

# Deep - Q Network Learning

Why Deep ?  $\Rightarrow$

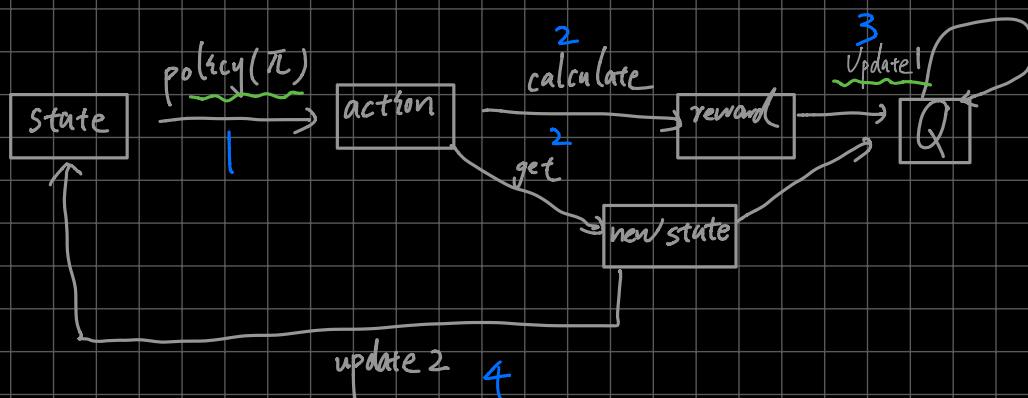
当 Action Space , state space 大到一定规模之后,  
Q-Learning 使用的 Q-table 方法便不能应对 (无法  
得到一张具有如此之大储存空间的 Q-table.)

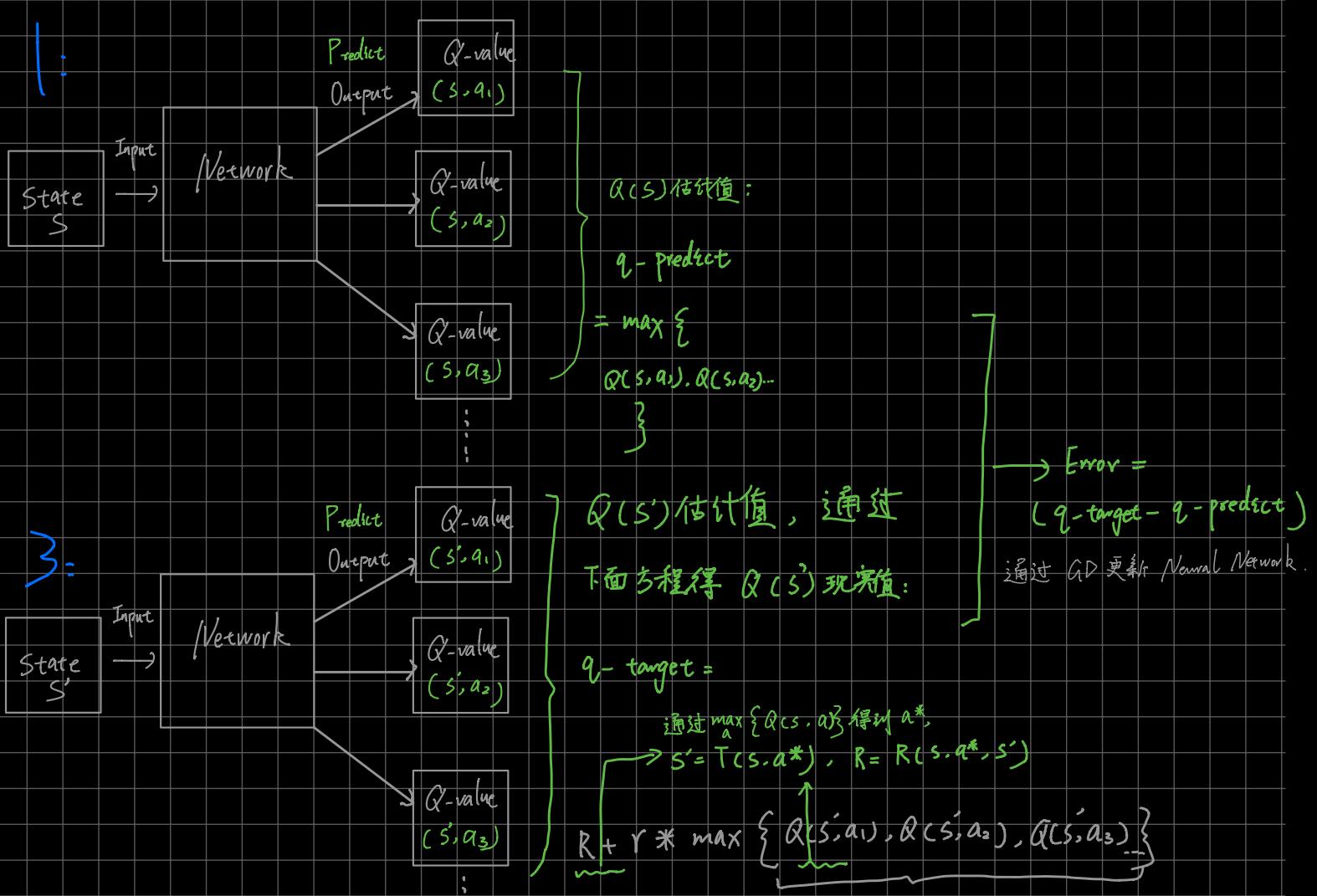
Solve: Deep Neural Network . ① 构建 NN , 使得  
其输入为  $s$  , 轮输出为  $A$  , 可以完成 choose-action  
② 构建 NN , 使得其输入为  $s-A$  , 轮输出为  $Q$  , 可  
以完成 Policy-update .

