

# Proximal Policy Optimization (PPO)

policy gradient進階版

default reinforcement learning algorithm at OpenAI



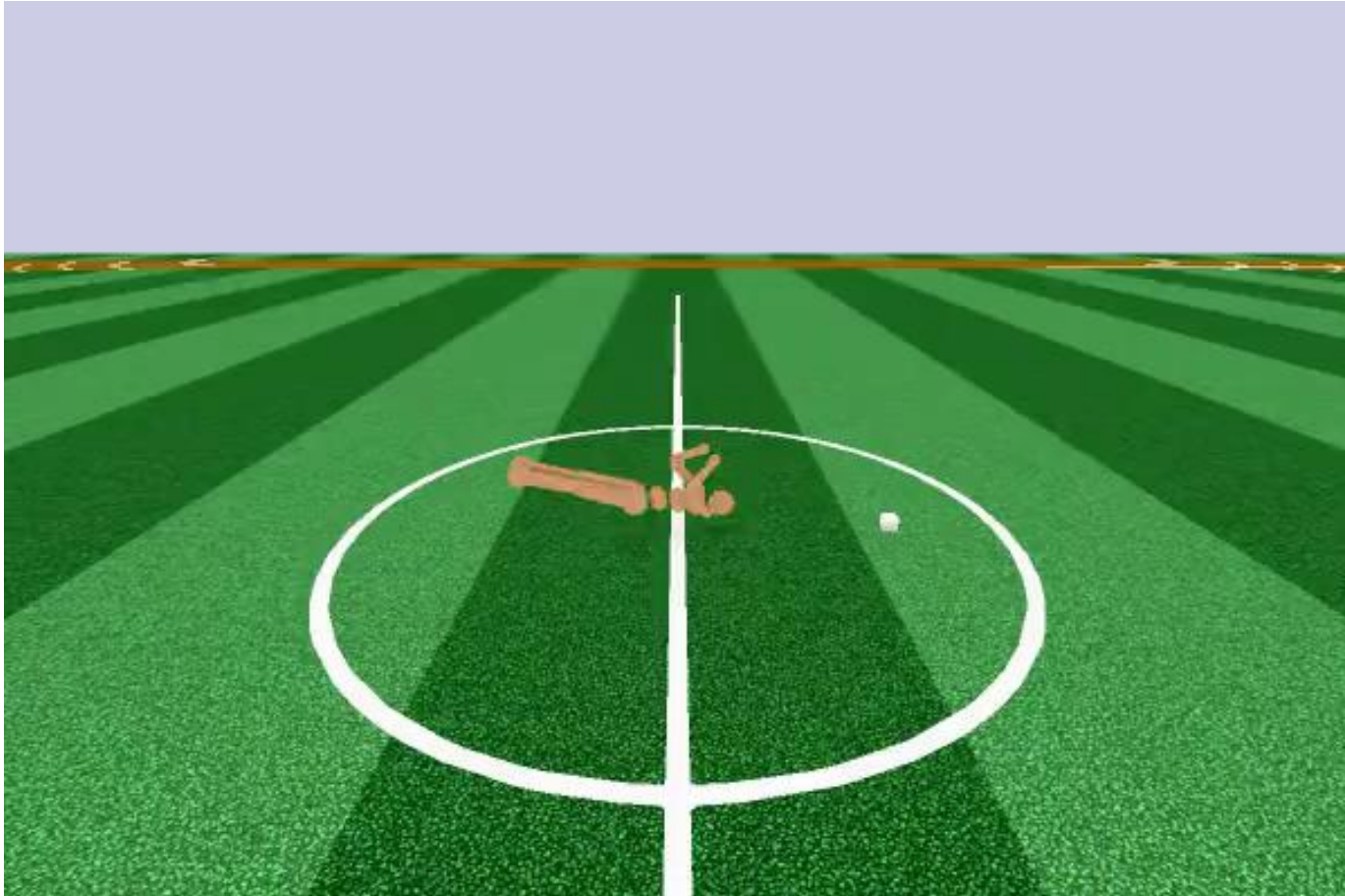
# DeepMind

<https://youtu.be/gn4nRCC9TwQ>



# OpenAI

<https://blog.openai.com/openai-baselines-ppo/>



# Policy Gradient (Review)

# Basic Components

事先給定好的  
You cannot control

訓練過程可以調動的



Actor

Env

Reward  
Function

Video  
Game



Get 20 scores when  
killing a monster

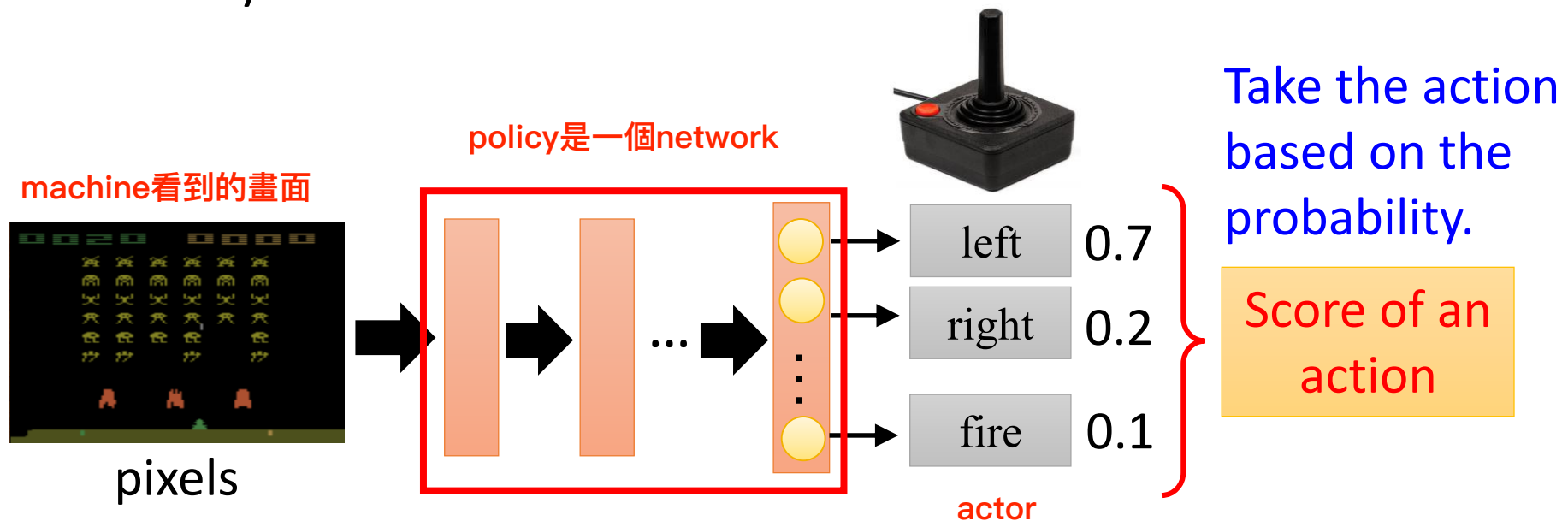
Go



The rule  
of GO

# Policy of Actor

- Policy  $\pi$  is a network with parameter  $\theta$ 
  - Input: the observation of machine represented as a vector or a matrix
  - Output: each action corresponds to a neuron in output layer



