


Chapter 10. Actor-Critic Methods.

actor : policy update

critic : policy evaluation.

- ① simple AC.
- ② A2C
- ③ off-policy AC.
- ④ DPG.

1. The simplest A-C

$\left\{ \begin{array}{l} \text{reinforce (MC)} \\ \text{Temporal-difference learning} \end{array} \right.$

At time step t

generate a_t following $\pi(a|s_t, \theta_t)$ observe r_{t+1}, s_{t+1}
 and then generate a_{t+1} following $\pi(a|s_{t+1}, \theta_t)$

Critics:

$$w_{t+1} = w_t + \alpha_w [r_{t+1} + \gamma q(s_{t+1}, a_{t+1}, w_t) - q(s_t, a_t, w_t)] \triangleright_w q(s_t, a_t, w_t)$$

Actor:

$$\theta_{t+1} = \theta_t + \alpha_\theta \triangleright_{\pi} |_{\pi(a|s_t, \theta_t)} q(s_t, a_t, w_{t+1})$$

policy gradient is on policy

2. Advantage AC. (A2C)

The gradient is $\nabla_{\theta} J(\theta) = E[X]$

where: $X(S, A) = \nabla_{\theta} \ln \pi(A|S, \theta_t) [q(S, A) - b(S)]$

$E[X]$ is invariant to $b(S)$

$\nabla \text{Var}[X]$ is not invariant to $b(S)$

$$\xrightarrow{\quad} \text{tr}[\text{Var}(X)] = E[X^T X] - \bar{x}^T \bar{x}$$

$$\bar{x} = E(X)$$

~~W W W~~ we hope $X(S, A)$ $\text{Var}(X)$ smallest.
 ↑
 It's not right
 the optimal baseline:

$$b^*(S) = \frac{E_{\text{A}\pi} [\|\nabla_{\theta} \ln \pi(A|S, \theta_t)\|^2 q(S, A)]}{E_{\text{A}\pi} [\|\nabla_{\theta} \ln \pi(A|S, \theta_t)\|^2]}$$

↓

$$b^*(S) \approx b(S) = E_{\text{A}\pi} [q(S, A)] = V_{\pi}(S)$$

the gradient-ascent algorithm is

$$\theta_{t+1} = \theta_t + \alpha E[\nabla_{\theta} \ln \pi(A|S, \theta_t) \delta_{\pi}(S, A)]$$

advantage function: $\delta_{\pi}(S, A) = q_{\pi}(S, A) - V_{\pi}(S)$

stochastic version:

$$\begin{aligned}\theta_{t+1} &= \theta_t + \alpha \nabla_{\theta} \ln \pi(a_t | s_t, \theta_t) \delta_{t+1}(s_t, a_t) \\ &= \theta_t + \alpha \underbrace{\frac{\delta_t(s_t, a_t)}{\pi(a_t | s_t, \theta_t)}}_{\text{Step-size}} \nabla_{\theta} \pi(a_t | s_t, \theta_t)\end{aligned}$$

$$s_t = q_t(s_t, a_t) - V_t(s_t) \rightarrow = r_{t+1} + \gamma V_t(s_{t+1}) - V_t(s_t)$$

↑ ↑
2 networks one network

A2C / TD actor-critic;

At step time t :

generate a_t following $\pi(a_t | s_t, \theta_t)$ and then observe r_{t+1}, s_{t+1}

TD error:

$$\delta_t = r_{t+1} + \gamma V(s_{t+1}, w_t) - V(s_t, w_t)$$

Critic:

$$w_{t+1} = w_t + \alpha_w \delta_t \nabla_w V(s_t, w_t)$$

Actor:

$$\theta_{t+1} = \theta_t + \alpha_{\theta} \delta_t \nabla_{\theta} \pi(a_t | s_t, \theta_t)$$

3. off policy actor-critic

Introduction: importance sampling

$$E_{X \sim p_0} [f(x)] = \sum_x p_0(x) f(x) = \sum_x p_1(x) \underbrace{\frac{p_0(x)}{p_1(x)} x}_{f(x)} = E_{X \sim p_1} [f(x)]$$

$$\bar{f} = \frac{1}{n} \sum_{i=1}^n f(x_i)$$

$\frac{p_0(x_i)}{p_1(x_i)}$ is called importance weight

The theorem of off-policy policy gradient.

