# Advantage Actor-Critic





CIDA Lab.



Presentation: Baekryun Seong

2024.07.10

Table of Contents	Tabl	le of	Con	tents
-------------------	------	-------	-----	-------

The Actor

The Critic

A2C Algorithm

Network Architecture

Neural Actor Critic

In A2C, we will use the policy gradient and a learned value function.

In REINFORCE, A policy is reinforced by reward, directly.

But in Actor-Critic, a policy is reinforced with a learned reinforcing signal.



# Ch 6.1 The Actor

In REINFORCE, we used the policy gradient of:  $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{t} [R_{t}(\tau) \nabla_{\theta} \log \pi_{\theta}(a_{t}|s_{t})].$ Now in actor-critic, we will use:  $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{t} [A_{t}^{\pi} \nabla_{\theta} \log \pi_{\theta}(a_{t}|s_{t})]$ where  $A_{t}^{\pi}$  is an advantage function.

# Ch 6.2 The Critic

Policy gradient is very useful, but it has large variance. So, we introduce an advantage function.

$$A^{\pi}(s_t, a_t) = Q^{\pi}(s_t, a_t) - V^{\pi}(s_t)$$

where

$$\mathbb{E}_{a\in A}[A^{\pi}(s_t,a)]=0.$$

Advantage function reinforces the actions with rewards over expectation.

	$Q^{\pi}(s,a)$	$V^{\pi}(s)$	$A^{\pi}(s,a)$
Case 1	110	100	10
Case 2	-90	-100	10

To get advantage function, we need to estimate  $Q^{\pi}, V^{\pi}$ . We generally estimate  $V^{\pi}$  first and infer  $Q^{\pi}$  from it and reward. There is two main reason:

1.  $Q^{\pi}$  is more complicated function.

2. It is hard to get  $V^{\pi}$  from  $Q^{\pi}$  because it requires summation over all states.

With approximated value function  $V^{\pi}$ , we get:

 $\begin{aligned} A_{NSTEP}^{\pi}A(s_{t},a_{t}) &= Q^{\pi}(s_{t},a_{t}) - V^{\pi}(s_{t}) \\ &\approx r_{t} + \gamma r_{t+1} + \gamma^{2}r_{t+2} + \dots + \gamma^{n}r_{t+n} + \gamma^{n+1}\hat{V}^{\pi}(s_{t+n+1}) - \hat{V}^{\pi}(s_{t}) \end{aligned}$ 

Like SARSA or TD( $\lambda$ ), we can approximate advantage utilizing different steps.

$$A_{GAE}^{\pi}(s_t, a_t) = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}$$
  
where  $\delta_t = r_t + \gamma V^{\pi}(s_{t+1}) - V^{\pi}(s_t)$ .  
 $\delta_t$  is usually called *TD error*.

 $\lambda$  regulates bias-variance tradeoff. For further explanation, see <u>this</u>.

Now we train advantage function using  $V^{\pi}$  target, n-step or GAE.

 $V^{\pi}$  target can be set differently. In one step,

 $V_{tar}^{\pi}(s) = r + \gamma \hat{V}^{\pi}(s;\theta).$ 

In n-step, MC and GAE, targets are:

$$V_{tar}^{\pi}(s_t) = r_t + \gamma r_{t+1} + \dots + \gamma^n r_{t+n} + \gamma^{n+1} \hat{V}^{\pi}(s_{t+n+1}),$$

$$V_{tar}^{\pi}(s_t) = \sum_{t'=t}^{T} \gamma^{t'-t} r_{t'},$$
$$V_{tar}^{\pi}(s_t) = A_{GAE}^{\pi}(s_t, a_t) + \hat{V}^{\pi}(s_t).$$

