

---

# Your Project Title (e.g. Replication: Policy Optimization with Demonstrations)

---

**Author names**

Department of Computer Science  
National Yang Ming Chiao Tung University  
{xxx, yyy, zzz}@nycu.edu.tw

## 1 Problem Overview

This paper proposed an algorithm called ACER that combines experience replay and TRPO-like gradient update.

The goal of the problem is to maximize the discounted return  $R_t = \sum_{i \geq 0} \gamma^i r_{t+i}$  in expectation

Notation: For the value function:  $V^\pi(x_t) = \mathbb{E}_{a_t} [Q^\pi(x_t, a_t) | x_t]$

And for the action-value function:  $Q^\pi(x_t, a_t) = \mathbb{E}_{x_{t+1:\infty}, a_{t+1:\infty}} [R_t | x_t, a_t]$

where the actions are determined by policy  $\pi$

ACER estimates its policy  $\pi_\theta(a_t|x_t)$  and value function  $V_{\theta_v}^\pi(x_t)$  with deep neural networks.

## 2 Background and The Algorithm

The original policy gradient is computed as follows:

$$g = E_{x_{0:\infty}, a_{0:\infty}} [\sum_{t \geq 0} A_t^\pi(x_t, a_t) \nabla_\theta \log \pi_\theta(a_t|x_t)]$$

Off-policy learning with experience replay may appear to be an obvious strategy for improving the sample efficiency of actor-critic methods. However, controlling the variance and stability of off-policy estimators is notoriously hard. Importance sampling is one of the most popular approaches for off policy learning. With off-policy sampling, the gradient should be approximate as (Eq 4):

$$g^{marg} = E_{x_t \sim \beta, a_t \sim \mu} [\rho_t \nabla_\theta \log \pi_\theta(a_t|x_t) Q^\pi(x_t, a_t)]$$

They adapted behavior policy  $\mu$  from memory experience as policy and marginal importance weight by importance sampling.

In the following subsection, they adopt the Retrace algorithm to estimate  $Q^\pi$ , propose an importance weight truncation technique to improve the stability of the off-policy actor critic, and introduce a computationally efficient trust region scheme for policy optimization.

Below is the retrace estimator (Eq 5):

$$Q^{\text{ret}}(x_t, a_t) = r_t + \gamma \bar{\rho}_{t+1} [Q^{\text{ret}}(x_{t+1}, a_{t+1}) - Q(x_{t+1}, a_{t+1})] + \gamma V(x_{t+1})$$

which is off-policy, low variance and is proven to converge to the value function of the target policy for any behavior policy. Here,  $\bar{\rho} = \min\{c, \rho_t\}$ , where  $\rho_t = \frac{\pi(a_t|x_t)}{\mu(a_t|x_t)}$ .

We compute Q in Eq(5) using the critic neural network prediction and substitute  $Q_\pi$  in Eq(4) with retrace estimator. Because Retrace is return-based, the advantages of this improvement are to reduce bias in the policy gradient, and to enable faster learning of the critic.

In order to avoid the high variance and bias, the paper decomposed the original equation and implement two methods. First, to deal with the high variance, it clipped the marginal importance weight ( $\rho_t$ ) into  $\bar{\rho}_t = \min(c, \rho_t)$ , so that the variance of the gradient estimate is bounded.

However with the gradient clipped, we need a term to preserve the relative magnitude of different gradient. Second, the paper added the correction term to ensure the estimate is unbiased:

$$\left[\frac{\rho_t(a) - c}{\rho_t(a)}\right]_+$$

So the gradient becomes:

$$g^{marg} = \mathbb{E}_{x_t} [\mathbb{E}_{a_t} [\bar{\rho}_t \nabla_{\theta} \log \pi_{\theta}(a_t | x_t) Q^{\pi}(x_t, a_t)] + \mathbb{E}_{a \sim \pi} \left( \left[\frac{\rho_t(a) - c}{\rho_t(a)}\right]_+ \nabla_{\theta} \log \pi_{\theta}(a | x_t) Q^{\pi}(x_t, a) \right)]$$

