

Distributed Proximal Policy Optimization (by Google Deep Mind)

Distributed: 与 A3C 算法一致，采用并行的分布式训练。以降低参数之间的相关性，提高 model 的收敛性和稳定性。

Proximal Policy: 采用加入 KL 惩罚项的方法，限制了 new policy 的更新幅度 \Rightarrow 降低了 model 对 step size / 敏感度 lr 提高了鲁棒性。

PPO: (by OpenAI) 近端策略优化

特点：现有最顶尖的 RL 算法之一 \Rightarrow

对比 ACER.
TRPO { ① 在众多问题上表现优异 (more general)
② 算法简单，易于实现、调试 (simpler to implement)

对比 A3C.
DDPG { ③ 采样效率更高 (better sample complexity)
④ 鲁棒性好 (Robustness)

前身：TRPO (Trust Region Policy Optimization)

同级算法：ACER (Actor Critic with Experience Replay)

缺：与 PPO 表现相近，但实现十分复杂，且泛化性差 (与参数共享、dropout 不兼容)

Why PPO? \Rightarrow

Policy Gradient 类的算法 (Simple PG, AC, A3C, DDPG...)

(i) 对超参数敏感，鲁棒性差。

① Learning rate

①, ②过大：Policy 跳动太大，难以收敛

② Step size.

①, ②太小：学习速度缓慢。

(ii) 采样效率低下，学习简单任务即需要百万级 iteration

Algorithm 1 Proximal Policy Optimization (adapted from [8])

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

empirical collection

 Run policy π_θ for T timesteps, collecting $\{s_t, a_t, r_t\}$

 Estimate advantages $\hat{A}_t = \sum_{t' > t} \gamma^{t'-t} r_{t'} - V_\phi(s_t)$ \Rightarrow Advantage = $(V_{t'} - V_t)$
 $\pi_{old} \leftarrow \pi_\theta$ \uparrow actor update steps

for $j \in \{1, \dots, M\}$ **do**

$J_{PPO}(\theta) = \sum_{t=1}^T \frac{\pi_\theta(a_t | s_t)}{\pi_{old}(a_t | s_t)} \hat{A}_t - \lambda \text{KL}[\pi_{old} | \pi_\theta]$ \Rightarrow KL penalty
 \Rightarrow 在迭代过程中保持不变，此为 PPO 核心思想

 Update θ by a gradient method w.r.t. $J_{PPO}(\theta)$

end for

for $j \in \{1, \dots, B\}$ **do**

$L_{BL}(\phi) = -\sum_{t=1}^T (\sum_{t' > t} \gamma^{t'-t} r_{t'} - V_\phi(s_t))^2$ \Rightarrow c-loss

 Update ϕ by a gradient method w.r.t. $L_{BL}(\phi)$

end for

if $\text{KL}[\pi_{old} | \pi_\theta] > \beta_{high} \text{KL}_{target}$ **then**

$\lambda \leftarrow \alpha \lambda$ \uparrow 权衡因子

else if $\text{KL}[\pi_{old} | \pi_\theta] < \beta_{low} \text{KL}_{target}$ **then**

$\lambda \leftarrow \lambda / \alpha$ \downarrow 调整因子

end if

λ \downarrow (scaling term)

end for

λ \downarrow Adaptive KL penalty
 \downarrow (因为 KL penalty 自适应能力差)

KL penalty :

$$KL(\theta) = \mathbb{E}_t \left[\frac{\pi_\theta(a_t | s_t)}{\pi_{old}(a_t | s_t)} \hat{A}_t - \lambda \cdot KL[\pi_{old} | \pi_\theta] \right]$$

Expect over T time steps.

in time step t

Update
Actor

Clipped surrogate objective :

