

Towards Understanding Asynchronous Advantage Actor-critic : Convergence and Linear Speedup

Hyelin Choi

Department of Mathematics
Sungkyunkwan University

May 9, 2024

Definition

- Advantage function $A_{\pi}(s, a) := Q_{\pi}(s, a) - V_{\pi}(s)$
- Initial state distribution η

Definition

- Discounted state visitation measure (induced by policy π)

$$d_{\pi}(s) := (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s \mid s_0 \sim \eta, \pi)$$

- State-action visitation distribution

$$d_{\pi}(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s \mid s_0 \sim \eta, \pi) \pi(a|s)$$

Definition

- KL divergence

$$\begin{aligned} D_{KL}(P\|Q) &= \sum_x P(x) \log \left(\frac{P(x)}{Q(x)} \right) \\ &= \sum_x P(x) \log P(x) - \sum_x P(x) \log Q(x) \\ &= \mathbb{E}_P [\log P(x)] - \mathbb{E}_P [\log Q(x)] \end{aligned}$$

Asynchronous Advantage Actor-critic

Actor-critic

- Actor network

Policy Gradient Method \longrightarrow Policy optimization $\pi(a|s)$

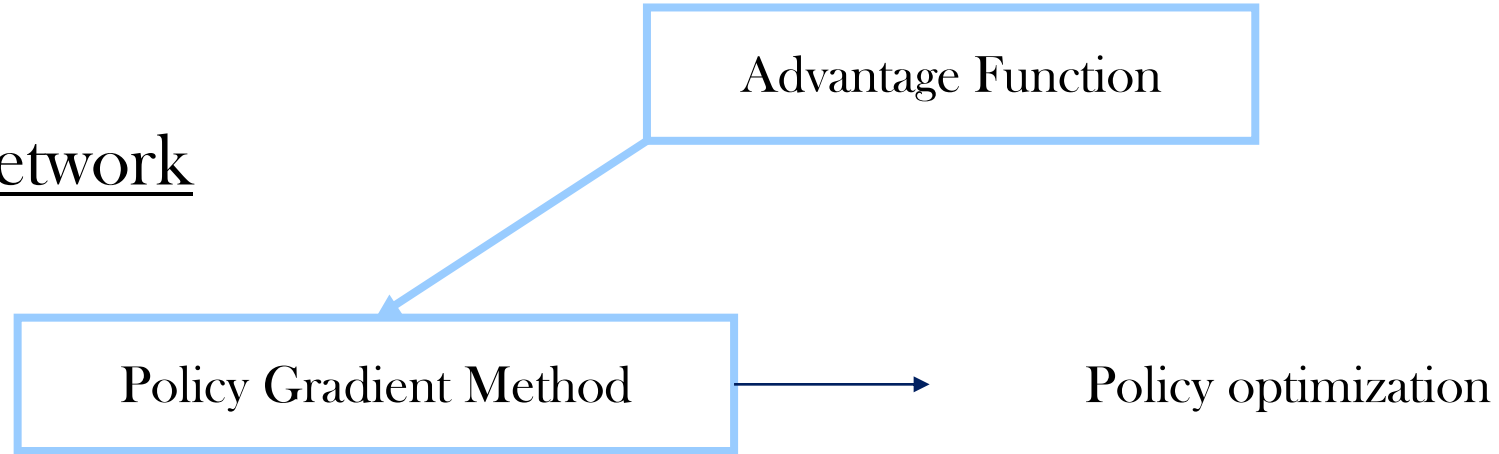
- Critic network

Temporal Difference Learning Algorithm \longrightarrow Policy evaluation $V(s)$

Asynchronous Advantage Actor-critic

A3C

- Actor network



- Critic network

Temporal Difference Learning Algorithm → Policy evaluation

- Policy Gradient Method

$$\max_{\theta \in \mathbb{R}^d} J(\theta) \text{ with } J(\theta) := (1 - \gamma) \mathbb{E}_{s \sim \eta} [V_{\pi_\theta}(s)]$$

$$\theta_{k+1} = \theta_k + \alpha \nabla J(\theta_k)$$

$$\text{where } \nabla J(\theta) = \mathbb{E}_{s, a \sim d_\theta} [\nabla \log \pi_\theta(s, a) Q_{\pi_\theta}(s, a)]$$

