

proof of (2):

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^H \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left( \sum_{t'=t}^H \gamma^{t'-t} r(s_{t'}, a_{t'}) - b(s_t) \right) \right]$$

This form shows the policy gradient as a sum over actions, each scaled by an advantage estimate: a discounted future reward sum minus a baseline.

### Maximize

We want to maximize expected return:

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)], \quad R(\tau) = \sum_{t=0}^H \gamma^t r(s_t, a_t), \quad \tau = (s_0, a_0, \dots, s_H, a_H).$$

### Use the log-derivative trick

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)] = \mathbb{E}_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log p_{\pi_{\theta}}(\tau) R(\tau)].$$

proof of the trick:

$$\nabla_{\theta} p(\tau; \theta) = p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta) \tag{1}$$

$$\nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)] = \nabla_{\theta} \int p_{\pi_{\theta}}(\tau) R(\tau) d\tau \tag{2}$$

$$= \int \nabla_{\theta} p_{\pi_{\theta}}(\tau) R(\tau) d\tau \tag{3}$$

$$= \int p_{\pi_{\theta}}(\tau) \nabla_{\theta} \log p_{\pi_{\theta}}(\tau) R(\tau) d\tau \tag{4}$$

$$= \mathbb{E}_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log p_{\pi_{\theta}}(\tau) R(\tau)] \tag{5}$$

$$= \nabla_{\theta} J(\theta). \tag{6}$$

### Decompose $\log p_{\pi_{\theta}}(\tau)$ :

$$p(\tau) = \rho_0(s_0) \prod_{t=0}^H \pi_{\theta}(a_t | s_t) P(s_{t+1} | s_t, a_t),$$

so

$$\nabla_{\theta} \log p_{\pi_{\theta}}(\tau) = \sum_{t=0}^H \nabla_{\theta} \log \pi_{\theta}(a_t | s_t).$$

### Plug in:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[ \left( \sum_{t=0}^H \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) R(\tau) \right] = \sum_{t=0}^H \mathbb{E}_{\tau} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau)].$$

### Reward-to-Go (Causality)

Only future rewards depend on  $a_t$ . Dropping past rewards: (mathematically if you write out the integral terms that do not depend on the integrator factor out and you are left with the gradient of integral of a probability which is the constant 1 and hence zero)

$$\nabla_{\theta} J(\theta) = \sum_{t=0}^H \mathbb{E}_{\tau} \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \underbrace{\sum_{t'=t}^H \gamma^{t'} r(s_{t'}, a_{t'})}_{\text{reward-to-go}} \right].$$

Factor out  $\gamma^t$ :

$$= \sum_{t=0}^H \mathbb{E}_{\tau} \left[ \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'=t}^H \gamma^{t'-t} r(s_{t'}, a_{t'}) \right].$$

**Add a Baseline (As to why you can do this and what a baseline is see pp 329 of Sutton)**

Subtract a baseline  $b(s_t)$ :

$$\nabla_{\theta} J(\theta) = \sum_{t=0}^H \mathbb{E}_{\tau} \left[ \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left( \sum_{t'=t}^H \gamma^{t'-t} r(s_{t'}, a_{t'}) - b(s_t) \right) \right].$$