How ?  $\Rightarrow$

DQN 仍采用与 Q-Learning 相同的 流程:

(除了上述 涉及 Q-table 的二  
部分由 NN 替代之外)





Key points?  $\Rightarrow$  Experience Replay + Fixed Q-targets  
 (DQN 两大利器)

Experience Replay (经验回放):

作用: 解决 Sample 特性 (High correlated and non-stationary.) 带来的难以收敛问题。

实质: 将 online  $\rightarrow$  semi-offline 用 memory 储存样本, 通过随机采样去除相关性。

Fixed Q-target

作用: Separate Target Network from Q-network. 使训练震荡发散可能性降低, 更加稳定。

震荡发散: 若 y 使用同一个 Q-network, 则使 Q 大的样本, y 也会大, 则震荡发散可能性增大。

$\Rightarrow$  用估计值去逼近一个不断更新的 target label, 训练是肯定不稳定的。在 ML

中, 训练时 label 是已知的, 因此是让估计值朝着固定的目标更新优化的。而 RL

没有 label, 因此我们训练时把“假定的”label 固定不变, 让估计值去逼近它。

这就是 Fixed Q-target 思想来源。

# DQN Algorithm :

**Algorithm 1:** deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta \Rightarrow eval\ network$

Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta \Rightarrow target\ network$

**For** episode = 1,  $M$  **do**

    Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$

**For**  $t = 1, T$  **do**

        With probability  $\varepsilon$  select a random action  $a_t$       **action selection**

        otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$

        Execute action  $a_t$  in **emulator** and observe reward  $r_t$  and image  $x_{t+1}$

        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $D$       **Experience Replay**

        Sample random minibatch of transitions  $(\phi_i, a_i, r_i, \phi_{i+1})$  from  $D$

        Set  $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$       **target Network**

        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to the **样本** network parameters  $\theta$       **更新 Q-network 权重**      **模型优化的目标**

        Every C steps reset  $\hat{Q} = Q$       **每隔C步更新target network**

**End For**

**End For**

Replay Memory ( $D$ ) : capacity  $N$

1: Initialization       $\left\{ \begin{array}{l} Eval\ Network\ (Q) : random\ weights\ \theta \\ Target\ Network\ (\hat{Q}) = \theta^- = \theta \end{array} \right.$

2: For episode in EPISODE :

$$s_1 = \{x_1\}$$

$$\phi_1 = \phi(s_1)$$

$$step\_counter = 0$$

3: while not Is-terminated :