- Policy Gradient Method

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [G_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)]$$

REINFORCE

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [Q_{\pi_{\theta}}(s, a) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)]$$

Q-value Actor-Critic

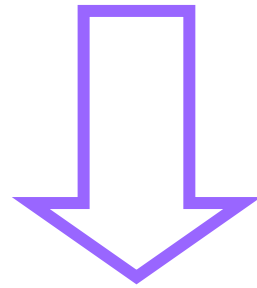
$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [A(s, a) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)]$$

Advantage Actor-Critic

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [(r + \gamma V_{\pi_{\theta}}(s') - V_{\pi_{\theta}}(s)) \nabla_{\theta} \log \pi_{\theta}(a | s)]$$

TD Actor-Critic

$$\nabla J(\theta) = \mathbb{E}_{s,a \sim d_\theta} [\nabla \log \pi_\theta(s, a) Q_{\pi_\theta}(s, a)]$$



$$\nabla J(\theta) = \mathbb{E}_{s,a \sim d_\theta} [\nabla \log \pi_\theta(s, a) (Q_{\pi_\theta}(s, a) - V_{\pi_\theta}(s))]$$

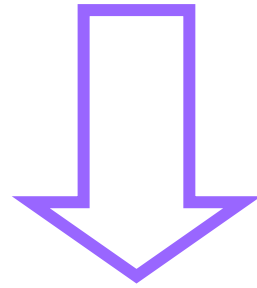
- Need to show: $\mathbb{E}_{s,a \sim d_\theta} [\nabla \log \pi_\theta(s, a) B(s)] = 0$

Asynchronous Advantage Actor-critic

Synchronous vs Asynchronous

- Actor step

$$\theta_{k+1} = \theta_k + \alpha \mathbb{E}_{s,a \sim d_\theta} [(Q_{\pi_\theta}(s, a) - V_{\pi_\theta}(s)) \underbrace{\nabla \log \pi_\theta(s, a)}_{\psi_\theta(s_k, a_k)}]$$



$$\theta_{k+1} = \theta_k + \alpha [(r(s_k, a_k, s_{k+1}) + \gamma \phi(s_{k+1})^T w_k - \phi(s_k)^T w_k) \psi_\theta(s_k, a_k) + \lambda \psi_\theta(x^p)]$$

■ Regularization

η_p : prior distribution of states

π_p : prior policy

Regularization term encourages π_θ to imitate π_p , incorporating prior knowledge into training process.

$$\begin{aligned} J_\lambda(\theta) &:= J(\theta) - \lambda \mathbb{E}_{s \sim \eta_p} [D_{KL}(\pi_p(\cdot|s) | \pi_\theta(\cdot|s))] \\ &= J(\theta) + \lambda R(\theta) \end{aligned}$$

- Bellman equation

$$V_{\pi_{\theta}}(s) = \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s), s' \sim \mathcal{P}(\cdot|s,a)} [r(s, a, s') + \gamma V_{\pi_{\theta}}(s')]$$

- Value Function Approximation

state feature mapping

$$V_{\pi_{\theta}}(s) \approx \hat{V}_w(s) = \phi(s)^T w$$

- KL divergence

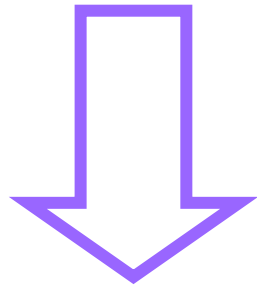
$$\begin{aligned} D_{KL}(P\|Q) &= \sum_x P(x) \log \left(\frac{P(x)}{Q(x)} \right) \\ &= \sum_x P(x) \log P(x) - \sum_x P(x) \log Q(x) \\ &= \mathbb{E}_P [\log P(x)] - \mathbb{E}_P [\log Q(x)] \end{aligned}$$

$$\psi_{\theta}(s, a) := \nabla \log \pi_{\theta}(s, a)$$

$$x^p := (s^p \sim \eta_p, a^p \sim \pi_p(\cdot | s_p))$$

- Unbiased Estimator

$$\nabla J_{\lambda}(\theta) = \nabla J(\theta) + \lambda \nabla R(\theta)$$

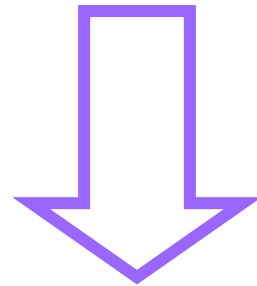


$$\hat{\nabla} J_{\lambda}(\theta) = \hat{\nabla} J(\theta) + \lambda \hat{\nabla} R(\theta)$$

