

chapter 13

Policy gradient method: action without calculating value function
parameterized policy:

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

Notation value function $\hat{v}(s, w)$

Gradient ascent, $\theta_{t+1} = \theta_t + \alpha \widehat{\nabla j(\theta_t)}$ → policy gradient method.

method learns both value & policy → Actor critic Algorithm.

$$\pi(a|s, \theta) = \frac{e^{h(s, a, \theta)}}{\sum_a e^{h(s, a, \theta)}} \quad // \quad h(s, a, \theta) = \text{can be function of } (x(s, a), \theta)$$

maybe $\approx \theta^T x(s, a)$

Two advantages: ① reach deterministic policy ② stochastic value consideration.

Policy gradient theorem:

$$J(\theta) = \mathbb{E}_{\pi_\theta}(s_0)$$

distribution of staying at state s .

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

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

Reinforce:

$$\theta_{t+1} = \theta_t + \alpha \sum_a q_\pi(s_t, a, w) \nabla \pi(a|s_t, \theta)$$

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

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

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

$a_t, \pi_t \sim \pi$
(sampling)

$$\theta_{t+1} = \theta_t + \alpha G_t \frac{\nabla \pi(a_t | s_t, \theta)}{\pi(a_t | s_t, \theta)}$$

$$= \theta_t + \alpha G_t \nabla \ln [\pi(a_t | s_t, \theta)]$$

Reinforce: Monte Carlo Policy - Gradient Control (episodic) π

Input: $\pi(a|s, \theta)$

$\alpha > 0$

$\theta \in \mathbb{R}^d$

Loop forever:

Generate $s_0, a_0, r_1, \dots, s_{T-1}, a_{T-1}, r_T$ following $\pi(\cdot | \cdot, \theta)$

Loop for each episode $t = 0, 1, \dots, T-1$

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

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

$$\nabla \ln \pi(a|s, \theta) = X(s, a) - \sum_b \pi(b|s, \theta) X(s, b)$$

Reinforce with baseline:

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

Algorithm:

Input: $\pi(a|s, \theta)$, $\hat{v}(s, w)$, $\alpha^v > 0$, $\alpha^w > 0$, $\theta \in \mathbb{R}^d$, $w \in \mathbb{R}^d$

Loop forever (for each episode):

Generate $s_0, a_0, r_1, \dots, s_{T-1}, a_{T-1}, r_T$ following $\pi(\cdot | \cdot, \theta)$

Loop for each step of episode: $t = 0, 1, \dots, T-1$

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

$$\delta \leftarrow G - \hat{v}(s_t, w)$$

$$w \leftarrow w + \alpha^w \delta \nabla \hat{v}(s_t, w)$$

$$\theta \leftarrow \theta + \alpha^v \gamma^t \delta \nabla \ln \pi(a_t | s_t, \theta)$$

Action critique method:

$$\begin{aligned} \theta_{t+1} &= \theta_t + \alpha \left(G_{t:t+1} - \hat{v}(s_t, \underline{w}) \right) \frac{\nabla \pi(a_t | s_t, \theta)}{\pi(a_t | s_t, \theta_t)} \\ &= \theta_t + \alpha \left(r_{t+1} + \gamma \hat{v}(s_{t+1}, \underline{w}) - \hat{v}(s_t, \underline{w}) \right) \frac{\nabla \pi(a_t | s_t, \theta)}{\pi(a_t | s_t, \theta_t)} \\ &= \theta_t + \alpha \delta_t \frac{\nabla \pi(a_t | s_t, \theta_t)}{\pi(a_t | s_t, \theta_t)} \end{aligned}$$

One step Action-critic:

input: $\pi(a|s, \theta)$, $\hat{v}(s, \underline{w})$, $\alpha^\theta > 0$, $\alpha^w > 0$

param: $\theta \in \mathbb{R}^d$, $\underline{w} \in \mathbb{R}^d$

Loop forever:

init s

$J \leftarrow 1$

Loop while s is not terminal:

$a \sim \pi(\cdot | s, \theta)$

Take action a , observe s', r

$\delta \leftarrow r + \gamma \hat{v}(s', \underline{w}) - \hat{v}(s, \underline{w})$ // s terminal then $\hat{v}(s', \underline{w}) = 0$

$\underline{w} \leftarrow \underline{w} + \alpha^w \delta \nabla \hat{v}(s, \underline{w})$

$\theta \leftarrow \theta + \alpha^\theta J \nabla \ln \pi(a | s, \theta)$

$J \leftarrow \gamma J$

$s \leftarrow s'$

eligibility trace (episode)

Loop forever:

init s , $z^\theta \leftarrow 0$, $z^w \leftarrow 0$, $J \leftarrow 1$

loop while s not terminal

+ action \rightarrow observe s', r find δ

$z^w \leftarrow \gamma z^w + \nabla \hat{v}(s, \underline{w})$

$z^\theta \leftarrow \gamma z^\theta + J \nabla \ln \pi(a | s, \theta)$

$\underline{w} \leftarrow \underline{w} + \alpha^w \delta z^w$

$\theta \leftarrow \theta + \alpha^\theta \delta z^\theta$

;

actor critic eligibility trace continuity

input: $\pi(a|s, \theta)$

input: $\hat{v}(s, w)$

$\lambda^w \in [0, 1], \lambda^a \in [0, 1], \alpha^w > 0, \alpha^a > 0, \alpha^p > 0$

init $\bar{r} \in \mathbb{R}$

init $w \in \mathbb{R}^d, \phi \in \mathbb{R}^d, s \in \mathcal{S}$

$z^w \leftarrow 0$ // eligibility trace vector

$z^a \leftarrow 0$

Loop forever:

$A \sim \pi(\cdot | s, \theta)$: take action

observe s', r

$\delta \leftarrow r - \bar{r} + \hat{v}(s', w) - \hat{v}(s, w)$

$\bar{r} \leftarrow \bar{r} + \alpha^p \delta$

$z^w \leftarrow \lambda^w z^w + \nabla \hat{v}(s, w)$

$z^a \leftarrow \lambda^a z^a + \nabla \ln \pi(A|s, \theta)$

$\theta \leftarrow \theta + \alpha^a \delta z^a$

$w \leftarrow w + \alpha^w \delta z^w$

$s \leftarrow s'$
