it}, g_{it})}{\partial x_i} \cdot \pi_i(v_i, \theta_{iw}) = 0 \quad \textcircled{1}$$

$f_{x_i}$  satisfies \textcircled{1} and is compatible with the policy parameterization in the sense that

$$\frac{\partial f_{x_i}(s_i, g_i)}{\partial x_i} = \frac{\partial \pi_i(g_i | s_i)}{\partial \theta_i} \frac{1}{\pi_i(g_i | s_i)} \quad \textcircled{2}$$

$$\textcircled{1} \textcircled{2} \Rightarrow \sum_{i=1}^N \sum_{s} d^{x_i}(s_i) \sum_g \pi_i(g_i | s_i) [q_{\pi_i}(s_i, g_i) - f_{x_i}(s_i, g_i)] \frac{\partial \pi_i(g_i | s_i)}{\partial \theta_i} \frac{1}{\pi_i(g_i | s_i)} \cdot \pi_i(v_i, \theta_{iw}) = 0$$

$$\sum_{i=1}^N \sum_{s} d^{x_i}(s_i) \sum_g \frac{\partial \pi_i(g_i | s_i)}{\partial \theta_i} [q_{\pi_i}(s_i, g_i) - f_{x_i}(s_i, g_i)] \cdot \pi_i(v_i, \theta_{iw}) = 0 \quad \textcircled{3}$$

which tells us that the error in  $f_{x_i}(s_i, g_i)$  is orthogonal to the gradient of the policy parameterization

$$\begin{aligned} \frac{\partial}{\partial \theta_i} J(\theta_i, \theta_{iw}) &= \sum_{i=1}^N \sum_{s} d^{x_i}(s_i) \sum_g \frac{\partial \pi_i(g_i | s_i)}{\partial \theta_i} q_{\pi_i}(s_i, g_i) \cdot \pi_i(v_i, \theta_{iw}) - \sum_{i=1}^N \sum_{s} d^{x_i}(s_i) \sum_g \frac{\partial \pi_i(g_i | s_i)}{\partial \theta_i} [q_{\pi_i}(s_i, g_i) - f_{x_i}(s_i, g_i)] \pi_i(v_i, \theta_{iw}) \\ &= \sum_{i=1}^N \sum_{s} d^{x_i}(s_i) \sum_g \frac{\partial \pi_i(g_i | s_i)}{\partial \theta_i} [q_{\pi_i}(s_i, g_i) - q_{\pi_i}(s_i, g_i) + f_{x_i}(s_i, g_i)] \pi_i(v_i, \theta_{iw}) \\ &= \sum_{i=1}^N \sum_{s} d^{x_i}(s_i) \sum_g \frac{\partial \pi_i(g_i | s_i)}{\partial \theta_i} f_{x_i}(s_i, g_i) \pi_i(v_i, \theta_{iw}) \\ &= \sum_{i=1}^N \sum_{s} d^{x_i}(s_i) \sum_g \pi_i(g_i | s_i) \underbrace{\frac{1}{\pi_i(g_i | s_i)} \frac{\partial \pi_i(g_i | s_i)}{\partial \theta_i}}_{\substack{\text{global system} \\ \mathbb{E}_{\pi}}} f_{x_i}(s_i, g_i) \pi_i(v_i, \theta_{iw}) \\ &= \underbrace{\mathbb{E}_{\pi}}_{\substack{\text{global system}}} \frac{\partial}{\partial \theta_i} \{ \ln [\pi_i(g_i | s_i)] f_{x_i}(s_i, g_i) \} \cdot \pi_i(v_i, \theta_{iw}) \end{aligned}$$

