

# SAMPLE EFFICIENT ACTOR-CRITIC WITH EXPERIENCE REPLAY

- 목적
  - stable
  - sample efficient
- 논문에서 제안한 여러가지 방법들
  - **truncated importance sampling with bias correction**
  - **stochastic dueling network architectures**
  - **a new trust region policy optimization method**
- Replay Buffer
  - DQN에서 처음 적용
  - sample correlation을 줄이기 위한 용도로 사용된 기술이지만, 실제로는 **sample efficiency**도 향상 시킨다

## Background and Problem Setup

$$\text{Max. } R_t = \mathbb{E}(\sum_{i \geq 0} \gamma^i r_{t+i})$$

$$Q^\pi(x_t, a_t) = \mathbb{E}_{x_{t+1}:\infty, a_{t+1}:\infty} [R_t | x_t, a_t]$$

$$V^\pi(x_t) = \mathbb{E}_{a_t} [Q^\pi(x_t, a_t) | x_t]$$

$$A^\pi(x_t, a_t) = Q^\pi(x_t, a_t) - V^\pi(x_t)$$

$$\mathbb{E}_{a_t} [A^\pi(x_t, a_t)] = 0$$

$$g = \mathbb{E}_{x_0:\infty, a_0:\infty} [\sum_{t \geq 0} A^\pi(x_t, a_t) \nabla_\theta \log \pi_\theta(a_t | x_t)] \quad (1)$$

### A3C

- trade-off bias and variance

$$\hat{g}^{\text{a3c}} = \sum_{t \geq 0} \left( \left( \sum_{i=0}^{k-1} \gamma^i r_{t+i} \right) + \gamma^k V_{\theta_v}^\pi(x_{t+k}) - V_{\theta_v}^\pi(x_t) \right) \nabla_\theta \log \pi_\theta(a_t | x_t) \quad (2)$$

### ACER

- A3C + serveral modification , new modules
- a single deep neural network to estimate the policy  $\pi_\theta(a_t | x_t)$  and value function  $V_{\theta_v}^\pi(x_t)$

# Discrete Actor Critic With Experience Replay

## off-policy learning with experience replay

- off-policy learning with experience replay은 actor-critics의 **sample efficiency** 를 향상 시킨다.
- 그러나 **off-policy**의 **variance**와 **stability**를 **control**하는 것은 어려운 일이다.
- **Importance sampling**은 off-policy learning의 가장 popular한 approach이다.

$$\hat{g}^{\text{imp}} = \left( \prod_{t=0}^K \rho_t \right) \sum_{t=0}^k \left( \sum_{i=0}^k \gamma^i r_{t+i} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | x_t) \quad (3)$$
$$\rho_t = \frac{\pi(a_t | x_t)}{\mu(a_t | x_t)}$$

- $\left( \prod_{t=0}^K \rho_t \right)$  은 **high variance**를 야기한다.

gradient를 approximation 하기 위해 **marginal value functions over the limiting distribution** 사용 해서 해결.

$$g^{\text{marg}} = \mathbb{E}_{x_t \sim \beta, a_t \sim \mu} [\rho_t \nabla_{\theta} \log \pi_{\theta}(x_t | x_t) Q^{\pi}(x_t, a_t)] \quad (4)$$

limiting distribution  $\beta(x) = \lim_{t \rightarrow \infty} p(x_t = x | x_0, \mu)$  with behavior policy  $\mu$

극한분포는 이산 또는 연속 시간 확률과정에서 시간이 무한대로 갈 때, 확률과정의 분포가 일정한 분포를 가지는 경우 이를 주어진 확률과정의 극한분포라고 한다.

여기서 중요한 점 두가지.

1.  $Q^u$  대신에  $Q^{\pi}$  를 사용했다. 따라서  $Q^{\pi}$  를 추정해야 한다.
2. importance weight의 product 가 없고 marginal importance weight  $\rho_t$  를 추정할 필요가 있다.

Off-Policy Actor-Critic 논문에서는  $R_t^{\lambda} = r_t + (1 - \lambda)\gamma V(x_{t+1}) + \lambda\gamma\rho_{t+1}R_{t+1}^{\lambda}$  라는 재귀식을 통해  $Q^{\pi}$  를 계산하는데  $\rho_t$  를 계속 곱해주기 때문에 학습이 불안정해 줄 수 있다.

