

$$\nabla_{\theta} J(\theta)$$

$$\xi = (s_0, a_0, s_1, a_1, s_2, \dots)$$

TRAJECTORY

$$J(\theta) = \mathbb{E}_{\xi \sim P_{\theta}(\xi)} R(\xi) = \sum_{\xi} P_{\theta}(\xi) R(\xi)$$

$$\nabla_{\theta} J(\theta) = \sum_{\xi} \nabla_{\theta} P_{\theta}(\xi) R(\xi)$$

$$P_{\theta}(\xi) = P(s_0) \pi_{\theta}(a_0 | s_0) P(s_1 | s_0, a_0) \pi_{\theta}(a_1 | s_1) P(s_2 | s_1, a_1) \dots$$

$$\log P_{\theta}(\xi) = \log P(s_0) + \log \pi_{\theta}(a_0 | s_0) + \log P(s_1 | s_0, a_0) + \dots$$

$$\nabla_{\theta} \log P_{\theta}(\xi) = \nabla_{\theta} \log P(s_0) + \nabla_{\theta} \log \pi_{\theta}(a_0 | s_0) + \nabla_{\theta} \log P(s_1 | s_0, a_0) + \dots$$

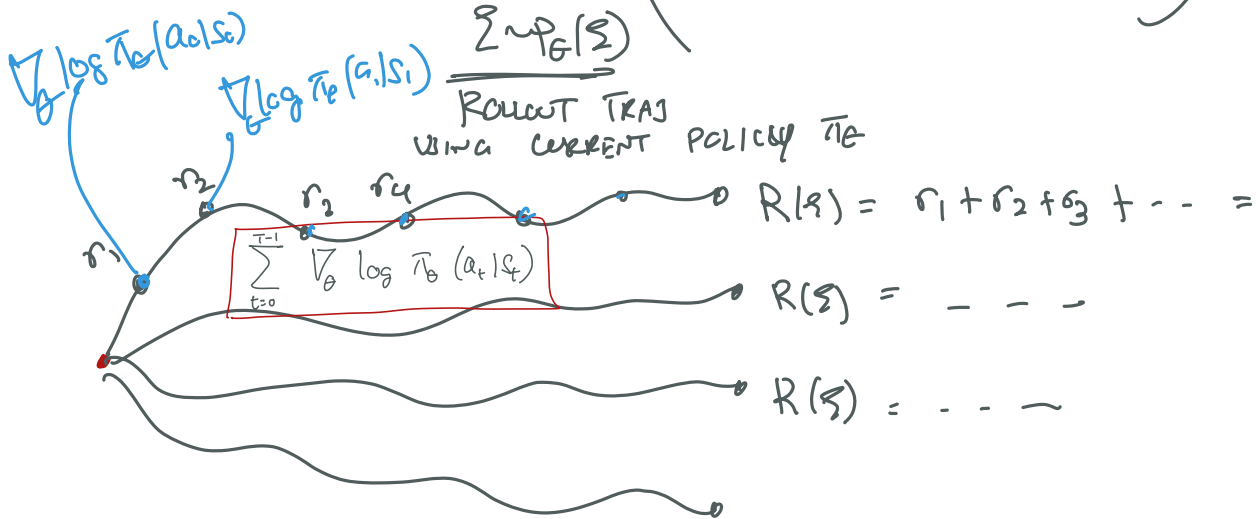
$$= \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

$$= \nabla_{\theta} \log P_{\theta}(\xi)$$

$$\nabla_{\theta} J(\theta) = \sum_{\xi} P_{\theta}(\xi) \left[\frac{\nabla_{\theta} P_{\theta}(\xi)}{P_{\theta}(\xi)} R(\xi) \right]$$

$$= \sum_{\xi} P_{\theta}(\xi) \left[\nabla_{\theta} \log P_{\theta}(\xi) R(\xi) \right]$$

$$= E \left(\boxed{\nabla_{\theta} \log P_{\theta}(\Sigma)} R(\Sigma) \right)$$



PERFORMANCE DIFFERENCE LEMMA

$$J(\pi_{\theta'}) - J(\pi_{\theta}) \geq 0$$

"NEW POLICY" "OLD POLICY"

$$= \sum_{t=0}^{T-1} E_{s_t \sim d_{\pi_{\theta'}}} \left[Q_{(s_t, \pi_{\theta'}(s_t))}^{\pi_{\theta}} - Q_{(s_t, \pi_{\theta}(s_t))}^{\pi_{\theta}} \right]$$

$\geq 0 \quad \nabla_{\pi_{\theta}}(s_t)$

$$= \sum_{t=0}^{T-1} E_{s_t \sim d_{\pi_{\theta'}}} \left[Q^{\pi_{\theta}}(\dots) - \gamma V^{\pi_{\theta}}(\cdot) \right]$$

$$A^{\pi_{\theta}}(s, \pi_{\theta'}(s))$$