# Example: Playing Video Game

Start with  
observation  $s_1$

Observation  $s_2$

Observation  $s_3$

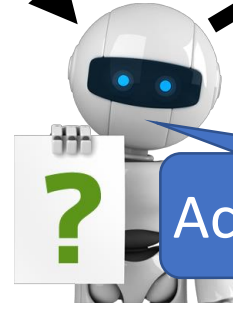


遊戲初始的畫面，根據  
NN policy決定actor



Obtain reward  
 $r_1 = 0$

Action  $a_1$ : "right"



Obtain reward  
 $r_2 = 5$

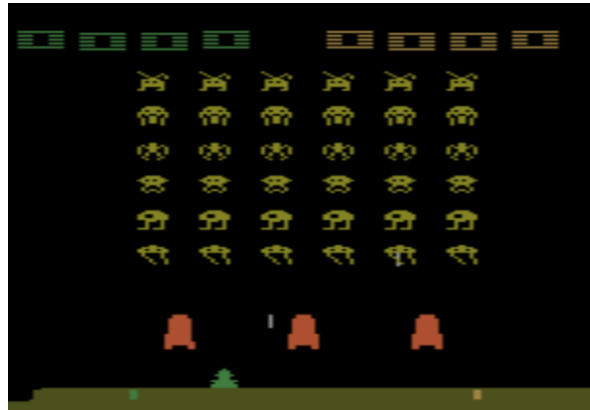
Action  $a_2$ : "fire"  
(kill an alien)

# Example: Playing Video Game

Start with  
observation  $s_1$



Observation  $s_2$



Observation  $s_3$



After many turns  
.....➡



Obtain reward  $r_T$

Action  $a_T$

This is an <sup>一場遊戲</sup>episode.

Total reward:

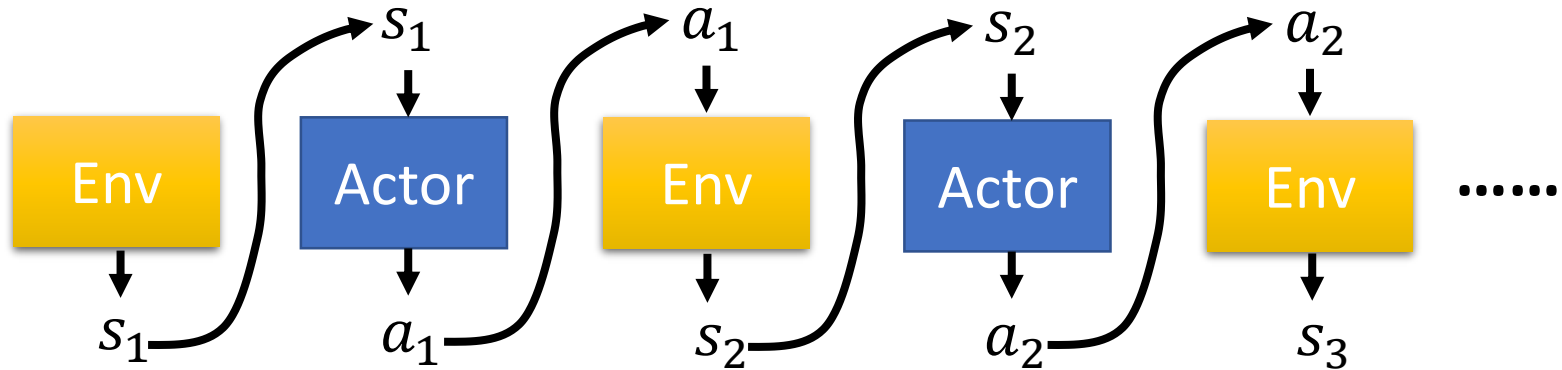
$$R = \sum_{t=1}^T r_t$$

We want the total  
reward be maximized.