## Multi-Step Estimation of The State-Action Value Function

- 이 논문에서는 Retrace(Munos et al., 2016 [Safe and Efficient Off-Policy Reinforcement Learning](#)) 방법을 사용해서  $Q^{\pi}(x_t, a_t)$  를 추정한다. ( $\rho$  대신에  $\bar{\rho}$  를 쓰는 것만으로도 variance가 낮아진다고 한다.)

$$Q^{\text{ret}}(x_t, a_t) = r - t + \gamma\bar{\rho}_{t+1}[Q^{\text{ret}}(x_{t+1}, a_{t+1}) - Q(x_{t+1}, a_{t+1})] + \gamma V(x_{t+1}) \quad (5)$$

$\bar{\rho}_t$  는 **truncated importance weight** 라 한다.  $\bar{\rho}_t = \min \{c, \rho_t\}$

$Q$  는  $Q^\pi$  의 current value estimate 고  $V(x) = \mathbb{E}_{a \sim \pi} Q(x, a)$

- **Retrace**는 **low variance**를 갖고 수렴이 보장된 **off-policy, return-based algorithm** 이다.
- $Q$  를 계산 하기 위해 discrete action space인 경우 "**two headed**"를 갖는 **convolutional neural network** 적용했다. (  $Q_{\theta_v}(x_t, a_t)$  와  $\pi_\theta(a_t|x_t)$  를 동시에 추정하기 위해 )
- Retrace 는 multistep returns를 사용하기 때문에 , **bias**를 줄인다.
- critic  $Q_{\theta_v}(x_t, a_t)$  를 학습하기 위해  $Q^{\text{ret}}(x_t, a_t)$  를 target으로 MSE 를 사용했고 parameter  $\theta_v$  를 업데이트 하기 위해 다음과 같은 standard gradient를 사용했다.

$$(Q^{\text{ret}}(x_t, a_t) - Q_{\theta_v}(x_t, a_t)) \nabla_{\theta_v} Q_{\theta_v}(x_t, a_t) \quad (6)$$

**The purpose of the multi-step estimator  $Q^{\text{ret}}$**

- to reduce bias in the policy gradient.
- to enable faster learning of the critic, hence further reducing bias.

## Importance Weight Truncation with Bias Correction

- (식 4)에서 marginal importance weight는 커질 수 있어서, instability를 야기한다. 즉 식 (3)에서 식 (4)로 넘어오며 importance weight에 대한 곱셈항을 제거했지만,  $\rho_t$  가 **unbounded**라는 사실은 변함이 없기 때문에 여전히 학습이 불안정해질 수 있는 요소가 남아 있다.
- high variance에 대해 safe-guard를 하기위해
  - importance weight를 truncate하고 다음과 같이  $g^{\text{marg}}$  를 correction term을 도입해서 나눠준다.

$$\begin{aligned} g^{\text{marg}} &= \mathbb{E}_{x_t, a_t} [\rho_t \nabla_{\theta} \log \pi_{\theta}(x_t|x_t) Q^{\pi}(x_t, a_t)] \\ &= \mathbb{E}_{x_t} \left[ \mathbb{E}_{a_t} [\bar{\rho}_t \nabla_{\theta} \log \pi_{\theta}(x_t|x_t) Q^{\pi}(x_t, a_t)] + \mathbb{E}_{a \sim \pi} \left( \left[ \frac{\rho_t(a) - c}{\rho_t(a)} \right]_+ \nabla_{\theta} \log \pi_{\theta}(x_t|x_t) Q^{\pi}(x_t, a) \right) \right] \quad (7) \end{aligned}$$

- (수식 7)의 앞의 부분은 the **importance weight** 를 **clipping** 하여 gradient estimate 의 variance가 bound되게 한다.
- (수식 7)의 뒷 부분 correction term은  $\rho_t(a) > c$  일 때 active된다.

corret term 의  $Q^{\pi}(x_t, a)$  은 neural network approximation  $Q_{\theta_v}(x_t, a)$  로 모델링 한다.

## Truncation with bias correction trick

- variance를 줄여 주기 위해 advantage 사용

$$\bar{g}^{\text{marg}} = \mathbb{E}_{x_t} \left[ \mathbb{E}_{a_t} [\bar{\rho}_t \nabla_{\theta} \log \pi_{\theta}(x_t | x_t) Q^{\text{ret}}(x_t, a_t)] + \mathbb{E}_{a \sim \pi} \left( \left[ \frac{\rho_t(a) - c}{\rho_t(a)} \right]_+ \nabla_{\theta} \log \pi_{\theta}(x_t | x_t) Q_{\theta_v}(x_t, a) \right) \right] \quad (8)$$

(식 8)은 Markov process의 stationary distribution에 대해 expectation을 포함하고 있는데 이것은 sampling trajectories로 approximation할 수 있다.

$$\begin{aligned} \hat{g}^{\text{acer}} &= \bar{\rho}_t \nabla_{\theta} \log \pi_{\theta}(x_t | x_t) [Q^{\text{ret}}(x_t, a_t) - V_{\theta_v}(x_t)] \\ &\quad + \mathbb{E}_{a \sim \pi} \left( \left[ \frac{\rho_t(a) - c}{\rho_t(a)} \right]_+ \nabla_{\theta} \log \pi_{\theta}(x_t | x_t) [Q_{\theta_v}(x_t, a) - V_{\theta_v}(x_t)] \right) \end{aligned} \quad (9)$$

