

Week 7: Policy Optimization II: Actor-Critic and Natural Policy Gradient

Bolei Zhou

The Chinese University of Hong Kong

October 10, 2020

Today's Plan

- ① Review on policy gradient
- ② Actor critic and advantage actor critic
- ③ Natural policy gradient
- ④ Sample code

Review on Policy Gradient

- ① policy optimization is to maximize $J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[R(\tau)]$
- ② The gradient is

$$\nabla_\theta J(\theta) = \nabla_\theta \sum_{\tau} P(\tau; \theta) R(\tau) = \mathbb{E}_{\tau} [R(\tau) \nabla_\theta \log P(\tau; \theta)]$$

- ① After decomposing τ we derived that

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau} \left[\left(\sum_{t=0}^{T-1} r_t \right) \left(\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) \right) \right]$$

- ② After applying the causality, we have

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} G_t \cdot \nabla_\theta \log \pi_\theta(a_t | s_t) \right]$$

Reduce Variance by a Baseline

- 1 The original update

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} \mathbf{G}_t \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- 1 $G_t = \sum_{t'=t}^{T-1} r_{t'}$ is the return for a Monte-Carlo trajectory which might have high variance
- 2 We subtract a baseline $b(s)$ from the actual return of the policy gradient to reduce variance

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} (\mathbf{G}_t - b(s_t)) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- 3 A good baseline is the expected return

$$b(s_t) = \mathbb{E}[r_t + r_{t+1} + \dots + r_{T-1}]$$

Reduce Variance by a Baseline

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} (G_t - b_{\mathbf{w}}(s_t)) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- ① Baseline $b(s)$ can reduce variance, without changing the expectation
- ② $b_{\mathbf{w}}(s)$ also has a parameter \mathbf{w} to learn so that we have two set of parameters \mathbf{w} and θ

Vanilla Policy Gradient Algorithm with Baseline

procedure POLICY GRADIENT(α)

Initialize policy parameters θ and baseline values $b(s)$ for all s , e.g. to 0

for iteration = 1, 2, ... **do**

Collect a set of m trajectories by executing the current policy π_θ

for each time step t of each trajectory $\tau^{(i)}$ **do**

Compute the *return* $G_t^{(i)} = \sum_{t'=t}^{T-1} r_{t'}$

Compute the *advantage estimate* $\hat{A}_t^{(i)} = G_t^{(i)} - b(s_t)$

Re-fit the baseline to the empirical returns by updating \mathbf{w} to minimize

$$\sum_{i=1}^m \sum_{t=0}^{T-1} \|b(s_t) - G_t^{(i)}\|^2$$

Update policy parameters θ using the policy gradient estimate \hat{g}

$$\hat{g} = \sum_{i=1}^m \sum_{t=0}^{T-1} \hat{A}_t^{(i)} \nabla_\theta \log \pi_\theta(a_t^{(i)} | s_t^{(i)})$$

with an optimizer like SGD ($\theta \leftarrow \theta + \alpha \cdot \hat{g}$) or Adam

return θ and baseline values $b(s)$

Reduce Variance Using a Critic

- 1 The update is

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^{T-1} \mathbf{G}_t \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- 2 In practice, G_t is a sample from Monte Carlo policy gradient, which is the unbiased but noisy estimate of $Q^{\pi_{\theta}}(s_t, a_t)$
- 3 Instead we can use a **critic** to estimate the action-value function,

$$Q_{\mathbf{w}}(s, a) \approx Q^{\pi_{\theta}}(s, a)$$

- 4 Then the update becomes

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^{T-1} Q_{\mathbf{w}}(s_t, a_t) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

Actor-Critic Policy Gradient

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^{T-1} Q_{\mathbf{w}}(s_t, a_t) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- ① It becomes Actor-Critic Policy Gradient
 - ① **Actor**: the policy function used to generate the action
 - ② **Critic**: the value function used to evaluate the reward of the actions
- ② Actor-critic algorithms maintain two sets of parameters
 - ① **Actor**: Updates policy parameters θ , in direction suggested by critic
 - ② **Critic**: Updates action-value function parameters \mathbf{w}

Estimate the Action-Value Function

- ① Estimating the critic is solving a familiar problem: policy evaluation
 - ① How good is policy π_θ under current parameter θ
- ② Policy evaluation was explored in previous lectures, e.g.
 - ① Monte-Carlo policy evaluation
 - ② Temporal-Difference learning
 - ③ Least-squares policy evaluation

