

Reinforcement Learning

Chapter 6: Actor Critic Methods

Ali Bereyhi

`ali.bereyhi@utoronto.ca`

Department of Electrical and Computer Engineering
University of Toronto

Fall 2025

Extension to Other PGMs

*In practice, all PGMs we studied in Chapter 5 can be implemented via the **actor-critic** idea through the following general framework*

Loop over the following three steps

- ① Specify **policy** and **value** networks
 - ↳ They could be either **separate** DNNs or DNNs with **shared layers**
- ② Set the policy and sample a batch of trajectories
- ③ Go over the batch for multiple epochs
 - ↳ After each **mini-batch** estimate the **policy** and **value** gradient
 - ↳ Update both **policy** and **value** networks after each **mini-batch**

Let's now look at the TRPO and PPO algorithms in **actor-critic** framework

TRPO: Actor-Critic

TRPO():

- 1: Initiate with θ and \mathbf{w} , as well as factor $\alpha < 1$ and learning rate β
- 2: **while** *interacting* **do**
- 3: Sample a *batch of trajectories* $S_0, A_0 \xrightarrow{R_1} \dots \xrightarrow{R_T} S_T$ by policy π_θ
- 4: **for** multiple epochs **do**
- 5: **for** samples in each *mini-batch* **do**
- 6: Compute *sample advantage* using $v_{\mathbf{w}}(\cdot)$
- 7: Update policy *gradient* $\hat{\nabla}_\theta$ and *value gradient* $\hat{\nabla}_{\mathbf{w}}$
- 8: **end for**
- 9: Compute a *Hessian estimator* $\hat{\mathbf{H}}$ and solve $\hat{\mathbf{H}}\mathbf{y} = \hat{\nabla}$ for \mathbf{y}
- 10: *Backtrack on a line* to find minimum i satisfying $\bar{D}_{\text{KL}}(\pi_{\theta'} \parallel \pi_\theta) \leq d_{\max}$

$$\theta' \leftarrow \theta + \alpha^i \sqrt{\frac{2d_{\max}}{\mathbf{y}^\top \hat{\mathbf{H}} \mathbf{y}}} \mathbf{y}$$

- 11: Update $\theta \leftarrow \theta'$ and $\mathbf{w} \leftarrow \mathbf{w} + \beta \hat{\nabla}_{\mathbf{w}}$
- 12: **end for**
- 13: **end while**

Recall: Use Importance Sampling

Attention: Importance Sampling

It is important to *remember* that in each iteration of PGM

we estimate the *policy* gradient via *importance sampling*

Let's denote the policy of current mini-batch with π_{θ} : in next *mini-batch* we compute the *policy* gradient as


$$\hat{\nabla}_{\theta} \leftarrow \hat{\nabla}_{\theta} + \sum_{t=0}^T U_t \frac{\nabla \pi_{\mathbf{x}} (A_t | S_t) |_{\mathbf{x}=\theta}}{\pi_{\theta} (A_t | S_t)}$$

with U_t being the sample advantage of policy π_{θ}

Recall: Use Importance Sampling

- + You mentioned this before! What is new about this?!
- Well! We should also consider it in our value estimation!

Denote the policy that we sampled with in line 3 with $\pi_{\theta_{\text{old}}}$: after multiple **mini-batches** we have an updated **policy** gradient π_{θ}

 To update the value network in this mini-batch, we consider the **Bellman equation** which says

$$v_{\pi_{\theta}}(S_t) = \mathbb{E}_{\pi_{\theta}} \{R_{t+1} + \gamma v_{\pi_{\theta}}(S_{t+1})\}$$

and set the labels for **value network training** as

$$\hat{v}_{\pi_{\theta}}(S_t) = R_{t+1} + \gamma v_{\pi_{\theta}}(S_{t+1})$$

But this is only valid if we had sampled the trajectory by π_{θ} !

Recall: Use Importance Sampling

- + Shall we use *importance sampling* here as well?!
- Sure!

