

# RL Tutorial 2

FIT 2017

# DQN Tricks

- Double DQN
- max操作导致“过优化”
- $Q^{\text{target}}(s, a)$ 是目标,  $Q^{\text{approx}}$ 是NN,  $Y$ 是误差,  $a' = \arg\max_a Q^{\text{target}}(s', a)$

$$\begin{aligned} Z &= r_{s,a} + \gamma \max_{a1} Q^{\text{approx}}(s', a1) - r_{s,a} + \gamma \max_{a2} Q^{\text{target}}(s', a2) \\ &= \gamma \max_{a1} Q^{\text{approx}}(s', a1) - \gamma \max_{a2} Q^{\text{target}}(s', a2) \\ &\geq \gamma Q^{\text{approx}}(s', a') - Q^{\text{target}}(s', a') = \gamma Y_{s',a'} \end{aligned}$$

# Double Q Network

- 于是DQN的估计不再是无偏的
- 解决办法：训两个Q网络，一个选择动作一个用于计算，交替更新

---

**Algorithm 1** Double Q-learning

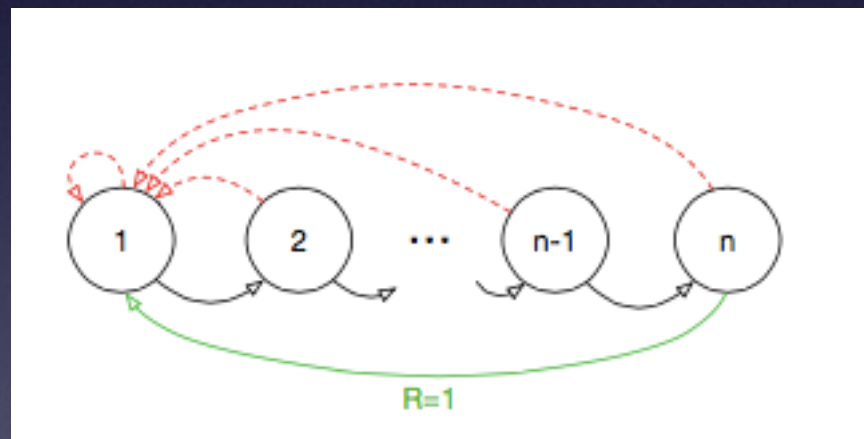
---

```
1: Initialize  $Q^A, Q^B, s$ 
2: repeat
3:   Choose  $a$ , based on  $Q^A(s, \cdot)$  and  $Q^B(s, \cdot)$ , observe  $r, s'$ 
4:   Choose (e.g. random) either UPDATE(A) or UPDATE(B)
5:   if UPDATE(A) then
6:     Define  $a^* = \arg \max_a Q^A(s', a)$ 
7:      $Q^A(s, a) \leftarrow Q^A(s, a) + \alpha(s, a) (r + \gamma Q^B(s', a^*) - Q^A(s, a))$ 
8:   else if UPDATE(B) then
9:     Define  $b^* = \arg \max_a Q^B(s', a)$ 
10:     $Q^B(s, a) \leftarrow Q^B(s, a) + \alpha(s, a) (r + \gamma Q^A(s', b^*) - Q^B(s, a))$ 
11:   end if
12:    $s \leftarrow s'$ 
13: until end
```