Action-Value Actor-Critic Algorithm

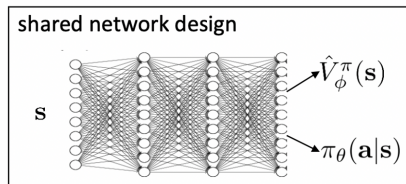
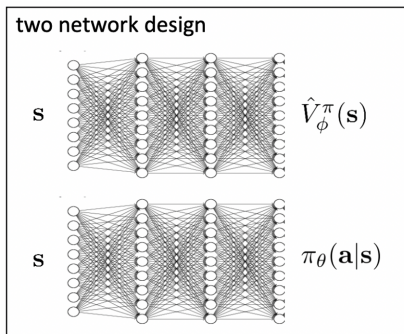
- ① Using a linear value function approximation: $Q_{\mathbf{w}}(s, a) = \psi(s, a)^T \mathbf{w}$
 - ① **Critic**: update \mathbf{w} by a linear TD(0)
 - ② **Actor**: update θ by policy gradient

Algorithm 1 Simple QAC

- 1: **for** each step **do**
 - 2: generate sample s, a, r, s', a' following π_{θ}
 - 3: $\delta = r + \gamma Q_{\mathbf{w}}(s', a') - Q_{\mathbf{w}}(s, a)$ #TD error
 - 4: $\mathbf{w} \leftarrow \mathbf{w} + \beta \delta \psi(s, a)$
 - 5: $\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(s, a) Q_{\mathbf{w}}(s, a)$
 - 6: **end for**
-

Actor-Critic Function Approximators

- 1 We can have two separate functions to approximate value function and policy function, or use a shared network design (feature extraction is shared but output two heads), as below:



Compatible Function Approximation

Theorem (Compatible Function Approximation Theorem)

If the following two conditions are satisfied:

- 1. Value function approximator is compatible to the policy*

$$\nabla_w Q_w(s, a) = \nabla_\theta \log \pi_\theta(s, a)$$

- 2. Value function parameters w minimize the mean-squared error*

$$\epsilon(w) = E_{\pi_\theta} [(Q^{\pi_\theta}(s, a) - Q_w(s, a))^2]$$

then the policy gradient is exact

$$\nabla_\theta J(\theta) = E_{\pi_\theta} [Q_w(s, a) \nabla_\theta \log \pi_\theta(s, a)]$$

RS Sutton, et al. Policy gradient methods for reinforcement learning with function approximation. NIPS'01

Compatible Function Approximation

- 1 Given a parameterized policy, compatible function approximation theorem can be used to derive an appropriate form of the value-function parameterization.
 - 1 For example, consider a policy as Gibbs distribution where ϕ_{sa} is the feature vector for state-action pair s,a

$$\pi_{\theta}(s, a) = \frac{e^{\theta^T \phi_{sa}}}{\sum_b e^{\theta^T \phi_{sb}}} \quad (1)$$

- 2 Then meeting the compatibility condition requires

$$\nabla_w Q_w(s, a) = \nabla_{\theta} \log \pi_{\theta}(s, a) = \phi_{sa} - \sum_b \pi_{\theta}(s, b) \phi_{sb}$$

- 3 so the compatible parameterization of $Q_w(s, a)$ is

$$Q_w(s, a) = w^T \left(\phi_{sa} - \sum_b \pi_{\theta}(s, b) \phi_{sb} \right) = w^T \nabla_{\theta} \log \pi_{\theta}(s, a)$$

Reduce the Variance of Actor-Critic by a Baseline

- 1 Recall Q-function / state-action-value function:

$$Q^{\pi,\gamma}(s, a) = \mathbb{E}_{\pi}[r_1 + \gamma r_2 + \dots | s_1 = s, a_1 = a]$$

- 2 State value function can serve as a great baseline

$$\begin{aligned} V^{\pi,\gamma}(s) &= \mathbb{E}_{\pi}[r_1 + \gamma r_2 + \dots | s_1 = s] \\ &= \mathbb{E}_{a \sim \pi}[Q^{\pi,\gamma}(s, a)] \end{aligned}$$

- 3 Advantage function: combining Q with baseline V

$$A^{\pi,\gamma}(s, a) = Q^{\pi,\gamma}(s, a) - V^{\pi,\gamma}(s)$$

- 4 Then the policy gradient becomes:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) A^{\pi,\gamma}(s, a)]$$

Advantage Actor-Critic

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A^{\pi_{\theta}}(s, a)]$$

- ① Critic should estimate the advantage function
- ② Use two function approximators and two sets of parameter vectors,

$$\begin{aligned} V_{\mathbf{v}}(s) &\approx V^{\pi}(s) \\ Q_{\mathbf{w}}(s, a) &\approx Q^{\pi}(s, a) \end{aligned}$$

