

# Stochastic Policy Gradient (SPG)

Policy Gradient

策略梯度

Deterministic Policy Gradient (DPG)

SPG formula:

$$\nabla_{\theta} J(\pi_{\theta}) = \int_S p^{\pi}(s) \int_A (\nabla_{\theta} \pi_{\theta}(a|s)) Q_{\pi}(s, a) da ds$$

policy gradient  
对状态  $s$  积分

like Bellman Equation

$$= E_{s \sim p^{\pi}, a \sim \pi_{\theta}} [\underbrace{\nabla_{\theta} \log(\pi_{\theta}(a|s)) \cdot Q^{\pi}(s, a)}_{\text{Bellman Equation}}]$$

SPG相比 DPG 多了一个  $\log()$ .

本质源于 SPG 需要加一层对策略  $\pi$  的期望，而 DPG 中  $\pi$  确定，为  $u$ ，因此不用。

数变换后得  $\log(u)$

DPG formula:

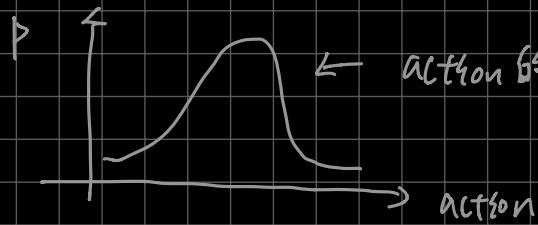
$$\nabla_{\theta} J(\pi_{\theta}) = \int_S p^u(s) \cdot \underbrace{(\nabla_{\theta} u_{\theta}(s)) \cdot \left( \nabla_a Q_u(s, a) \Big|_{a=u_{\theta}(s)} \right)}_{\text{在确定型策略中, } a=u_{\theta}(s)} ds$$

在确定型策略中， $a = u_{\theta}(s)$ ，  
 $\pi$  的确定形式

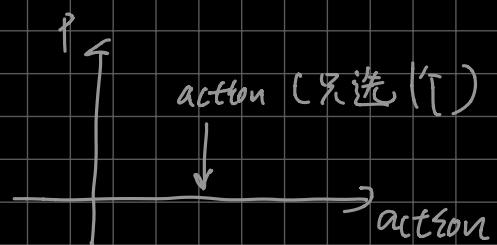
因此没有了内层的积分

$$= E_{s \sim p^u} [(\nabla_{\theta} u_{\theta}(s)) \cdot (\nabla_a Q^u(s, a) \Big|_{a=u_{\theta}(s)})]$$

SPG:



action 的概率分布



action (只选一个)

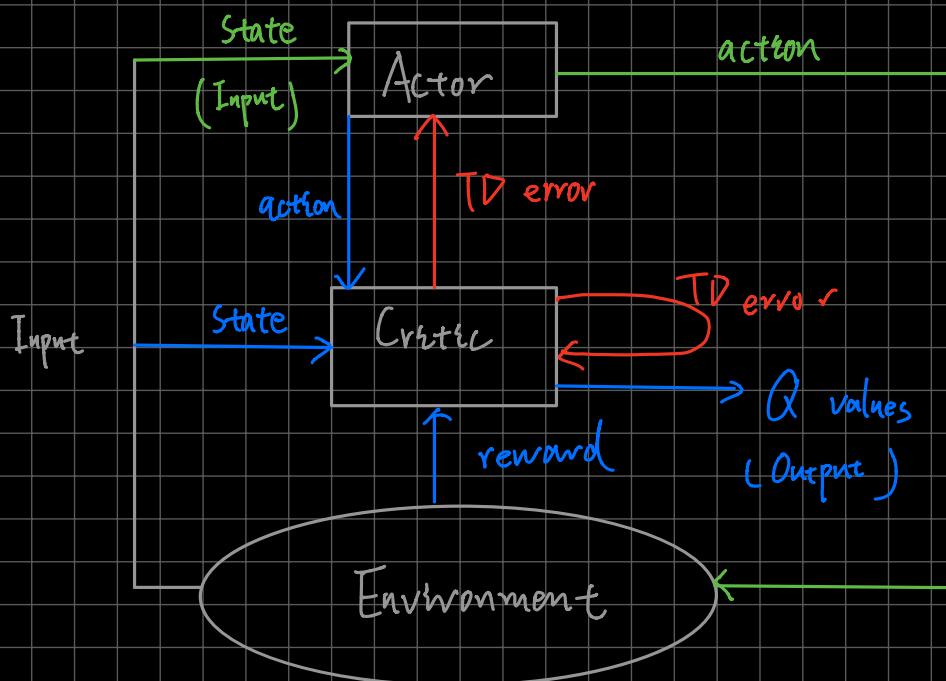
# ACCL Actor - Critic Algorithm

Actor: Policy Q-network, used to choose action

Critic: Value Q-network, used to estimate action

is good or not. And generates TD-error to guide updating.