---

# Prioritized Replay

- 使更Exciting的样本更容易被采样
- 对于奖励稀疏的问题效果很好，例如



- 维护优先队列
- 显然不能仅仅是贪婪的采样



# Prioritized Replay

- 按照rank-based或者proportional
- 等于改变了样本分布，需要修偏

---

**Algorithm 1** Double DQN with proportional prioritization

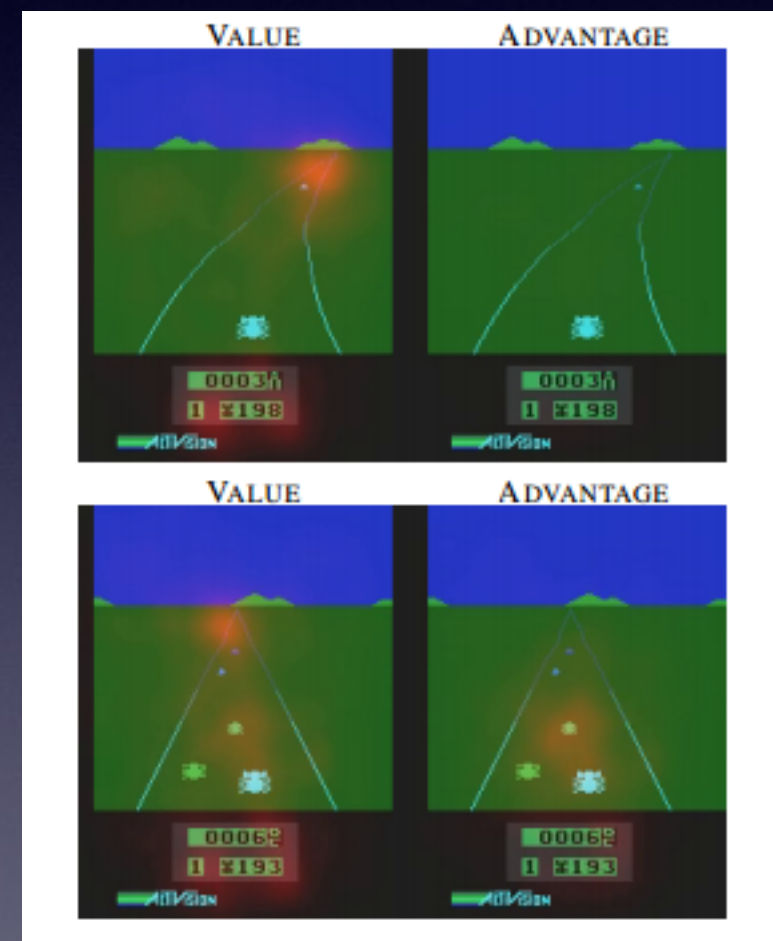
---

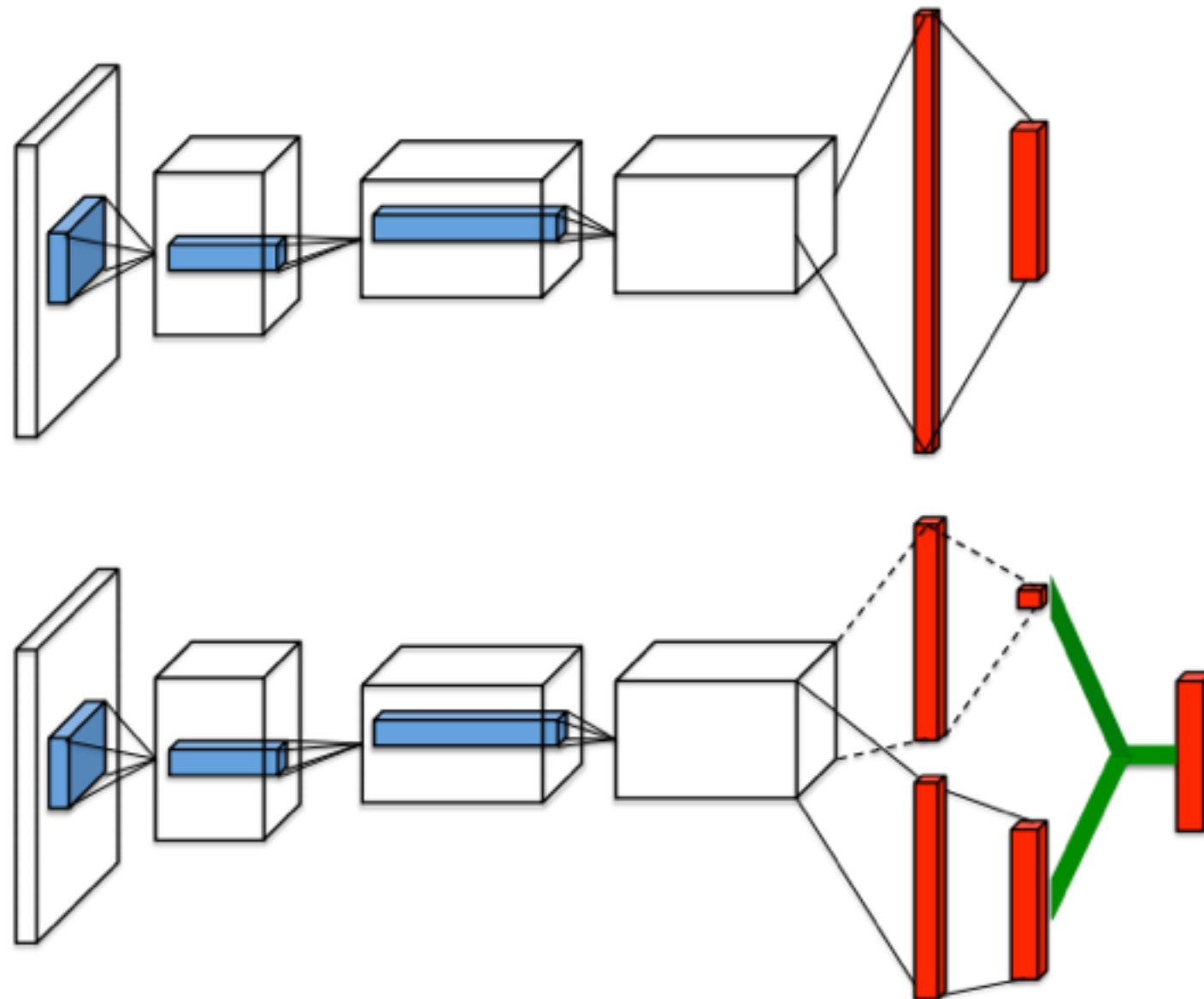
```
1: Input: minibatch  $k$ , step-size  $\eta$ , replay period  $K$  and size  $N$ , exponents  $\alpha$  and  $\beta$ , budget  $T$ .
2: Initialize replay memory  $\mathcal{H} = \emptyset$ ,  $\Delta = 0$ ,  $p_1 = 1$ 
3: Observe  $S_0$  and choose  $A_0 \sim \pi_\theta(S_0)$ 
4: for  $t = 1$  to  $T$  do
5:   Observe  $S_t, R_t, \gamma_t$ 
6:   Store transition  $(S_{t-1}, A_{t-1}, R_t, \gamma_t, S_t)$  in  $\mathcal{H}$  with maximal priority  $p_t = \max_{i < t} p_i$ 
7:   if  $t \equiv 0 \pmod K$  then
8:     for  $j = 1$  to  $k$  do
9:       Sample transition  $j \sim P(j) = p_j^\alpha / \sum_i p_i^\alpha$ 
10:      Compute importance-sampling weight  $w_j = (N \cdot P(j))^{-\beta} / \max_i w_i$ 