(April 28, 2019). Convergence  $\theta_i, \theta_{iw}$ .  $J(\theta_i, \theta_{iw}) = \sum_{i=1}^N \sum_{s=1}^{T_i} d(s_i) \sum_g \pi_i(g_i | s_i) q_{\pi_i}(s_i, g_i) \cdot \bar{\pi}_i(v_i, \theta_{iw})$

from April 19, 2019, we know:  $\nabla_{\theta_i} J(\theta_i, \theta_{iw}) = \sum_{i=1}^N \sum_{s=1}^{T_i} d(s_i) \sum_g \pi_i(g_i | s_i) \underbrace{\frac{\partial}{\partial \theta_i} \{ \ln[\pi_i(g_i | s_i)] q_{\pi_i}(s_i, g_i) \}}_{\mathbb{E}_{\pi}} \cdot \bar{\pi}_i(v_i, \theta_{iw})$

$$\nabla_{\theta_{iw}} J(\theta_i, \theta_{iw}) = \sum_{i=1}^N \sum_{s=1}^{T_i} d(s_i) \sum_g \pi_i(g_i | s_i) \underbrace{[q_{\pi_i}(s_i, g_i) \cdot \nabla_{\theta_{iw}} \pi_i(v_i, \theta_{iw})]}_{\mathbb{E}_{\pi}}$$

Based on Mar. 25, 2019.

Dimitri Bertsekas, John Tsitsiklis, Neuro-Dynamic Programming, 1996 Pg 9-94-96 Gradient Methods

From the first order Taylor expansion:  $J(\theta_{i+1}, \theta_{iw}) = J(\theta_i, \theta_{iw}) + \gamma_t \nabla_{\theta_i} J(\theta_i, \theta_{iw}) S_t + O(\gamma_t)$ .

$\underbrace{\nabla_{\theta_i} J(\theta_i, \theta_{iw}) S_t}_{\text{descent direction}} < 0$

and

and

$$J(\theta_i, \theta_{iw+1}) = J(\theta_i, \theta_{iw}) + \gamma_{wt} \nabla_{\theta_{iw}} J(\theta_i, \theta_{iw}) S_{wt} + O(\gamma_{wt})$$

A technical condition on the direction  $S_t, S_{wt}$ ,

$$\begin{aligned} \theta_i & \left\{ \begin{array}{l} C_1 \|\nabla_{\theta_i} J(\theta_i, \theta_{iw})\|^2 \leq -\nabla_{\theta_i} J(\theta_i, \theta_{iw}) S_t \\ \|S_t\| \leq C_2 \|\nabla_{\theta_i} J(\theta_i, \theta_{iw})\| \end{array} \right. \\ (3.36) \quad \theta_{iw} & \left\{ \begin{array}{l} C_{w1} \|\nabla_{\theta_{iw}} J(\theta_i, \theta_{iw})\|^2 \leq -\nabla_{\theta_{iw}} J(\theta_i, \theta_{iw}) S_{wt} \\ \|S_{wt}\| \leq C_{w2} \|\nabla_{\theta_{iw}} J(\theta_i, \theta_{iw})\| \end{array} \right. \end{aligned}$$

① proposition 3.4  
Constant Stepsize  $|\theta_i| \quad g(\xi) = J(\theta_i + \xi z, \theta_{iw}) \Rightarrow \frac{d}{d\xi} g(\xi) = z^T \nabla J(\theta_i + \xi z, \theta_{iw})$

$$\begin{aligned} J(\theta_i + z, \theta_{iw}) - J(\theta_i, \theta_{iw}) &= J(\theta_i + 1z, \theta_{iw}) - J(\theta_i + 0z, \theta_{iw}) \\ &= g(1) - g(0) \\ &= \int_0^1 \frac{d}{d\xi} g(\xi) d\xi = \int_0^1 z^T \nabla J(\theta_i + \xi z, \theta_{iw}) d\xi = \int_0^1 d\xi [z^T \nabla J(\theta_i, \theta_{iw}) + z^T \nabla J(\theta_i + \xi z, \theta_{iw}) - z^T \nabla J(\theta_i, \theta_{iw})] \\ &\leq \int_0^1 z^T \nabla J(\theta_i, \theta_{iw}) d\xi + \left| \int_0^1 z^T [\nabla J(\theta_i + \xi z, \theta_{iw}) - \nabla J(\theta_i, \theta_{iw})] d\xi \right| \\ &\leq z^T \nabla J(\theta_i, \theta_{iw}) \int_0^1 d\xi + \int_0^1 \|z\| \cdot \|\nabla J(\theta_i + \xi z, \theta_{iw}) - \nabla J(\theta_i, \theta_{iw})\| d\xi \\ &\hookrightarrow \leq L \cdot \|\theta_i + \xi z - \theta_i\| \quad (3.37) \end{aligned}$$