$$\text{ratio}(\theta) = \frac{\pi_\theta(a_t | s_t)}{\pi_{old}(a_t | s_t)}$$

趋势数, 0.1 ~ 0.2

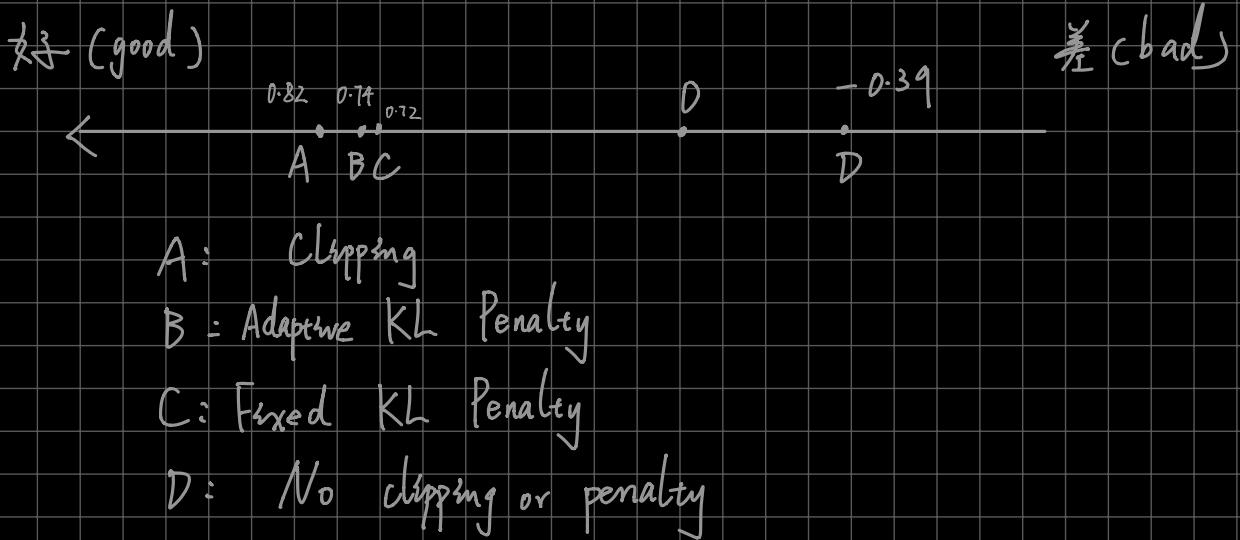
$$CLIP(\theta) = \mathbb{E}_t \left[\min(V_t(\theta) \cdot \hat{A}_t, \text{clip}(\text{ratio}(\theta), 1-\varepsilon, 1+\varepsilon) \cdot \hat{A}_t) \right]$$

△ CLIP method is better than KL method, 因为 它与 GD 兼容

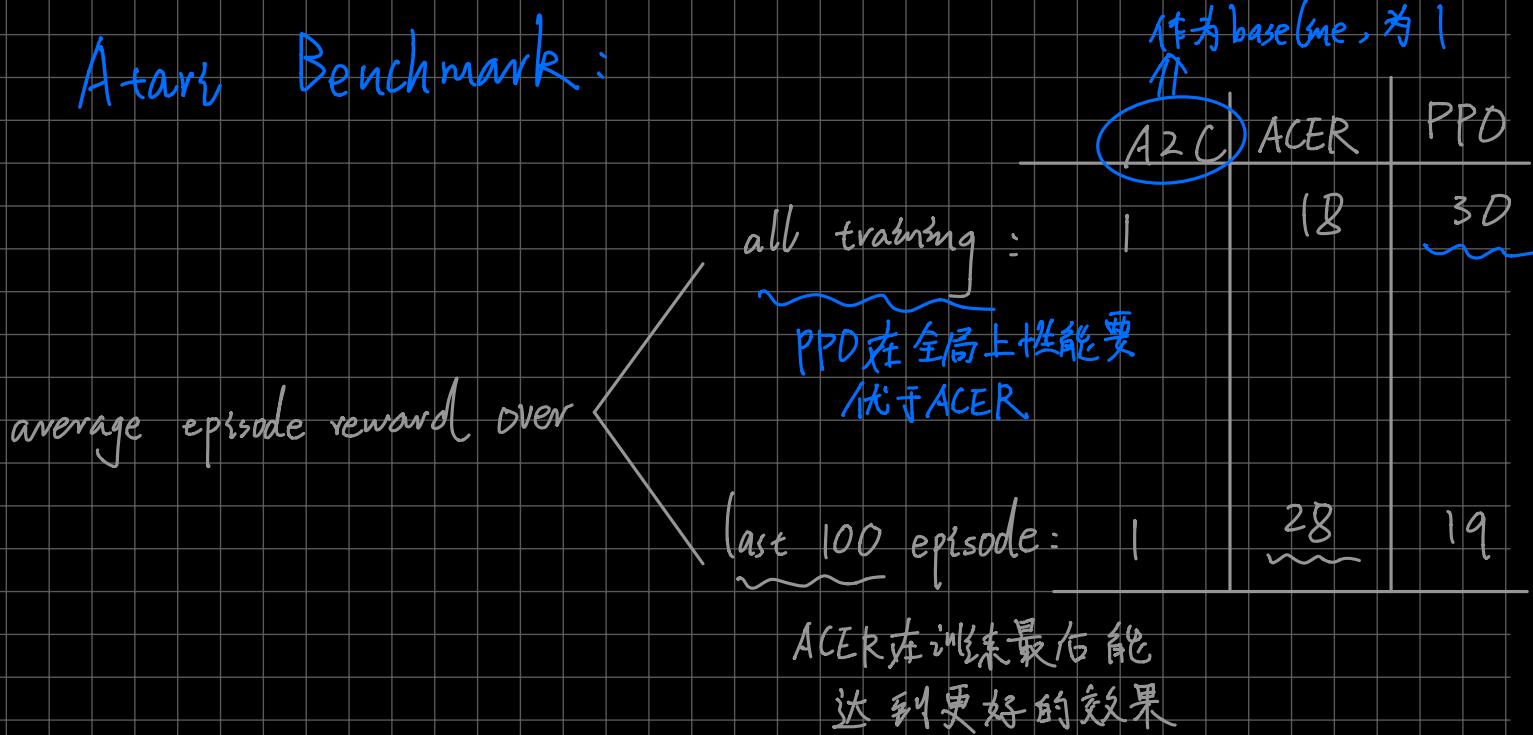
② 简化了实现 \Rightarrow 移除了 KL 惩罚项 (和 adaptive KL 部分)