actor存在的目的就是想辦法提高R



# Actor, Environment, Reward

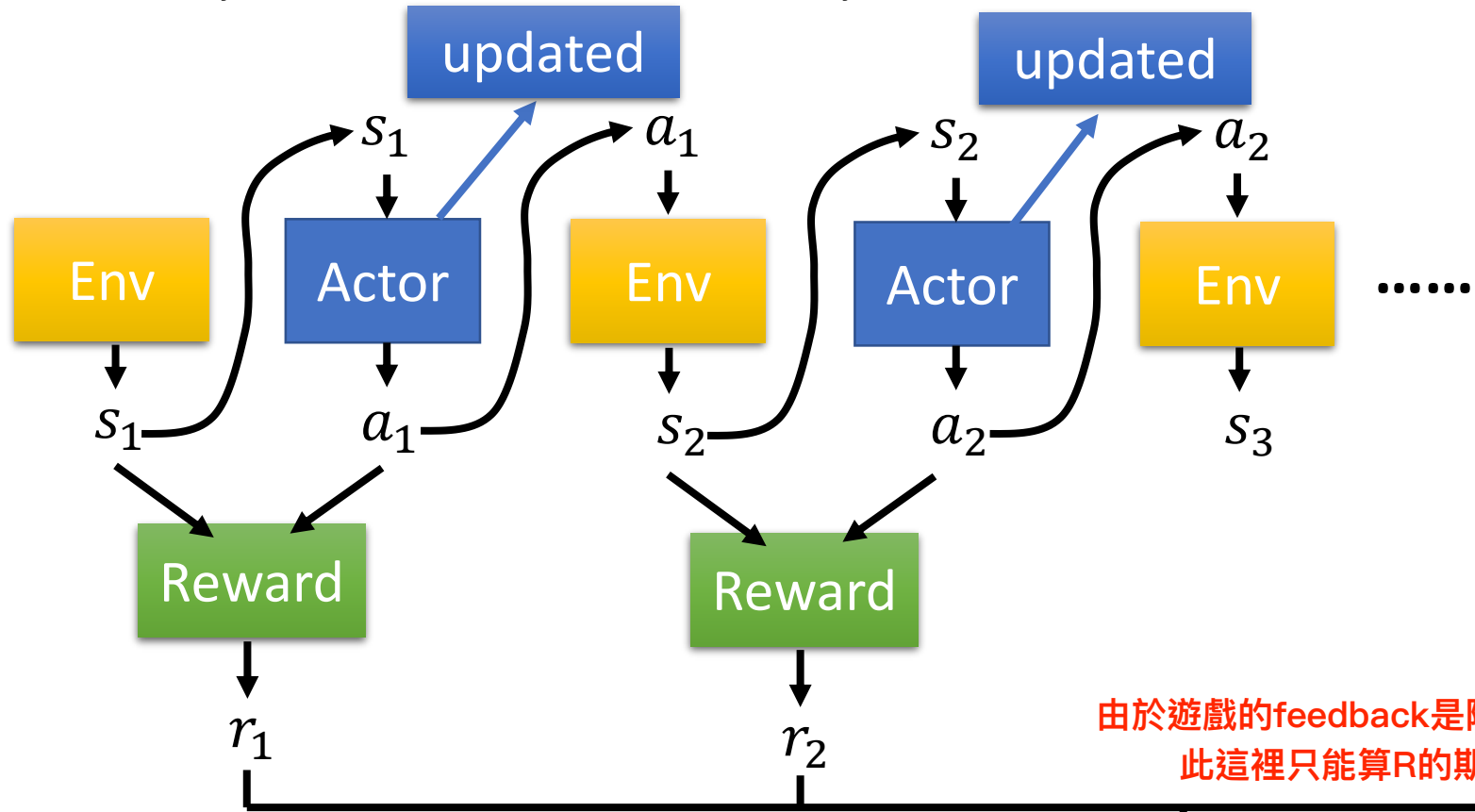


trajectory發生的機率  
**Trajectory**  $\tau = \{s_1, a_1, s_2, a_2, \dots, s_T, a_T\}$   
 $p_{\theta}(\tau)$

$$= p(s_1)p_{\theta}(a_1|s_1)p(s_2|s_1, a_1)p_{\theta}(a_2|s_2)p(s_3|s_2, a_2) \dots$$

$$= p(s_1) \prod_{t=1}^T \overset{\text{取決於自己學出來的 NN}}{p_{\theta}(a_t|s_t)} \underset{\text{取決於遊戲環境}}{p(s_{t+1}|s_t, a_t)}$$

# Actor, Environment, Reward




由於遊戲的feedback是隨機的，因此這裡只能算R的期望值

Expected Reward

看到期望值就想到窮舉~

$$\bar{R}_\theta = \sum_{\tau} R(\tau) p_\theta(\tau) = E_{\tau \sim p_\theta(\tau)} [R(\tau)]$$

$$R(\tau) = \sum_{t=1}^T r_t$$


我們要做的事情就是maximum expected reward

# Policy Gradient

$$\bar{R}_\theta = \sum_{\tau} R(\tau) p_\theta(\tau) \quad \nabla \bar{R}_\theta = ?$$

$$\nabla \bar{R}_\theta = \sum_{\tau} R(\tau) \nabla p_\theta(\tau) = \sum_{\tau} R(\tau) p_\theta(\tau) \frac{\nabla p_\theta(\tau)}{p_\theta(\tau)}$$