4:  $a_t = choose\_action(s_t) \rightarrow \varepsilon\text{-greedy} :$

$$a_t = \begin{cases} \text{random choice (ACTION)} \\ \operatorname{argmax}_a Q(\phi(s_t), a; \theta) \end{cases}$$

5:  $r_t, x_{t+1} = \text{get-env-feedback}(s_t, a_t)$

$\xrightarrow{\text{reward bias}} \Rightarrow \text{Dualing Network}$

6:  $s_{t+1} = s_t$

$$\phi_{ct+1} = \phi(s_{t+1})$$

transition

Store  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $D$

7: Experience Replay

Sample transition batches from  $D$

randomly

8: Policy update:

$\Rightarrow$  prioritized replay

$$y_t = \begin{cases} r_j & \Rightarrow s_{t+1} = \text{"terminated"} \\ r_t + \gamma \cdot \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$$

next state

$\hookrightarrow$  upward bias  $\Rightarrow$

9: Gradient descent:

Double Q-Network

$$\text{error} = \left( \underbrace{y_t - \hat{Q}(\phi_t, a_t; \theta)}_{q\text{-target}} \right)^2$$

$\hat{Q}$  - predict

$$\theta = \theta - \alpha \frac{\partial \text{error}}{\partial \theta}$$

10: Step-counter  $t = 1$

if step-counter  $= C$ :

$$\theta^- = \theta \longrightarrow \hat{Q} = Q$$

# Advanced DQN

Double Q-Network : solve upward-bias problem

DQN 三大改进

Prioritized Experience Replay : solve random choose mini-batch

① 变相增加样本 ② 独立于训练过程中状态的影响

Dueling Network : solve reward bias problem

$Q(s, a)$   $\begin{cases} \rightarrow V(s) \text{ (标量)} \\ \rightarrow A(s, a) : \text{Advantage} \text{ (矢量)} \end{cases}$

## Double DQN:

Problem:

$$q\text{-target}_t = R_{t+1} + \gamma \cdot \max_a \hat{Q}(S_{t+1}, a; \theta_t^-)$$

$\downarrow$

target-net

Upward-bias: 由于  $\hat{Q}(S_{t+1}, a; \theta_t^-)$  本身就有误差

而估计  $q\text{-target}$  时采用 max policy, 因此每次向着最大误差的方向改进神经网络, 导致 Overestimate.

Solve  $\Rightarrow$

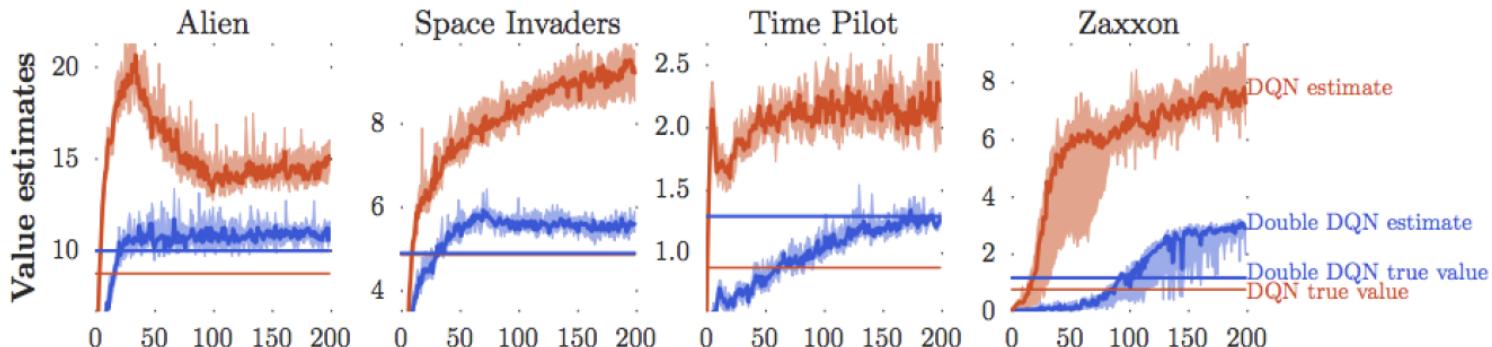
$$q\text{-target}_t = R_{t+1} + \gamma \cdot \hat{Q}(S_{t+1}, \underbrace{\arg\max_a Q(S_{t+1}, a; \theta_t^-)}_{\downarrow}; \theta_t^-)$$

通过 eval-net 估计 action values, 通过 np.argmax (action values, axis=1)

选出 action index, 再通过 target-net 估计  $q\text{-target}$

Function: Solve upward bias problem, 使得 action values 的估计值下降, 接近中心.

