

# **Proximal Policy Optimization (PPO)**

Konpat Preechakul  
Chulalongkorn University  
November 2019

# Policy gradient fails

*A, Q*

✗ Critic is not oracle, it has its flaws

✗ Critic is prone to forgetting

✗ How to know when to trust? *confidence*

✗ **How to limit the trust?**

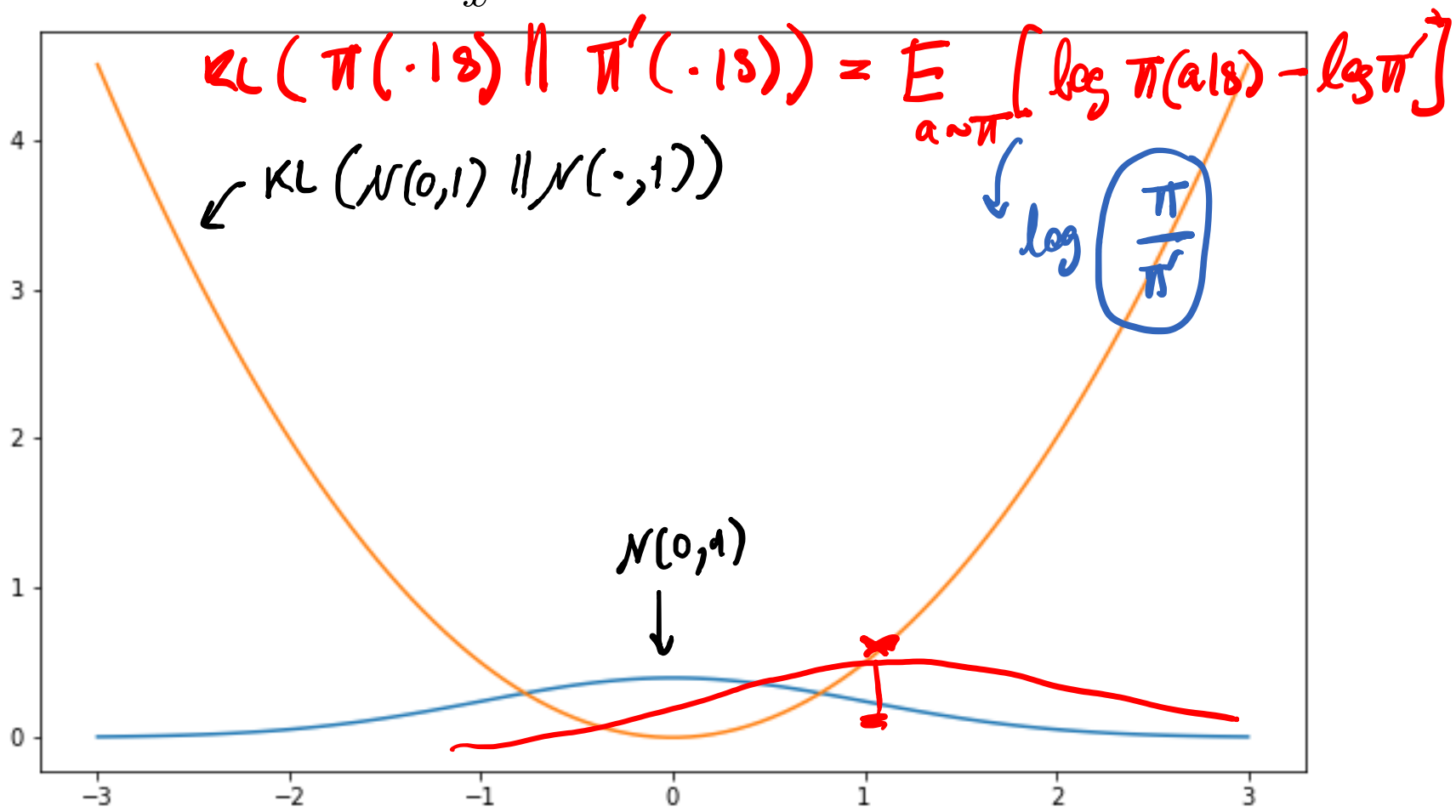
✗ Not updating too much ...