11:      Compute TD-error  $\delta_j = R_j + \gamma_j Q_{\text{target}}(S_j, \arg \max_a Q(S_j, a)) - Q(S_{j-1}, A_{j-1})$ 
12:      Update transition priority  $p_j \leftarrow |\delta_j|$ 
13:      Accumulate weight-change  $\Delta \leftarrow \Delta + w_j \cdot \delta_j \cdot \nabla_\theta Q(S_{j-1}, A_{j-1})$ 
14:    end for
15:    Update weights  $\theta \leftarrow \theta + \eta \cdot \Delta$ , reset  $\Delta = 0$ 
16:    From time to time copy weights into target network  $\theta_{\text{target}} \leftarrow \theta$ 
17:  end if
18:  Choose action  $A_t \sim \pi_\theta(S_t)$ 
19: end for
```

---

# Dueling Network

- $Q(s, a) = V(s) + A(s, a)$
- 右图中前方没有车时
- Action随意





# Dueling Network

# 基于价值的方法

- Q-Learning
- 估计状态-动作对的价值, 选取 $\text{argmax}$



# 基于策略的方法

- NN的输出直接就是动作

$$a = \pi(a|s, \theta)$$

# 基本思想

- 存在一个策略  $\pi(a|s; \theta)$
- 如果我们知道每个状态下正确的行动  $a^*$
- 转化为炼丹: 
$$\max_{\theta} \sum_{n=1}^N \log p(a_n^* | s_n; \theta)$$
- 然而我们并不知道  $a^*$