Notice that the sampling of (s, a) pair is under the effect of the behavior policy  $\mu$ . However the truncation has can cause the bias effect on the actual gradient, For example, there will be no difference for  $\rho$  with value 10000 or value 1 if c is with value 1. so they introduced a bias term to retain the relative effect of  $\rho$  on the gradient. The weight of the bias term is designed as:

$$\frac{\rho_t(a) - c}{\rho_t(a)}$$

To milden the effect of large  $\rho$  value causing large variance on gradient. Then, using the measure called "truncation with bias correction trick" to approximated the  $Q^{\pi}(x, a_t)$  with  $Q_{\theta_v}(x_t, a)$ . Also, approximate the expectation of Markov process by sampling trajectories from the behavior  $\mu$ . Last, subtracted the baseline  $V_{\theta_v}(x_t)$  to reduce variance. Hence, we can rewrite the off-policy ACER gradient as

$$\hat{g}_t^{acer} = \bar{\rho}_t \nabla_{\theta} \log \pi_{\theta}(a_t | x_t) [Q^{ret}(x_t, a_t) - V_{\theta_v}(x_t)] + \mathbb{E}_{a \sim \pi} \left( \left[\frac{\rho_t(a) - c}{\rho_t(a)}\right]_+ \nabla_{\theta} \log \pi_{\theta}(a | x_t) (Q_{\theta_v}(x_t, a) - V_{\theta_v}(x_t)) \right)$$

As the above, it turned to updated actor critic by Q estimates, when  $c = 0$ . And it turned to updated policy gradient by Retrace, when  $c = \infty$ .

In order to reduce high variance of the policy updates in actor-critic, they solve the problem by retraining stepsize. However, the existing method of TRPO is time consuming, they introduced a new TRPO method. The main difference between TRPO and improved TRPO is that they update the policy by running average of past policies which is the average policy network. There are two separate parts of average policy network output, one is the probability distribution  $f$  and the other is the parameters  $\phi_{\theta}$  of distribution  $f$ . The relationship between these two terms is:  $\phi_{\theta} : \pi(\cdot | x) = f(\cdot | \phi_{\theta}(x))$ . Each episode, they update policy by  $\phi_{\theta}$ . The following equation is the ACER policy gradient with respect to  $\phi$ :

$$\begin{aligned} \hat{g}_t^{acer} = & \bar{\rho}_t \nabla_{\phi_{\theta}(x_t)} \log f(a_t | \phi_{\theta}(x)) [Q^{ret}(x_t, a_t) - V_{\theta_v}(x_t)] \\ & + \mathbb{E}_{a \sim \pi} \left( \left[\frac{\rho_t(a) - c}{\rho_t(a)}\right]_+ \nabla_{\phi_{\theta}(x_t)} \log f(a_t | \phi_{\theta}(x)) [Q_{\theta_v}(x_t, a) - V_{\theta_v}(x_t)] \right) \end{aligned}$$

As for the trust region update, in the first stage, they solved optimization problem by KL divergence :

