

# Week 6: Policy Optimization II

Bolei Zhou

*UCLA*

February 10, 2026

# Announcement

- ① Assignment 2 due by last week, and Assignment 3 is out
  - ① [https://github.com/ucla-rlcourse/cs260r-assignment-2026winter/  
blob/main/assignment3/assignment3.ipynb](https://github.com/ucla-rlcourse/cs260r-assignment-2026winter/blob/main/assignment3/assignment3.ipynb)
  - ② Implement TD3 and PPO

# Review of Policy Gradient

- ① Softmax policy: weight actions using a linear combination of features  
 $\phi(s, a)^T \theta$

$$\pi_\theta(s, a) = \frac{\exp^{\phi(s, a)^T \theta}}{\sum_{a'} \exp^{\phi(s, a')^T \theta}}$$

- ② Gaussian policy:
  - ① Mean is a linear combination of state features  $\mu(s) = \phi(s)^T \theta$
  - ② Variance may be fixed  $\sigma^2$  or can also be parameterized
  - ③ Policy is Gaussian, the continuous  $a \sim \mathcal{N}(\mu(s), \sigma^2)$
- ③ Stochastic policy can be considered as a probability (discrete action) or a probability density function (continuous action)
- ④ Above are the most common stochastic policy parameterizations to use

# Review of Policy Gradient

- ① In policy optimization, for the policy function  $\pi_\theta$  the objective is to maximize

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[R(\tau)]$$

- ①  $R(\tau)$  could be any reasonable reward function on  $\tau$
- ② So the gradient is

$$\begin{aligned}\nabla_\theta J(\theta) &= \nabla_\theta \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \mathbb{E}_\tau [R(\tau) \nabla_\theta \log P(\tau; \theta)]\end{aligned}$$

- ① **log-derivative trick:**  $\log \nabla_\theta f_\theta(x) = f_\theta(x) \nabla_\theta \log f_\theta(x)$
- ② Trajectory distribution:  $P(\tau; \theta) = \mu(s_0) \prod_{t=0}^{T-1} \pi_\theta(a_t | s_t) p(s_{t+1} | s_t, a_t)$
- ③ **cancelling trick:**  $\nabla_\theta \log P(\tau; \theta) = \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t)$

# Review of Policy Gradient

- ① After decomposing  $\tau$  we derived that

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[ \left( \sum_{t=0}^{T-1} r_t \right) \left( \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \right]$$

- ② After applying the causality, we have

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} G_t \cdot \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- ①  $G_t = \sum_{t'=t}^{T-1} r_{t'}$  is the return for a trajectory at step  $t$

# Summary of Policy Gradient

- ① The policy gradient has many forms

actual return

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) G_t] \text{ — REINFORCE} \\ &= \mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) Q^w(s, a)] \text{ — Q Actor-Critic} \\ &\quad \text{q function approximator (boot strapped estimate)} \\ &= \mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) A^w(s, a)] \text{ — Advantage Actor-Critic} \\ &\quad \text{q - v (more than average)} \\ &= \mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) \delta] \text{ — TD Actor-Critic}\end{aligned}$$

these use MC/TD learning  
to estimate Q/A

- ② Critic uses policy evaluation (e.g. MC or TD learning) to estimate  $Q^{\pi}(s, a)$ ,  $A^{\pi}(s, a)$ , or  $V^{\pi}(s, a)$

# Two Approaches to Reinforcement Learning

- ① Value-based RL: solve RL through dynamic programming
  - ① Classic RL and control theory
  - ② Representative algorithms: Deep Q-learning and its variant
  - ③ Representative researchers: Richard Sutton (no more than 20 pages on PG out of the 500-page textbook), David Silver, from DeepMind
- ② Policy-based RL: solve RL mainly through learning
  - ① Machine learning and deep learning
  - ② Representative algorithms: PG, and its variants TRPO, PPO
  - ③ Representative researchers: Pieter Abbeel, Sergey Levine, John Schulman, from OpenAI, Berkeley

# This Week's Plan: More Recent Advances

Two lines of work on policy optimization:

- ① Stochastic policy-based approaches: Policy Gradient → TRPO → ACKTR → PPO
- ② Deterministic policy-based (Value-based) approaches: Q-learning → DDPG → TD3
- ③ SAC: an algorithm that optimizes a stochastic policy in an off-policy way, forming a bridge between stochastic policy optimization and DDPG-style approaches

# More Recent Advances for Policy Optimization

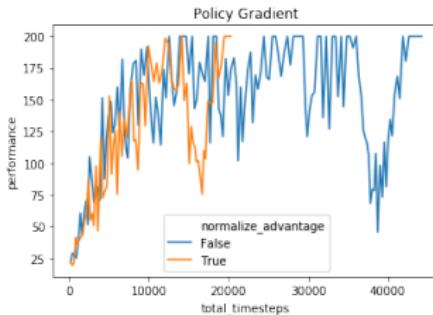
- ① Policy Gradient → TRPO → ACKTR → PPO → SAC
  - ① **TRPO**: Trust region policy optimization. Schulman, L., Moritz, Jordan, Abbeel. 2015
  - ② **ACKTR**: Scalable trust-region method for deep reinforcement learning using Kronecker-factored approximation. Y. Wu, E. Mansimov, S. Liao, R. Grosse, and J. Ba. 2017
  - ③ **PPO**: Proximal policy optimization algorithms. Schulman, Wolski, Dhariwal, Radford, Klimov. 2017
  - ④ **SAC**: Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor, Haarnoja et al., 2018
- ② Q-learning → DQN → DDPG → TD3
  - ① **DQN**: Human-level control through deep reinforcement learning, V Mnih et al., 2015
  - ② **DDPG**: Deterministic Policy Gradient Algorithms, Silver et al., 2014
  - ③ **TD3**: Addressing Function Approximation Error in Actor-Critic Methods, Fujimoto et al., 2018

# Problems with Policy Gradient

- ① Poor sample efficiency as PG is on-policy learning,

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

- ② Large policy update or improper step size destroy the training
  - ① This is different from supervised learning where the learning and data are independent
  - ② In RL, step too far  $\rightarrow$  bad policy  $\rightarrow$  bad data collection
  - ③ We may not be able to recover from a bad policy, which collapses the training



# Policy Gradient with Importance Sampling

- ① We can turn PG into off-policy learning using importance sampling
- ② **Importance sampling:** IS calculates the expected value of  $f(x)$  where  $x$  has a data distribution  $p$  we want to reuse old data and transitions
- ③ we can sample data from another distribution  $q$  and use the probability ratio between  $p$  and  $q$  to re-calibrate the result

$$\mathbb{E}_{x \sim p}[f(x)] = \mathbb{E}_{x \sim q} \left[ \frac{p(x)}{q(x)} f(x) \right]$$

$f(x)$  is like a reward function

$p(x)/q(x)$  we know the value,  
its just  $p$  is hard to sample from?