舉例來說在GAN來說，這裡的R就是D

$R(\tau)$  do not have to be differentiable

It can even be a black box.

$$= \sum_{\tau} R(\tau) p_\theta(\tau) \nabla \log p_\theta(\tau)$$

$$\nabla f(x) = f(x) \nabla \log f(x)$$

$$= E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)] \approx \frac{1}{N} \sum_{n=1}^N R(\tau^n) \nabla \log p_\theta(\tau^n)$$

唯一需要算gradient的部分

$$= \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p_\theta(a_t^n | s_t^n)$$

# Policy Gradient

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)]$$

用train好的參數跟遊戲玩一場，並且記錄下trajectory

Given policy  $\pi_\theta$

$\tau^1:$	$(s_1^1, a_1^1)$	$R(\tau^1)$
	$(s_2^1, a_2^1)$	$R(\tau^1)$
	$\vdots$	$\vdots$
$\tau^2:$	$(s_1^2, a_1^2)$	$R(\tau^2)$
	$(s_2^2, a_2^2)$	$R(\tau^2)$
	$\vdots$	$\vdots$

only used once

Update  
Model

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta = \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p_\theta(a_t^n | s_t^n)$$

Data  
Collection

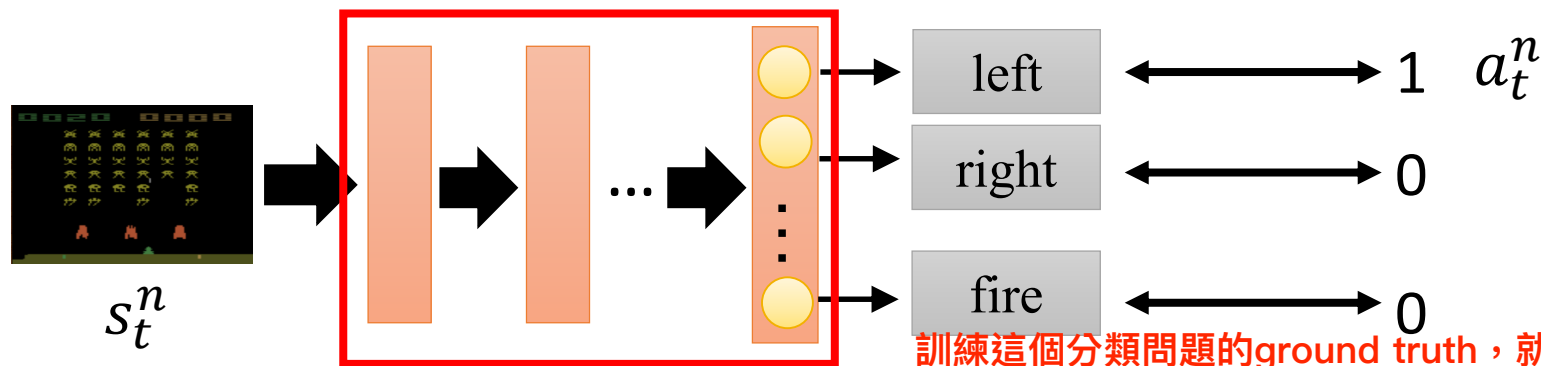
$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

## Implementation

$$\nabla \bar{R}_\theta = \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p_\theta(a_t^n | s_t^n)$$

Consider as classification problem

$$s_t^n \quad a_t^n \quad R(\tau^n)$$



訓練這個分類問題的ground truth，就採用 sampling，sample到哪個就當他GT label

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \log p_\theta(a_t^n | s_t^n) \xrightarrow{\text{TF, pyTorch ...}} \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \nabla \log p_\theta(a_t^n | s_t^n)$$

一整場遊戲只有一個大R

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \log p_\theta(a_t^n | s_t^n) \xrightarrow{\quad} \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p_\theta(a_t^n | s_t^n)$$

在R L 中只是多一個weight

# Tip 1: Add a Baseline

這樣會造成大家都提升，只是有大有小

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

It is possible that  $R(\tau^n)$  is always positive.

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p_\theta(a_t^n | s_t^n)$$

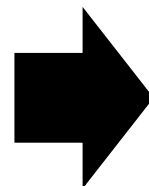
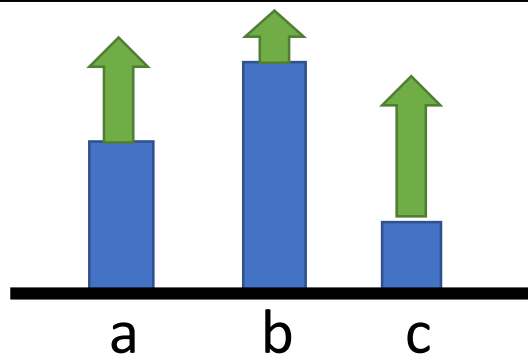
因此這邊做了一點感變：不要讓reward總是正的

大家同減去baseline  $b$

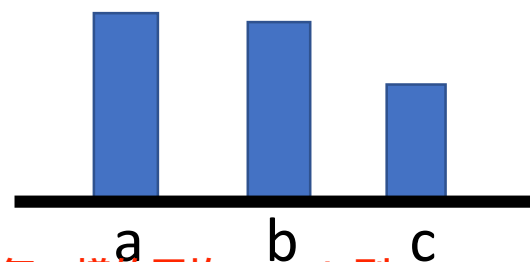
$$b \approx E[R(\tau)]$$

mean: 每次episode計算一次

Ideal case



It is probability ...