- KL divergence:

→ 
$$\text{KL}(P\|Q) = \sum_x P(x) [\log P(x) - \log Q(x)]$$

# KL Divergence in picture

$$\text{KL}(P\|Q) = \sum_x P(x) [\log P(x) - \log Q(x)]$$



# Trust region policy optimization

$$\rightarrow d^* = \underset{d}{\operatorname{argmax}} \ J(\theta + d) \quad \text{s.t.} \quad \text{KL}(\theta \parallel \theta + d) = c$$

- A constrained optimization

- We relax it using Lagrangian:

$$\star \mathcal{L}(d, \lambda) = \underline{J(\theta + d)} + \underline{\lambda}(\underline{\text{KL}(\theta \parallel \theta + d)} - c)$$

- Optimal  $d$  is at the critical point

$$\nabla_{d, \lambda} \mathcal{L} = 0$$

# Policy improvement guarantee

$$\mathcal{A}_\pi(\pi') = \mathbb{E}_{a_t \sim \pi} \left[ \frac{\pi'(a_t|s_t)}{\pi(a_t|s_t)} \gamma^t A^\pi(s_t, a_t) \right]$$

$$J(\theta') - J(\theta) = \sum_t \mathbb{E}_{s_t \sim P'(s_t)} [\mathcal{A}_\pi(\pi')]$$

Lower bound:

$$J(\theta') - J(\theta) = \sum_t \mathbb{E}_{s_t \sim P(s_t)} [\mathcal{A}_\pi(\pi')] - \sum_t \epsilon t \mathcal{O}\left(\frac{r_{\max}}{1-\gamma}\right)$$

$\uparrow$   
 $\mathbb{E}_{s, a \sim \pi} \left[ \frac{\pi'(a|s)}{\pi(a|s)} A^\pi(s, a) \right] \approx \text{KL}(\pi' \parallel \pi)$

s.t.  $\sqrt{\frac{1}{2} \text{KL}^{\max}(\pi \parallel \pi')} = \epsilon$

# Constrained optimization is hard

→ • Using a fixed constant is not likely to work:

$$\uparrow \bar{J}(\theta + d) = \underline{J}(\theta + d) + \beta \text{KL}(\theta || \theta + d)$$

Handwritten red annotations: an arrow points to  $\bar{J}$ , another to  $\underline{J}$ , a circle around  $\beta$  with an arrow pointing to it, and a bracket under  $\text{KL}(\theta || \theta + d)$  with an arrow pointing to it.

✗ Lagrangian involves exotic terms like “inverse”

$$\propto \left( \nabla_{\theta}^2 \text{KL} \right)^{-1} \nabla_{\theta} J(\theta)$$

Handwritten red annotations: a circle around  $\left( \nabla_{\theta}^2 \text{KL} \right)^{-1}$  with an arrow pointing to it, and the word "exotic" written below it.

**Is there an easy way to  
constrain KL?**



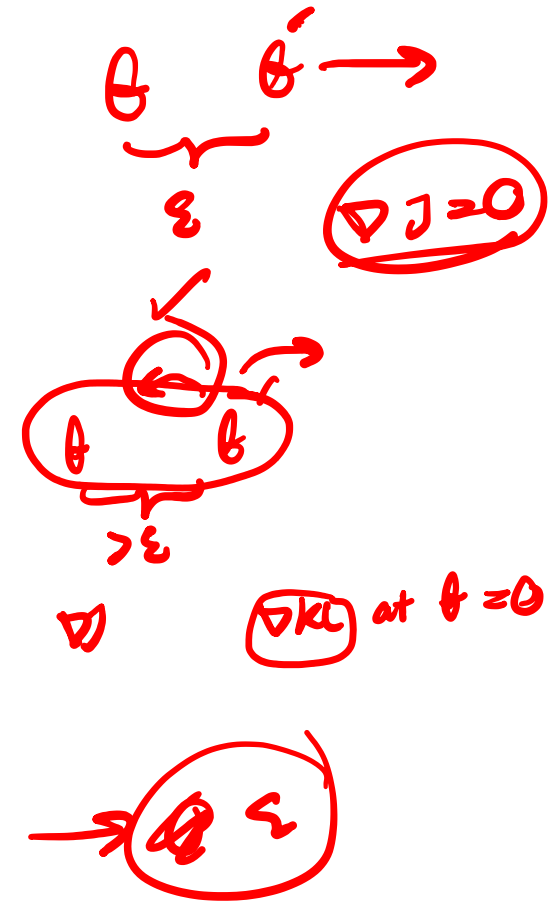
# Goal

$$KL(\theta \parallel \theta')$$

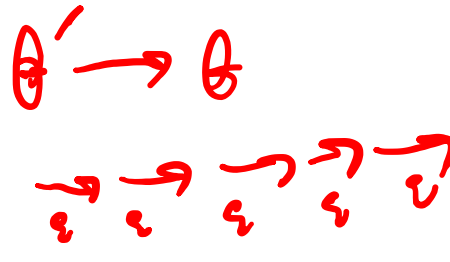
- Design objective function  $J$
- That has “zero” gradient
- When the constraint is breached

if  $KL > \epsilon \rightarrow \nabla J = 0$   
 else  $\nabla J$

- Optimize normally ...