- ③ Update both parameters by TD learning or MC

Advantage Actor-Critic

- ① For the true value function $V^{\pi_{\theta}}(s)$, the TD error $\delta^{\pi_{\theta}}$

$$\delta^{\pi_{\theta}} = r(s, a) + \gamma V^{\pi_{\theta}}(s') - V^{\pi_{\theta}}(s)$$

is an unbiased estimate of the advantage function

$$\begin{aligned}\mathbb{E}_{\pi_{\theta}}[\delta^{\pi_{\theta}} | s, a] &= \mathbb{E}_{\pi_{\theta}}[r + \gamma V^{\pi_{\theta}}(s') | s, a] - V^{\pi_{\theta}}(s) \\ &= Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s) \\ &= A^{\pi_{\theta}}(s, a)\end{aligned}$$

- ② So we can use the TD error to compute the policy gradient

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) \delta^{\pi_{\theta}}]$$

- ③ It requires only **one set of critic parameter κ** that is estimated by TD

$$\delta_v = r + \gamma V_{\kappa}(s') - V_{\kappa}(s)$$

N-step estimators

- 1 We used the Monte-Carlo estimates of the reward
- 2 We can also use TD methods for the policy gradient update, or any intermediate blend between TD and MC methods:
- 3 Consider the following n -step returns for $n = 1, 2, \infty$

$$n = 1(TD) \quad G_t^{(1)} = r_{t+1} + \gamma v(s_{t+1})$$

$$n = 2 \quad G_t^{(2)} = r_{t+1} + \gamma r_{t+2} + \gamma^2 v(s_{t+2})$$

$$n = \infty(MC) \quad G_t^{(\infty)} = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{T-t-1} r_T$$

- 4 Then the advantage estimators become

$$\hat{A}_t^{(1)} = r_{t+1} + \gamma v(s_{t+1}) - v(s_t)$$

$$\hat{A}_t^{(2)} = r_{t+1} + \gamma r_{t+2} + \gamma^2 v(s_{t+2}) - v(s_t)$$

$$\hat{A}_t^{(\infty)} = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{T-t-1} r_T - v(s_t)$$

$\hat{A}_t^{(1)}$ has low variance and high bias. $\hat{A}_t^{(\infty)}$ has high variance but low bias

Critic at Different Time-Scales

- ① Critic can estimate linear value function $V_{\kappa}(s) = \psi(s)^T \kappa$ from many targets as follows
- ② For MC, the update is

$$\Delta\kappa = \alpha(\textcolor{red}{G}_t - V_{\kappa}(s))\psi(s)$$

- ③ For TD(0), the update is

$$\Delta\kappa = \alpha(\textcolor{red}{r} + \gamma V_{\kappa}(\textcolor{red}{s}') - V_{\kappa}(s))\psi(s)$$

- ④ For k-step return, the update is

$$\Delta\kappa = \alpha\left(\sum_{i=0}^k \gamma^i \textcolor{red}{r}_{t+i} + \gamma^k V_{\kappa}(\textcolor{red}{s}_{t+k}) - V_{\kappa}(s_t)\right)\psi(s)$$

Actors at Different Time-Scales

- 1 The policy gradient can also be estimated at many time-scales

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A^{\pi_{\theta}}(s, a)]$$

- 2 Monte-Carlo Actor-Critic policy gradient uses error from complete return

$$\Delta \theta = \alpha (G_t - V_{\kappa}(s_t)) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$$

- 3 TD Difference Actor-Critic policy gradient uses TD error

$$\Delta \theta = \alpha (r + \gamma V_{\kappa}(s_{t+1}) - V_{\kappa}(s_t)) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$$

- 4 k-step return Actor-Critic policy gradient uses the k-step return error

$$\Delta \theta = \alpha \left(\sum_{i=0}^k \gamma^i r_{t+i} + \gamma^k V_{\kappa}(s_{t+k}) - V_{\kappa}(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$$

Summary of Policy Gradient Algorithms

- 1 The policy gradient has many forms

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) G_t] \text{ — REINFORCE}$$

$$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q^{\mathbf{w}}(s, a)] \text{ — Q Actor-Critic}$$

$$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A^{\mathbf{w}}(s, a)] \text{ — Advantage Actor-Critic}$$

$$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta] \text{ — TD Actor-Critic}$$

- 2 Critic uses policy evaluation (e.g. MC or TD learning) to estimate $Q^{\pi}(s, a)$, $A^{\pi}(s, a)$, or $V^{\pi}(s, a)$

Natural Policy Gradient

- ① Policy gradient is the steepest ascend in **parameter space** with Euclidean metric:

$$d^* = \nabla_{\theta} J(\theta) = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \max J(\theta + d), s.t. \|d\| \leq \epsilon$$