sample from  $q(x)$  which is easier to sample from

- ② Code example of IS: link
- ③ Using importance sampling in policy objective

$$J(\theta) = \mathbb{E}_{s,a \sim \pi_\theta} [r(s,a)] = \mathbb{E}_{s,a \sim \hat{\pi}} \left[ \frac{\pi_\theta(s,a)}{\hat{\pi}(s,a)} r(s,a) \right]$$

this is our current policy

$\hat{\pi}$  is a behavior policy

# Increasing the Robustness with Trust Regions

- ① The behavior policy could be the old policy directly, thus we can have a surrogate objective function

expectation can be estimated with samples

$$\theta = \arg \max_{\theta} J(\theta; \theta_{old}) = \arg \max_{\theta} \mathbb{E}_{\pi_{\theta_{old}}} \left[ \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} R_t \right]$$

- ② The estimate might be excessively large when  $\pi_{\theta}/\pi_{\theta_{old}}$  is too large
- ③ Solution: to limit the difference between subsequent policies
- ④ For instance, use the Kullbeck-Leibler (KL) divergence to measure the distance between two policies

works because policies are distributions over actions

$$KL(\pi_{\theta_{old}} || \pi_{\theta}) = \sum_a \pi_{\theta_{old}}(a|s) \log \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)}$$

this is like cross entropy loss

# Increasing the Robustness with Trust Regions

- ① Thus our objective with trust region becomes, to maximize

$$J(\theta; \theta_{old}) = \mathbb{E}_t \left[ \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} R_t \right]$$

$$\text{subject to } KL(\pi_{\theta_{old}}(\cdot | s_t) || \pi_\theta(\cdot | s_t)) \leq \delta$$

limit updates so that the KL divergence doesn't change too much

- ② In the trust region, we limit our parameter search within a region controlled by  $\delta$ . This is the intuition behind algorithms TRPO and PPO



Gradient ascend



Trust region

## Prior Work: Natural Policy Gradient

- ① Policy gradient is the steepest ascend in **parameter space** with Euclidean metric:

$$d^* = \nabla_{\theta} J(\theta) = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \max J(\theta + d), \text{s.t. } \|d\| \leq \epsilon$$

constrain updates to parameters by \epsilon  
this is on the parameters

this may change agent behavior a lot

- ① Drawback: it is sensitive to the parameterization of the policy function
- ② Natural (policy) gradient is the steepest ascend in **distribution space (policy output)** with KL-divergence as constraint

$$d^* = \arg \max J(\theta + d), \text{s.t. } KL[\pi_{\theta} || \pi_{\theta+d}] = c \quad (1)$$

this is on the output of the distribution

- ① Fixing KL-divergence as a constant  $c$  makes sure we move along the distribution space with constant speed, regardless the curvature (thus robust to the reparametrization of the model as we only care about the distribution induced by the parameter)

# KL-Divergence as the Distance between Two Distributions

- ① KL-divergence measures the “closeness” of two distributions.

$$KL[\pi_\theta || \pi_{\theta'}] = E_{\pi_\theta}[\log \pi_\theta] - E_{\pi_\theta}[\log \pi_{\theta'}]$$

- ② Although the KL-divergence is non-symmetric and thus not a true metric, we can use it anyway. This is because as  $d$  goes to zero, KL-divergence is asymptotically symmetric. So, within a local neighbourhood, KL-divergence is approximately symmetric
- ③ We can prove that the second-order Taylor expansion of KL-divergence is

F is the fischer information metric (second derivative of KL?)

$$KL[\pi_\theta || \pi_{\theta+d}] \approx \frac{1}{2} d^T F d$$

where  $F$  is the **Fisher Information Matrix** as second order derivative of KL divergence  $E_{\pi_\theta}[\nabla \log \pi_\theta \nabla \log \pi_\theta^T]$

# Natural Policy Gradient

$$d^* = \arg \max J(\theta + d), \text{ s.t. } KL[\pi_\theta || \pi_{\theta+d}] = c$$

- ① Write the above as Lagrangian form, with the  $J(\theta + d)$  approximated by its first order Taylor expansion and the constraint KL-divergence approximated by its second order Taylor expansion

$$\begin{aligned} d^* &= \arg \max_d J(\theta + d) - \lambda(KL[\pi_\theta || \pi_{\theta+d}] - c) \\ &\approx \arg \max_d J(\theta) + \nabla_\theta J(\theta)^T d - \frac{1}{2}\lambda d^T F d + \lambda c \end{aligned}$$

its because  $KL(P||P) = 0$ , so we need to take the second derivative so that we get the curvature

- ② To solve this maximization we set its derivative wrt.  $d$  to zero, then we have the **natural policy gradient**:  $d = \frac{1}{\lambda} F^{-1} \nabla_\theta J(\theta)$

# Natural Policy Gradient

- ① Natural policy gradient is a **second-order** optimization that is more accurate and works regardless of how the model is parameterized (model invariant).

$$\theta_{t+1} = \theta_t + F^{-1} \nabla_\theta J(\theta)$$

we have to calculate KL divergence (second order is much slower)  
moreover inverse is terrible for compute

- ① where  $F = E_{\pi_\theta(s,a)}[\nabla \log \pi_\theta(s, a) \nabla \log \pi_\theta(s, a)^T]$  is the Fisher information matrix, also as the second-order derivative of the KL-divergence
- ②  $F$  measures the curvature of the policy(distribution) relative to the model parameter  $\theta$
- ② Natural policy gradient produces the same policy change **regardless of the model parameterization**.
- ③ Idea behind TRPO. A detailed read on natural gradient:  
<https://agustinus.kristia.de/blog/natural-gradient/>

# Back to TRPO

- ① The objective of TRPO is to maximize

$$J(\theta; \theta_{old}) = \mathbb{E}_t \left[ \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} R_t \right]$$

$$\text{subject to } KL(\pi_{\theta_{old}}(\cdot | s_t) || \pi_\theta(\cdot | s_t)) \leq \delta$$

- ② In the trust region, we limit our parameter search within a region controlled by  $\delta$ . This is the intuition behind algorithms TRPO and PPO



Gradient ascend



Trust region

# Natural Policy Gradient part of TRPO

- ① Following Taylor's series expansion on both terms above up to the second-order
- ② After some derivations we can have

$$J(\theta; \theta_t) \approx g^T (\theta - \theta_t)$$

$$KL(\theta || \theta_t) \approx \frac{1}{2} (\theta - \theta_t)^T H (\theta - \theta_t)$$

where  $g = \nabla_{\theta} J_{\theta_t}(\theta)$  and  $H = \nabla_{\theta}^2 KL(\theta || \theta_t)$  and  $\theta_t$  is the old policy parameter

- ③ Then the objective becomes:

$$\theta_{t+1} = \arg \max_{\theta} g^T (\theta - \theta_t) \text{ s.t. } \frac{1}{2} (\theta - \theta_t)^T H (\theta - \theta_t) \leq \delta$$