## Efficient Trust Region Policy Optimization

- The policy updates of actor-critic methods do often **exhibit high variance**
- To ensure stability, we must **limit the per-step changes to the policy**.
- **TRPO**
  - requires repeated computation of Fisher-vector products for each update. (prohibitively expensive in large domains)
- **average policy network**
  - a **running average of past policies**.
  - forces the updated policy **to not deviate far from this average**.
- policy network를 distribution  $f$  와 이 distribution의 statistics  $\phi_{\theta}(x)$  를 generate 하는 deep neural network 로 나눈다. 즉  $f$  가 주어지면 policy는  $\phi_{\theta} : \pi(\cdot | x) = f(\cdot | \phi_{\theta}(x))$  에 의해 characterized 된다.
  - 예)  $f$  는 statistics로 probability vector  $\phi_{\theta}(x)$  를 갖는 categorical distribution으로 선택할 수 있다.
- $\theta : \theta_a \leftarrow \alpha \theta_a + (1 - \alpha) \theta$

$$\begin{aligned} \hat{g}^{\text{acer}} &= \bar{\rho}_t \nabla_{\phi_{\theta}(x_t)} \log f(a_t | \phi_{\theta_t}(x)) [Q^{\text{ret}}(x_t, a_t) - V_{\theta_v}(x_t)] \\ &\quad + \mathbb{E}_{a \sim \pi} \left( \left[ \frac{\rho_t(a) - c}{\rho_t(a)} \right]_+ \nabla_{\phi_{\theta_t}(x_t)} \log f(a_t | \phi_{\theta_t}(x)) [Q_{\theta_v}(x_t, a) - V_{\theta_v}(x_t)] \right) \end{aligned} \quad (10)$$

- averaged policy network 가 있을 때, 제안된 trust region 업데이트는 두 단계를 거친다.
  - 선형화된 KL divergence 제약식을 갖는 optimization 문제를 푼다

$$\begin{aligned}
& \underset{z}{\text{minimize}} && \frac{1}{2} \|\hat{g}^{\text{acer}} - z\|_2^2 \\
& \text{subject to} && \nabla_{\phi_\theta(x_t)} D_{\text{KL}}[f(\cdot|\theta_a(x_t)) \| f(\cdot|\phi_\theta(x_t))]^T z \leq \delta \quad (11)
\end{aligned}$$

- 제약식이 선형이기 때문에, overall optimization problem 은 simple quadratic programming problem으로 reduce 할 있는데, 이것의 solution은 KKT codition을 사용한 closed 형태로 쉽게 derived할 수 있다.

$$z^* = \hat{g}_t^{\text{acer}} - \max \left\{ 0, \frac{k^T \hat{g}_t^{\text{acer}} - \delta}{\|k\|_2^2} \right\} k \quad (12)$$

## ACER Pseudo-Code for Discrete Actions

## A ACER PSEUDO-CODE FOR DISCRETE ACTIONS

---

### Algorithm 1 ACER for discrete actions (master algorithm)

---

// Assume global shared parameter vectors  $\theta$  and  $\theta_v$ .

// Assume ratio of replay  $r$ .

**repeat**

    Call ACER on-policy, Algorithm 2.

$n \leftarrow \text{Poisson}(r)$

**for**  $i \in \{1, \dots, n\}$  **do**

        Call ACER off-policy, Algorithm 2.

**end for**

**until** Max iteration or time reached.

---



---

### Algorithm 2 ACER for discrete actions

---

Reset gradients  $d\theta \leftarrow 0$  and  $d\theta_v \leftarrow 0$ .

Initialize parameters  $\theta' \leftarrow \theta$  and  $\theta'_v \leftarrow \theta_v$ .

**if not** On-Policy **then**

    Sample the trajectory  $\{x_0, a_0, r_0, \mu(\cdot|x_0), \dots, x_k, a_k, r_k, \mu(\cdot|x_k)\}$  from the replay memory.

**else**

    Get state  $x_0$

**end if**

**for**  $i \in \{0, \dots, k\}$  **do**

    Compute  $f(\cdot|\phi_{\theta'}(x_i))$ ,  $Q_{\theta'_v}(x_i, \cdot)$  and  $f(\cdot|\phi_{\theta_a}(x_i))$ .

**if** On-Policy **then**

        Perform  $a_i$  according to  $f(\cdot|\phi_{\theta'}(x_i))$

        Receive reward  $r_i$  and new state  $x_{i+1}$

$\mu(\cdot|x_i) \leftarrow f(\cdot|\phi_{\theta'}(x_i))$

**end if**

$\bar{\rho}_i \leftarrow \min \left\{ 1, \frac{f(a_i|\phi_{\theta'}(x_i))}{\mu(a_i|x_i)} \right\}$ .