# 基本思想

- 和上次DQN时一样，我们可以收集agent和环境互动的数据，从而得到一个状态-行动-收益序列

$$\tau = ((s_0, a_0, r_0), \dots, (s_{T-1}, a_{T-1}, r_{T-1}), s_T)$$

- 例如，我们令R是序列收益之和

- 那么期望

$$\hat{E} = R \cdot \prod_{t=0}^{T-1} \pi(a_t | s_t; \theta)$$

# 基本思想

$$\hat{E} = R \cdot \prod_{t=0}^{T-1} \pi(a_t|s_t; \theta)$$

$$\hat{g} = \hat{E}' = R \cdot \nabla_{\theta} \left( \prod_{t=0}^{T-1} \pi(a_t|s_t; \theta) \right)$$

- 通常使用log likelihood

$$R \cdot \nabla_{\theta} \left( \sum_{t=0}^{T-1} \log \pi(a_t|s_t; \theta) \right)$$

- 使用梯度下降来优化刚才的期望



# 策略梯度

- 即策略期望总收益的梯度，一般记作  $\nabla_{\theta} J(\theta)$
- 可以证明  $E[\hat{g}]$  是对它的无偏估计
- 换言之，确定型策略梯度算法给出的期望恰好就是策略梯度 (DPG原始论文证明)
- 优势：处理连续场景

# DPG

**Theorem 1** (Deterministic Policy Gradient Theorem). *Suppose that the MDP satisfies conditions A.1 (see Appendix; these imply that  $\nabla_{\theta}\mu_{\theta}(s)$  and  $\nabla_a Q^{\mu}(s, a)$  exist and that the deterministic policy gradient exists. Then,*

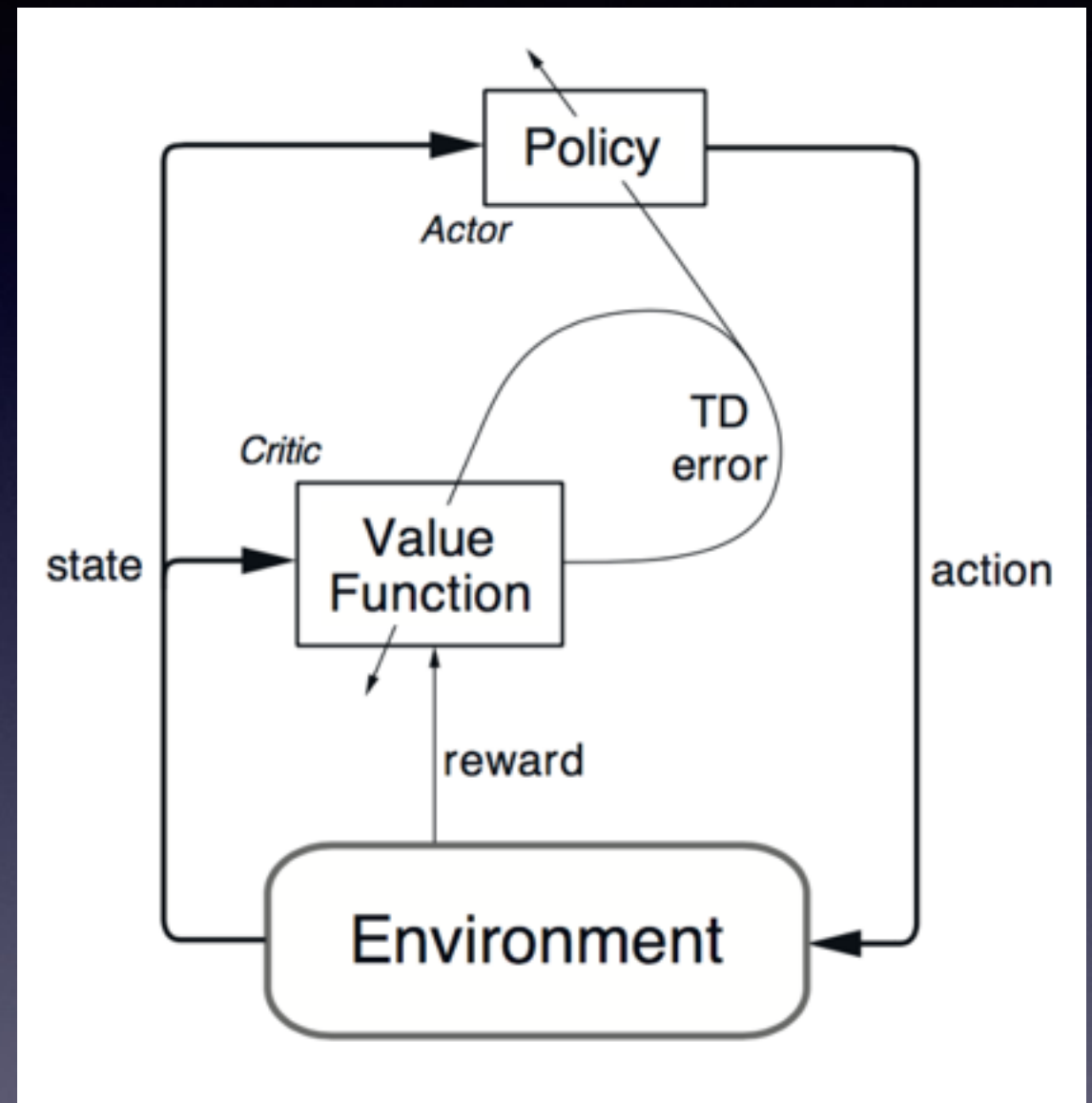
$$\begin{aligned}\nabla_{\theta}J(\mu_{\theta}) &= \int_{\mathcal{S}} \rho^{\mu}(s) \nabla_{\theta}\mu_{\theta}(s) \nabla_a Q^{\mu}(s, a)|_{a=\mu_{\theta}(s)} \mathrm{d}s \\ &= \mathbb{E}_{s \sim \rho^{\mu}} \left[ \nabla_{\theta}\mu_{\theta}(s) \nabla_a Q^{\mu}(s, a)|_{a=\mu_{\theta}(s)} \right] \quad (9)\end{aligned}$$