- ① Drawback: it is sensitive to the parameterization of the policy function
- ② Steepest ascend in **distribution space (policy output)** with KL-divergence as constraint

$$d^* = \arg \max J(\theta + d), s.t. KL[\pi_{\theta} || \pi_{\theta+d}] = c \quad (2)$$

- ① Fixing KL-divergence as a constant c make sure we move along the distribution space with constant speed, regardless the curvature (thus robust to the reparametrization of the model as we only care about the distribution induced by the parameter)

KL-Divergence as Metric for Two Distributions

- 1 KL-divergence measure the “closeness” of two distributions.

$$KL[\pi_{\theta}||\pi_{\theta'}] = E_{\pi_{\theta}}[\log \pi_{\theta}] - E_{\pi_{\theta}}[\log \pi_{\theta'}]$$

- 2 Although as KL-divergence is non-symmetric and thus not a true metric, we can use it anyway. This is because as d goes to zero, KL-divergence is asymptotically symmetric. So, within a local neighbourhood, KL-divergence is approximately symmetric
- 3 We can prove that the second-order Taylor expansion of KL-divergence is

$$KL[\pi_{\theta}||\pi_{\theta+d}] \approx \frac{1}{2}d^T F d$$

where F is the **Fisher Information Matrix** as second order derivative of KL divergence $E_{\pi_{\theta}}[\nabla \log \pi_{\theta} \nabla \log \pi_{\theta}^T]$

Natural Policy Gradient

$$d^* = \arg \max J(\theta + d), s.t. KL[\pi_\theta || \pi_{\theta+d}] = c$$

- 1 Write the above as Lagrangian form, with the $J(\theta + d)$ approximated by its first order Taylor expansion and the constraint KL-divergence approximated by its second order Taylor expansion

$$\begin{aligned} d^* &= \arg \max_d J(\theta + d) - \lambda(KL[\pi_\theta || \pi_{\theta+d}] - c) \\ &\approx \arg \max_d J(\theta) + \nabla_\theta J(\theta)^T d - \frac{1}{2} \lambda d^T F d + \lambda c \end{aligned}$$

- 2 To solve this maximization we set its derivative wrt. d to zero, then we have the **natural policy gradient**: $d = \frac{1}{\lambda} F^{-1} \nabla_\theta J(\theta)$

Natural Policy Gradient

- ① Natural policy gradient is a **second-order** optimization that is more accurate and works regardless of how the model is parameterized (model invariant).

$$\theta_{t+1} = \theta_t + \alpha \mathbf{F}^{-1} \nabla_{\theta} J(\theta)$$

- ① where $\mathbf{F} = E_{\pi_{\theta}(s,a)}[\nabla \log \pi_{\theta}(s,a) \nabla \log \pi_{\theta}(s,a)^T]$ is the Fisher information matrix, also as the second-order derivative of the KL-divergence
- ② \mathbf{F} measures the curvature of the policy(distribution) relative to the model parameter θ
- ② Natural policy gradient produces the same policy change **regardless of the model parameterization**.
- ③ Idea behind TRPO. A detailed read on natural gradient: <https://wiseodd.github.io/techblog/2018/03/14/natural-gradient/>

Code example

- 1 REINFORCE code on CartPole:
`https://github.com/cuhkrlcourse/RLEXample/blob/master/policygradient/reinforce.py`
- 2 Policy Gradient on Pong : `https://github.com/cuhkrlcourse/RLEXample/blob/master/policygradient/pg-pong-pytorch.py`
- 3 Policy Gradient with Baseline on Pong:
`https://github.com/cuhkrlcourse/RLEXample/blob/master/policygradient/pgb-pong-pytorch.py`
- 4 Actor Critic on Pong: `https://github.com/cuhkrlcourse/RLEXample/blob/master/policygradient/ac-pong-pytorch.py`

Next Week: State of the Arts on Policy Optimization

① Policy Gradient→TRPO→ACKTR→PPO

- ① **TRPO**: Trust region policy optimization. Schulman, L., Moritz, Jordan, Abbeel. 2015
- ② **ACKTR**: Scalable trust-region method for deep reinforcement learning using Kronecker-factored approximation. Y. Wu, E. Mansimov, S. Liao, R. Grosse, and J. Ba. 2017
- ③ **PPO**: Proximal policy optimization algorithms. Schulman, Wolski, Dhariwal, Radford, Klimov. 2017

② Q-learning→DDPG→TD3→SAC

- ① **DDPG**: Deterministic Policy Gradient Algorithms, Silver et al. 2014
- ② **TD3**: Addressing Function Approximation Error in Actor-Critic Methods, Fujimoto et al. 2018
- ③ **SAC**: Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor, Haarnoja et al. 2018