Experiment by OpenAI:

Continuous control benchmark:



Atari Benchmark:



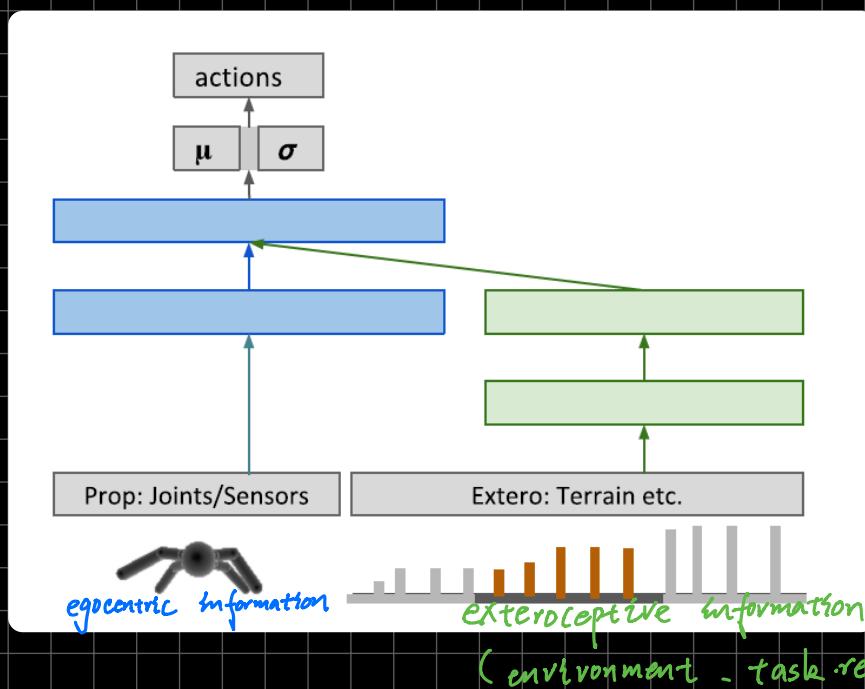
Distributed PPO

data collection - gradient calculation are distributed over workers

Important points :

- Pull : Averaging gradients from workers
- Push: Push parameters synchronously
- Normalization :

Network Architecture :



$$\left. \begin{array}{l} \text{observation: } s_t \Rightarrow \frac{s_t - \bar{w}}{\sigma} \\ \text{reward: } r_t \Rightarrow \frac{r_t}{\sigma} \\ \text{advantage: } A_t \Rightarrow \text{per-batch norm} \end{array} \right\} \begin{array}{l} \text{Time} \\ \text{steps} \end{array}$$

Aim to achieve a separation of basic locomotion skills and terrain perception - navigation

Algorithm 2 Distributed Proximal Policy Optimization (chief)

```

for i ∈ {1, …, N} do
    for j ∈ {1, …, M} do
        update actor
        update critic
        Episode number
        global update steps
        number of workers
        threshold
        Wait until at least W - D gradients wrt. θ are available
        average gradients and update global θ
    end for
    for j ∈ {1, …, B} do
        update actor
        update critic
        Episode number
        global update steps
        number of workers
        threshold
        Wait until at least W - D gradients wrt. ϕ are available
        average gradients and update global ϕ
    end for
end for

```

Algorithm 3 Distributed Proximal Policy Optimization (worker)