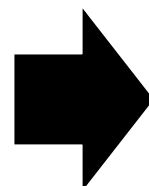
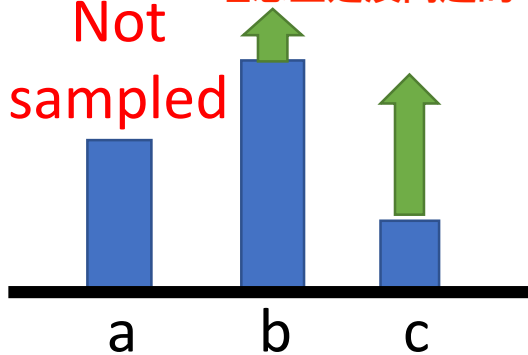


理想上是沒問題的，但是前提是每一樣的平均sample到

Not  
sampled

Sampling

.....



The probability of the actions not sampled will decrease.

譬如上面的例子，沒被sample到反而下降

# Tip 2: Assign Suitable Credit

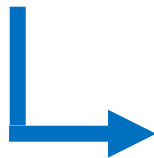
理想上這個問題只要sample夠多就可以被解決  
但實作上無法sample所有情況，因此要涉及合理的credit

$\times 3$	$\times -2$	$\times -2$	$\times -7$	$\times -2$	$\times -2$
$(s_a, a_1)$	$(s_b, a_2)$	$(s_c, a_3)$	$(s_a, a_2)$	$(s_b, a_2)$	$(s_c, a_3)$
+5	+0	-2	-5	+0	-2
$R = +3$			$R = -7$		

同一場episode都是weighted  
by same weight

但這是不公平的，因為說不定  
只有某個action需要做調整

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (\cancel{R(t^n)} - b) \nabla \log p_\theta(a_t^n | s_t^n)$$



$$\sum_{t'=t}^{T_n} r_{t'}^n$$

只看這個時間點以後所得到的reward

# Tip 2: Assign Suitable Credit

實際上在算R的時候需要跟環境  
interaction，因此based on theta

Advantage  
Function  $A^\theta(s_t, a_t)$

相對其他action的好壞

How good it is if we take  $a_t$  other  
than other actions at  $s_t$ .

Estimated by “critic” (later)

Can be state-dependent

b 可以是state dependent，  
一個NN的output

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (\underbrace{R(t^n)}_{\text{red box}} - b) \nabla \log p_\theta(a_t^n | s_t^n)$$

$$\sum_{t'=t}^{T_n} r_{t'}^n$$

$$\sum_{t'=t}^{T_n} \underbrace{\gamma^{t'-t}}_{\text{decay factor}} r_{t'}^n$$

Add discount factor

$$\gamma < 1$$

代表說距離越遠影響越小



# From on-policy to off-policy

Using the experience more than once

# On-policy v.s. Off-policy

- On-policy: The agent learned and the agent interacting with the environment is the same.
- Off-policy: The agent learned and the agent interacting with the environment is different.



阿光下棋

自己下自己學 : on policy



佐為下棋、阿光在旁邊看

在旁邊看著學 : off policy

每次更新完參數後又要跟環  
境大量互動 collect data  
這樣太花時間了！

# On-policy $\rightarrow$ Off-policy

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)]$$

- Use  $\pi_\theta$  to collect data. When  $\theta$  is updated, we have to sample training data again.
- Goal: Using the sample from  $\pi_{\theta'}$  to train  $\theta$ .  $\theta'$  is fixed, so we can re-use the sample data.

## Importance Sampling

從p這個distribution sample x

$$E_{x \sim p} [f(x)] \approx \frac{1}{N} \sum_{i=1}^N f(x^i)$$

$x^i$  is sampled from  $p(x)$

We only have  $x^i$  sampled from  $q(x)$

但我們無法從p sample data，只能從另一個distribution sample (q)

$$= \int f(x) p(x) dx = \int f(x) \frac{p(x)}{q(x)} q(x) dx = E_{x \sim q} [f(x) \frac{p(x)}{q(x)}]$$

Importance weight

# Issue of Importance Sampling

mean一樣

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

但variance不一定一樣

$$\text{Var}_{x \sim p}[f(x)] \quad \text{Var}_{x \sim q}\left[f(x) \frac{p(x)}{q(x)}\right]$$

$$\text{VAR}[X]$$

$$= E[X^2] - (E[X])^2$$

$$\text{Var}_{x \sim p}[f(x)] = E_{x \sim p}[f(x)^2] - (E_{x \sim p}[f(x)])^2$$

$$\text{Var}_{x \sim q}\left[f(x) \frac{p(x)}{q(x)}\right] = E_{x \sim q}\left[\left(f(x) \frac{p(x)}{q(x)}\right)^2\right] - \left(E_{x \sim q}\left[f(x) \frac{p(x)}{q(x)}\right]\right)^2$$