- ④ This is a quadratic equation and can be solved analytically:

$$\theta_{t+1} = \theta_t + \sqrt{\frac{2\delta}{g^T H^{-1} g}} H^{-1} g$$

# Natural Policy Gradient part of TRPO

- ① Natural gradient is the steepest ascent direction with respect to the Fisher information

$$\theta_{t+1} = \theta_t + \sqrt{\frac{2\delta}{g^T H^{-1} g}} H^{-1} g$$

- ②  $H$  is the Fisher Information Matrix (FIM) which can be computed explicitly as

$$H = \nabla_\theta^2 KL(\pi_{\theta_t} || \pi_\theta) = E_{a,s \sim \pi_{\theta_t}} \left[ \nabla_\theta \log \pi_\theta(a, s) \nabla_\theta \log \pi_\theta(a, s)^T \right]$$

- ③ Learning rate ( $\delta$ ) can be thought of as choosing a step size that is normalized **with respect to the change in the policy**
- ④ This is beneficial because it means that any parameter update won't significantly change the output of the policy network

# Natural Policy Gradient part of TRPO

"update policy in a way that is invariant to how parameters are represented"

---

## Algorithm 1 Natural Policy Gradient

---

Input: initial policy parameters  $\theta_0$

for  $k = 0, 1, 2, \dots$  do

mc like policy -> get the rollouts

Collect set of trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Form sample estimates for this is the critic neural network (to calculate the value function, which we then use to calculate the q-value function)

- policy gradient  $\hat{g}_k$  (using advantage estimates) steepest ascent in the parameter space
- and KL-divergence Hessian / Fisher Information Matrix  $\hat{H}_k$

Compute Natural Policy Gradient update:

$$\theta_{k+1} = \theta_k + \sqrt{\frac{2\delta}{\hat{g}_k^T \hat{H}_k^{-1} \hat{g}_k}} \hat{H}_k^{-1} \hat{g}_k$$

This is the Natural Gradient

This is the Weight using the KL Divergence

end for

---

- ① Sham Kakade. "A Natural Policy Gradient." NIPS 2001

# Trust Region Policy Optimization (TRPO)

- ① FIM and its inverse are very expensive to compute
- ② TRPO estimates the term  $x = H^{-1}g$  by solving the following linear equation  $Hx = g$
- ③ Consider the optimization problem for a quadratic equation

Solving  $Ax = b$  is equivalent to

$$x = \arg \min_x f(x) = \frac{1}{2}x^T Ax - b^T x$$

since  $f'(x) = Ax - b = 0$

- ④ Thus we can optimize the quadratic equation as

$$\min_x \frac{1}{2}x^T Hx - g^T x$$

- ⑤ Use conjugate gradient method to solve it. It is very similar to the gradient ascent but can be done in fewer iterations

# Trust Region Policy Optimization (TRPO)

- ① Resulting algorithm is a refined version of natural policy gradient

---

**Algorithm 3** Trust Region Policy Optimization

---

Input: initial policy parameters  $\theta_0$

**for**  $k = 0, 1, 2, \dots$  **do**

    Collect set of trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

    Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

    Form sample estimates for

- policy gradient  $\hat{g}_k$  (using advantage estimates)

- and KL-divergence Hessian-vector product function  $f(v) = \hat{H}_k v$

    Use CG with  $n_{cg}$  iterations to obtain  $x_k \approx \hat{H}_k^{-1} \hat{g}_k$

    Estimate proposed step  $\Delta_k \approx \sqrt{\frac{2\delta}{x_k^T \hat{H}_k x_k}} x_k$

    Perform backtracking line search with exponential decay to obtain final update

$$\theta_{k+1} = \theta_k + \alpha^j \Delta_k$$

**end for**

---

get an approximation of  $\hat{H}_k^{-1} \hat{g}_k$  with Conjugate gradient method because its faster

# Appendix of Trust Region Policy Optimization (TRPO)

- ① Schulman, et al. ICML 2015: a lot of proofs
- ② The appendix A of the TRPO paper provides a 2-page proof that establishes the guaranteed monotonic improvement that **the policy update in each TRPO iteration creates a better policy**

$$J(\pi_{t+1}) - J(\pi_t) \geq M_t(\pi_{t+1}) - M(\pi_t)$$

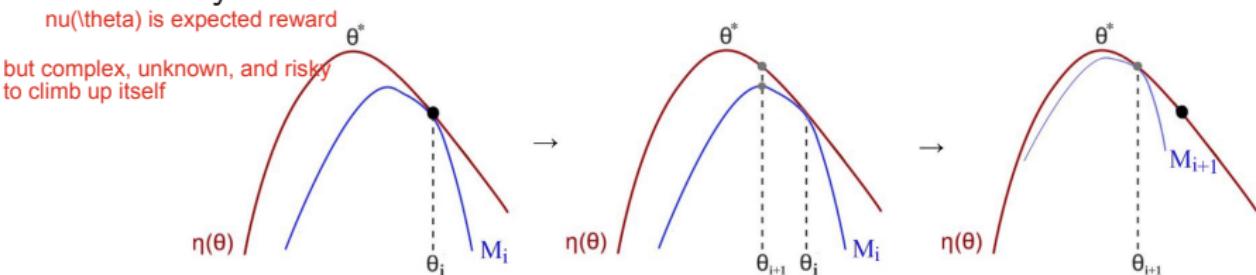
where  $M_t(\pi) = L_{\pi_t}(\pi) - \alpha D_{KL}(\pi_t, \pi)$

- ③ Thus by maximizing  $M_t$  at each iteration, we guarantee that the true objective  $J$  is non-decreasing

# Appendix of Trust Region Policy Optimization (TRPO)

maximize a surrogate function?

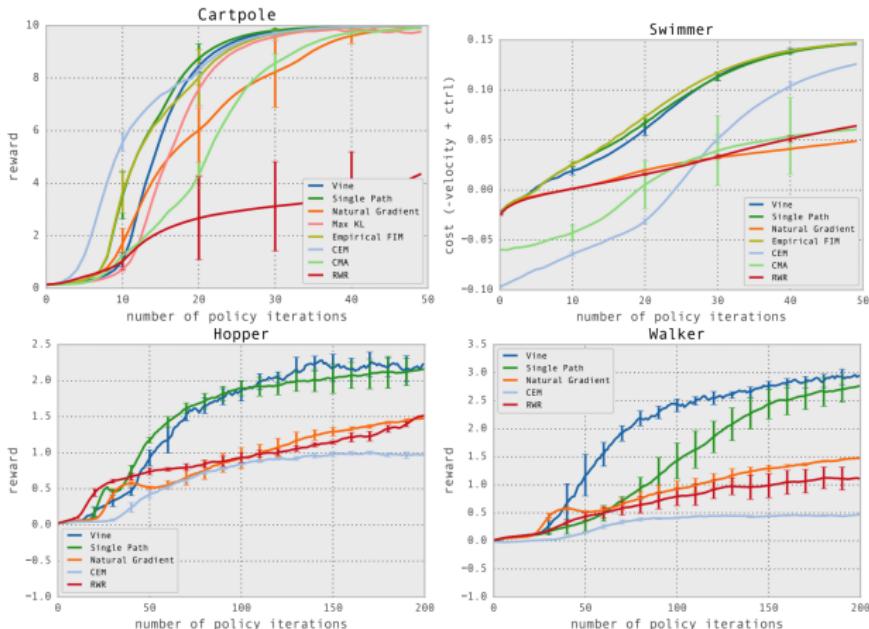
- ① It is a type of Minorize-Maximization (MM) algorithm which is a class of methods that includes expectation maximization
- ② The MM algorithm achieves this iteratively by maximizing surrogate function (the blue line below) approximating the expected reward locally.



