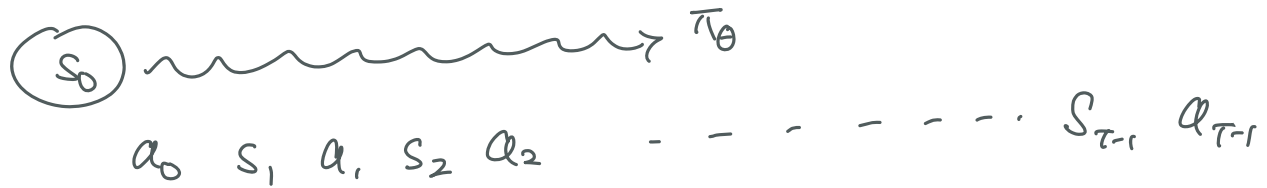


LET'S SAY WE HAVE A POLICY $\pi_\theta(a|s)$

ROLLOUT π_θ FROM START STATE s_0

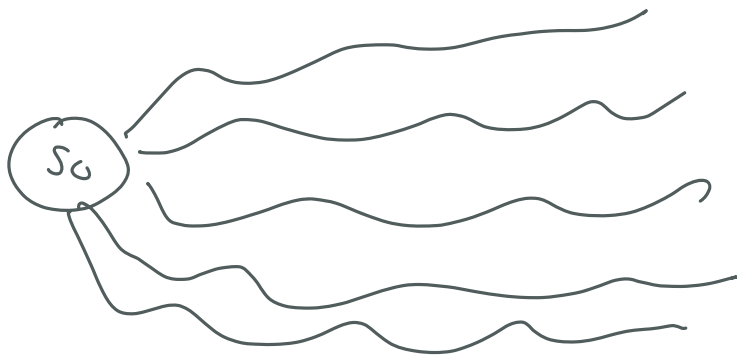


$$\xi = (s_0, a_0, s_1, a_1, \dots, s_{T-1}, a_{T-1})$$

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

EXPECTED TOTAL REWARD

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



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

NAIVE APPROACH:

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

APPLY CHAIN RULE.

$$\begin{aligned} \nabla_\theta P_\theta(\xi) &= P(s_0) \nabla_\theta \pi_\theta(a_0 | s_0) \overset{\text{UNKNOWN}}{\boxed{P(s_1 | s_0, a_0)}} \dots \\ &+ P(s_0) \pi_\theta(a_0 | s_0) \overset{\text{UNKNOWN}}{\boxed{P(s_1 | s_0, a_0)}} \nabla_\theta \pi_\theta(a_1 | s_1) \dots \\ &+ \dots \end{aligned}$$

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

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

$$\nabla_{\theta} \log P_{\theta}(\xi) = 0 + \nabla_{\theta} \log \pi_{\theta}(a_0 | s_0) + 0 + \nabla_{\theta} \log \pi_{\theta}(a_1 | s_1) + \dots$$

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

$$\frac{1}{P_{\theta}(\xi)} \nabla_{\theta} P_{\theta}(\xi)$$

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

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

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

APPROXIMATE THE EXPECTATION BY SAMPLING TRAJECTORIES

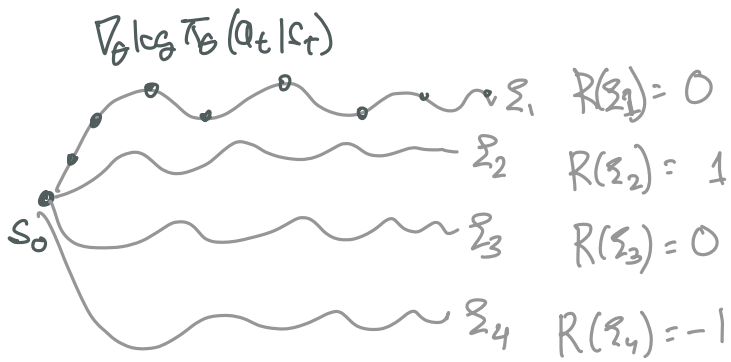
$$\left\{ \xi_i \right\}_{i=1}^N$$

BY ROLLING OUT POLICY π_{θ} IN REAL WORLD

$$\tilde{\nabla}_{\theta} J(\theta) = \frac{1}{N} \sum_{i=1}^N \left[\nabla_{\theta} \log P_{\theta}(\xi_i) R(\xi_i) \right]$$

Plugging in (1)

$$= \frac{1}{N} \sum_{i=1}^N \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t=0}^{T-1} r(s_t, a_t) \right]$$



$$E_{\substack{s \sim d^{\pi_{\theta}} \\ a \sim \pi_{\theta}}} \left[\nabla_{\theta} \log \pi_{\theta}(a | s) \cdot Q^{\pi_{\theta}}(s, a) \right]$$