**end for**

$Q^{ret} \leftarrow \begin{cases} 0 & \text{for terminal } x_k \\ \sum_a Q_{\theta'_v}(x_k, a) f(a|\phi_{\theta'}(x_k)) & \text{otherwise} \end{cases}$

**for**  $i \in \{k-1, \dots, 0\}$  **do**

$Q^{ret} \leftarrow r_i + \gamma Q^{ret}$

$V_i \leftarrow \sum_a Q_{\theta'_v}(x_i, a) f(a|\phi_{\theta'}(x_i))$

    Computing quantities needed for trust region updating:

$$\begin{aligned} g &\leftarrow \min \{c, \rho_i(a_i)\} \nabla_{\phi_{\theta'}(x_i)} \log f(a_i|\phi_{\theta'}(x_i)) (Q^{ret} - V_i) \\ &\quad + \sum_a \left[ 1 - \frac{c}{\rho_i(a)} \right]_+ f(a|\phi_{\theta'}(x_i)) \nabla_{\phi_{\theta'}(x_i)} \log f(a|\phi_{\theta'}(x_i)) (Q_{\theta'_v}(x_i, a) - V_i) \\ k &\leftarrow \nabla_{\phi_{\theta'}(x_i)} D_{KL} [f(\cdot|\phi_{\theta_a}(x_i)) \| f(\cdot|\phi_{\theta'}(x_i))] \end{aligned}$$

Accumulate gradients wrt  $\theta'$ :  $d\theta' \leftarrow d\theta' + \frac{\partial \phi_{\theta'}(x_i)}{\partial \theta'} \left( g - \max \left\{ 0, \frac{k^T g - \delta}{\|k\|_2^2} \right\} k \right)$

Accumulate gradients wrt  $\theta'_v$ :  $d\theta_v \leftarrow d\theta_v + \nabla_{\theta'_v} (Q^{ret} - Q_{\theta'_v}(x_i, a))^2$

Update Retrace target:  $Q^{ret} \leftarrow \bar{\rho}_i (Q^{ret} - Q_{\theta'_v}(x_i, a_i)) + V_i$

**end for**

Perform asynchronous update of  $\theta$  using  $d\theta$  and of  $\theta_v$  using  $d\theta_v$ .

Updating the average policy network:  $\theta_a \leftarrow \alpha \theta_a + (1 - \alpha) \theta$

---

## RESULTS ON ATARI

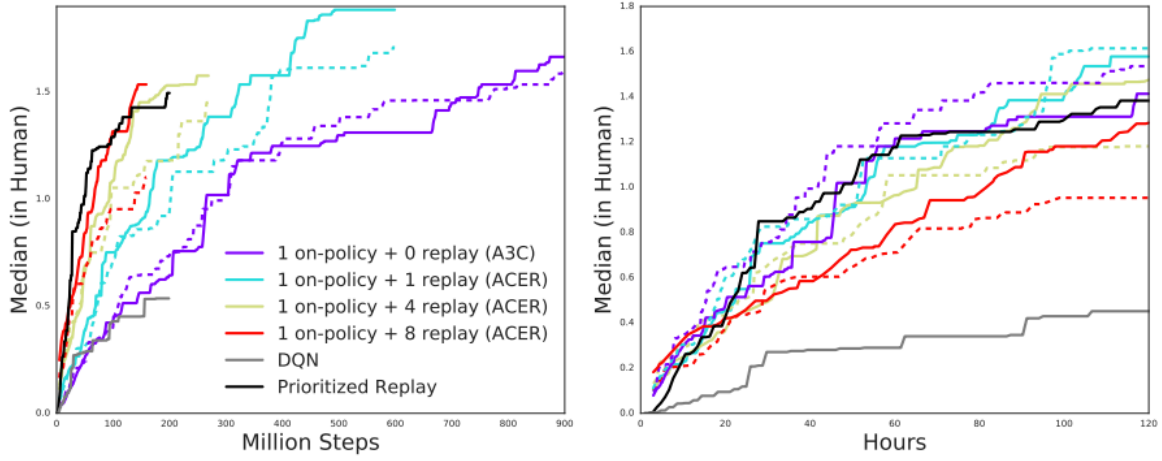


Figure 1: ACER improvements in sample (LEFT) and computation (RIGHT) complexity on Atari. On each plot, the median of the human-normalized score across all 57 Atari games is presented for 4 ratios of replay with 0 replay corresponding to on-policy A3C. The colored solid and dashed lines represent ACER with and without trust region updating respectively. The environment steps are counted over all threads. The gray curve is the original DQN agent (Mnih et al., 2015) and the black curve is one of the Prioritized Double DQN agents from Schaul et al. (2016).

