

POLICY GRADIENT (RECREM)

$$J(\theta) = E \left[\sum_{t=0}^{T-1} r(s_t, a_t) \right]$$

$a_t \sim \pi_{\theta}(\cdot | s_t)$
 $s_{t+1} \sim P(\cdot | s_t, a_t)$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} E \left[\sum_{t=0}^{T-1} r(s_t, a_t) \right]$$

$a_t \sim \pi_{\theta}(\cdot | s_t)$
 $s_{t+1} \sim P(\cdot | s_t, a_t)$

$$= \nabla_{\theta} \sum_{s_0, a_0, s_1, a_1, \dots} P(s_0) \pi_{\theta}(a_0 | s_0) P(s_1 | s_0, a_0) \pi_{\theta}(a_1 | s_1) \left[\sum_{t=0}^{T-1} r(s_t, a_t) \right]$$

All trajectories

BAD IDEA: TRY CHAIN RULE.

$$\nabla_{\theta} \sum_{s_0, a_0, s_1, a_1, \dots} P(s_0) \nabla_{\theta} \pi_{\theta}(a_0 | s_0) P(s_1 | s_0, a_0) \dots$$

+

$$\sum_{s_0, a_0, s_1, a_1} P(s_0) \pi_{\theta}(a_0 | s_0) P(s_1 | s_0, a_0) \nabla_{\theta} \pi_{\theta}(a_1 | s_1) \dots$$

MUCH BETTER IDEA

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

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E} \left[\sum_{t=0}^{T-1} r(s_t, a_t) \right]$$

$a_t \sim \pi_{\theta}(\cdot | s_t)$
 $s_{t+1} \sim \mathcal{P}(\cdot | s_t, a_t)$

$$= \nabla_{\theta} \left(\sum_{\xi} \frac{P_{\theta}(\xi)}{P_{\theta}(\xi)} \right) \left[\sum_{t=0}^{T-1} r(s_t, a_t) \right]$$

$$= \mathbb{E}_{\xi \sim P_{\theta}(\xi)} \left[\frac{\nabla_{\theta} P_{\theta}(\xi)}{P_{\theta}(\xi)} \left[\sum_{t=0}^{T-1} r(s_t, a_t) \right] \right]$$

Just follow using current policies

$$= \mathbb{E}_{\xi \sim P_{\theta}(\xi)} \left[\nabla_{\theta} \log P_{\theta}(\xi) \left[\sum_{t=0}^{T-1} r(s_t, a_t) \right] \right]$$

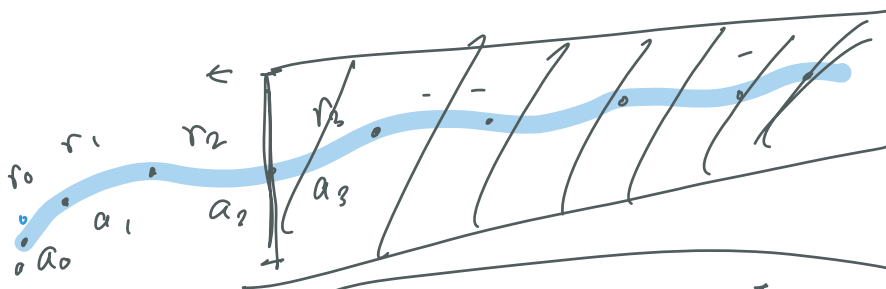
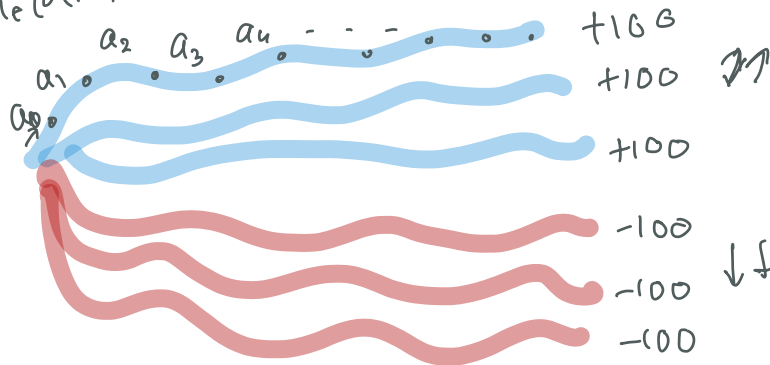
$$\nabla_{\theta} \left[\log \left(\cancel{P(s_0)} \cdot \pi_{\theta}(a_0 | s_0) \cdot P(s_1 | s_0, a_0) \cdot \pi_{\theta}(a_1 | s_1) \cdot \dots \right) \right]$$

$$\nabla_{\theta} \left[\log \cancel{P(s_0)} + \log \pi_{\theta}(a_0 | s_0) + \log \cancel{P(s_1 | s_0, a_0)} + \log \pi_{\theta}(a_1 | s_1) + \dots \right]$$

$$\nabla_{\theta} J(\theta) = E_{\pi_{\theta}(\xi)} \left[\nabla_{\theta} \log P_{\theta}(\xi) \left[\sum_{t=0}^{T-1} r(s_t, a_t) \right] \right]$$

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

① Roll out $\pi_{\theta}(a_t | s_t)$



$$\nabla_{\theta} \log \pi_{\theta}(a_3 | s_3) Q^{\pi_{\theta}}(s_3, a_3)$$