- 确定性策略梯度
- 按照Q最大的方向调整策略梯度

# Actor-Critic

- 求解策略梯度
- 两个网络
- Actor是策略，参数为 $u$ ， $\pi(s; u)$  输出action
- Critic是值函数，参数为 $w$ ， $Q(s, a; w)$  输出Q值
- Actor为Critic决定 $a$
- Critic为Actor提供Loss function，评判结果好坏

- Critic: Q-Learning Methods
- Actor: 用Q值计算梯度更新策略





- We use a **critic** to estimate the action-value function,

$$Q_w(s, a) \approx Q^{\pi_\theta}(s, a)$$

- Actor-critic algorithms maintain *two* sets of parameters

**Critic** Updates action-value function parameters  $w$

**Actor** Updates policy parameters  $\theta$ , in direction suggested by critic

- Actor-critic algorithms follow an *approximate* policy gradient

$$\nabla_\theta J(\theta) \approx \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)]$$

$$\Delta\theta = \alpha \nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)$$

而不再是  $Q^{\pi_\theta}(s, a)$

- Using linear value fn approx.  $Q_w(s, a) = \phi(s, a)^\top w$

**Critic** Updates  $w$  by linear TD(0)

**Actor** Updates  $\theta$  by policy gradient

**function** QAC

  Initialise  $s, \theta$

  Sample  $a \sim \pi_\theta$

**for** each step **do**

    Sample reward  $r \sim \mathcal{R}_s^a$ ; sample transition  $s' \sim \mathcal{P}_{s'}^a$ .

    Sample action  $a' \sim \pi_\theta(s', a')$

$\delta = r + \gamma Q_w(s', a') - Q_w(s, a)$

$\theta = \theta + \alpha \nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)$

$w \leftarrow w + \beta \delta \phi(s, a)$

$a \leftarrow a', s \leftarrow s'$

**end for**

**end function**

### 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**  
        Select action  $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$  noise sample 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  $R$   
        Sample a random minibatch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $R$  experience replay  
        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  
        Update the actor policy using the sampled gradient:  
            Critic target Network  $\nabla_{\theta^\mu} \mu|_{s_i} \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$  update Actor Network  
            Actor target Network  
        Update the target networks:  
            update Actor & Critic Target Network  
             $\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$   
             $\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$   
    **end for**  
**end for**

# DDPG

Blue tricks

- The **policy gradient** has many equivalent forms

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \mathbf{v}_t] && \text{REINFORCE} \\
 &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q^w(s, a)] && \text{Q Actor-Critic} \\
 &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A^w(s, a)] && \text{Advantage Actor-Critic} \\
 &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta] && \text{TD Actor-Critic} \\
 &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta \mathbf{e}] && \text{TD}(\lambda) \text{ Actor-Critic} \\
 G_{\theta}^{-1} \nabla_{\theta} J(\theta) &= \mathbf{w} && \text{Natural Actor-Critic}
 \end{aligned}$$

