

## Core assumptions

1.  $\pi(\cdot | s) : \mathbb{R}^n \times \mathbb{R}^d \rightarrow [0, 1]^m$  is continuously differentiable and Lipschitz smooth w.r.t. L2 norm  
 $\Rightarrow \|D_\theta \pi(a | s)\|_2 \leq B_\pi$
2.  $f(s, a) : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$  is cont. differentiable and Lipschitz smooth in both arguments w.r.t. L2 norm  
 $\Rightarrow \|D_s f(s, a)\|_2 \leq B_f \quad \|D_a f(s, a)\|_2 \leq B_f$
3.  $r(s, a) : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$  is cont. differentiable and Lipschitz smooth in both arguments w.r.t. L2 norm  
 $\Rightarrow \|D_s r(s, a)\|_2 \leq B_r \quad \|D_a r(s, a)\|_2 \leq B_r$
4. Zeroth-order gradients use the baseline  
 $\hat{D}_\theta^{[0]} J(\theta) = \sum_{n=1}^N (r(s_n, a_n) - b) D_\theta \log \pi_\theta(a_n | s_n)$

In the paper we define and bound

$$B = \|\mathbb{E}[\hat{D}_\theta^{[0]} J(\theta)] - \mathbb{E}[\hat{D}_\theta^{[0]} J(\theta)]\|$$

as pointed out by Reviewer 6 LUU,  
 $B = 0$ . Instead we want to bound  
the sample error

$$B = \|\hat{D}_\theta^{[0]} J(\theta) - \bar{D}_\theta^{[0]} J(\theta)\|_2$$

$$= \left\| \frac{1}{N} \sum_{i=1}^N \hat{D}_\theta^{[0]} J_i(\theta) - \frac{1}{N} \sum_{i=1}^N \bar{D}_\theta^{[0]} J_i(\theta) \right\|_2$$

$$= \frac{1}{N} \left\| \sum_{i=1}^N (\hat{D}_\theta^{[0]} J_i(\theta) - \bar{D}_\theta^{[0]} J_i(\theta)) \right\|_2 \quad (1)$$

1. drop sample subscript for simplicity

2. assume  $a \sim \pi_\theta(\cdot | s)$  for easier derivation  
where  $\text{sg}(\cdot)$  is stop gradient operation (Ma et al. 2024)

Aside: Why do this?

Example with  $r(s_3, a_3)$ . The assumption crosses out the red terms!

$$\begin{aligned} D_\theta r(s_3, a_3) &= D_{a_3} r(s_3, a_3) D_\theta \pi(a_3 | s_3) \\ &\quad + D_{a_3} r(s_3, a_3) P_{s_2} f(s_2, a_2) D_{a_2} f(s_2, a_2) D_\theta \pi(a_2 | s_2) \\ &\quad + D_{a_3} r(s_3, a_3) D_{a_2} f(s_2, a_2) D_\theta \pi(a_2 | s_2) \\ &\quad + D_{a_3} r(s_3, a_3) P_{s_2} f(s_2, a_2) D_\theta \pi(a_2 | s_2) D_{a_2} f(s_2, a_2) D_\theta \pi(a_2 | s_2) \\ &\quad + D_{s_3} r(s_3, a_3) D_{s_2} f(s_2, a_2) D_{a_2} f(s_2, a_2) D_\theta \pi(a_2 | s_2) \\ &\quad + D_{s_3} r(s_3, a_3) D_{a_2} f(s_2, a_2) D_\theta \pi(a_2 | s_2) \\ &\quad + D_{s_3} r(s_3, a_3) P_{a_2} f(s_2, a_2) D_\theta \pi(a_2 | s_2) \end{aligned}$$

Issues wrt Reviewer 4 (E239)

1. Expanding of  $\hat{D}_\theta^{[0]} J(\theta)$  is incorrect
2. We implicitly assumed  $a_t \sim \pi_\theta(\cdot | s_t)$  but never stated it

$$\hat{D}_\theta^{[0]} J(\theta) = \hat{D}_\theta^{[0]} J(\theta)$$

$$= D_\theta \sum_{n=1}^N r(s_n, a_n) - \sum_{n=1}^N (r(s_n, a_n) - b) D_\theta \log \pi_\theta(a_n | s_n)$$