## Continuous Actor Critic with Experience Replay

### Policy Evaluation

- Retrace는  $Q_{\theta_v}$  를 학습하기 위한 target를 제시하지만  $V_{\theta_v}$  에 대해서는 target를 제시하지 않는다.
- $Q_{\theta_v}$  가 주어졌을 때  $V_{\theta_v}$  를 계산하기 위해서 importance sampling을 사용하지만 이 추정치는 **high variacne**를 갖는다.
- **Stochastic Dueling Networks(SDN)**
  - $V^\pi$  와  $Q^\pi$  off-policy를 추정하기 위해 사용된 Dueling network 에 영감을 받음.
  - 매 time step 마다 SDN은  $Q^\pi$  에 대해  $\tilde{Q}_{\theta_v}$  로 **stochastic** 추정하고,  $V^\pi$  에 대해  $V_{\theta_v}$  **deterministic** 추정한다.

$$\tilde{Q}_{\theta_v}(x_t, a_t) \sim V_{\theta_v}(x_t) + A_{\theta_v}(x_t, a_t) - \frac{1}{n} \sum_{i=1}^n A_{\theta_v}(x_t, u_i), \quad \text{and} \quad u_i \sim \pi_\theta(\cdot | x_t) \quad (13)$$

여기서  $n$  은 parameter 다.

- $\mathbb{E}_{a \sim \pi(\cdot | x_t)} [\mathbb{E}_{u_1: n \sim \pi(\cdot | x_t)} (\tilde{Q}_{\theta_v}(x_t, a_t))] = V_{\theta_v}(x_t)$
- $\tilde{Q}_{\theta_v}$  를 학습함으로써  $V^\pi$  에 대해 학습할 수 있다.  $Q^\pi$  를  $\mathbb{E}_{u_1: n \sim \pi(\cdot | x_t)} (\tilde{Q}_{\theta_v}(x_t, a_t)) = Q^\pi(x_t, a_t)$  과 같이 완벽하게 학습했다고 가정하면  $V_{\theta_v}(x_t) = \mathbb{E}_{a \sim \pi(\cdot | x_t)} [\mathbb{E}_{u_1: n \sim \pi(\cdot | x_t)} (\tilde{Q}_{\theta_v}(x_t, a_t))] = \mathbb{E}_{a \sim \pi(\cdot | x_t)} [Q^\pi(x_t, a_t)] = V^\pi(x_t)$
- 그래서  $\tilde{Q}_{\theta_v}(x_t, a_t)$  에 대한 taget은  $V_{\theta_v}$  를 업데이트에 할 때 오류가 같이 전파된다.

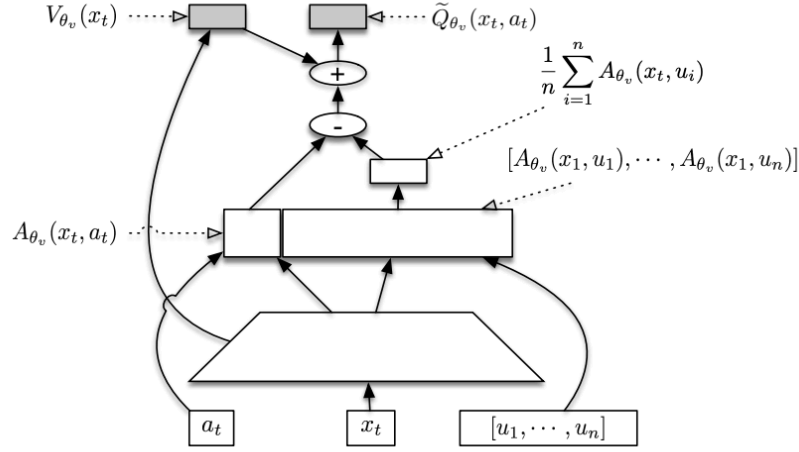


Figure 2: A schematic of the Stochastic Dueling Network. In the drawing,  $[u_1, \dots, u_n]$  are assumed to be samples from  $\pi_\theta(\cdot|x_t)$ . This schematic illustrates the concept of SDNs but does not reflect the real sizes of the networks used.

- SDN에 덧붙여서  $V^\pi$  를 추정하기 위해 다음과 같은 novel target을 만들었다.

$$V^{\text{target}}(x_t) = \min \left\{ 1, \frac{\pi(a_t|x_t)}{\mu(a_t|x_t)} \right\} (Q^{\text{ret}}(x_t, a_t) - Q_{\theta_v}(x_t, a_t)) + V_{\theta_v}(x_t) \quad (14)$$

- 마지막으로 continuous domain 에서  $Q^{\text{ret}}$  를 추정하기 위해, 조금 다른 truncated importance weights  $\bar{\rho}_t = \min \left\{ 1, \left( \frac{\pi(a_t|x_t)}{\mu(a_t|x_t)} \right)^{\frac{1}{d}} \right\}$  여기서  $d$  는 action space의 dimensionality이다.