<http://blog.csdn.net/>

- Each leads a stochastic gradient ascent algorithm
- Critic uses **policy evaluation** (e.g. MC or TD learning) to estimate  $Q^{\pi}(s, a)$ ,  $A^{\pi}(s, a)$  or  $V^{\pi}(s)$

# Policy Gradient的变种



# A3C

- Asynchronous Advantage Actor-Critic
- 异步：用多线程替代Experience Replay

- 梯度  $\nabla_{\theta'} \log \pi(a_t | s_t; \theta') A(s_t, a_t; \theta, \theta_v)$

- 选择Advantage最大而非Q最大

- A是n-step TD Error  $\sum_{i=0}^{k-1} \gamma^i r_{t+i} + \gamma^k V(s_{t+k}; \theta_v) - V(s_t; \theta_v)$



- ▶ Estimate state-value function

$$V(s, \mathbf{v}) \approx \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \dots | s]$$

- ▶ Q-value estimated by an  $n$ -step sample

$$q_t = r_{t+1} + \gamma r_{t+2} \dots + \gamma^{n-1} r_{t+n} + \gamma^n V(s_{t+n}, \mathbf{v})$$

- ▶ Actor is updated towards target

$$\frac{\partial l_u}{\partial \mathbf{u}} = \frac{\partial \log \pi(a_t | s_t, \mathbf{u})}{\partial \mathbf{u}} (q_t - V(s_t, \mathbf{v}))$$

- ▶ Critic is updated to minimise MSE w.r.t. target

$$l_v = (q_t - V(s_t, \mathbf{v}))^2$$

# A3C流程

---

**Algorithm S2** Asynchronous n-step Q-learning - pseudocode for each actor-learner thread.

---

```
// Assume global shared parameter vector  $\theta$ .
// Assume global shared target parameter vector  $\theta^-$ .
// Assume global shared counter  $T = 0$ .
Initialize thread step counter  $t \leftarrow 1$ 
Initialize target network parameters  $\theta^- \leftarrow \theta$ 
Initialize thread-specific parameters  $\theta' = \theta$ 
Initialize network gradients  $d\theta \leftarrow 0$ 
repeat
  Clear gradients  $d\theta \leftarrow 0$ 
  Synchronize thread-specific parameters  $\theta' = \theta$ 
   $t_{start} = t$ 
  Get state  $s_t$ 
  repeat
    Take action  $a_t$  according to the  $\epsilon$ -greedy policy based on  $Q(s_t, a; \theta')$ 
    Receive reward  $r_t$  and new state  $s_{t+1}$ 
     $t \leftarrow t + 1$ 
     $T \leftarrow T + 1$ 
  until terminal  $s_t$  or  $t - t_{start} == t_{max}$ 
   $R = \begin{cases} 0 & \text{for terminal } s_t \\ \max_a Q(s_t, a; \theta^-) & \text{for non-terminal } s_t \end{cases}$ 
  for  $i \in \{t - 1, \dots, t_{start}\}$  do
     $R \leftarrow r_i + \gamma R$ 
    Accumulate gradients wrt  $\theta'$ :  $d\theta \leftarrow d\theta + \frac{\partial (R - Q(s_i, a_i; \theta'))^2}{\partial \theta'}$ 
  end for
  Perform asynchronous update of  $\theta$  using  $d\theta$ .
  if  $T \bmod I_{target} == 0$  then
     $\theta^- \leftarrow \theta$ 
  end if
until  $T > T_{max}$ 
```

训练数据(s,a,r)序列生成, 最大t\_max个

n-step TD error

累积"batch"内的梯度

异步更新网络参数

低频率同步target网络参数

# 具体算法

# Examples

- <https://github.com/dennybritz/reinforcement-learning> 各类资源、习题与答案
- <https://github.com/yanpanlau/DDPG-Keras-Torcs> 300行DDPG训练TORCS自动驾驶 (Keras+TF)
- <https://github.com/NVlabs/GA3C> 核弹厂的A3C

# Thanks!

Q&A?