instead we build  $M_i$  (blue curve), the surrogate function such that its always  $\leq$  than  $\eta$

$$M = L(\theta) = -C \text{KL}(\theta_{\text{old}} || \theta_{\text{new}})$$

# Result and Demo of TRPO



- 1 Demo video is at <https://www.youtube.com/embed/KJ15iGGJFvQ>

# Limitations of TRPO

## ① Scalability issue for TRPO

- ① Computing  $H$  every time for the current policy model is expensive
- ② It requires a large batch of rollouts to approximate  $H$ .

$$H = E_{s,a \sim \pi_{\theta_t}} \left[ (\nabla_{\theta} \log \pi_{\theta}(s, a))(\nabla_{\theta} \log \pi_{\theta}(s, a))^T \right]$$

- ③ Conjugate Gradient(CG) makes implementation more complicated
- ② TRPO does not work well on some of the tasks compared to DQN

	B. Rider	Breakout	Enduro	Pong	Q*bert	Seaquest	S. Invaders
Random Human (Mnih et al., 2013)	354 7456	1.2 31.0	0 368	-20.4 -3.0	157 18900	110 28010	179 3690
Deep Q Learning (Mnih et al., 2013)	4092	168.0	470	20.0	1952	1705	581
UCC-I (Guo et al., 2014)	5702	380	741	21	20025	2995	692
TRPO - single path	1425.2	10.8	534.6	20.9	1973.5	1908.6	568.4
TRPO - vine	859.5	34.2	430.8	20.9	7732.5	788.4	450.2

# ACKTR: Calculating Natural Gradient with KFAC

- ① Y. Wu, et al. "Scalable trust-region method for deep reinforcement learning using Kronecker-factored approximation". NIPS 2017.
- ② ACKTR speeds up the optimization by reducing the complexity of calculating the inverse of the Fisher Information Matrix (FIM) using the Kronecker-factored approximation curvature (K-FAC).

$$F = E_{x \sim \pi_{\theta_t}} \left[ (\nabla_{\theta} \log \pi_{\theta}(x))^T (\nabla_{\theta} \log \pi_{\theta}(x)) \right]$$

calculate Natural Gradient with a different approximation method

- ① It is replaced as layer-wise calculation

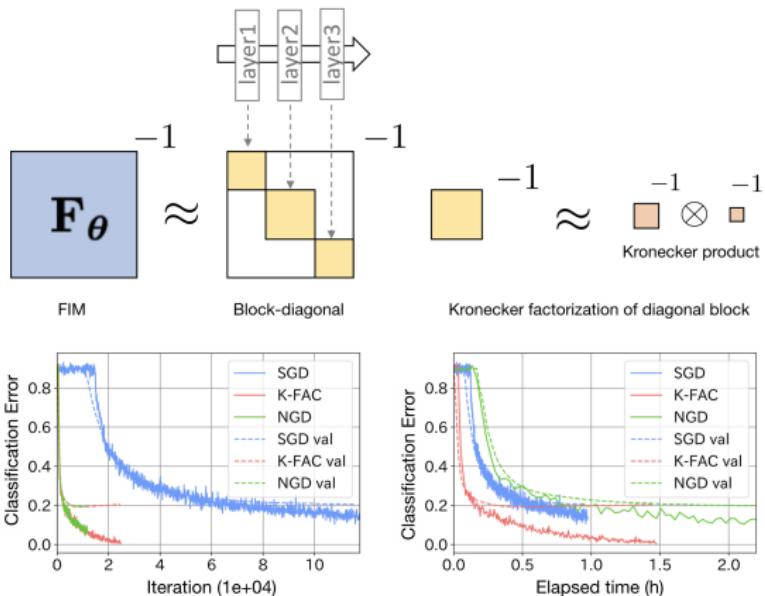
$$\begin{aligned} F_{\ell} &= \mathbb{E}[\text{vec}\{\nabla_W L\} \text{vec}\{\nabla_W L\}^T] = \mathbb{E}[aa^T \otimes \nabla_s L (\nabla_s L)^T] \\ &\approx \mathbb{E}[aa^T] \otimes \mathbb{E}[\nabla_s L (\nabla_s L)^T] := A \otimes S := \hat{F}_{\ell}, \end{aligned}$$

where  $A$  denotes  $\mathbb{E}[aa^T]$  and  $S$  denotes  $\mathbb{E}[\nabla_s L (\nabla_s L)^T]$ .

# Optimizing Neural Networks with Kronecker-factorized Approximate Curvature.

- ① Martens et al. ICML'15: <https://arxiv.org/pdf/1503.05671.pdf>
- ② Stochastic Gradient Descend (SGD) is the first-order optimization, which is widely used for network optimization; Natural Gradient Descend converges faster in terms of iterations than first-order method, by taking into consideration the curvature of the loss function. However it involves computing the inverse of the Fisher Information Matrix  $\theta_{t+1} = \theta_t + \alpha F^{-1} \nabla_\theta J(\theta)$ , where  $F = E_{\pi_\theta(s,a)}[\nabla \log \pi_\theta(s, a) \nabla \log \pi_\theta(s, a)^T]$
- ③ K-FAC (Kronecker-factorized Approximate Curvature) approximates Fisher Information Matrix for large-scale neural network optimization

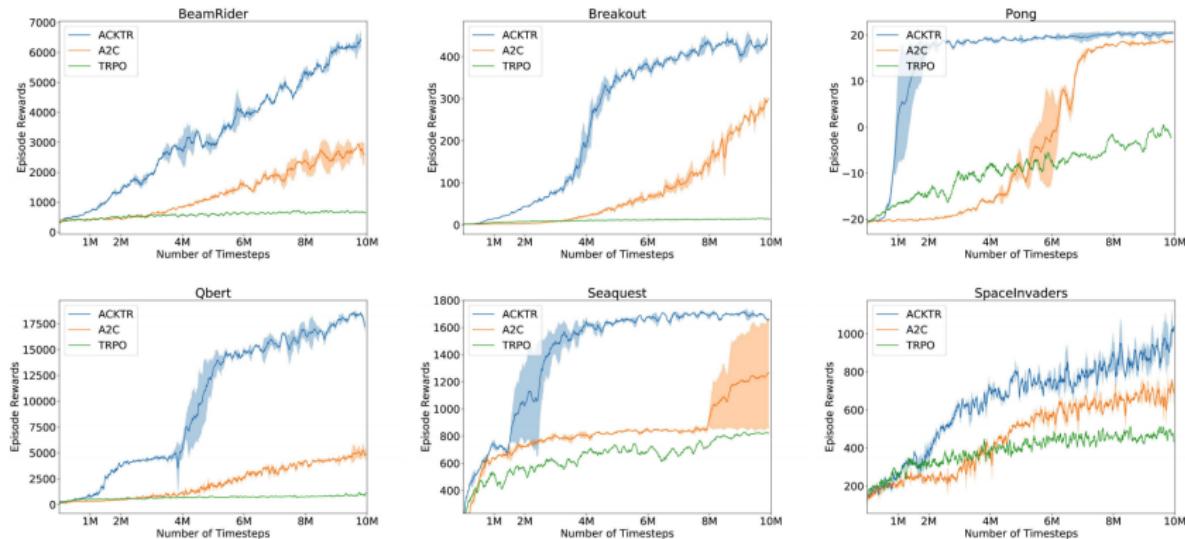
# Optimizing Neural Networks with Kronecker-factorized Approximate Curvature.