By *importance sampling* we could say

$$\begin{aligned} v_{\pi_{\theta}}(S_t) &= \mathbb{E}_{\pi_{\theta}} \{R_{t+1} + \gamma v_{\pi_{\theta}}(S_{t+1})\} \\ &= \mathbb{E}_{\pi_{\theta_{\text{old}}}} \left\{ (R_{t+1} + \gamma v_{\pi_{\theta}}(S_{t+1})) \frac{\pi_{\theta}(A_t|S_t)}{\pi_{\theta_{\text{old}}}(A_t|S_t)} \right\} \end{aligned}$$

and now we can compute value estimators from our sample trajectories as

$$\hat{v}_{\pi_{\theta}}(S_t) = \left(R_{t+1} + \gamma v_{\pi_{\theta_{\text{old}}}}(S_{t+1}) \right) \frac{\pi_{\theta}(A_t|S_t)}{\pi_{\theta_{\text{old}}}(A_t|S_t)}$$

Recall: Use Importance Sampling

This means that the advantage should be computed as

$$\begin{aligned}
 U_t &= \underbrace{\left(R_{t+1} + \gamma v_{\pi_{\theta_{\text{old}}}}(S_{t+1}) \right)}_{\text{samples in batch}} \underbrace{\frac{\pi_{\theta}(A_t|S_t)}{\pi_{\theta_{\text{old}}}(A_t|S_t)}}_{\text{importance sampling}} - v_{\mathbf{w}}(S_t) \\
 &\approx \left(R_{t+1} + \gamma v_{\pi_{\theta_{\text{old}}}}(S_{t+1}) - v_{\pi_{\theta_{\text{old}}}}(S_t) \right) \frac{\pi_{\theta}(A_t|S_t)}{\pi_{\theta_{\text{old}}}(A_t|S_t)} \\
 &= U_t^{\text{old}} \frac{\pi_{\theta}(A_t|S_t)}{\pi_{\theta_{\text{old}}}(A_t|S_t)}
 \end{aligned}$$

This is *the correct way* of estimating *advantage*: other implementations that ignore this will have *bias* especially with *too much epochs*

Recall: Use Importance Sampling

Consequently, the value gradient is updated as

$$\begin{aligned}\hat{\nabla}_{\mathbf{w}} &\leftarrow \hat{\nabla}_{\mathbf{w}} + \frac{1}{T} \sum_{t=0}^{T-1} U_t \nabla v_{\mathbf{w}}(S_t) \\ &\leftarrow \hat{\nabla}_{\mathbf{w}} + U_t^{\text{old}} \frac{\pi_{\boldsymbol{\theta}}(A_t | S_t)}{\pi_{\boldsymbol{\theta}_{\text{old}}}(A_t | S_t)} \nabla v_{\mathbf{w}}(S_t)\end{aligned}$$

at the end of each trajectory

- + Shall we also consider this when we *update the policy*?
- Of course!

Recall: Use Importance Sampling

In a given *mini-batch* we update the *policy* gradient as

$$\begin{aligned}\hat{\nabla}_{\theta} &\leftarrow \hat{\nabla}_{\theta} + \sum_{t=0}^T U_t \frac{\nabla \pi_{\mathbf{x}}(A_t | S_t) |_{\mathbf{x}=\theta}}{\pi_{\theta}(A_t | S_t)} \\ &\leftarrow \hat{\nabla}_{\theta} + \sum_{t=0}^T U_t^{\text{old}} \frac{\pi_{\theta}(A_t | S_t)}{\pi_{\theta_{\text{old}}}(A_t | S_t)} \frac{\nabla \pi_{\mathbf{x}}(A_t | S_t) |_{\mathbf{x}=\theta}}{\pi_{\theta}(A_t | S_t)} \\ &\leftarrow \hat{\nabla}_{\theta} + \sum_{t=0}^T U_t^{\text{old}} \frac{\nabla \pi_{\mathbf{x}}(A_t | S_t) |_{\mathbf{x}=\theta}}{\pi_{\theta_{\text{old}}}(A_t | S_t)}\end{aligned}$$

which is now consistent with *importance sampling*

TRPO: Actor-Critic

TRPO_AC() :