$$\theta_{k+1} = \theta_k + \alpha \left[\underbrace{\left(r(s_k, a_k, s_{k+1}) + \gamma \phi(s_{k+1})^T w_k - \phi(s_k)^T w_k \right) \psi_\theta(s_k, a_k)}_{v(x_k, \theta_k, w_k): \text{unbiased estimator of } \nabla J(\theta)} + \underbrace{\lambda \psi_\theta(x^p)}_{\text{Unbiased estimator of } \nabla R(\theta)} \right]$$

$v(x_k, \theta_k, w_k)$: unbiased estimator of $\nabla J(\theta)$

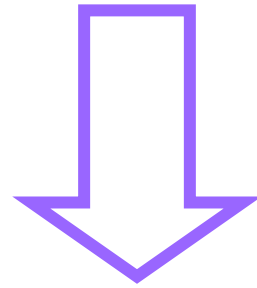
Unbiased estimator of $\nabla R(\theta)$



$$\theta_{k+1} = \theta_k + \alpha [v(x_k, \theta_k, w_k) + \lambda \psi_\theta(x^p)]$$

- Critic step

TD(0) algorithm



$$w_{k+1} = w_k + \beta \hat{\delta}(x_k, w_k) \nabla \hat{V}_{w_k}(s_k)$$

- TD(0) algorithm

$$V(\boldsymbol{s}_t) \leftarrow V(\boldsymbol{s}_t) + \beta [r_{t+1} + \gamma V(\boldsymbol{s}_{t+1}) - V(\boldsymbol{s}_t)]$$

- TD(0) algorithm

$$V(\boldsymbol{s}_t) \leftarrow V(\boldsymbol{s}_t) + \beta \mathbb{E}_{\boldsymbol{s}, \boldsymbol{a}, \boldsymbol{s}'} [r(\boldsymbol{s}, \boldsymbol{a}, \boldsymbol{s}') + \gamma V_{\pi_{\theta}}(\boldsymbol{s}') - V_{\pi_{\theta}}(\boldsymbol{s})]$$

- Value Function Approximation

$$V_{\pi_{\theta}}(s) \approx \hat{V}_w(s) = \phi(s)^T w$$

- TD error

$$\hat{\delta}(x_k, w_k) := r(s_k, a_k, s_{k+1}) + (\gamma \phi(s_{k+1}) - \phi(s_k))^T w_k$$

Π_{R_w} is a projection operator that projects a vector to a l_2 norm ball with radius R_w .

It prevents the actor and critic updates from going too far in the wrong direction.

Critic step

$$w_{k+1} = \Pi_{R_w} (w_k + \beta g(x_k, w_k))$$

Actor step

$$\theta_{k+1} = \theta_k + \alpha [v(x_k, \theta_k, w_k) + \lambda \psi_\theta(x^p)]$$

Critic step $\omega_{k+1} = \Pi_{R_\omega} \left(\omega_k + \beta g(x_{(k)}, \omega_{k-\tau_k}) \right)$

Actor step $\theta_{k+1} = \theta_k + \alpha \left(v(\hat{x}_{(k)}, \theta_{k-\tau_k}, \omega_{k-\tau_k}) + \lambda \psi_{\theta_{k-\tau_k}}(x_{(k)}^p) \right)$

- τ_k : the **delay** in the k th actor and critic updates

Algorithm 1 A3C: each worker's view.