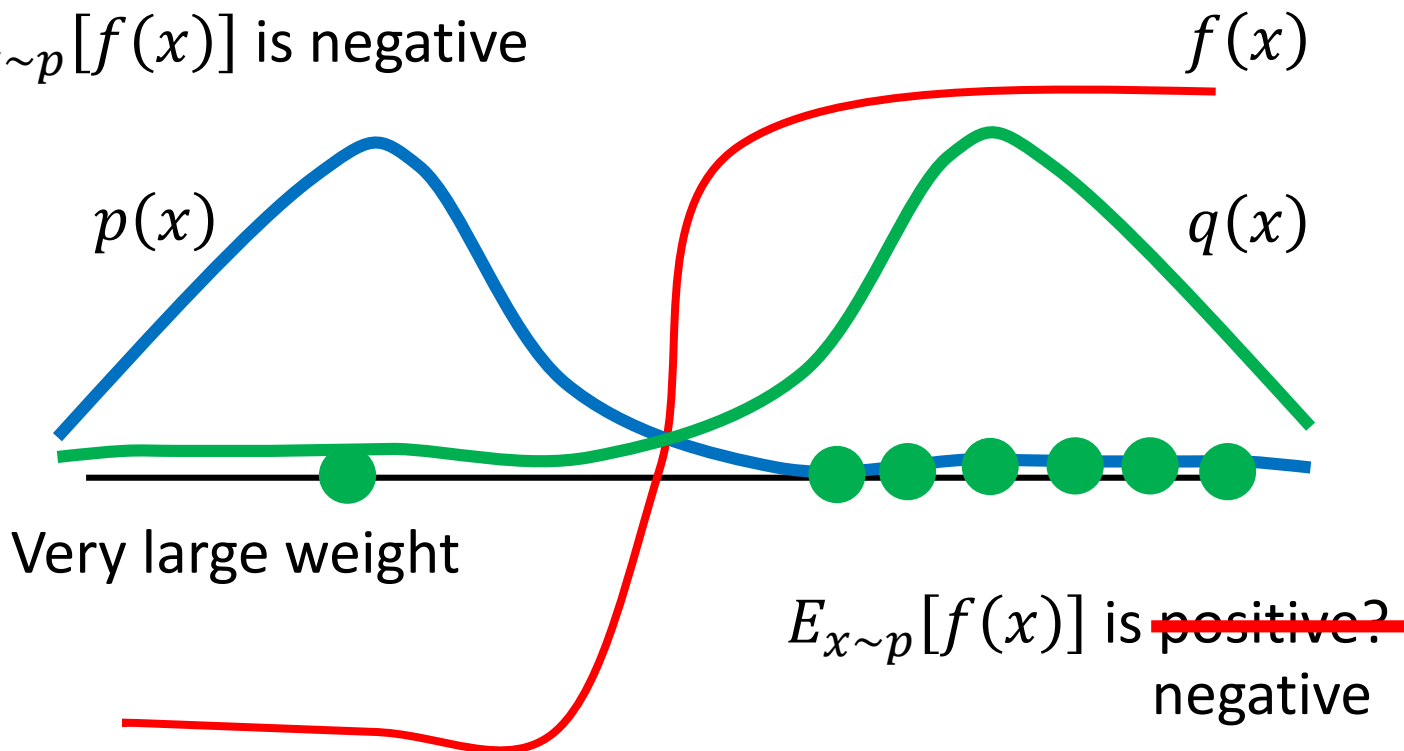
$$= E_{x \sim p}\left[f(x)^2 \frac{p(x)}{q(x)}\right] - (E_{x \sim p}[f(x)])^2$$

# Issue of Importance Sampling

sample不夠多次，等式無法成立

$$E_{x \sim p}[f(x)] = E_{x \sim q}[f(x) \frac{p(x)}{q(x)}]$$

$E_{x \sim p}[f(x)]$  is negative



# On-policy $\rightarrow$ Off-policy

theta' 只負責做demo，負責互動後取出資料

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)]$$

- Use  $\pi_\theta$  to collect data. When  $\theta$  is updated, we have to sample training data again.
- Goal: Using the sample from  $\pi_{\theta'}$  to train  $\theta$ .  $\theta'$  is fixed, so we can re-use the sample data.

$$\nabla \bar{R}_\theta = E_{\tau \sim p_{\theta'}(\tau)} \left[ \overset{\text{important weight}}{\frac{p_\theta(\tau)}{p_{\theta'}(\tau)}} R(\tau) \nabla \log p_\theta(\tau) \right]$$

- Sample the data from  $\theta'$ .
- Use the data to train  $\theta$  many times.