Actor, Critic to use GD来  
update, 而 Actor 产生的 TD-error 相当于表明此次  
的 ascent 是否正确 (若正石角, 就多 ascent,  
错误就少 ascent)



Why  
↓  
propose AC algorithm?  $\Rightarrow$   
(Google DeepMind)

- {
- DQN, Q-learning, SARSA (SARSA- $\lambda$ ) 等 Value based method 无法处理连续动作空间的情形。
  - 传统的 Policy Gradient 为回合更新制，model 学习效率低下，学习速度缓慢。
- Combine: AC algorithm

优: AC 算法结合了 Value-based 与 Policy based 算法的优势，可处理连续、离散问题，又可单步更新，提升学习效率。

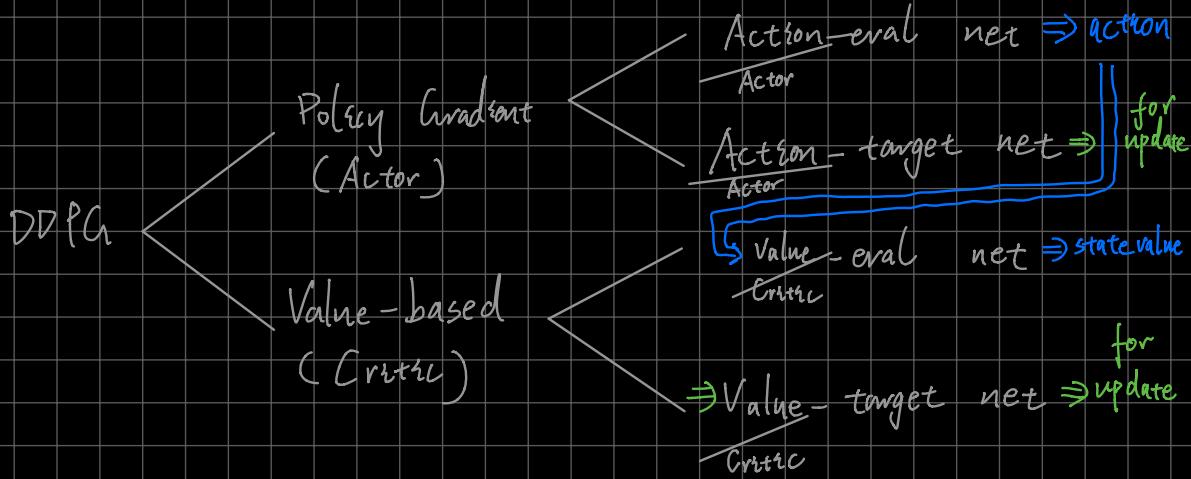
缺: 难收敛  $\Rightarrow$  收敛性取决于 critic 的价值判断，而 critic 依赖的 actor 不断更新，导致更难收敛。

How to solve?  $\Rightarrow$  Deep Deterministic Policy Gradient  
采用 DQN 相似策略，用 memory 储存 transitions，并通过采样来降低数据间相关性，提高收敛性、稳定性

△ 参数更新的相关性越低，算法收敛性越高

DDPG, A3C 均是为降低 AC 算法参数间相关性做的改进（采样、并行更新）

# Famous DPG Algorithm: DDPG (Deep DPG)



DDPG is like DQN + AC. 前半部分(①)很像DQN，后半部分②像AC。

DDPG实际上算 Value-based：因为 policy 更新的目标实际上是最优的 Q。

## Algorithm 1 DDPG algorithm

Randomly initialize critic network  $Q(s, a|\theta^Q)$  and actor  $\mu(s|\theta^\mu)$  with weights  $\theta^Q$  and  $\theta^\mu$ .

