

* Actor

Reminder: causal policy-gradient estimator

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left(\sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \left(\sum_{t'=t}^{T-1} r(s_{i,t'}, a_{i,t'}) \right)\end{aligned}$$

We can subtract a function of state without adding bias
base, $b(s_t)$

$$\nabla_{\theta} J(\theta) = \mathbb{E} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left(\sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) - b(s_{t'}) \right) \right]$$

Where $b(s_t)$ is baseline,

$$b(s_t) = \mathbb{E} \left[\sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right]$$

Intuition: We only prioritize actions that do better
than average

$b(s_t)$

Actor-Critic

We can briefly write down overall transition,

$$\left(\sum_{t=t}^{T-1} r_t \right)$$

$$\begin{aligned} \nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \frac{\hat{Q}^{\pi}(s_{i,t}, a_{i,t})}{\hat{V}_{\phi}^{\pi}(s_{i,t})} \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) (r(s_{i,t}, a_{i,t}) + \gamma \hat{V}_{\phi}^{\pi}(s_{i,t+1})) \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) (r(s_{i,t}, a_{i,t}) + \gamma \hat{V}_{\phi}^{\pi}(s_{i,t+1})) \\ &\quad \quad \quad \hat{Q}_{\phi}^{\pi}(s_{i,t}, a_{i,t}) - \hat{V}_{\phi}^{\pi}(s_{i,t}) \end{aligned}$$

$$\hat{A}^{\pi}(s_{i,t}, a_{i,t}) = \hat{Q}_{\phi}^{\pi}(s_{i,t}, a_{i,t}) - \hat{V}_{\phi}^{\pi}(s_{i,t})$$

↑ Advantage function.

+ lower Variance (due to Critic)

- not unbiased (since Critic is not perfect)

Policy Gradient

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left(\sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) - b(s_t) \right)$$

+ no bias

- higher variance