**Figure:** The comparison of training of ConvNet for CIFAR-10 dataset. Solid line: train, dashed line: validation. (data iteration vs. time) . Read more on K-FAC tutorial

# Performance of ACKTR

## ① Performance of ACKTR



## ② Introductory link:

<https://openai.com/index/openai-baselines-acktr-a2c/>

# Proximal Policy Optimization (PPO)

- ① The loss function in the Natural Policy Gradient and TRPO

$$\begin{aligned} & \text{maximize}_{\theta} \mathbb{E}_t \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} A_t \right] \\ & \text{subject to } \mathbb{E}_t [KL[\pi_{\theta_{old}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)]] \leq \delta \\ & \quad \text{Hard constraint, hard to implement in practice} \end{aligned}$$

how good the new policy is while following data from the old policy

- ② It can be also written as an unconstrained form,

$$\text{maximize}_{\theta} \mathbb{E}_{\pi_{\theta_{old}}} \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} A_t \right] - \beta \mathbb{E}_{\pi_{\theta_{old}}} [KL[\pi_{\theta_{old}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)]]$$

use this penalty instead to enforce the hard constraint

- ③ PPO: include the adaptive KL Penalty, so the optimization will have better insurance that we are optimizing within a trust region

# PPO with adaptive KL penalty

---

**Algorithm 4** PPO with Adaptive KL Penalty

---

Input: initial policy parameters  $\theta_0$ , initial KL penalty  $\beta_0$ , target KL-divergence  $\delta$

**for**  $k = 0, 1, 2, \dots$  **do**

    Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

    Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

    Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}(\theta) - \beta_k \bar{D}_{KL}(\theta || \theta_k)$$

    by taking  $K$  steps of minibatch SGD (via Adam)

**if**  $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \geq 1.5\delta$  **then**

$$\beta_{k+1} = 2\beta_k$$

**else if**  $\bar{D}_{KL}(\theta_{k+1} || \theta_k) \leq \delta / 1.5$  **then**

$$\beta_{k+1} = \beta_k / 2$$

**end if**

**end for**

---

- ① Same as or better performance than TRPO, but can be solved much faster by first-order optimization (SGD)

# PPO with clipping

- ① Let  $r_t(\theta)$  denote the probability ratio  $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$  so the different surrogate objectives are:
  - ① PG without trust region:  $L_t(\theta) = r_t(\theta)\hat{A}_t$
  - ② KL constraint:  $L_t(\theta) = r_t(\theta)\hat{A}_t$  s.t.  $KL[\pi_{\theta_{old}}, \pi_\theta] \leq \delta$
  - ③ KL penalty:  $L_t(\theta) = r_t(\theta)\hat{A}_t - \beta KL[\pi_{\theta_{old}}, \pi_\theta]$
- ② A new objective function to clip the estimated advantage function if the new policy is far away from the old policy ( $r_t$  is too large)

$$L_t(\theta) = \min \left( r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t \right)$$

- ① If the probability ratio between the new policy and the old policy falls outside the range  $(1 - \epsilon)$  and  $(1 + \epsilon)$ , the advantage function will be clipped.
- ②  $\epsilon$  is set to 0.2 for the experiments in the PPO paper.

# How the clipping works

- ① Clipping serves as a regularizer by removing incentives for the policy to change dramatically

$$L^{CLIP}(\theta) = \min \left( \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{A}_t, \text{clip} \left( \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t \right)$$

how much a ratio has changed  
must stay within this interval

- ② When the advantage is positive, we encourage the action thus  $\pi_\theta(a|s)$  increases,

we don't want to increase the Probability of taking the action (L) by too much

$$L(\theta; \theta_{old}) = \min \left( \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}, (1 + \epsilon) \right) \hat{A}_t$$

- ③ When the advantage is negative, we discourage the action thus  $\pi_\theta(a|s)$  decreases,

$$L(\theta; \theta_{old}) = \max \left( \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}, (1 - \epsilon) \right) \hat{A}_t$$

# PPO with clipping

---

**Algorithm 5** PPO with Clipped Objective

---

Input: initial policy parameters  $\theta_0$ , clipping threshold  $\epsilon$

**for**  $k = 0, 1, 2, \dots$  **do**

    Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

    Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

    Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}^{CLIP}(\theta)$$

    by taking  $K$  steps of minibatch SGD (via Adam), where

$$\mathcal{L}_{\theta_k}^{CLIP}(\theta) = \mathbb{E}_{\tau \sim \pi_k} \left[ \sum_{t=0}^T \left[ \min(r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^{\pi_k}) \right] \right]$$

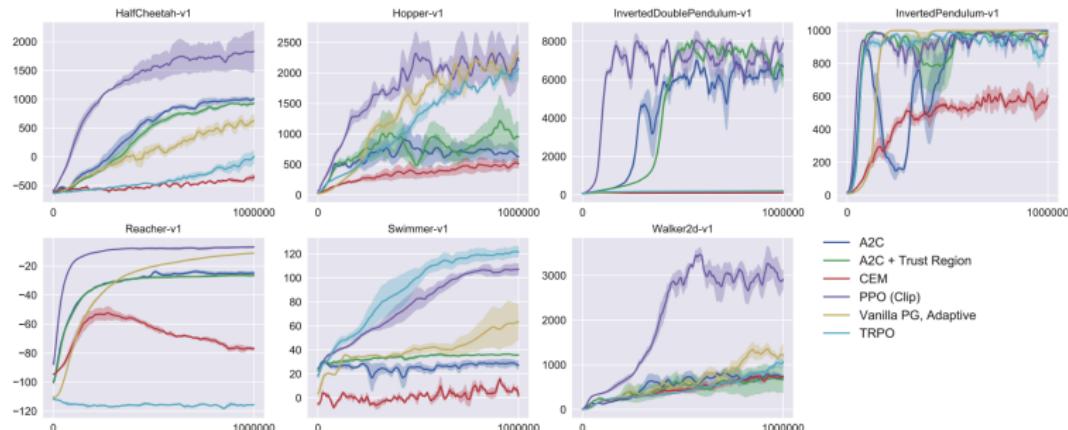
**end for**

---

- ① PPOs have the stability and reliability of trust-region methods but are much simpler to implement, requiring only few lines of code change to a vanilla policy gradient implementation

# Result of PPO

- ① Result on the continuous control tasks (MuJoCo)  
<https://gym.openai.com/envs/mujoco>



- ② Demo of PPO at  
<https://openai.com/index/openai-baselines-ppo/>
- ③ Emergence of Locomotion Behaviours in Rich Environments by DeepMind (Distributed PPO):  
[https://www.youtube.com/embed/hx\\_bgoTF7bs](https://www.youtube.com/embed/hx_bgoTF7bs)