- 1: **Global initialize:** Global counter $k=0$, initial θ_0, ω_0 in the shared memory.
 - 2: **Worker initialize:** Counter $t=0$. Sample $s_0 \sim \eta, \hat{s}_0 \sim \eta$.
 - 3: **for** $t = 0, 1, 2, \dots$ **do**
 - 4: Read θ, ω in the shared memory.
 - 5: **option 1 (i.i.d. sampling):**
 - 6: $x_t = (s_t \sim \mu_{\pi_{\theta_t}}, a_t \sim \pi_{\theta_t}(\cdot|s_t), s'_t \sim \mathcal{P}(\cdot|s_t, a_t)).$
 - 7: $\hat{x}_t = (\hat{s}_t \sim d_{\pi_{\theta_t}}, \hat{a}_t \sim \pi_{\theta_t}(\cdot|\hat{s}_t), \hat{s}'_t \sim \mathcal{P}(\cdot|\hat{s}_t, \hat{a}_t)).$
 - 8: **option 2 (Markovian sampling):**
 - 9: $x_t = (s_t, a_t \sim \pi_{\theta}(\cdot|s_t), s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t)).$
 - 10: $\hat{x}_t = (\hat{s}_t, \hat{a}_t \sim \pi_{\theta}(\cdot|\hat{s}_t), \hat{s}'_{t+1} \sim \mathcal{P}(\cdot|\hat{s}_t, \hat{a}_t)).$
 - 11: With probability γ : $\hat{s}_{t+1} = \hat{s}'_{t+1}$; Otherwise: $\hat{s}_{t+1} \sim \eta$.
 - 12: Compute $g(x_t, \omega) = \delta(x_t, \omega) \nabla_{\omega} \hat{V}_{\omega}(s_t)$.
 - 13: Compute $v(\hat{x}_t, \theta, \omega) = \hat{\delta}(\hat{x}_t, \omega) \psi_{\theta}(\hat{s}_t, \hat{a}_t)$.
 - 14: Compute $\psi_{\theta}(x_t^p)$ with $x_t^p = (s_t^p \sim \eta_p, a_t^p \sim \pi_p(\cdot|s_t^p)).$
 - 15: In the shared memory, perform update (13).
 - 16: **end for**
-

$$\omega_{k+1} = \Pi_{R_{\omega}} (\omega_k + \beta g(x_{(k)}, \omega_{k-\tau_k})) \quad (13a)$$

$$\theta_{k+1} = \theta_k + \alpha (v(\hat{x}_{(k)}, \theta_{k-\tau_k}, \omega_{k-\tau_k}) + \lambda \psi_{\theta_{k-\tau_k}}(x_{(k)}^p)) \quad (13b)$$

option 2 (Markovian sampling):

$x_t = (s_t, a_t \sim \pi_\theta(\cdot|s_t), s_{t+1} \sim \mathcal{P}(\cdot|s_t, a_t))$. Critic

$\hat{x}_t = (\hat{s}_t, \hat{a}_t \sim \pi_\theta(\cdot|\hat{s}_t), s'_{t+1} \sim \mathcal{P}(\cdot|\hat{s}_t, \hat{a}_t))$. Actor

With probability γ : $\hat{s}_{t+1} = s'_{t+1}$; Otherwise: $\hat{s}_{t+1} \sim \eta$

- In Markovian sampling case, we maintain separate Markov chains for actor and critic.
- For **Critic**, we generate samples following the original transition kernel \mathbf{P} .
- For **Actor**, we generate samples following a transition kernel $\hat{\mathbf{P}} = \gamma\mathbf{P} + (1 - \gamma)\eta$.
If the actor's chain evolves under \mathbf{P} like critic, asymptotically the initial distribution η is forgotten, which will introduce asymptotic error.

A3C

option 1 (i.i.d. sampling):

$$x_t = (s_t \sim \mu_{\pi_{\theta_t}}, a_t \sim \pi_{\theta_t}(\cdot|s_t), s'_t \sim \mathcal{P}(\cdot|s_t, a_t)). \quad \text{Critic}$$

$$\hat{x}_t = (\hat{s}_t \sim d_{\pi_{\theta_t}}, \hat{a}_t \sim \pi_{\theta_t}(\cdot|\hat{s}_t), \hat{s}'_t \sim \mathcal{P}(\cdot|\hat{s}_t, \hat{a}_t)). \quad \text{Actor}$$

- $\mu_{\pi_{\theta}}$ is the stationary distribution of the Markov chain with transition distribution P and π_{θ} .
- $d_{\pi_{\theta}}$ is the discounted state visitation measure induced by policy π_{θ} .

■ Advantage

Training time roughly reduces linearly as the number of workers increases.