---

Importance  
Sampling

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

# On-policy $\rightarrow$ Off-policy

Gradient for update

$$\nabla f(x) = f(x) \nabla \log f(x)$$

A: 這個state  $s_t$ ，採取action  $a_t$ ，互動完後得到的cumulative reward減去baseline

$$= E_{(s_t, a_t) \sim \pi_{\theta}} [A^{\theta}(s_t, a_t) \nabla \log p_{\theta}(a_t^n | s_t^n)]$$

$$A^{\theta'}(s_t, a_t)$$

This term is from sampled data.

$$= E_{(s_t, a_t) \sim \pi_{\theta'}} \left[ \frac{P_{\theta}(s_t, a_t)}{P_{\theta'}(s_t, a_t)} \cancel{A^{\theta}(s_t, a_t)} \nabla \log p_{\theta}(a_t^n | s_t^n) \right]$$

遊戲出現 $s_t$ 的機率，太難計算了乾

脆忽略他，說服自己不重要

$$= E_{(s_t, a_t) \sim \pi_{\theta'}} \left[ \frac{p_{\theta}(a_t | s_t)}{p_{\theta'}(a_t | s_t)} \cancel{\frac{p_{\theta}(s_t)}{p_{\theta'}(s_t)}} A^{\theta'}(s_t, a_t) \nabla \log p_{\theta}(a_t^n | s_t^n) \right]$$

network output

probability

demo  $\theta'$

$$J^{\theta'}(\theta) = E_{(s_t, a_t) \sim \pi_{\theta'}} \left[ \frac{p_{\theta}(a_t | s_t)}{p_{\theta'}(a_t | s_t)} A^{\theta'}(s_t, a_t) \right] \quad \text{When to stop?}$$

objective function

但這裡有一個限制， $\theta'$ 與 $\theta$ 不能差太多，不然結果會太過不一致

# Add Constraint

穩紮穩打，步步為營



# PPO / TRPO

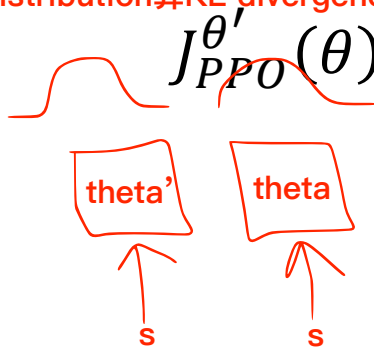
$\theta$  cannot be very different from  $\theta'$   
Constraint on behavior not parameters

## Proximal Policy Optimization (PPO)

對兩個action (network output)  
distribution算KL divergence

regularization, 讓 $\theta'$  與  $\theta$ 不要差太多

$$\nabla f(x) = f(x) \nabla \log f(x)$$


$$J_{PPO}^{\theta'}(\theta) = J^{\theta'}(\theta) - \beta KKL(\theta, \theta')$$

$$J^{\theta'}(\theta) = E_{(s_t, a_t) \sim \pi_{\theta'}} \left[ \frac{p_{\theta}(a_t | s_t)}{p_{\theta'}(a_t | s_t)} A^{\theta'}(s_t, a_t) \right]$$

## TRPO (Trust Region Policy Optimization)

TRPO是將KL divergence變成一個constraint

$$J_{TRPO}^{\theta'}(\theta) = E_{(s_t, a_t) \sim \pi_{\theta'}} \left[ \frac{p_{\theta}(a_t | s_t)}{p_{\theta'}(a_t | s_t)} A^{\theta'}(s_t, a_t) \right]$$

$$KL(\theta, \theta') < \delta$$

# PPO algorithm

- Initial policy parameters  $\theta^0$
- In each iteration
  - Using  $\theta^k$  <sup>using for demo</sup> to interact with the environment to collect  $\{s_t, a_t\}$  and compute advantage  $A^{\theta^k}(s_t, a_t)$
  - Find  $\theta$  optimizing  $J_{PPO}(\theta)$

$$J^{\theta^k}(\theta) \approx \sum_{(s_t, a_t)} \frac{p_{\theta}(a_t|s_t)}{p_{\theta^k}(a_t|s_t)} A^{\theta^k}(s_t, a_t)$$

$$J_{PPO}^{\theta^k}(\theta) = J^{\theta^k}(\theta) - \beta KL(\theta, \theta^k)$$

Update parameters  
several times

adaptive KL divergence

- If  $KL(\theta, \theta^k) > KL_{max}$ , increase  $\beta$
- If  $KL(\theta, \theta^k) < KL_{min}$ , decrease  $\beta$  <sub>反之</sub>

當KL太大，表示最後一項沒有發揮作用，調大weight

Adaptive  
KL Penalty

# PPO algorithm

$$J_{PPO}^{\theta^k}(\theta) = J^{\theta^k}(\theta) - \beta K L(\theta, \theta^k)$$

$$J^{\theta^k}(\theta) \approx \sum_{(s_t, a_t)} \frac{p_{\theta}(a_t | s_t)}{p_{\theta^k}(a_t | s_t)} A^{\theta^k}(s_t, a_t)$$