Initialize target network  $Q'$  and  $\mu'$  with weights  $\theta^{Q'} \leftarrow \theta^Q, \theta^{\mu'} \leftarrow \theta^\mu$

Initialize replay buffer  $R$

**for** episode = 1, M **do**

    Initialize a random process  $\mathcal{N}$  for action exploration

    Receive initial observation state  $s_1$

**for** t = 1, T **do**      noise sample       $\Rightarrow$  解决连续空间 action 探索能力差的问题

        Select action  $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$  according to the current policy and exploration noise

        Execute action  $a_t$  and observe reward  $r_t$  and observe new state  $s_{t+1}$

        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $R$

            experience replay

        Sample a random minibatch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $R$

Critic-eval network

        Set  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$

        Update critic by minimizing the loss:  $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$

update Critic Network  
action eval network

        Update the actor policy using the sampled gradient:

it's update

为 Q-network, 不是 target

Critic target Network

Actor target Network

$$\nabla_{\theta^\mu} \mu|_{s_i} \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$$

→ 往最大 Q 值的 action 方向  
更新 policy

update Actor Network

        Update the target networks:

soft-update 软更新

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$

$$\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$$

update Actor & Critic Target Network

TQ factor

**end for**

**end for**

# A3C (Asynchronous Advantage Actor-Critic) Algorithm

Brief Introduction: 将 AC 分布于多个 threads 中同步训练。

Advantages:

① 提升训练速度

② 减弱训练相关性，加速收敛

Features:

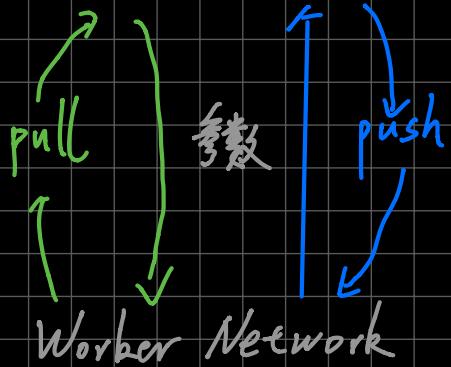
① Global Network and Worker Network

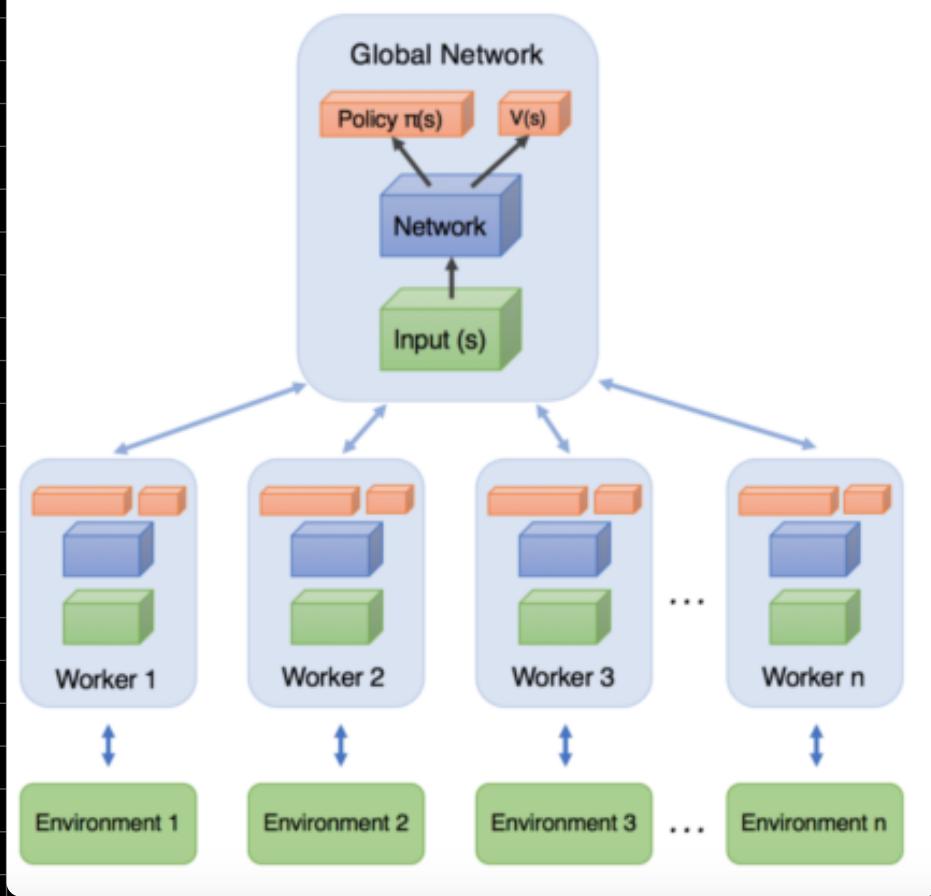
have same structure.

② Global Network stores parameters,

while Worker Network updates parameters  
through training.