## Trust Region Updating

- continuous action space에서  $g_t^{\text{acer}}$  를 유도하기 위해 stochastic dueling network에 대해 ACER policy gradient를 고려해 보자.  $\phi$  에 대해서

$$g_t^{\text{acer}} = \mathbb{E}_{x_t} \left[ \mathbb{E}_{a_t} \left[ \bar{\rho}_t \nabla_{\phi_{\theta}(x_t)} \log f(a_t|\phi_{\theta}(x_t)) (Q^{\text{opc}}(x_t, a_t) - V_{\theta_v}(x_t)) \right] + \mathbb{E}_{a \sim \pi} \left( \left[ \frac{\rho_t(a) - c}{\rho_t(a)} \right]_+ (\tilde{Q}_{\theta_v}(x_t, a) - V_{\theta_v}(x_t)) \nabla_{\phi_{\theta}(x_t)} \log f(a|\phi_{\theta}(x_t)) \right) \right] \quad (15)$$

- (식 15)에서는  $Q^{\text{ret}}$  대신에  $Q^{\text{opc}}$  를 사용했다.
- $Q^{\text{opc}}$  는 truncated importance ratio를 1로 대체 한다는 것을 제외하고는 Retrace 와 같다. (Appendix B 참조)
- Observation  $x_t$  가 주어졌을 때 다음과 같은 Monte Carlo approximation을 얻기 위해  $a'_t \sim \pi_\theta(\cdot|x_t)$  로 샘플링을 한다.

$$\hat{g}_t^{\text{acer}} = \bar{\rho}_t \nabla_{\phi_{\theta}(x_t)} \log f(a_t|\phi_{\theta}(x_t)) (Q^{\text{opc}}(x_t, a_t) - V_{\theta_v}(x_t)) + \left[ \frac{\rho_t(a'_t) - c}{\rho_t(a'_t)} \right]_+ (\tilde{Q}_{\theta_v}(x_t, a'_t) - V_{\theta_v}(x_t)) \nabla_{\phi_{\theta}(x_t)} \log f(a'_t|\phi_{\theta}(x_t)) \quad (16)$$

- $f$  와  $\hat{g}_t^{\text{acer}}$  가 주어 졌을 때 update를 완성하기 위해 "Discrete Actor Criti With Experience Replay - Efficient Trust Region Policy Optimization"에서 설명한 step을 따른다.



## $Q(\lambda)$ with Off-Policy Correctors

$$Q^{\text{opc}}(x_t, a_t) = r_t + \gamma [Q^{\text{opc}}(x_{t+1}, a_{t+1}) - Q(x_{t+1}, a_{t+1})] + \gamma V(x_{t+1}) \quad (21)$$

## Algorithm ACER for Continuous Actions

---

### Algorithm 3 ACER for Continuous Actions

---

Reset gradients  $d\theta \leftarrow 0$  and  $d\theta_v \leftarrow 0$ .  
Initialize parameters  $\theta' \leftarrow \theta$  and  $\theta'_v \leftarrow \theta_v$ .  
Sample the trajectory  $\{x_0, a_0, r_0, \mu(\cdot|x_0), \dots, x_k, a_k, r_k, \mu(\cdot|x_k)\}$  from the replay memory.  
**for**  $i \in \{0, \dots, k\}$  **do**  
    Compute  $f(\cdot|\phi_{\theta'}(x_i))$ ,  $V_{\theta'_v}(x_i)$ ,  $\tilde{Q}_{\theta'_v}(x_i, a_i)$ , and  $f(\cdot|\phi_{\theta_a}(x_i))$ .  
    Sample  $a'_i \sim f(\cdot|\phi_{\theta'}(x_i))$   
     $\rho_i \leftarrow \frac{f(a_i|\phi_{\theta'}(x_i))}{\mu(a_i|x_i)}$  and  $\rho'_i \leftarrow \frac{f(a'_i|\phi_{\theta'}(x_i))}{\mu(a'_i|x_i)}$   
     $c_i \leftarrow \min \left\{ 1, (\rho_i)^{\frac{1}{d}} \right\}$ .  
**end for**  
 $Q^{\text{ret}} \leftarrow \begin{cases} 0 & \text{for terminal } x_k \\ V_{\theta'_v}(x_k) & \text{otherwise} \end{cases}$   
 $Q^{\text{opc}} \leftarrow Q^{\text{ret}}$   
**for**  $i \in \{k-1, \dots, 0\}$  **do**  
     $Q^{\text{ret}} \leftarrow r_i + \gamma Q^{\text{ret}}$   
     $Q^{\text{opc}} \leftarrow r_i + \gamma Q^{\text{opc}}$   
    Computing quantities needed for trust region updating:  

$$g \leftarrow \min \{c, \rho_i\} \nabla_{\phi_{\theta'}(x_i)} \log f(a_i|\phi_{\theta'}(x_i)) (Q^{\text{opc}}(x_i, a_i) - V_{\theta'_v}(x_i))$$

$$+ \left[ 1 - \frac{c}{\rho'_i} \right]_+ (\tilde{Q}_{\theta'_v}(x_i, a'_i) - V_{\theta'_v}(x_i)) \nabla_{\phi_{\theta'}(x_i)} \log f(a'_i|\phi_{\theta'}(x_i))$$

$$k \leftarrow \nabla_{\phi_{\theta'}(x_i)} D_{KL} [f(\cdot|\phi_{\theta_a}(x_i)) \| f(\cdot|\phi_{\theta'}(x_i))]$$
  