(cont'd)

$$\text{let's consider } \nabla_{\theta_i} J(\theta_i, \theta_{iw}) = \sum_{i=1}^N \sum_S d^{(S)}(S_i) \sum_S \pi_i(g_i | S_i) \frac{\partial}{\partial \theta_i} \left\{ \ln [\pi_i(g_i | S_i)] q_{\pi_i}(S_i, g_i) \right\} \cdot \pi_i(v_i, \theta_{iw})$$

$\nabla_{\theta_i} J(\theta_i, \theta_{iw})$  satisfy Lipschitz continuity

$\nabla J(\theta_i + \xi z, \theta_{iw})$  satisfy Lipschitz continuity, too.

Then,

$$\begin{aligned} J(\theta_i + z, \theta_{iw}) - J(\theta_i, \theta_{iw}) &\leq z^T \nabla J(\theta_i, \theta_{iw}) \int_0^1 ds + \int_0^1 \|z\| \cdot \|\nabla J(\theta_i + \xi z, \theta_{iw}) - \nabla J(\theta_i, \theta_{iw})\| d\xi \\ &\leq z^T \nabla J(\theta_i, \theta_{iw}) + \|z\| \int_0^1 L \cdot \|\nabla J(\theta_i + \xi z, \theta_i)\| d\xi \\ &= z^T \nabla J(\theta_i, \theta_{iw}) + \|z\| \int_0^1 L \cdot \|z\| d\xi \\ &= z^T \nabla J(\theta_i, \theta_{iw}) + \|z\| \cdot \|z\| \cdot L \int_0^1 \xi d\xi = z^T \nabla J(\theta_i, \theta_{iw}) + \frac{1}{2} L \|z\|^2 \quad (3.39) \end{aligned}$$

$$\Rightarrow J(\theta_i + z, \theta_{iw}) - J(\theta_i, \theta_{iw}) \leq z^T \nabla J(\theta_i, \theta_{iw}) + \frac{1}{2} L \|z\|^2 \quad (3.39)$$

$$\begin{aligned} \downarrow J(\theta_{it} + \gamma s_t, \theta_{iw}) - J(\theta_{it}, \theta_{iw}) &\leq \nabla J(\theta_{it}, \theta_{iw})^T s_t + \frac{1}{2} L \|s_t\|^2 = \gamma \nabla J(\theta_{it}, \theta_{iw})^T s_t + \frac{1}{2} \gamma^2 L \|s_t\|^2 \\ &\quad \text{in direction } \nabla J(\theta_{it}, \theta_{iw}) \\ &\quad [ \theta_{it}: \text{sequence generated by gradient method } \theta_{it+1} = \theta_{it} + \gamma s_t ] \quad - C_1 \|\nabla J(\theta_{it}, \theta_{iw})\|^2 \quad C_2 \|\nabla J(\theta_{it}, \theta_{iw})\| \quad (3.36) \end{aligned}$$

$$\Rightarrow J(\theta_{it}, \theta_{iw}) - J(\theta_{it} + \gamma s_t, \theta_{iw}) \geq \gamma C_1 \|\nabla J(\theta_{it}, \theta_{iw})\|^2 - [C_2 \|\nabla J(\theta_{it}, \theta_{iw})\|]^2 \frac{L}{2} \gamma^2$$

$$= C_1 \gamma \|\nabla J(\theta_{it}, \theta_{iw})\|^2 - \frac{1}{2} \gamma^2 L C_2 \|\nabla J(\theta_{it}, \theta_{iw})\|^2$$

$$= \frac{1}{2} \gamma L^2 \left( \frac{2C_1}{LC_2} - \gamma \right) \|\nabla J(\theta_{it}, \theta_{iw})\|^2 \quad (3.40)$$

$$\Rightarrow 0 < \gamma < \frac{2C_1}{LC_2} \quad (3.38)$$

$\geq 0$

$\Rightarrow J(\theta_{it}, \theta_{iw})$  is monotonically nonincreasing

$\Rightarrow J(\theta_{it}, \theta_{iw}) \rightarrow -\infty$  or finite value, in  $\Theta$ ,  $J(\theta_{it}, \theta_{iw}) \rightarrow 0$ , so 3.40 implies  $\nabla J(\theta_{it}, \theta_{iw}) \rightarrow 0$

