

Policy Gradient Method

Action-value Methods

Learned the values of actions and then selected actions based on their estimated action values

Learn a parameterized policy

Select actions without consulting a value function

$\theta \in \mathbb{R}^{d'}$: policy's parameter vector

$$\pi(a|s, \theta) = P_{\pi} \{ A_t = a | S_t = s, \theta_t = \theta \}$$

$w \in \mathbb{R}^d$: value function's weight vector

$$\hat{V}(s, w)$$

"Learn the policy parameter based on the gradient of some scalar performance measure $J(\theta)$ w.r.t. the policy parameter"

$$\Theta_{t+1} = \Theta_t + \alpha \widehat{\nabla J(\Theta_t)}$$

$\widehat{\nabla J(\Theta_t)} \in \mathbb{R}^d$: Stochastic estimate whose expectation approximates the gradient of the performance measure w.r.t. its argument Θ_t

{ Learn value function : Action-value Method
 Learn policy : policy gradient Method
 Learn value function & policy : actor-critic Method

Softmax in action preference:

$$\pi(a|s, \theta) = \frac{e^{h(s,a,\theta)}}{\sum_b e^{h(s,b,\theta)}}$$

$$h(s,a,\theta) = \theta^T x(s,a)$$

① deterministic policy,

② action selection with arbitrary probabilities

The Policy Gradient Theorem

"How performance is affected by the policy parameter that does not involve derivatives of the state distribution"

Value of the start state under the parameterized policy

performance: $J(\theta) = V_{\pi_\theta}(s_0)$ average reward rate

$$\nabla J(\theta) = \sum_{s'} \eta(s') \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_{\pi}(s, a)$$

$$\propto \sum_s \mu(s) \sum_a q_{\pi}(s, a) \nabla \pi(a|s, \theta)$$

↑
proportional to

REINFORCE: MC Policy Gradient

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_{\pi}(s, a) \nabla \pi(a|s, \theta)$$

$$= \mathbb{E} \left[\sum_a q_{\pi}(s_t, a) \nabla \pi(a|s_t, \theta) \right]$$

$$= \mathbb{E}_{\pi} \left[G_t \frac{\nabla \pi(A_t | s_t, \theta)}{\pi(A_t | s_t, \theta)} \right]$$

$$\theta_{t+1} \doteq \theta_t + \alpha \sum_a \hat{q}(S_t, a, w) \nabla \pi(a|S_t, \theta)$$

*
$$\theta_{t+1} = \theta_t + \alpha \underbrace{G_t}_{\substack{\nabla \ln \pi(A_t | S_t, \theta_t) \\ \pi(A_t | S_t, \theta_t)}}$$

" use the complete return from time t , which includes all future rewards up until the end of the episode. "

REINFORCE: MC Policy-gradient Control (episodic)
for π :

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$$

$$\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t | S_t, \theta)$$

REINFORCE Proof:

$$\nabla J(\theta) = \mathbb{E}_{\pi} \left[\sum_a \pi(a|s_t, \theta) q_{\pi}(s_t, a) \frac{\nabla \pi(a|s_t, \theta)}{\pi(a|s_t, \theta)} \right]$$

$$= \mathbb{E}_{\pi} \left[q_{\pi}(s_t, A_t) \frac{\nabla \pi(A_t|s_t, \theta)}{\pi(A_t|s_t, \theta)} \right]$$

$$= \mathbb{E}_{\pi} \left[G_t \frac{\nabla \pi(A_t|s_t, \theta)}{\pi(A_t|s_t, \theta)} \right]$$

REINFORCE with Baseline (unbiased)

"include a comparison of the action value to an arbitrary baseline $b(s)$ "

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a (\gamma \pi(s, a) - b(s)) \nabla \pi(a|s, \theta)$$

$$\sum_a b(s) \nabla \pi(a|s, \theta) = b(s) \nabla \sum_a \pi(a|s, \theta) = b(s) \nabla 1 = 0$$

$$\theta_{t+1} \doteq \theta_t + \alpha (G_t - b(s_t)) \frac{\nabla \pi(A_t | s_t, \theta_t)}{\pi(A_t | s_t, \theta_t)}$$

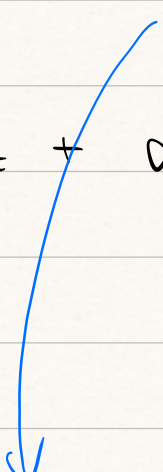
$$\begin{cases} G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k \\ \delta \leftarrow G - \hat{v}(S_t, w) \\ w \leftarrow w + \alpha^w \delta \nabla \hat{v}(S_t, w) \\ \theta \leftarrow \theta + \alpha^\theta \gamma^t \delta \nabla \ln \pi(A_t | S_t, \theta) \end{cases}$$

$$\alpha^w = \frac{0.1}{\mathbb{E} [\| \nabla \hat{v}(S_t, w) \|^2_\mu]}$$

Actor-critic Methods

$$\Theta_{t+1} \doteq \Theta_t + \alpha (G_{t:t+1} - \hat{V}(S_t, w)) \frac{\nabla \pi(A_t | S_t, \Theta_t)}{\pi(A_t | S_t, \Theta_t)}$$

$$= \Theta_t + \alpha (R_{t+1} + \gamma \hat{V}(S_{t+1}, w) - \hat{V}(S_t, w)) \frac{\nabla \pi(A_t | S_t, \Theta_t)}{\pi(A_t | S_t, \Theta_t)}$$

$$= \Theta_t + \alpha \delta \frac{\nabla \pi(A_t | S_t, \Theta_t)}{\pi(A_t | S_t, \Theta_t)}$$


One-step actor-critic methods replace the full return of REINFORCE with the one-step return (and use a learned state value function as the baseline)