    Accumulate gradients wrt  $\theta$ :  $d\theta \leftarrow d\theta + \frac{\partial \phi_{\theta'}(x_i)}{\partial \theta'} \left( g - \max \left\{ 0, \frac{k^T g - \delta}{\|k\|_2^2} \right\} k \right)$   
    Accumulate gradients wrt  $\theta'_v$ :  $d\theta_v \leftarrow d\theta_v + (Q^{\text{ret}} - \tilde{Q}_{\theta'_v}(x_i, a_i)) \nabla_{\theta'_v} \tilde{Q}_{\theta'_v}(x_i, a_i)$   

$$d\theta_v \leftarrow d\theta_v + \min \{1, \rho_i\} \left( Q^{\text{ret}}(x_t, a_i) - \tilde{Q}_{\theta'_v}(x_t, a_i) \right) \nabla_{\theta'_v} V_{\theta'_v}(x_i)$$
  
    Update Retrace target:  $Q^{\text{ret}} \leftarrow c_i \left( Q^{\text{ret}} - \tilde{Q}_{\theta'_v}(x_i, a_i) \right) + V_{\theta'_v}(x_i)$   
    Update Retrace target:  $Q^{\text{opc}} \leftarrow \left( Q^{\text{opc}} - \tilde{Q}_{\theta'_v}(x_i, a_i) \right) + V_{\theta'_v}(x_i)$   
**end for**  
Perform asynchronous update of  $\theta$  using  $d\theta$  and of  $\theta_v$  using  $d\theta_v$ .  
Updating the average policy network:  $\theta_a \leftarrow \alpha \theta_a + (1 - \alpha) \theta$