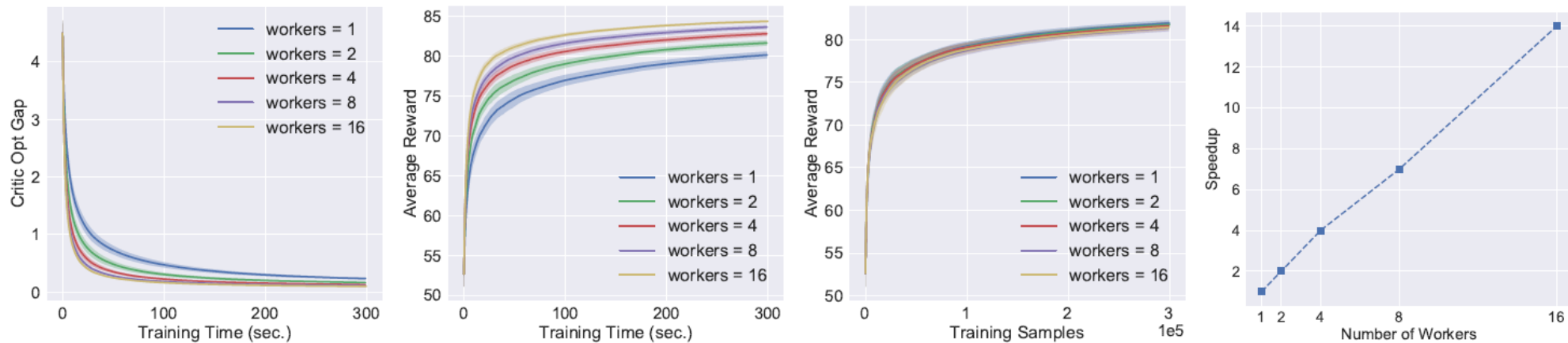


Figure 3: Convergence results of A3C with i.i.d. sampling in synthetic environment.

A3C

- Advantage

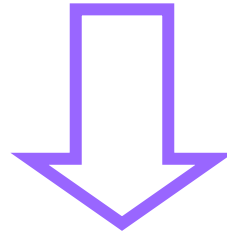
It operates in both discrete and continuous action spaces.

- Disadvantage

The delayed parameters introduce extra error.

Critic step $w_{k+1} = \Pi_{R_w} (w_k + \beta g(x_k, w_k))$

Actor step $\theta_{k+1} = \theta_k + \alpha [v(x_k, \theta_k, w_k) + \lambda \psi_\theta(x^p)]$



- error

$$g(x, \omega_{k-\tau_k}) - g(x, \omega_k)$$

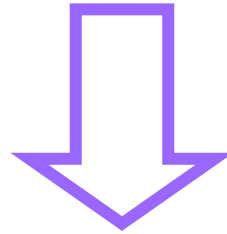
$$v(x, \theta_{k-\tau_k}, \omega_{k-\tau_k}) - v(x, \theta_k, \omega_k)$$

Critic step $\omega_{k+1} = \Pi_{R_\omega} (\omega_k + \beta g(x_{(k)}, \omega_{k-\tau_k}))$

Actor step $\theta_{k+1} = \theta_k + \alpha (v(\hat{x}_{(k)}, \theta_{k-\tau_k}, \omega_{k-\tau_k}) + \lambda \psi_{\theta_{k-\tau_k}}(x_{(k)}^p))$

Critic step $w_{k+1} = \Pi_{R_w} (w_k + \beta g(x_k, w_k))$

Actor step $\theta_{k+1} = \theta_k + \alpha [v(x_k, \theta_k, w_k) + \lambda \psi_\theta(x^p)]$



- error

$$g(x, \omega_{k-\tau_k}) - g(x, \omega_k)$$

$$v(x, \theta_{k-\tau_k}, \omega_{k-\tau_k}) - v(x, \theta_k, \omega_k)$$

Critic step $\omega_{k+1} = \Pi_{R_\omega} (\omega_k + \beta g(x_{(k)}, \omega_{k-\tau_k}))$

Actor step $\theta_{k+1} = \theta_k + \alpha (v(\hat{x}_{(k)}, \theta_{k-\tau_k}, \omega_{k-\tau_k}) + \lambda \psi_{\theta_{k-\tau_k}}(x_{(k)}^p))$

Thank you for listening