# Code of PPO

- ① Paper link of PPO: <https://arxiv.org/abs/1707.06347>
- ② Code example: <https://github.com/ucla-rlcourse/DeepRL-Tutorials/blob/master/14.PPO.ipynb>
- ③ Very interesting one: **The 37 Implementation Details of Proximal Policy Optimization (ICLR'22 blog track)**
  - ① <https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/>
  - ② Devil is in the details

# State of the Art on Policy Optimization

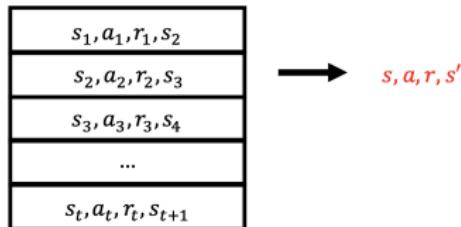
- ① **TRPO**: Schulman, L., Moritz, Jordan, Abbeel (2015). Trust region policy optimization
  - ① comment: Solid math proofs and guarantee, but hard to follow
- ② **ACKTR**: Y. Wu, E. Mansimov, S. Liao, R. Grosse, and J. Ba (2017). Scalable trust-region method for deep reinforcement learning using Kronecker-factored approximation.
  - ① comment: numeric optimization-based improvement, scalable to real-problems
- ③ **PPO**: Schulman, Wolski, Dhariwal, Radford, Klimov (2017). Proximal policy optimization algorithms
  - ① comment: Easy to read, elegant design of loss function, easy to implement, on-policy method

# Q-learning → DQN → DDPG → TD3

- ① **DQN**: Human-level control through deep reinforcement learning, V Mnih et al., Nature 2015
- ② **DDPG**: Deterministic Policy Gradient Algorithms, Silver et al. ICML 2014
- ③ **TD3**: Addressing Function Approximation Error in Actor-Critic Methods, Fujimoto et al. ICML 2018

# Review of DQN: Experience Replay

- To reduce the correlations among samples, store transition  $(s_t, a_t, r_t, s_{t+1})$  in replay memory  $\mathcal{D}$



- To perform experience replay, repeat the following
  - sample an experience tuple from the dataset:  $(s, a, r, s') \sim \mathcal{D}$
  - compute the target value for the sampled tuple:  $r + \gamma \max_{a'} \hat{Q}(s', a', \mathbf{w})$
  - use stochastic gradient descent to update the network weights

$$\Delta \mathbf{w} = \alpha \left( r + \gamma \max_{a'} \hat{Q}(s', a', \mathbf{w}) - Q(s, a, \mathbf{w}) \right) \nabla_{\mathbf{w}} \hat{Q}(s, a, \mathbf{w})$$

# Review of DQN: Fixed Targets

- ① To help improve stability, fix the target weights used in the target calculation for multiple updates
- ② Let a different set of parameter  $\mathbf{w}^-$  be the set of weights used in the target, and  $\mathbf{w}$  be the weights that are being updated
- ③ To perform experience replay with fixed target, repeat the following
  - ① sample an experience tuple from the dataset:  $(s, a, r, s') \sim \mathcal{D}$
  - ② compute the target value for the sampled tuple:  
 $r + \gamma \max_{a'} \hat{Q}(s', a', \mathbf{w}^-)$
  - ③ use stochastic gradient descent to update the network weights

$$\Delta \mathbf{w} = \alpha \left( r + \gamma \max_{a'} \hat{Q}(s', a', \mathbf{w}^-) - \hat{Q}(s, a, \mathbf{w}) \right) \nabla_{\mathbf{w}} \hat{Q}(s, a, \mathbf{w})$$

# Deep Deterministic Policy Gradient (DDPG)

this is deterministic

- ① Motivation: how to extend DQN to the environment with continuous action space?  
think driving a car for action space
- ② DDPG is very similar to DQN, which can be considered as a continuous action space version of DQN

check all actions and choose the one with the highest value

$$\text{DQN} : a^* = \arg \max_a Q^*(s, a)$$

$\mu_\theta$  is the action - outputs a specific action

$$\text{DDPG} : a^* = \arg \max_a Q^*(s, a) \approx \arg \max_a Q_\phi(s, \mu_\theta(s))$$

$\theta_\phi$  is the critic (computes how good an action is)

so this is like a regression?

- ① a deterministic policy  $\mu_\theta(s)$  directly gives the action that maximizes  $Q_\phi(s, a)$
- ② as action  $a$  is continuous we assume Q-function  $Q_\phi(s, a)$  is differentiable with respect to  $a$

1. input state  $S$
2. pass it to policy ( $\mu_\theta$ ) to get action
3. pass it to state  $S$  and action  $a$  into the Q-function ( $Q_\phi$ )
4. output the estimate value

# Deep Deterministic Policy Gradient (DDPG)

- ① Thus the DDPG has the following objective:

$$\textbf{Q-target: } y(r, s') = r + \gamma Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s'))$$

$$\textbf{Q-function update: } \min E_{s, r, s', d \sim \mathcal{D}} [(y(r, s') - Q_{\phi}(s, a))^2]$$

$$\textbf{policy update: } \max_{\theta} E_{s \sim \mathcal{D}} [Q_{\phi}(s, \mu_{\theta}(s))]$$

- ② Reply buffer and target networks for both value network and policy network are the same as DQN
- ③ DDPG example code (using the sample codebase for TD3):  
<https://github.com/sfujim/TD3/blob/master/DDPG.py>

# Twin Delayed DDPG (TD3)

- ① One drawback of DDPG is that the learned Q-function sometimes dramatically overestimates Q-values, which then leads to the breaking in training

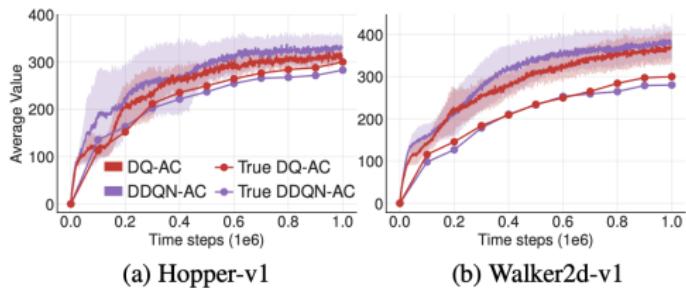


Figure 2. Measuring overestimation bias in the value estimates of actor critic variants of Double DQN (DDQN-AC) and Double Q-learning (DQ-AC) on MuJoCo environments over 1 million time steps.