April 28, 2019

Page 2

② Minimising Step size  $\|\theta\|$

$$(3.39) \quad \begin{cases} J(\theta_i + z, \theta_{iw}) - J(\theta_i, \theta_{iw}) \leq z^T \nabla J(\theta_i) + \frac{1}{2} L \|z\|^2 \\ z = \gamma_t s_t \end{cases} \quad (3.36) \quad \begin{cases} C_1 \|\nabla J(\theta_i, \theta_{iw})\|^2 \leq -[\nabla J(\theta_i, \theta_{iw})^T s_t] \\ \|s_t\| \leq C_2 \|\nabla J(\theta_i, \theta_{iw})\| \end{cases}$$

$\leftarrow$  descent direction

$$\Rightarrow J(\theta_{it+1}, \theta_{iw}) \leq J(\theta_{it}, \theta_{iw}) + (\gamma_t s_t)^T \nabla J(\theta_i) + \frac{1}{2} L \|s_t\|^2 = J(\theta_{it}, \theta_{iw}) + \gamma_t \left( \frac{1}{2} \gamma_t L \|s_t\|^2 - |\nabla J(\theta_{it}, \theta_{iw})^T s_t| \right)$$

$$\leq J(\theta_{it}, \theta_{iw}) - \gamma_t (C_1 - \frac{1}{2} \gamma_t C_2 L) \|\nabla J(\theta_{it}, \theta_{iw})\|^2$$

$\downarrow \gamma_t > 0$

when  $t > \bar{t}$

$$\Rightarrow J(\theta_{it+1}, \theta_{iw}) \leq J(\theta_{it}, \theta_{iw}) - \gamma_t C_1 \|\nabla J(\theta_{it}, \theta_{iw})\|^2 \quad (3.42)$$

$\forall t \geq \bar{t}$ , By adding 3.42 over all  $t \geq \bar{t}$ , we obtain

$$\sum_{t=\bar{t}}^{\infty} \gamma_t \|\nabla J(\theta_{it}, \theta_{iw})\|^2 \leq J(\theta_{i\bar{t}}, \theta_{iw}) - \lim_{t \rightarrow \infty} J(\theta_{it}, \theta_{iw}) < \infty$$

$\downarrow$  V.S. contradict

$$\sum_{t=0}^{\infty} \gamma_t = \infty, \quad \|\nabla J(\theta_{it}, \theta_{iw})\|^2 > \epsilon \quad \text{for all } t \geq \bar{t}$$

$\left\{ \begin{array}{l} \lim_{t \rightarrow \infty} \|\nabla J(\theta_{it}, \theta_{iw})\| = 0 \\ \liminf_{t \rightarrow \infty} \|\nabla J(\theta_{it}, \theta_{iw})\| > 0 \end{array} \right.$

Then, to show  $\lim_{t \rightarrow \infty} \|\nabla J(\theta_{it}, \theta_{iw})\| = 0$  反证法:

Assume the contrary, that is,  $\limsup_{t \rightarrow \infty} \|\nabla J(\theta_{it}, \theta_{iw})\| > 0$ .

Then,  $\exists \epsilon > 0$ , let  $\begin{cases} \|\nabla J(\theta_{it}, \theta_{iw})\| < \frac{\epsilon}{2}, \text{ for inf many } t \\ \|\nabla J(\theta_{it}, \theta_{iw})\| > \epsilon, \text{ for inf many } t \\ \frac{\epsilon}{2} \leq \|\nabla J(\theta_{it}, \theta_{iw})\| < \epsilon, \text{ if } t < i < i(t) \end{cases}$

since  $\|\nabla J(\theta_{i+1}, \theta_{iw})\| - \|\nabla J(\theta_i, \theta_{iw})\| \leq \|\nabla J(\theta_{i+1}, \theta_{iw}) - \nabla J(\theta_i, \theta_{iw})\| \leq L \|\theta_{i+1} - \theta_i\|$