$$\theta$$





# First attempt

$\min, \max(-, -)$

$$\nabla J \leftarrow$$

$$\nabla J_{KL}$$



if  $\boxed{KL(\theta \parallel \theta') > \epsilon}$

if  $\nabla J: \theta \rightarrow \leftarrow \theta': \nabla J$

else:  $\nabla J = 0$

else:  $\nabla J$

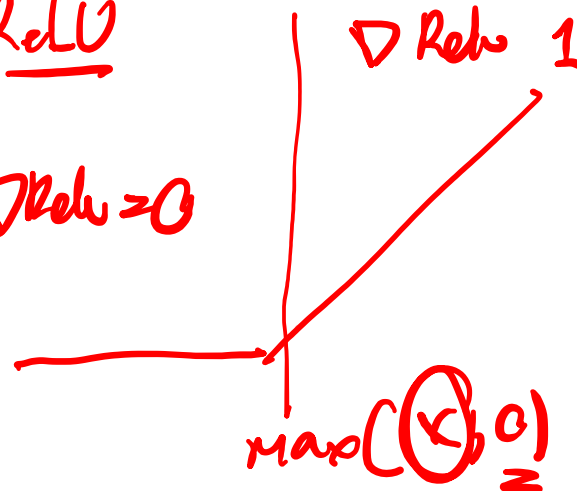
$$\stackrel{\text{KL}}{\leftarrow} \bar{e}^k \nabla J$$

ReLU

$\nabla \text{ReLU} = 0$

$\nabla \text{ReLU} = 1$

$\max(\underbrace{x}_z, 0)$



# Possible pseudocode

$\theta' \leftarrow \theta$

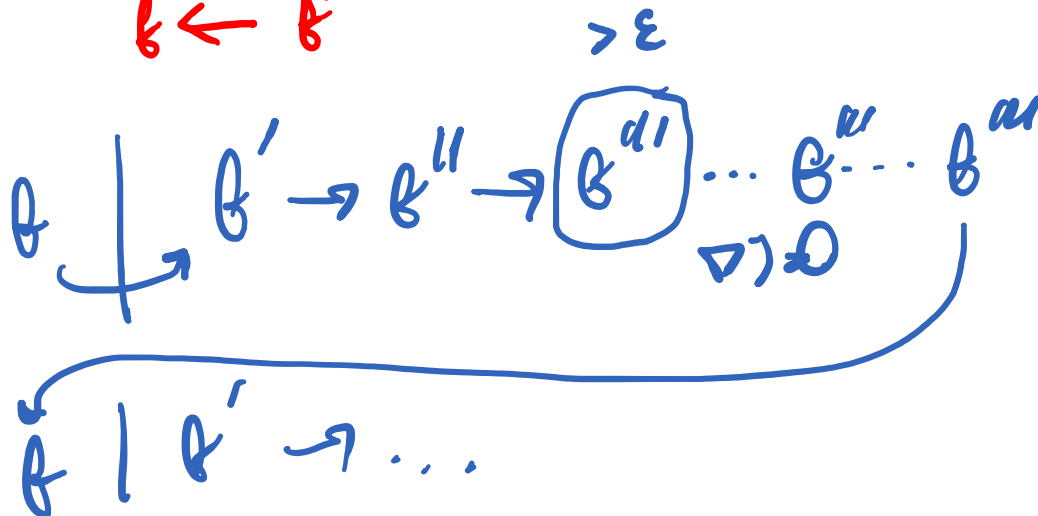
for until satisfied:

collect data using  $\pi$  for  $t$  steps  $\overbrace{(s_0, a_0, r_1, \dots)}^t$

{ for  $i$  to  $K$ :

$\theta' \leftarrow \theta' + \alpha \nabla J_K(\theta')$

$\theta \leftarrow \theta'$



# Is there a one-liner?

$$J(\theta) = E_{s,a \sim \pi} \left[ \frac{\pi'(a|s)}{\pi(a|s)} \cdot A^\pi(s,a) \right]$$

$$\pi' = \pi : \quad \frac{\pi'}{\pi} = 1$$

↓ ↓ ↓ ↓ ↓ ↓

$$\pi' \neq \pi : \quad \frac{\pi'}{\pi} \neq 1$$

→

$$1 - \epsilon \leq \frac{\pi'}{\pi} \leq 1 + \epsilon$$

↑  
0.2

one liner clip, clamp

$$\rightarrow \text{clip} \left( \frac{\pi'}{\pi}, 1 - \epsilon, 1 + \epsilon \right) \cdot A^\pi(s,a)$$

$$J(\theta)_{\text{PPO}} = E_{s,a \sim \pi} \left[ \min \left( \frac{\pi'}{\pi} \cdot A^\pi, \text{clip}(\dots) \cdot A^\pi \right) \right] \quad \underline{\underline{\text{PPO}}}$$