- ① true value is estimated using the average discounted return over 1000 episodes following the current policy

# Twin Delayed DDPG (TD3)

## ① TD3 introduces three critical designs

- ① **Clipped Double-Q Learning.** TD3 learns two Q-functions instead of one (hence “twin”), and uses the smaller of the two Q-values to form the targets in the Bellman error loss functions.
- ② **“Delayed” Policy Updates.** TD3 updates the policy (and target networks) less frequently than the Q-function. The paper recommends one policy update for every two Q-function updates.
- ③ **Target Policy Smoothing.** TD3 adds noise to the target action, to make it harder for the policy to exploit Q-function errors by smoothing out Q along changes in action.  
*prevents overfitting*

# Twin Delayed DDPG (TD3)

aimed to fix over estimation bias (agent thinks actions are better than they actually are)

- ① TD3 concurrently learns two Q-functions,  $Q_{\phi_1}$  and  $Q_{\phi_2}$ , by mean square Bellman error minimization
- ② Both Q-functions use a single target, calculated using whichever of the two Q-functions gives a smaller value as the following **Q-target**:

$$y(r, s') = r + \gamma \min_{i=1,2} Q_{\phi_{i,targ}}(s', a_{TD3}(s'))$$

standard Q-learning typically overestimates a Q-value

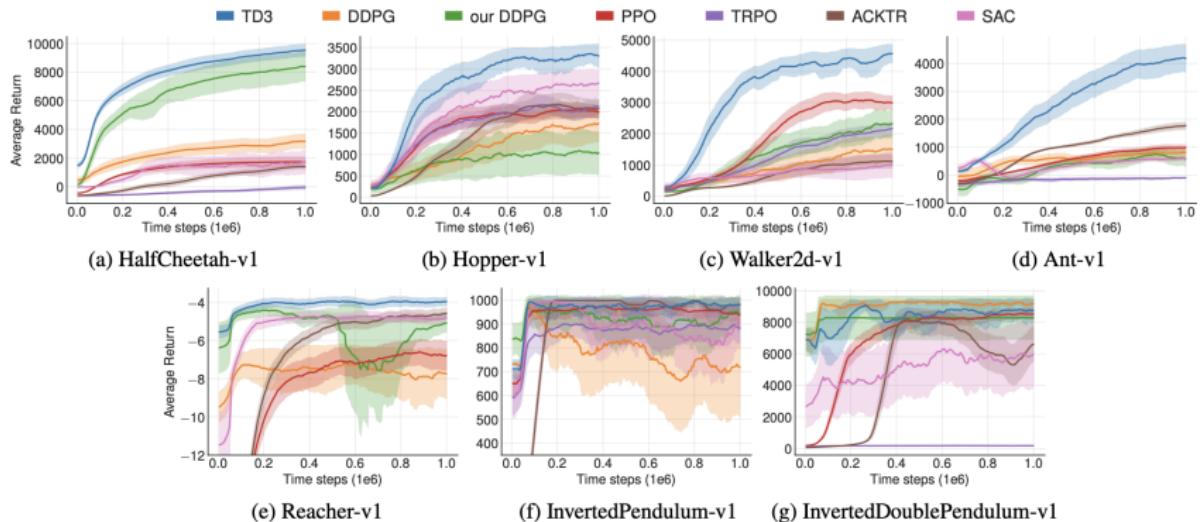
- ③ Target policy smoothing works as follows

$$a_{TD3}(s') = \text{clip}(\mu_{\theta,targ}(s') + \text{clip}(\epsilon, -c, c), a_{low}, a_{high}), \epsilon \sim \mathcal{N}(0, \sigma)$$

- ① It serves as a regularizer: to avoid the incorrect sharp peak for some action output by the policy

# Twin Delayed DDPG (TD3)

## ① Result and comparison



# Twin Delayed DDPG (TD3)

- ① TD3 paper: Fujimoto, et al. Addressing Function Approximation Error in Actor-Critic Methods. ICML'18:  
<https://arxiv.org/pdf/1802.09477.pdf>
- ② Author's Pytorch implementation (**very clean implementation!**):  
<https://github.com/sfujim/TD3/>
- ③ Have a comparison to see the difference between DDPG and TD3
  - ① <https://github.com/sfujim/TD3/blob/master/DDPG.py>
  - ② <https://github.com/sfujim/TD3/blob/master/TD3.py>
- ④ Authors are from neither Berkeley/OpenAI nor DeepMind!
  - ① Scott Fujimoto, Herke van Hoof, David Meger from McGill University

## Soft Actor-Critic (SAC)

- ① SAC optimizes a stochastic policy in an off-policy way, which unifies stochastic policy optimization and DDPG-style approaches
  - ② SAC incorporates **entropy regularization**
  - ③ Entropy is a quantity which measures how random a random variable is,  $H(P) = E_{x \sim P}[-\log P(x)]$
  - ④ Entropy-regularized RL: the policy is trained to maximize a trade-off between expected return and entropy, a measure of randomness in the policy

# Soft Actor-Critic (SAC)

- ① Value function is changed to include the entropy bonus from every timestep

$$V^\pi(s) = E_{\tau \sim \pi} \left[ \sum_t \gamma^t (R(s_t, a_t) + \alpha H(\pi(\cdot | s_t))) \mid s_0 = s \right]$$

- ② The recursive Bellman equation for the entropy-regularized  $Q^\pi$

$$\begin{aligned} Q^\pi(s, a) &= E_{s' \sim P, a' \sim \pi} [R(s, a) + \gamma(Q^\pi(s', a') + \alpha H(\pi(\cdot | s')))] \\ &= E_{s' \sim P, a' \sim \pi} [R(s, a) + \gamma(Q^\pi(s', a') - \alpha \log \pi(a' | s'))] \end{aligned}$$

- ③ So the sample update for Q function is

$$Q^\pi(s, a) \approx r + \gamma(Q^\pi(s', \hat{a}') - \alpha \log \pi(\hat{a}' | s')), \quad \hat{a}' \sim \pi(\cdot | s').$$

# Soft Actor-Critic (SAC)

we want to be able to find all paths that are good to the goal by adding the entropy term

- ① SAC concurrently learns a policy  $\pi_\theta$  and two Q-functions  $Q_{\phi_1}$  and  $Q_{\phi_2}$
- ② Like in TD3, the shared target makes use of the clipped double-Q trick and both Q-functions are learned with MSBE minimization

$$L(\phi_i, \mathcal{D}) = E[(Q_\phi(s, a) - y(r, s', d))^2] \quad (3)$$

$$y(r, s') = r + \gamma \left( \min_{j=1,2} Q_{\phi_{targ,j}}(s', \hat{a}') - \alpha \log \pi_\theta(\hat{a}'|s') \right), \quad (4)$$

$$\hat{a}' \sim \pi_\theta(\cdot|s'). \quad \begin{matrix} \text{bonus for acting more randomly} \\ \backslash \text{alpha controls randomness temp} \end{matrix} \quad (5)$$

- ③ Policy is learned through maximizing the expected future return plus expected future entropy, thus maximize  $V^\pi(s)$

$$\begin{aligned} V^\pi(s) &= E_{a \sim \pi}[Q^\pi(s, a)] + \alpha H(\pi(\cdot|s)) \\ &= E_{a \sim \pi}[Q^\pi(s, a) - \alpha \log \pi(a|s)]. \end{aligned}$$