Lipschitz continuity

$$= \gamma_t L \|s_t\| \leq \gamma_t L C_2 \|\nabla J(\theta_{it}, \theta_{iw})\|$$

all  $t \leq \bar{t}$  so that  $\gamma_t < 1$

$\Rightarrow$  也即是在  $\forall t \in \mathbb{R}$  时, 可使:

$$\|\nabla J(\theta_{it+1}, \theta_{iw})\| \leq 2 \|\nabla J(\theta_{it}, \theta_{iw})\|$$

$$\therefore \|S_t\| \leq C_2 \|\nabla J(\theta_{i,t}, \theta_{iw})\|$$

$$\left\{ \|\nabla J(\theta_i, \theta_{iw}) - \nabla J(\bar{\theta}_i, \theta_{iw})\| \leq L \|\theta_i - \bar{\theta}_i\|; \text{ Lipschitz condition (3.41)} \right\}$$

$$\therefore \frac{\epsilon}{2} \leq \|\nabla J(\theta_{i,t}, \theta_{iw})\| - \|\nabla J(\theta_i, \theta_{iw})\| \leq \|\nabla J(\theta_{i,t}, \theta_{iw}) - \nabla J(\theta_{i,t}, \theta_{iw})\|$$

$$\leq L \|\theta_{i,t+1} - \theta_{i,t}\| \leq L \sum_{i=t}^{i(t)-1} \gamma_i \|S_i\|$$

$$= L C_2 \sum_{i=t}^{i(t)-1} \gamma_i \|\nabla J(\theta_{i,t}, \theta_{iw})\| \leq L C_2 \in \sum_{i=t}^{i(t)-1} \gamma_i$$

$$\Rightarrow \frac{1}{2L C_2} \leq \sum_{i=t}^{i(t)-1} \gamma_i \quad (3.43) \quad \text{If } \sum_{i=t}^{i(t)-1} \gamma_i \geq \frac{1}{2L C_2} > 0 \leftarrow$$

$$\|\nabla J(\theta_{i,t+1}, \theta_{iw})\| \leq 2 \|\nabla J(\theta_{i,t}, \theta_{iw})\|$$

$$\|\nabla J(\theta_{i,t+1}, \theta_{iw})\| \geq \frac{\epsilon}{2}$$

$$(3.42) \quad J(\theta_{i,t+1}, \theta_{iw}) \leq J(\theta_{i,t}, \theta_{iw}) - \gamma_t C \|\nabla J(\theta_{i,t}, \theta_{iw})\|^2$$

$$J(\theta_{i,t+1}, \theta_{iw}) \leq J(\theta_{i,t}, \theta_{iw}) - C \left( \frac{\epsilon}{4} \right)^2 \sum_{i=t}^{i(t)-1} \gamma_i$$

$\therefore \Omega, J(\theta_{it}, \theta_{iw})$  converge to a finite value,

$$\lim_{\substack{t \rightarrow \infty \\ t \in \mathbb{N}}} \sum_{i=t}^{i(t)-1} \gamma_i = 0 \leftarrow$$

$$\text{So, } \lim_{t \rightarrow \infty} \nabla J(\theta_{i,t}, \theta_{iw}) = 0.$$

Finally, if  $\bar{\theta}_i$  is a limit point of  $\theta_{i,t}$ , then  $J(\theta_{i,t}, \theta_{iw})$  converges to the finite value  $J(\bar{\theta}_i, \theta_{iw})$ .