β is the behavior policy generate episodes

$$J(\theta) = \sum_{s \in S} d_\beta(s) V_\pi(s) =$$

d_β — stationary distribution under policy β .

$$\nabla_\theta J(\theta) = E_{\text{sample}} \left[\frac{\pi(a_t | s_t, \theta)}{p(a_t | s_t)} \nabla_\theta \ln \frac{\pi(a_t | s_t, \theta) q_{\pi}(s_t, a_t)}{\downarrow q_{\pi}(s_t, a_t) - V_\pi(s_t)} \right]$$

$$\theta_{t+1} = \theta_t + \alpha_\theta \frac{\pi(a_t | s_t, \theta_t)}{p(a_t | s_t)} \nabla_\theta \ln \pi(a_t | s_t, \theta_t) (q_{\pi}(s_t, a_t) - V_\pi(s_t))$$

$$\theta_{t+1} = \theta_t + \alpha \left(\frac{\pi(s_t, \theta_t)}{p(a_t | s_t)} \right) \nabla_{\theta} \pi(a_t | s_t, \theta_t)$$

off-policy sad

given policy $\beta(a|s)$ target policy: $\pi(a|s, \theta_0)$

A Value function $V(s, w_0)$ where w_0 is the initial parameter vector

① generate a_t following $\beta(s_t)$

observe r_{t+1}, s_{t+1}

② TD error

$$e_t = r_{t+1} + \gamma V(s_{t+1}, w_t) - V(s_t, w_t)$$

③ Critic:

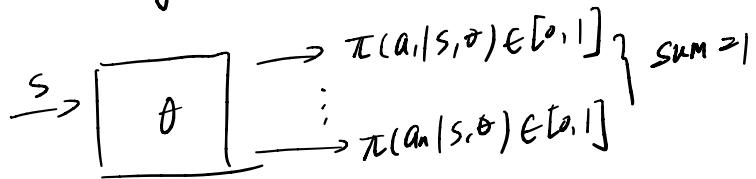
$$w_{t+1} = w_t + \alpha_w \frac{\pi(a_t | s_t, \theta_t)}{p(a_t | s_t)} e_t \nabla_w V(s_t, w_t)$$

④ Actor:

$$\theta_{t+1} = \theta_t + \alpha_{\theta} \frac{\pi(a_t | s_t, \theta_t)}{p(a_t | s_t)} e_t \nabla_{\theta} \ln \pi(a_t | s_t, \theta_t)$$

4. Deterministic actor-critic (DPG)

when action is infinite. — deterministic
finite — stochastic



deterministic policy is $\underline{a} = \mu(s, \theta) = u(s)$

$$\mu: S \rightarrow A$$



$\mu(s, \theta)$ in short $u(s)$

① policy gradient:

$$J(\theta) = \sum_{s \in S} d(s) V_u(s)$$

$$\nabla_\theta J(\theta) = E_{S \sim P_\theta} \left[\nabla_\theta \mu(s) (\nabla_\theta q_\mu(s, a)) \mid a = \mu(s) \right]$$

$$\theta_{t+1} = \theta_t + \alpha \mathbb{E}_{s \sim p_a} [\nabla_\theta U(s) (\nabla_a q_u(s, a)) \mid a = u(s)]$$

$$\theta_{t+1} = \theta_t + \alpha \nabla_\theta U(s_t) (\nabla_a q_u(s_t, a)) \mid a = u(s_t)$$

behavior policy: $\beta(a|s)$ A deterministic target policy: $u(s, \theta_0)$ A value function: $U(s, w)$

TD error:

$$\delta_t = r_{t+1} + \gamma q(s_{t+1}, u(s_{t+1}, \theta_t), w_t) - q(s_t, a_t, w_t)$$

Critic:

$$w_{t+1} = w_t + \alpha_w \delta_t \nabla_w q(s_t, a_t, w_t)$$

Actor:

$$\theta_{t+1} = \theta_t + \alpha \nabla_\theta U(s_t, \theta_t) (\nabla_a q(s_t, a|u(s_t))) \mid a = u(s_t)$$