# Pseudocode

---

**Algorithm 1** PPO, Actor-Critic Style

---

```
for iteration=1, 2, ... do
  for actor=1, 2, ..., N do
    Run policy  $\pi_{\theta_{\text{old}}}$  in environment for  $T$  timesteps
    Compute advantage estimates  $\hat{A}_1, \dots, \hat{A}_T$ 
  end for
  { Optimize surrogate  $L$  wrt  $\theta$ , with  $K$  epochs and minibatch size  $M \leq NT$ 
  {  $\theta_{\text{old}} \leftarrow \theta$ 
end for
```

*Handwritten notes:*

- $\}$  explore (next to the actor loop)
- $\uparrow \nabla J_{\text{PPO}}$  (pointing to the optimization step)
- $\theta, \theta_{\text{old}}$  (with a brace underneath and  $\leftarrow$  below it)

---

# Results

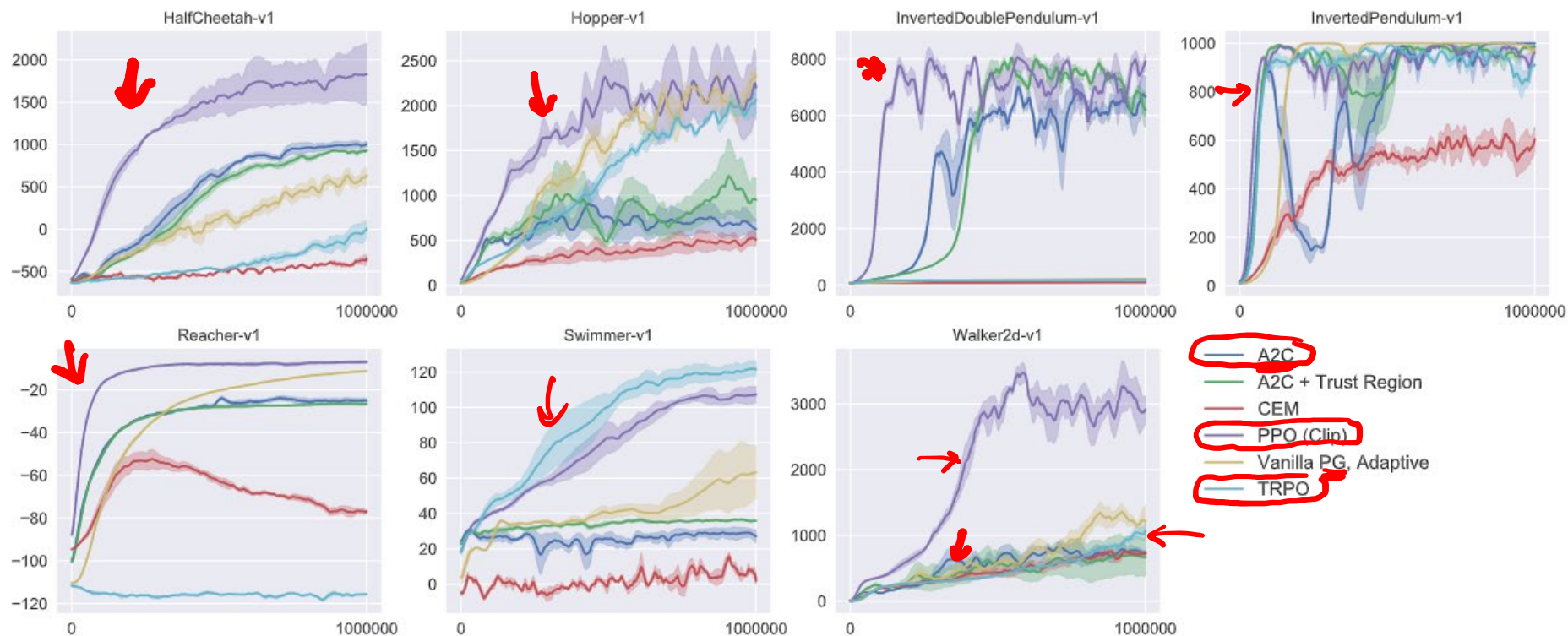


Figure 3: Comparison of several algorithms on several MuJoCo environments, training for one million timesteps.