

Policy Gradient Method

(1)

$J(\theta)$, $\pi: S \rightarrow a$ instead we will use

$\pi_\theta: S \rightarrow p(a|s)$ (stochastic policy)

$$\nabla_\theta J = \nabla_\theta \mathbb{E}_{p(\xi|\theta)} [R(\xi)]$$

$$[\xi := (s_0, a_0, s_1, a_1, \dots)]$$

$$= \nabla_\theta \sum_{\xi} p(\xi|\theta) R(\xi)$$

[where \mathcal{I} is the space of all possible trajectories]

[Note that $p(\xi|\theta) = p(s_0) \cdot p(a_0|s_0) \cdot p(s_1|a_0, s_0) \cdot \pi(a_1|s_1) \cdot p(s_2|s_1, a_1) \dots$]

$$= \sum_{\xi} \nabla_\theta p(\xi|\theta) \cdot R(\xi) \quad (\text{as } \nabla_\theta R(\xi) = 0)$$

$$= \sum_{\xi} \frac{p(\xi|\theta) \nabla_\theta p(\xi|\theta) \cdot R(\xi)}{p(\xi|\theta)}$$

$$= \mathbb{E}_{p(\xi|\theta)} \left[\frac{\nabla_\theta p(\xi|\theta) \cdot R(\xi)}{p(\xi|\theta)} \right]$$

$$= \mathbb{E}_{p(\xi|\theta)} \nabla_{\theta} \log p(\xi|\theta) \cdot R(\xi)$$

$$= \mathbb{E}_{p(\xi|\theta)} \left[\nabla_{\theta} \left[\sum_t \log p(s_t|s_{t-1}, a_{t-1}) + \sum_t \log \pi(a_t|s_t) \right] R(\xi) \right]$$

0 (NOT dependent on θ)

↑
(Dependent on θ)

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{p(\xi|\theta)} \left[\sum_t \nabla_{\theta} \log \pi(a_t|s_t) \cdot R(\xi) \right]$$

Policy Gradient Theorem.

③

Reinforce (Version 1)

① Run simulator n times with π_θ

$$\textcircled{2} \quad \tilde{\nabla} J = \frac{1}{N} \sum_{i=1}^N \left[\sum_t \nabla \log \pi(a_t^i | s_t^i) \right] R(\xi_t^i)$$

$$\textcircled{3} \quad \theta_{\text{new}} \leftarrow \theta_{\text{old}} + \alpha \tilde{\nabla} J.$$

Reinforce (Version 2)

Rewriting step ②. *what happened to other terms?*

$$\tilde{\nabla} J = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \left[\nabla \log \pi(a_t | s_t) \right] \underbrace{\sum_{t=t+1}^T r(s_t)}_V$$

$$Q^\pi(s_t, a_t)$$

(in expectation)

↙

Can get this by
regression on rollout states
and hence
lower variance.

(better than a bunch of rollouts as
lower variance and leverage nearby states)

$$\text{So } \nabla J = \frac{1}{N} \sum_{i=1}^N \left[\sum_t \nabla_{\theta} \log(\pi_{\theta}(a_t^i | s_t^i)) \cdot \tilde{Q}^{\pi}(s_t^i, a_t^i) \right] \quad (9)$$

Advantages over discrete gradient.

Here we understand policy as function of θ .
 θ has undra causal information.

Downsides

- ① Have to have stochastic policy -
- ② Has to be differentiable.
- ③ Be able to differentiate.
- ④ Step-size parameter.

You can rewrite ② in REINFORCE as

$$\nabla J = \mathbb{E}_{P(s|\pi_{\theta})} \mathbb{E}_{\pi(a|s)} \left[\nabla \log \pi(a|s) \cdot Q^{\pi}(s,a) \right]$$