---

## Reulsts on MuJoCo

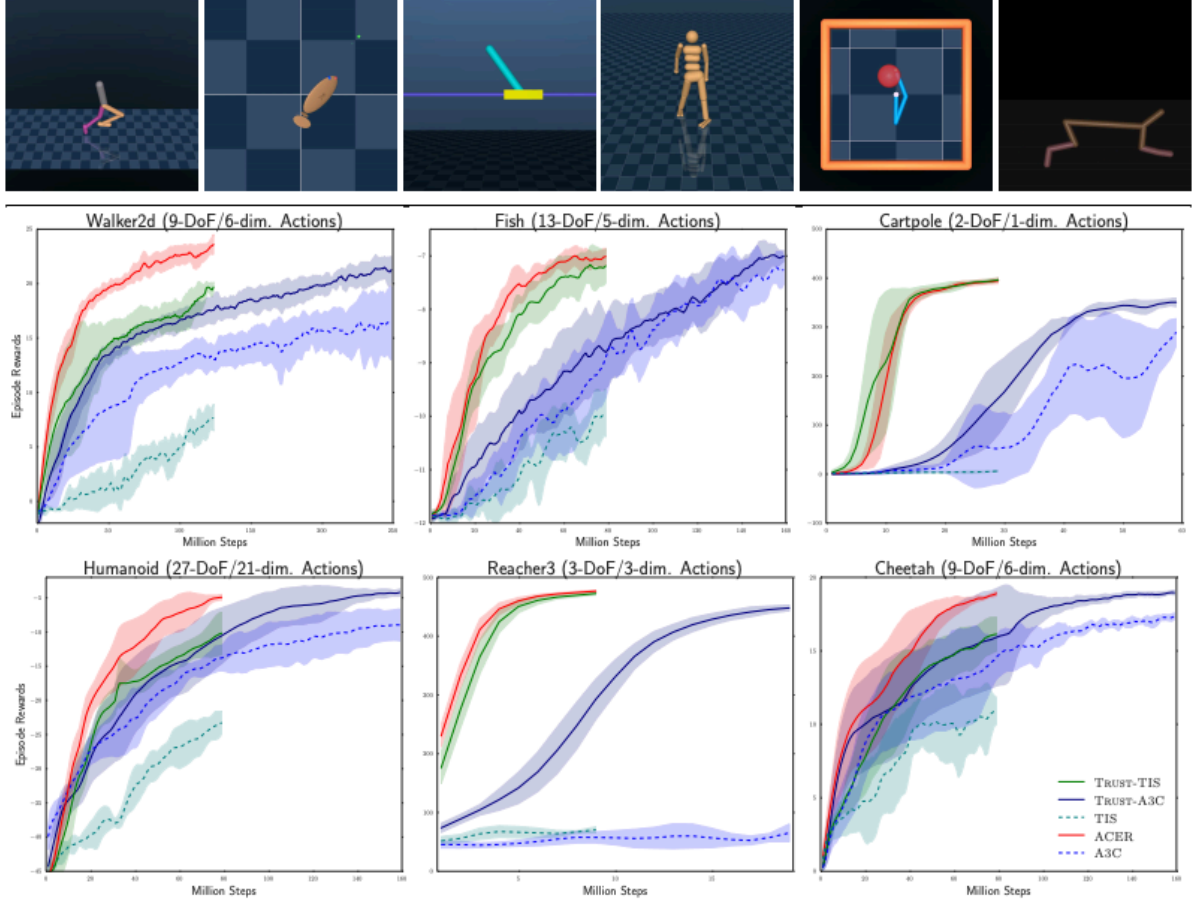


Figure 3: **[TOP]** Screen shots of the continuous control tasks. **[BOTTOM]** Performance of different methods on these tasks. ACER outperforms all other methods and shows clear gains for the higher-dimensionality tasks (humanoid, cheetah, walker and fish). The proposed trust region method by itself improves the two baselines (truncated importance sampling and A3C) significantly.

## Theoretical Analysis

- Retrace 가 이 논문에서 진전된 an application of the importance weight truncation와 bias correction trick 로 해석될 수 있음을 증명한다.
- 다음 수식을 고려해 보자

$$Q^\pi(x_t, a_t) = \mathbb{E}_{x_{t+1}a_{t+1}}[r_t + \gamma \rho_{t+1} Q^\pi(x_{t+1}, a_{t+1})] \quad (17)$$

- (식 17)을 얻기 위해 weight truncation 과 bias correction을 적용한다면

$$Q^\pi(x_t, a_t) = \mathbb{E}_{x_{t+1}a_{t+1}} \left[ r_t + \gamma \rho_{t+1} Q^\pi(x_{t+1}, a_{t+1}) + \gamma \mathbb{E}_{a \sim \pi} \left( \left[ \frac{\rho_{t+1}(a) - c}{\rho_{t+1}(a)} \right]_+ Q^\pi(x_{t+1}, a) \right) \right] \quad (18)$$

- (식 18) 에서  $Q^\pi$  를 recursively 하게 expanding 함으로써  $Q^\pi(x, a)$  는 다음과 같다.

$$Q^\pi(x, a) = \mathbb{E}_\mu \left[ \sum_{t \geq 0} \gamma^t \left( \prod_{i=1}^t \bar{\rho}_i \right) \left( r_t + \gamma \mathbb{E}_{b \sim \pi} \left( \left[ \frac{\rho_{t+1}(b) - c}{\rho_{t+1}(b)} \right]_+ Q^\pi(x_{t+1}, b) \right) \right) \right] \quad (19)$$

- expectation  $\mathbb{E}_\mu$  는  $\mu$  로 generate 한 actions을 취하는  $x$  에서 시작하는 trajectories에 취한다.

- $Q^\pi$  를 사용할 수 없을 때, current estimate  $Q$  로 대체한다.

$$\mathcal{B}Q(x, a) = \mathbb{E}_\mu \left[ \sum_{t \geq 0} \gamma^t \left( \prod_{i=1}^t \bar{\rho}_i \right) \left( r_t + \gamma \mathbb{E}_{b \sim \pi} \left( \left[ \frac{\rho_{t+1}(b) - c}{\rho_{t+1}(b)} \right]_+ Q(x_{t+1}, b) \right) \right) \right] \quad (20)$$

- 다음 명제는  $\mathcal{B}$  이 unique fiexe point  $Q^\pi$  로 contraction operator 라는 것을 보여준다.

**Proposition 1.** The operator  $\mathcal{B}$  is a contraction operator such that  $\|\mathcal{B}Q - Q^\pi\|_\infty \leq \gamma \|Q - Q^\pi\|_\infty$  and  $\mathcal{B}$  is equivalent to Retrace.

(Appendix C)

- Finally,  $\mathcal{B}$ , and therefore Retrace, generalizes both the Bellman operator  $\mathcal{T}^\pi$  and importance sampling.
- Specifically, when  $c = 0$ ,  $\mathcal{B} = \mathcal{T}^\pi$  and when  $c = \infty$ ,  $\mathcal{B}$  recovers importance sampling(Appendix C).

## Concluding Remarks

- continuous 과 discrete action spaces로 확장한 **a stable off-policy actor critic** 소개
- 다음 기법 사용
  - **truncated importance sampling with bias correction**
  - **stochastic dueling network architectures**
  - **a new trust region policy optimization method**