# Ch 6.3 A2C Algorithm

### Ch 6.3 A2C Algorithm

#### Algorithm 6.1 A2C algorithm

- 1: Set  $\beta \ge 0$  # entropy regularization weight
- 2: Set  $\alpha_A \ge 0 \ \#$  actor learning rate
- 3: Set  $\alpha_C \ge 0$  # critic learning rate
- 4: Randomly initialize the actor and critic parameters  $\theta_A$ ,  $\theta_C^4$
- 5: **for**  $episode = 0 \dots MAX\_EPISODE$  **do**
- 6: Gather and store data  $(s_t, a_t, r_t, s'_t)$  by acting in the environment using  $\hookrightarrow$  the current policy
- 7: **for** t = 0...T **do**
- 8: Calculate predicted V-value  $\hat{V}^{\pi}(s_t)$  using the critic network  $\theta_C$
- 9: Calculate the advantage  $\hat{A}^{\pi}(s_t, a_t)$  using the critic network  $\theta_C$
- 10: Calculate  $V_{tar}^{\pi}(s_t)$  using the critic network  $\theta_C$  and/or trajectory data
- 11: Optionally, calculate entropy  $H_t$  of the policy distribution, using  $\hookrightarrow$  the actor network  $\theta_A$ . Otherwise, set  $\beta = 0$

#### 12: end for

Calculate value loss, for example using MSE: 13:  $L_{\rm val}(\theta_C) = \frac{1}{T} \sum_{t=0}^{T} (\hat{V}^{\pi}(s_t) - V_{\rm tar}^{\pi}(s_t))^2$ 14: Calculate policy loss: 15:  $L_{\text{pol}}(\theta_A) = \frac{1}{T} \sum_{t=0}^{T} (-\hat{A}^{\pi}(s_t, a_t) \log \pi_{\theta_A}(a_t \mid s_t) - \beta H_t)$ 16: Update critic parameters, for example using SGD:<sup>5</sup> 17:  $\theta_C = \theta_C + \alpha_C \nabla_{\theta_C} L_{\text{val}}(\theta_C)$ 18: Update actor parameters, for example using SGD: 19:  $\theta_A = \theta_A + \alpha_A \nabla_{\theta_A} L_{\text{pol}}(\theta_A)$ 20: 21: end for

### Ch 6.5 Network Architecture

### Ch 6.5 Network Architecture

In A2C, we have to train two different networks: actor and critic.

But we can assume that two networks can share some information.

Sharing parameters makes the actor network to utilize state representation of critic net.

But it makes training very unstable, and scale problems may appear.



Figure 6.1 Actor-Critic network architectures: shared vs. separate networks

### Appendix Actor-Critic in Neuroscience

The learning rule of real neurons is generally called Hebbian learning. The most famous learning rule of Hebb is Spike-Timing Dependent Plasticity. (STDP) Researchers found that TD error exists in brain: *TD Reward Prediction Error*. (TD RPE) This is called *RPE hypothesis of dopamine neuron activity*. Dopamine level in brain is very similar to TD error:  $R_t + \gamma V(S_{t+1}) - V(S_t)$ . It is very interesting: reinforcement learning have started from dynamic programming, but it also appears in neuroscience.

Researchers also suggests the structure of actor-critic in real brain.

### Neural Actor-Critic



Figure 15.5: Actor-critic ANN and a hypothetical neural implementation. a) Actor-critic algorithm as an ANN. The actor adjusts a policy based on the TD error  $\delta$  it receives from the critic; the critic adjusts state-value parameters using the same  $\delta$ . The critic produces a TD error from the reward signal, R, and the current change in its estimate of state values. The actor does not have direct access to the reward signal, and the critic does not have direct access to the reward signal, and the critic algorithm. The actor and the value-learning part of the critic are respectively placed in the dorsal and ventral subdivisions of the striatum. The TD error is transmitted by dopamine neurons located in the VTA and SNpc to modulate changes in synaptic efficacies of input from cortical areas to the ventral and dorsal striatum. Adapted from Frontiers in Neuroscience, vol. 2(1), 2008, Y. Takahashi, G. Schoenbaum, and Y. Niv, Silencing the critics: Understanding the effects of cocaine sensitization on dorsolateral and ventral striatum in the context of an Actor/Critic model.

#### Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An introduction. MIT press.