# Soft Actor-Critic (SAC)

- ① We reparameterize  $a$  as sample from a squashed Gaussian policy

$$\hat{a}_\theta(s, \epsilon) = \tanh(\mu_\theta(s) + \sigma_\theta(s) \odot \epsilon), \epsilon \sim \mathcal{N}(0, I).$$

mean and standard deviation as predicted by the network

- ② Reparameterization trick allows us to rewrite the expectation over actions (which contains a pain point: the distribution depends on the policy parameters) into an expectation over noise (which removes the pain point: the distribution now has no dependence on parameters):

$$\begin{aligned} & E_{a \sim \pi_\theta}[Q^{\pi_\theta}(s, a) - \alpha \log \pi_\theta(a|s)] \\ &= E_{\epsilon \sim \mathcal{N}}[Q^{\pi_\theta}(s, \hat{a}_\theta(s, \epsilon)) - \alpha \log \pi_\theta(\hat{a}_\theta(s, \epsilon)|s)] \end{aligned}$$

now expectation is independent of the parameters?

sample many times and take the average

- ③ The policy is thus optimized as

$1/N \sum_k Q_{\{\pi_0\}}(s, \bar{a}, s_{-1}, \epsilon_k) ??$

$$\max_{\theta} E_{s \sim \mathcal{D}, \epsilon \sim \mathcal{N}} \left[ \min_{j=1,2} Q_{\phi_j}(s, \hat{a}_\theta(s, \epsilon)) - \alpha \log \pi_\theta(\hat{a}_\theta(s, \epsilon)|s) \right] \quad (6)$$

# Brief Intro on Reparameterization Trick

used in variational autoencoders

- ① Let's say we want to compute the gradient of the expected value of the function  $\nabla_{\theta} E_{x \sim p_{\theta}(x)}[f(x)]$ , the difficulty is that  $x$  depends on the distribution with parameter  $\theta$
- ② We can rewrite the samples of the distribution  $p_{\theta}$  in terms of a random variable  $\epsilon$  sampled from distribution  $q$  which is independent of  $\theta$

context: we have a neural net that outputs the parameters of a distribution, which we sample to get a value  $f(x)$ , but we can't back propagate because we are sampling

$$\epsilon \sim q \quad q \text{ is some fixed dist indep from } N \quad (7)$$

$$x = g_{\theta}(\epsilon) \quad \text{this uses our parameters to calculate } x \quad (8)$$

$$\nabla_{\theta} E_{x \sim p_{\theta}(x)}[f(x)] = \nabla_{\theta} E_{\epsilon \sim q}[f(g_{\theta}(\epsilon))] \quad (9)$$

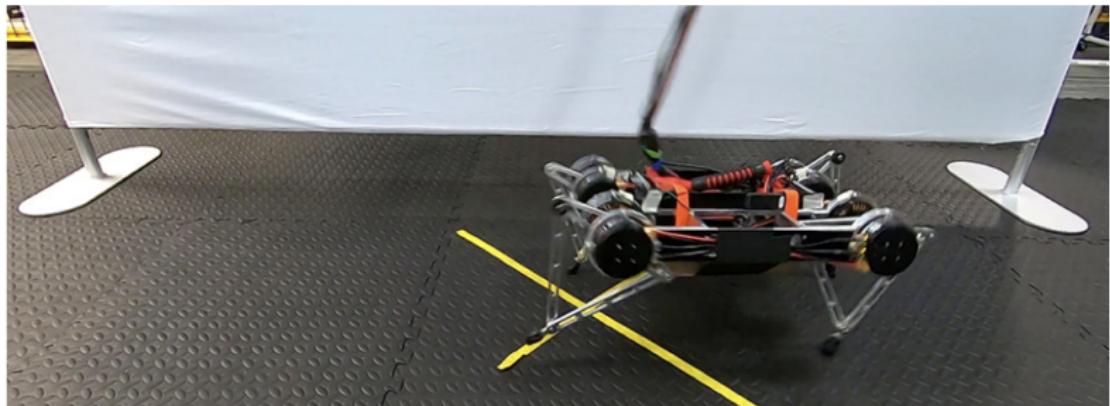
$$= E_{\epsilon \sim q}[\nabla_{\theta} f(g_{\theta}(\epsilon))] \quad (10)$$

since  $x$  no longer depends on sampling from  $p_{\theta}$ , we can move the gradient inside because of linearity and no dependence on  $q$  and  $\theta$

- ① where  $x$  is reparameterized as a function of  $\epsilon$  and the stochasticity of  $p_{\theta}$  is replaced by the distribution  $q$ , which could be just a standard Gaussian  $\mathcal{N}(0, 1)$ , thus  $g_{\theta}(\epsilon) = \mu_{\theta} + \epsilon \sigma_{\theta}$  where  $\epsilon \sim \mathcal{N}(0, 1)$
- ③ On REINFORCE and Reparameterization: <http://stillbreeze.github.io/REINFORCE-vs-Reparameterization-trick/>

# Soft Actor-Critic (SAC)

- ① Sample code on SAC: <https://github.com/pranz24/pytorch-soft-actor-critic/blob/master/sac.py>
- ② SAC is known as SOTA for robot learning:
  - ① Learning to Walk in the Real World with Minimal Human Effort.  
<https://arxiv.org/pdf/2002.08550.pdf>
  - ② <https://www.youtube.com/embed/cwyiq6dCg0c>



# Summary of the SOTAs

- ① Policy Gradient→TRPO→ACKTR→PPO
  - ① Stochastic policy thus output probability over discrete actions
  - ② Start with policy gradient and importance sampling for off-policy learning
- ② Q-learning→DDPG→TD3
  - ① Deterministic policy thus output continuous action spaces.
  - ② Start with Bellman equation, which doesn't care which transition tuples are used, or how the actions were selected, or what happens after a given transition
  - ③ Optimal Q-function should satisfy the Bellman equation for all possible transitions, so very easy for off-policy learning
- ③ SAC
  - ① SAC optimizes a stochastic policy in an off-policy way, which unifies stochastic policy optimization and DDPG-style approaches
  - ② off-policy method, high sample efficiency, it incorporates the clipped double-Q trick like TD3, and due to the inherent stochasticity of the policy in SAC, it also winds up benefiting from something like target policy smoothing.

# Tutorial and Toolkit of SOTA RL algorithms

- ① They become handy for your research project
- ② OpenAI Spinning Up: Nice implementations and summary of the algorithms from OpenAI
  - ① <https://spinningup.openai.com/>



- ③ Stable-baseline3 in PyTorch:
  - ① <https://github.com/DLR-RM/stable-baselines3>
- ④ CleanRL:
  - ① <https://github.com/vwxyzjn/cleanrl>
  - ② High-quality single file implementation of deep RL algorithms

# Great tutorial video on how to make deep RL work

- ① John Schulman (OpenAI research scientist): Nuts and Bolts of Deep RL Experimentation, August 26, 2017
  - ① Slide: <http://joschu.net/docs/nuts-and-bolts.pdf>
  - ② Video: <https://www.youtube.com/watch?v=8EcdaCk9KaQ>
- ② His more recent talk on Reinforcement Learning from Human Feedback: video:  
[https://www.youtube.com/watch?v=hhILw5Q\\_UFg](https://www.youtube.com/watch?v=hhILw5Q_UFg)