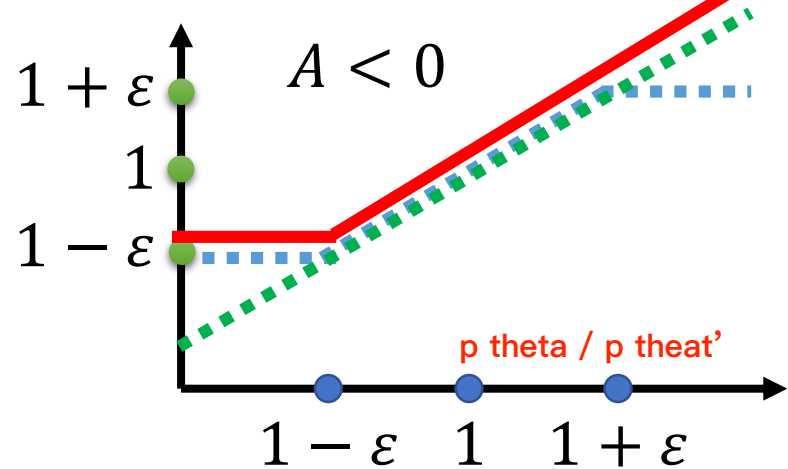
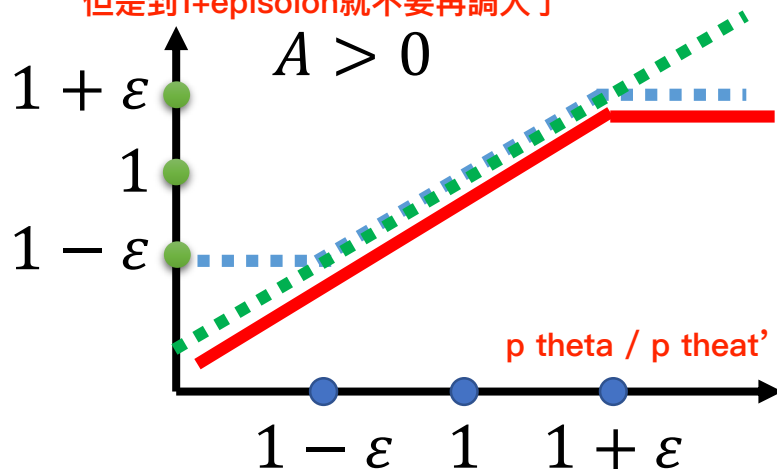
# PPO2 algorithm

$$J_{PPO2}^{\theta^k}(\theta) \approx \sum_{(s_t, a_t)} \min \left( \frac{p_{\theta}(a_t | s_t)}{p_{\theta^k}(a_t | s_t)} A^{\theta^k}(s_t, a_t), \right.$$

$$\left. \text{clip} \left( \frac{p_{\theta}(a_t | s_t)}{p_{\theta^k}(a_t | s_t)}, 1 - \varepsilon, 1 + \varepsilon \right) A^{\theta^k}(s_t, a_t) \right)$$

floor                      ceil  
反之

當今天actor是正的，我們就盡量調大p theta，  
但是到1+episilon就不要再調大了



# Experimental Results

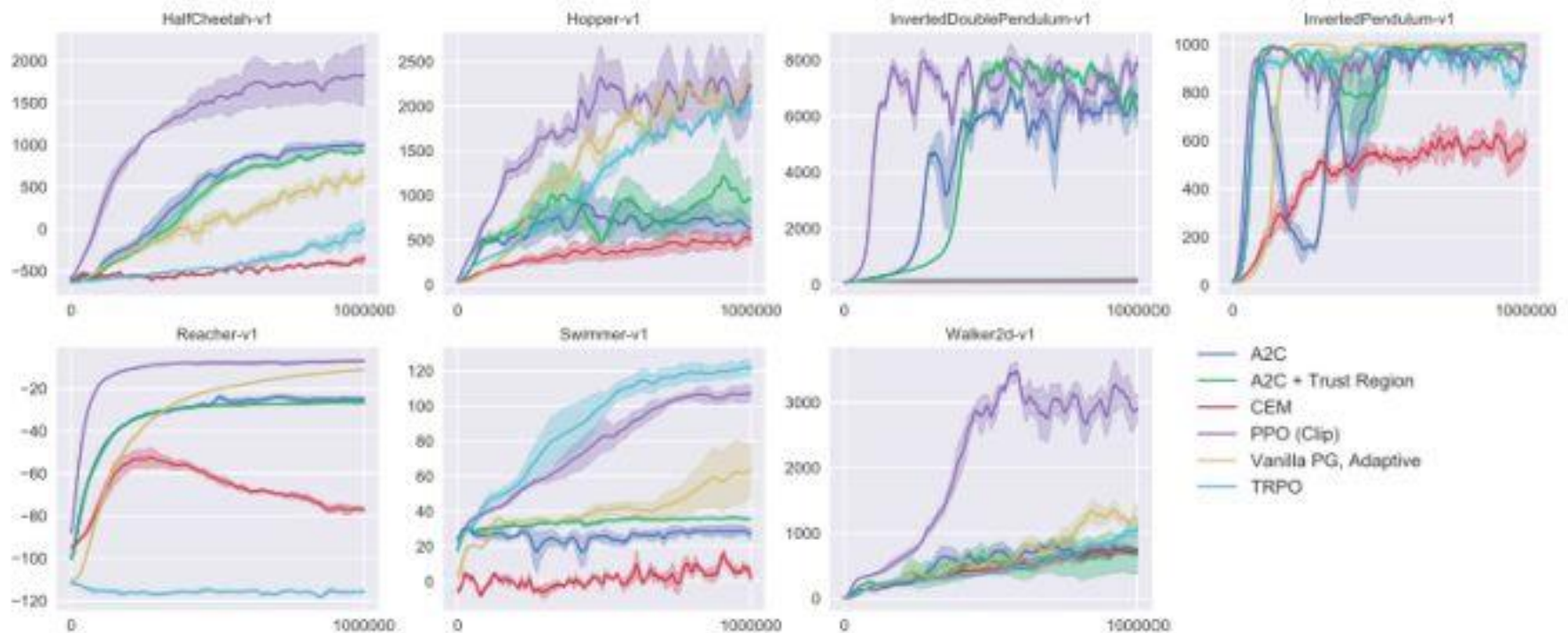


Figure 3: Comparison of several algorithms on several MuJoCo environments, training for one million timesteps.