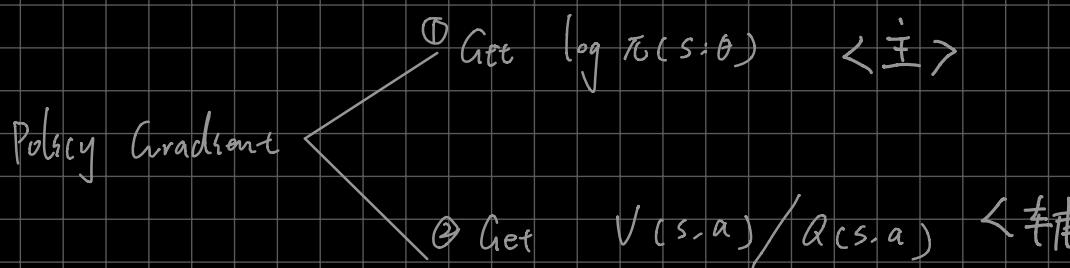
③ Global Network





# Comparison between Policy Gradient methods:

$$\text{Gradient} = \frac{1}{N_s} \sum_s [\nabla_{\theta} \log \pi(s; \theta) \cdot \nabla_{\theta} V(s, a)]$$



① Get  $\log \pi(s; \theta)$  (i) Policy Gradient of discrete problem (CartPole)  
 tf.nn.sparse-softmax-cross-entropy with logits  
 $(\text{logits} = \text{act-outs} - \text{label} = \text{self}.tf-as)$   
 From: PG network  
 output choose-action()  
 $\Downarrow$

$$\text{act-probs} = \text{tf.nn.softmax}(\text{act-outs})$$

action = np.random.choice(range(act-probs.shape[1]), p=act-probs)  
 $\Downarrow$  通过 action 的概率分布选择。  
 $\text{range}(1)$

(ii) AC of discrete problem (CartPole)

$$-\text{tf.log}(\underbrace{\text{self}.act-probs}_{[0, self.a]} \underbrace{[0, self.a]}_{\Downarrow})$$

From: Actor output choose-action()  
 $\Downarrow$

与 (i) 中相同

(iii) AC of continuous problem (Pendulum)

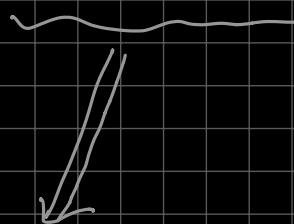
$$\underbrace{\text{self}.normal-dist.}_{\Downarrow} \log-prob(\underbrace{\text{self}.a}_{\Downarrow})$$

tf.distributions.Normal(mu, sigma)

choose-action()

$\text{tf.squeeze}()$ :  $\text{mu} * 2$        $\xrightarrow{\text{sigma} + 0 \sim 1}$

From: Actor Network output



$\text{tf.clip-by-value}(\text{self-normal-dist.sample}(1), \underline{\text{low}}, \overline{\text{high}})$

从高斯分布中采样 (action)  $\xrightarrow{\text{self.action-bound}}$

(iv)  $\nabla_{\theta^u} \mathcal{U}(s; \theta_u)_{si} = 1$

$\text{choose-action}()$ , 通过 deterministic method

选取  $\Rightarrow$

$\text{self.sess.run}(\text{self.a}, \text{feed\_dict})[0]$

可作为 action value 的表示

② Get  $V(s,a) / \text{td\_error} / Q(s,a)$

(i) Policy Gradient of discrete problem (CartPole)

discount-and-norm ( $\text{self} \cdot \text{ep} - \bar{V}$ )

experience reward in  
one turn (回合刷新)

(ii) AC of discrete problem (CartPole)

td-error (From Critic)

$$R + \gamma \cdot V_- - V$$

From: env reward

Critic network output  $\Rightarrow$

$$\begin{aligned} V_- &= \text{input} \\ V &= \text{output} \end{aligned}$$

(iii) AC of continuous problem (Pendulum)

td-error (From Critic)

$$\text{tf-reduce-mean}(R + \gamma \cdot V_- - V)$$

about discrete problem 不同之处

## (iv) DDPG (Pendulum)

+ f. reduce-mean ( $Q$ )  $\Rightarrow$  (From: Critic / eval)  
 continuous problem 要加  $\frac{1}{N}$  原先为:  $\frac{1}{N} \cdot \sum_i \nabla_a Q \cdot \nabla_{\theta U}$

△ DDPG 的 Critic 输出为 Q-value , 而  
 AC 的 Critic 输出为 V-value

$$Q(s, a; \theta) \quad V(s; \theta)$$

更新均靠 Critic network output it 算  
 loss 进行更新.

$$\frac{1}{N} \cdot \sum_{i=1}^s \text{td-error}^2$$