

Policy Gradient

$$\max_{\theta} \mathbb{E} \left[\sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \mid \pi_{\theta} \right] \rightarrow \text{objective fcn } J(\theta)$$

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

estimate this gradient.

* Likelihood Ratio Policy

Define state-action trajectory, τ , as:

$$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T)$$

Write the objective function, $J(\theta)$

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}(\tau)} \left[\sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \mid \pi_{\theta} \right]$$

$$= \sum_{\tau} \underbrace{P(\tau; \theta)}_{\uparrow \text{prob}} \cdot \underbrace{R(\tau)}_{\uparrow \text{R.V.}}$$

$P(\tau; \theta)$ is the probability of τ under π_{θ} :

$$P(\tau; \theta) = \mu(s_0) \prod_{t=0}^{T-1} [\pi(a_t \mid s_t; \theta) P(s_{t+1} \mid s_t, a_t)]$$

and the associated reward, $R(\tau)$, is defined as,

$$R(\tau) = \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t)$$

The goal is to find

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

$$= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau)$$

$$= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau)$$

$$= \sum_{\tau} P(\tau; \theta) R(\tau) \underbrace{\frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)}}_{\text{likelihood ratio}}$$

$$= \sum_{\tau} \underbrace{P(\tau; \theta)}_{\text{Prob}} R(\tau) \underbrace{\nabla_{\theta} \log P(\tau; \theta)}_{\text{RV}} \quad \textcircled{1}$$

$$= \mathbb{E}_{\tau} [R(\tau) \nabla_{\theta} \log P(\tau; \theta)] \quad \textcircled{2}$$

$$\approx \frac{1}{N} \sum_{i=1}^N R(\tau_i) \nabla_{\theta} \log P(\tau_i; \theta)$$

$$\nabla_{\theta} \log P(\tau; \theta) = \frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)}$$

$$\mathbb{E} \rightarrow \text{PMF} \times \text{RV}$$

The last term is called "Monte Carlo Sampling" to compute expectation.

Computing the gradient yields

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N R(\tau_i) \nabla_{\theta} \log P(\tau_i; \theta)$$

Next page

$$\begin{aligned}
 \nabla_{\theta} \log P(\tau_i; \theta) &= \nabla_{\theta} \log \left[\underbrace{\mu(s_0)}_{\text{Initial state dist}} \prod_{t=0}^{\tau-1} \underbrace{\pi_{\theta}(a_t|s_t)}_{\text{policy}} \underbrace{P(s_{t+1}|s_t, a_t)}_{\text{trans. prob}} \right] \\
 &= \nabla_{\theta} \left[\log \mu(s_0) + \sum_{t=0}^{\tau-1} \log \pi_{\theta}(a_t|s_t) + \log P(s_{\tau}|s_{\tau-1}, a_{\tau-1}) \right] \\
 &= \sum_{t=0}^{\tau-1} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \\
 &\quad \underbrace{\hspace{10em}}_{\text{No dynamics required}}
 \end{aligned}$$

We need a policy, $\pi_{\theta}(a_t|s_t)$ to be stochastic

Summary

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\
 &\approx \frac{1}{N} \sum_{i=1}^N \underbrace{R(\tau_i)} \underbrace{\nabla_{\theta} \log P(\tau_i; \theta)} \\
 &= \frac{1}{N} \sum_{i=1}^N \underbrace{\sum_{t=0}^{\tau-1} r(s_t, a_t)}_{s_t, a_t} \underbrace{\sum_{t=0}^{\tau-1} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)}
 \end{aligned}$$

Causality