- 1: Initiate with $\theta = \theta_{\text{old}}$ and \mathbf{w} , as well as factor $\alpha < 1$ and learning rate β
- 2: **while** *interacting* **do**
- 3: Sample a *batch of trajectories* $S_0, A_0 \xrightarrow{R_1} \dots \xrightarrow{R_T} S_T$ by policy $\pi_{\theta_{\text{old}}}$
- 4: Compute *sample advantage* U_t^{old} using $v_{\mathbf{w}}(\cdot)$
- 5: **for** multiple epochs in *each mini-batch* **do**
- 6: Compute gradient estimators $(\hat{\nabla}_{\theta}, \hat{\nabla}_{\mathbf{w}})$ from U_t^{old} via *importance sampling*
- 7: Compute a *Hessian estimator* $\hat{\mathbf{H}}$ and solve $\hat{\mathbf{H}}\mathbf{y} = \hat{\nabla}_{\theta}$ for \mathbf{y}
- 8: *Backtrack on a line* to find minimum i satisfying $\bar{D}_{\text{KL}}(\pi_{\theta'} \parallel \pi_{\theta}) \leq d_{\text{max}}$

$$\theta' \leftarrow \theta + \alpha^i \sqrt{\frac{2d_{\text{max}}}{\mathbf{y}^T \hat{\mathbf{H}} \mathbf{y}}} \mathbf{y}$$

- 9: Update $\theta \leftarrow \theta'$ and $\mathbf{w} \leftarrow \mathbf{w} + \beta \hat{\nabla}_{\mathbf{w}}$
- 10: **end for**
- 11: Update $\pi_{\theta} \leftarrow \pi_{\theta_{\text{old}}}$
- 12: **end while**

PPO: Actor-Critic

In PPO, we maximize in each iteration the *restricted surrogate*

$$\tilde{\mathcal{L}}(\pi_{\mathbf{x}}) = \mathbb{E}_{\pi_{\theta}} \left\{ \min \left\{ U_t \frac{\pi_{\mathbf{x}}(A_t|S_t)}{\pi_{\theta}(A_t|S_t)}, \ell_{\varepsilon}(U_t) \right\} \right\}$$

Similar to TRPO, we can estimate *restricted surrogate* via *importance sampling*

$$\begin{aligned} \tilde{\mathcal{L}}(\pi_{\mathbf{x}}) &= \text{mean} \left[\sum_t \min \left\{ U_t \frac{\pi_{\mathbf{x}}(A_t|S_t)}{\pi_{\theta}(A_t|S_t)}, \ell_{\varepsilon}(U_t) \right\} \right] \\ &= \text{mean} \left[\sum_t \min \left\{ U_t^{\text{old}} \frac{\pi_{\mathbf{x}}(A_t|S_t)}{\pi_{\theta_{\text{old}}}(A_t|S_t)}, \ell_{\varepsilon}(U_t) \right\} \right] \end{aligned}$$

where the *mean* is computed over the sample trajectories of a *mini-batch*

PPO Algorithm: Actor-Critic

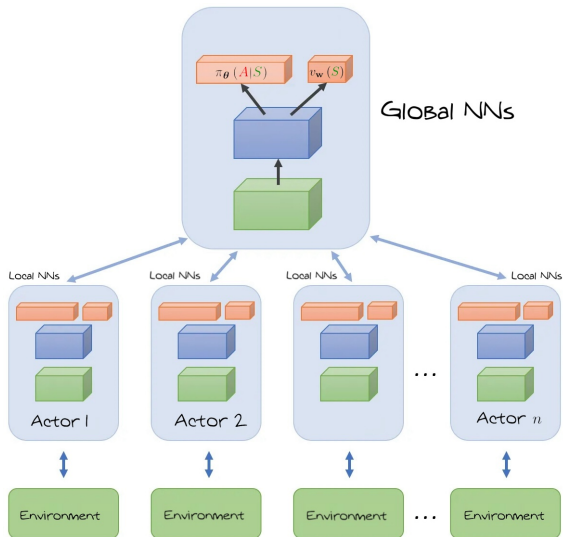
PPO_AC() :