Thus we have  $\nabla J(\theta_{i,t}, \theta_{iw}) \rightarrow 0$ , implying that  $\nabla J(\bar{\theta}_i, \theta_{iw}) = 0$ . Q.E.D.

contradicting.  
矛盾

$$(3.36.) \quad \theta_i : \begin{cases} C_1 \| \nabla_{\theta_i} J(\theta_{it}, \theta_{iw}) \|^2 \leq - \nabla_{\theta_i} J^T(\theta_{it}, \theta_{iw}) s_t \\ \| s_t \| \leq C_2 \| \nabla_{\theta_i} J(\theta_{it}, \theta_{iw}) \| \end{cases}$$

$$\nabla_{\theta_i} J(\theta_i, \theta_{iw}) = \sum_{i=1}^N \sum_S d^{\pi_i}(s_i) \sum_g \pi_i(g; s_i) \frac{\partial}{\partial \theta_i} \{ \ln [\pi_i(g; s_i)] q_{\pi_i}(s_i, g_i) \} \cdot \pi_i(v_i, \theta_{iw})$$

$$\text{technical condition : } \begin{cases} C_1 \left\| \sum_{i=1}^N \sum_S d^{\pi_i}(s_i) \sum_g \pi_i(g; s_i) \frac{\partial}{\partial \theta_{it}} \{ \ln [\pi_i(g; s_i)] q_{\pi_i}(s_i, g_i) \} \cdot \pi_i(v_i, \theta_{iw}) \right\|^2 \leq \\ - \sum_{i=1}^N \sum_S d^{\pi_i}(s_i) \sum_g \pi_i(g; s_i) \frac{\partial}{\partial \theta_{it}} \{ \ln [\pi_i(g; s_i)] q_{\pi_i}(s_i, g_i) \} \cdot \pi_i(v_i, \theta_{iw}) \cdot s_t \\ \| s_t \| \leq C_2 \left\| \sum_{i=1}^N \sum_S d^{\pi_i}(s_i) \sum_g \pi_i(g; s_i) \frac{\partial}{\partial \theta_{it}} \{ \ln [\pi_i(g; s_i)] q_{\pi_i}(s_i, g_i) \} \cdot \pi_i(v_i, \theta_{iw}) \right\| \end{cases}$$

$$(3.36) \quad \theta_{iw} : \begin{cases} C_{w1} \| \nabla_{\theta_{iw}} J(\theta_i, \theta_{iw}) \|^2 \leq - \nabla_{\theta_{iw}} J^T(\theta_i, \theta_{iw}) s_{wt} \\ \| s_{wt} \| \leq C_{w2} \| \nabla_{\theta_{iw}} J(\theta_i, \theta_{iw}) \| \end{cases}$$

$$\nabla_{\theta_{iw}} J(\theta_i, \theta_{iw}) = \sum_{i=1}^N \sum_S d^{\pi_i}(s_i) \sum_g \pi_i(g; s_i) [q_{\pi_i}(s_i, g_i) \nabla_{\theta_{iw}} \pi_i(v_i, \theta_{iw})]$$

$$\text{technical condition : } \begin{cases} C_{w1} \left\| \sum_{i=1}^N \sum_S d^{\pi_i}(s_i) \sum_g \pi_i(g; s_i) [q_{\pi_i}(s_i, g_i) \nabla_{\theta_{iw}} \pi_i(v_i, \theta_{iw})] \right\|^2 \leq \\ - \sum_{i=1}^N \sum_S d^{\pi_i}(s_i) \sum_g \pi_i(g; s_i) [q_{\pi_i}(s_i, g_i) \nabla_{\theta_{iw}} \pi_i(v_i, \theta_{iw})] s_{wt} \\ \| s_{wt} \| \leq C_{w2} \left\| \sum_{i=1}^N \sum_S d^{\pi_i}(s_i) \sum_g \pi_i(g; s_i) [q_{\pi_i}(s_i, g_i) \nabla_{\theta_{iw}} \pi_i(v_i, \theta_{iw})] \right\| \end{cases}$$

For  $\boxed{\theta_{iw}}$ ,  $\lim_{t \rightarrow \infty} \nabla J(\theta_i, \theta_{iw}) = 0$  proof is similar to  $\boxed{\theta_i}$ ,  $\lim_{t \rightarrow \infty} \nabla J(\theta_{it}, \theta_{iw}) = 0$  when  $\theta_i$  or  $\theta_{iw}$  satisfies (3.36)  
 $\nabla J(\theta_i, \theta_{iw})$  satisfies Lipschitz continuity