- Q value is usually over-estimated



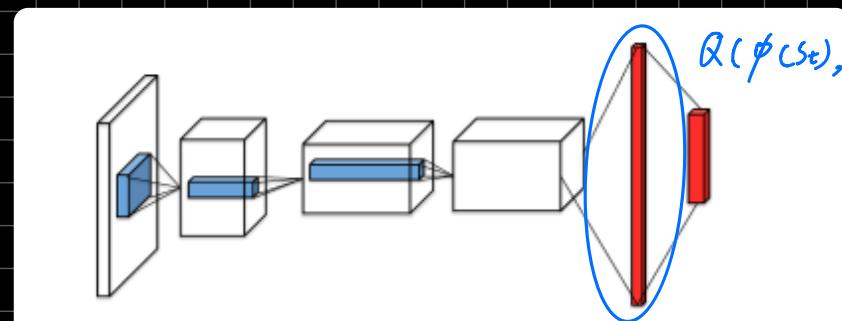
## Dueling DQN:

Problem: Reward bias. 在游戏中，存在很多对 action 不敏感的状态 ( $s \rightarrow s'$  不受采取 action 的影响)。在这些状态下计算 action value func 意义没有 state value func 大。因此若  $\text{state value} \gg \text{action value}$ , 此时 Q value 取决于 state value

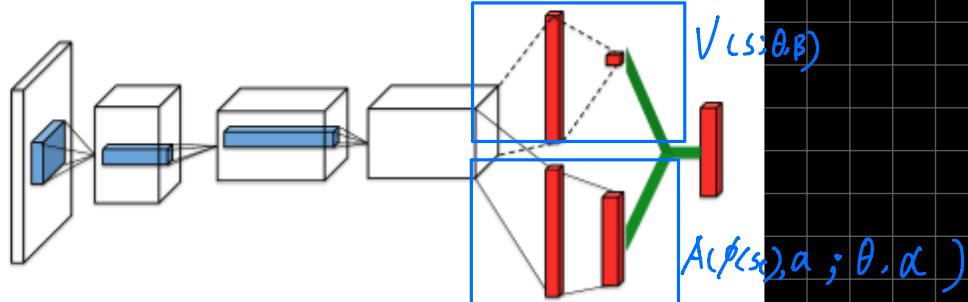
要使用 Q value 来选取 best action, 需去除 state value 的影响, 这样才能体现 action value 的作用。因此采用 Dueling network.

Solve  $\Rightarrow$

Original:



Dueling:



Original DQN:  $Q = \underbrace{Q}_{\text{eval-net}}(\phi(s_t), a; \theta)$

Dueling DQN: eval-net, predict the value of each action, and choose the max Q value action

$$(Q(s_t, a; \theta, \alpha - \beta)) = \underbrace{V(s; \theta, \beta)}_{\text{Value}} + \underbrace{A(s, a; \theta, \alpha)}_{\text{advantage of: action from current state}}$$

predicted Q value

$\Rightarrow$  实际应用中，采取 average 操作，使 advantage 更加稳定，便于优化。

$$Q(s_t, a; \theta, \alpha - \beta) = V(s; \theta, \beta) + \underbrace{[A(s, a; \theta, \alpha) - \frac{1}{|A|} \cdot \sum_a A(s, a; \theta, \alpha)]}_{\downarrow}$$

self.V + self.A - np.reduce\_mean(self.A, axis=1, keepdims=True)

Function: ① 提高 DQN 收敛速度

② 降低 model variance，使得 model 范化性  
能更好，更稳定。