- 1: Initiate with $\theta = \theta_{\text{old}}$ and \mathbf{w} , as well as learning rates α and β
- 2: **while** *interacting* **do**
- 3: Sample a *batch of trajectories* $S_0, A_0 \xrightarrow{R_1} \dots \xrightarrow{R_T} S_T$ by policy $\pi_{\theta_{\text{old}}}$
- 4: Compute *sample advantage* U_t^{old} using $v_{\mathbf{w}}(\cdot)$
- 5: **for** multiple epochs in *each mini-batch* **do**
- 6: Compute *value* gradient estimator $\hat{\nabla}_{\mathbf{w}}$ from U_t^{old} via *importance sampling*
- 7: Compute the *restricted surrogate*

$$\tilde{\mathcal{L}}(\pi_{\mathbf{x}}) = \text{mean} \left[\sum_t \min \left\{ U_t^{\text{old}} \frac{\pi_{\mathbf{x}}(A_t|S_t)}{\pi_{\theta_{\text{old}}}(A_t|S_t)}, \ell_{\varepsilon}(U_t) \right\} \right]$$

- 8: Update $\theta \leftarrow \theta + \alpha \nabla \tilde{\mathcal{L}}(\pi_{\mathbf{x}})|_{\mathbf{x}=\theta}$ and $\mathbf{w} \leftarrow \mathbf{w} + \beta \hat{\nabla}_{\mathbf{w}}$
- 9: **end for**
- 10: Update $\pi_{\theta} \leftarrow \pi_{\theta_{\text{old}}}$
- 11: **end while**

Distributed Actor-Critic



Distributed Setting: Asynchronous vs Synchronous

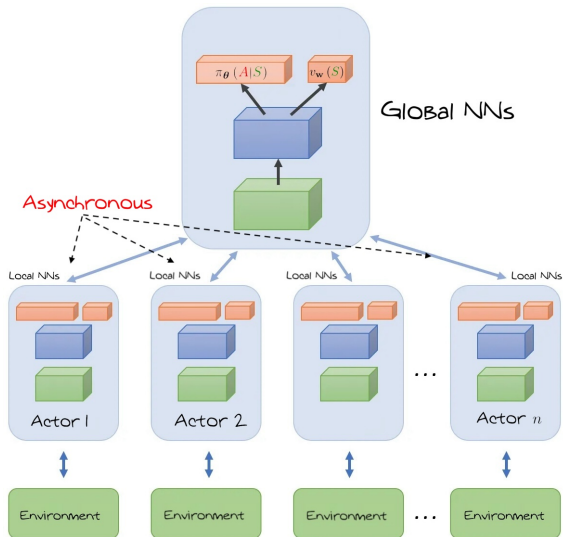
We can implement **actor-critic** approaches in a **distributed** fashion

- Multiple actors few samples with local policy and value network
 - ↳ Each actor focuses on a **specific part of environment**
- They share their gradient estimators with the **server**
- **Server** treats the collected estimators as of a large **mini-batch**
 - ↳ It updates its networks on this large **mini-batch**
- All actors update their local networks every time **server** shares its networks

We can implement this setting

- **Synchronous**
 - ↳ This is basically the **same** as what we do in A2C, TRPO or PPO
- **Asynchronous**
 - ↳ **Server** does **not** wait for all actors to send their estimators
 - ↳ It uses **what it has** every couple of rounds and remaining in next rounds

A3C: Asynchronous A2C



Some Final Remarks

It turns out that *asynchronous* update can negatively impact *convergence*

- A3C is hence not really extended to other PGMs
- In practice we usually implement *actor-critic* approaches in *synchronous distributed* form

-
- + Is that it? Are we free to go now?!
 - Pretty much Yes! Just we may further take a look at *deterministic policy gradient* approaches as well
 - + Is it a new set of approaches?!
 - **No!** It's a specific form of *actor-critic* methods that are *better compatible* with *continuous actions*