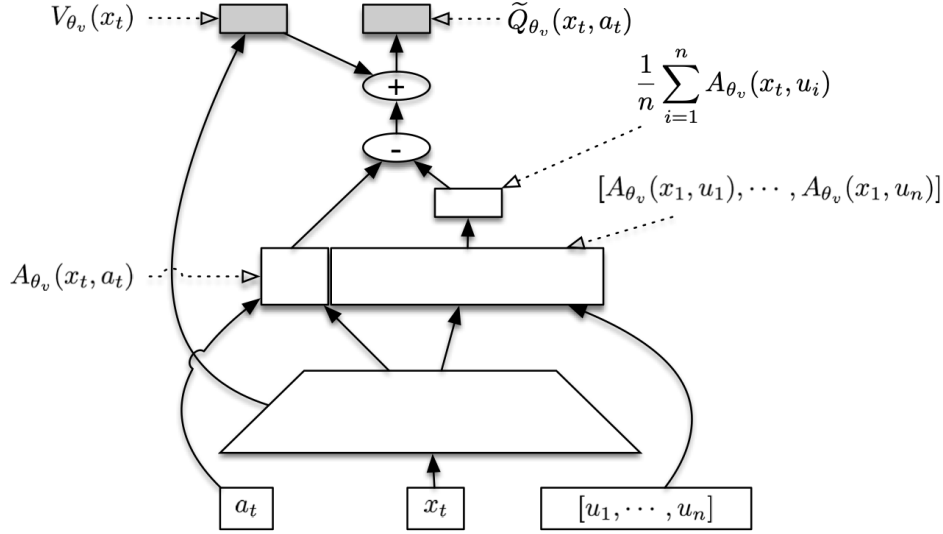
$$\begin{aligned} \underset{z}{\text{minimize}} \quad & \frac{1}{2} \|\hat{g}_t^{acer} - z\|_2^2 \\ \text{subject to} \quad & \nabla_{\phi_{\theta}(x_t)} D_{KL} [f(\cdot | \phi_{\theta_a}(x_t)) \| f(\cdot | \phi_{\theta}(x_t))]^T z \leq \delta \end{aligned}$$

Because the constraint is linear, they don't need to calculate the Hessian. With the advantages, it can be solved by the KKT condition easily. Letting  $k = \nabla_{\phi_{\theta}(x_t)} D_{KL} [f(\cdot | \phi_{\theta_a}(x_t)) \| f(\cdot | \phi_{\theta}(x_t))]$ , the solution is :

$$z^* = \hat{g}_t^{acer} - \max \left\{ 0, \frac{k^T \hat{g}_t^{acer} - \delta}{\|k\|_2^2} \right\} k$$

If the constraint is satisfied, the gradient would not change.

Figure 1: The Stochastic Dueling Network takes in a state and an action, then output the corresponding value and the action-value with the help of value and advantages



In the second stage, the trust region step is conducted in the space of the statistics of the distribution  $f$  rather than the space of the policy parameters. It can successfully avoid additional back-propagation by the policy network.

For the continuous case, since it is unpractical to compute value by summing up the action-values, we should deterministically get output from the critic network. To achieve this, the paper uses a Stochastic Dueling Network(SDN) to estimate both  $V^\pi$  and  $Q^\pi$  while maintaining consistency between two estimates.(see Figure 1)

The corresponding estimate for  $Q^\pi$  of SDN is computed as follows:

$$\tilde{Q}_{\theta_v}(x_t, a_t) \sim V_{\theta_v}(x_t) + A_{\theta_v}(x_t, a_t) - \frac{1}{n} \sum_{i=1}^n A_{\theta_v}(x_t, u_i), \text{ where } u_i \sim \pi(\cdot|x_t)$$

To train the network, the paper as well uses  $Q^{ret}$  for target action-value. For the state value  $V^\pi$ , they uses the following target:

$$V^{target}(x_t) = \min\{1, \frac{\pi(a_t|x_t)}{\mu(a_t|x_t)}\}(Q^{ret}(x_t, a_t) - V_{\theta_v}(x_t))$$

, which is derived by applying the truncation and bias correction trick like what have been done on the policy gradient update.

For continuous trust region updating, the paper replaced  $Q^{ret}$  with  $Q^{opc}$ , which is basically  $Q^{ret}$  without truncated importance ratio. Though might resulting in a less stable learning,  $Q^{opc}$  could better utilize the returns since its not truncated. The choice is made since by experiment they found out that  $Q^{opc}$  often leads to faster learning on continuous cases.

### 3 Detailed Implementation

Please explain your implementation in detail. You may do this with the help of pseudo code or a figure of system architecture. Please also highlight which parts of the algorithm lead to the most difficulty in your implementation.

## **4 Empirical Evaluation**

Please showcase your empirical results in this section. Please clearly specify which sets of experiments of the original paper are considered in your report. Please also report the corresponding hyperparameters of each experiment.

## **5 Conclusion**

Please provide succinct concluding remarks for your report. You may discuss the following aspects:

- The potential future research directions
- Any technical limitations
- Any latest results on the problem of interest

## **References**