Empirical Collection

update worker

Actor

update worker

Critic

Adaptive KL

for $i \in \{1, \dots, N\}$ do number of data points per worker
 for $w \in \{1, \dots, T/K\}$ do number of time steps (因为有多个 worker, 因此一次 empirical collection 降为 K 即可)
 Run policy π_θ for K timesteps, collecting $\{s_t, a_t, r_t\}$ for $t \in \{(i-1)K, \dots, iK-1\}$
 Estimate return $\hat{R}_t = \sum_{t=(i-1)K}^{iK-1} \gamma^{t-(i-1)K} r_t + \gamma^K V_\phi(s_{iK})$
 Estimate advantages $\hat{A}_t = \hat{R}_t - V_\phi(s_t)$
 Store partial trajectory information
 end for
 $\pi_{old} \leftarrow \pi_\theta$
 for $m \in \{1, \dots, M\}$ do 与 PPO一样, 在更新时使用所有的 data point 计算 $J(\theta)$
 $J_{PPO}(\theta) = \sum_{t=1}^T \frac{\pi_\theta(a_t|s_t)}{\pi_{old}(a_t|s_t)} \hat{A}_t - \lambda \text{KL}[\pi_{old}|\pi_\theta] - \xi \max(0, \text{KL}[\pi_{old}|\pi_\theta] - 2\text{KL}_{target})^2$
 if $\text{KL}[\pi_{old}|\pi_\theta] > 4\text{KL}_{target}$ then KL penalty item
 break and continue with next outer iteration $i+1$ threshold
 end if
 Compute $\nabla_\theta J_{PPO}$
 send gradient wrt. to θ to chief
 wait until gradient accepted or dropped; update parameters
 end for
 for $b \in \{1, \dots, B\}$ do Early Stopping
 $L_{BL}(\phi) = -\sum_{t=1}^T (\hat{R}_t - V_\phi(s_t))^2 \Rightarrow \text{critic loss}$
 Compute $\nabla_\phi L_{BL}$
 send gradient wrt. to ϕ to chief
 wait until gradient accepted or dropped; update parameters
 end for
 if $\text{KL}[\pi_{old}|\pi_\theta] > \beta_{high} \text{KL}_{target}$ then
 $\lambda \leftarrow \bar{\alpha} \lambda$
 else if $\text{KL}[\pi_{old}|\pi_\theta] < \beta_{low} \text{KL}_{target}$ then additional penalty term
 $\lambda \leftarrow \lambda / \bar{\alpha}$
 end if
 end for
 \Rightarrow wait for global update finished, then update local parameters.
 λ is shared among workers. (But λ update is based on local data)

Implementation

Hyper Parameters:

N : number of episode T : timesteps

γ : decay rate $\begin{cases} M: \text{action update steps} \\ B: \text{critic update steps} \end{cases}$ λ : KL factor
 $\begin{cases} \text{Blow-Bhigh} \\ \alpha \end{cases}$

for $i \in \{1, \dots, N\}$ do \Rightarrow training episode

Empirical Collection \Rightarrow For T timesteps, collect $\{s_t, a_t, r_t\}$ from interaction with RL environment

Calculate \hat{A}_t for $t \in \{1, \dots, T\}$

$$\hat{A}_t = \sum_{t' > t} \gamma^{t'-t} \cdot r_{t'} - \underbrace{V_\phi(s_t)}_{\Rightarrow \text{critic network}}$$

$\pi_{old} \leftarrow \pi_\theta \Rightarrow$ update π_{old}

for $j \in \{1, \dots, M\}$ do \Rightarrow action update step

$$J(\theta) = \sum_{t=1}^T \frac{\pi_\theta(a_t | s_t)}{\pi_{old}(a_t | s_t)} \cdot \hat{A}_t - \lambda \cdot KL(\pi_\theta | \pi_{old})$$

Calculate $\nabla_\theta J(\theta)$ and update $J(\theta)$

for $j \in \{1, \dots, B\}$ do \Rightarrow critic update step

$$L(\phi) = - \sum_{i=1}^T (\hat{A}_t - V_\phi(s_t))^2$$

Calculate $\nabla_\phi L(\phi)$ and update $L(\phi)$

If $KL(\pi_\theta | \pi_{old}) > \beta_{high} \cdot KL_{target}$:

$$\lambda = \lambda \cdot \alpha$$

Else if $KL(\pi_0 | \pi_{old}) < \beta_{low} \cdot KL_{target}$:

$$\lambda = \lambda / \alpha$$

可能原因：

① env. seed 无法复现，每次会不同 ✓
且达到的 reward 上限变化

② env. seed 无设置 X

③ KL-divergence 不满足交换律：X

$KL(\text{old-}p_i | p_i)$ ✓

$KL(p_i | \text{old-}p_i)$ X

④ Kernel-initializer 与 Constant initializer 影响不大 X

⑤ update 中 adv 是否进行了 norm 操作 = \Rightarrow ✓