$$= \sum_{n=1}^N D_{a_n} r(s_n, a_n) D_\theta \pi(a_n | s_n) + \left( \sum_{n=1}^N \left( \prod_{i=1}^{n-1} D_{s_i} f(s_i, a_i) \right) D_{a_n} f(s_n, a_n)^\top D_\theta \pi(a_n | s_n) \right)^\top P_{s_n} r(s_n, a_n)$$

$$- \sum_{n=1}^N (r(s_n, a_n) - b) D_\theta \log \pi_\theta(a_n | s_n)$$

$$= \sum_{n=1}^N D_\theta \pi_\theta(a_n | s_n)^\top \left( D_{a_n} r(s_n, a_n) - (r(s_n, a_n) - b) \pi_\theta(a_n | s_n)^{o-1} \right)$$

$$+ \left( \sum_{n=1}^N \left( \prod_{i=1}^{n-1} D_{s_i} f(s_i, a_i) \right) D_{a_n} f(s_n, a_n)^\top D_\theta \pi(a_n | s_n) \right)^\top P_{s_n} r(s_n, a_n)$$

$$\leq \sum_{n=1}^N D_\theta \pi_\theta(a_n | s_n)^\top \left( D_{a_n} r(s_n, a_n) - (D_\theta r(s_n, a_n) - \frac{B_r}{2} \pi_\theta(a_n | s_n)) \right)$$

$$+ \left( \sum_{n=1}^N \left( \prod_{i=1}^{n-1} D_{s_i} f(s_i, a_i) \right) D_{a_n} f(s_n, a_n)^\top D_\theta \pi(a_n | s_n) \right)^\top P_{s_n} r(s_n, a_n) \quad \text{apply (2)}$$

$$= \frac{B_r}{2} \sum_{n=1}^N D_\theta \pi_\theta(a_n | s_n)^\top \pi_\theta(a_n | s_n) + \left( \sum_{n=1}^N \left( \prod_{i=1}^{n-1} D_{s_i} f(s_i, a_i) \right) D_{a_n} f(s_n, a_n)^\top D_\theta \pi(a_n | s_n) \right)^\top P_{s_n} r(s_n, a_n) \quad (3)$$

Plug (3) into (1)

$$B \leq \|\hat{D}_\theta^{[0]} J(\theta) - \bar{D}_\theta^{[0]} J(\theta)\|$$

$$= \left\| \frac{B_r}{2} \sum_{n=1}^N D_\theta \pi_\theta(a_n | s_n)^\top \pi_\theta(a_n | s_n) + \left( \sum_{n=1}^N \left( \prod_{i=1}^{n-1} D_{s_i} f(s_i, a_i) \right) D_{a_n} f(s_n, a_n)^\top D_\theta \pi(a_n | s_n) \right)^\top P_{s_n} r(s_n, a_n) \right\|_2$$

$$\leq \left\| \frac{B_r}{2} \sum_{n=1}^N D_\theta \pi_\theta(a_n | s_n)^\top \pi_\theta(a_n | s_n) \right\|_2$$

$$+ \left\| \left( \sum_{n=1}^N \left( \prod_{i=1}^{n-1} D_{s_i} f(s_i, a_i) \right) D_{a_n} f(s_n, a_n)^\top D_\theta \pi(a_n | s_n) \right)^\top P_{s_n} r(s_n, a_n) \right\|_2$$

$$\leq \frac{1}{2} K B_r B_\pi + \left\| \left( \sum_{n=1}^N \left( \prod_{i=1}^{n-1} D_{s_i} f(s_i, a_i) \right) D_{a_n} f(s_n, a_n)^\top D_\theta \pi(a_n | s_n) \right)^\top P_{s_n} r(s_n, a_n) \right\|_2$$

$$\leq \frac{1}{2} K B_r B_\pi + (K-1) B_r B_\pi B_g^{K-1}$$

$$\leq K B_r B_\pi \left( \frac{1}{2} + B_g^{K-1} \right)$$

Critically the most important term above is  $B_g^{K-1}$  which implies that non-smooth dynamics contribute the most to this upper bound

## Summary of changes:

1. Need to work in L2 norms for Lipschitz property to work

2. Drop  $|r(s, a)| \leq 1$  assumption

3. Assume action-less baseline

$$b = \sum_{n=1}^N r(s_n, 0)$$

4. Assume  $a \sim \pi_\theta(\cdot | sg(s))$

5. More explicitly note that

$r(s, a)$  and  $f(s, a)$  are Lipschitz in both arguments