

## Lecture 23 - April 11, 2023

- Deep Q-Network (DQN) → Finite action
  - Deep Q-Network (DQN)
  - Double DQN
  - Prioritized DQN
  - Dueling DQN
- Deep Policy Gradients (DPG) → Large/continuous Action
  - REINFORCE
  - REINFORCE with Baseline
  - Advantage Actor Critic (A2C)
  - Deep Deterministic Policy Gradient (DDPG)

Project 3 → Due April 17

HWS → Due April 18

TA's office hour:

Wednesdays, 2pm-3pm (in-person)

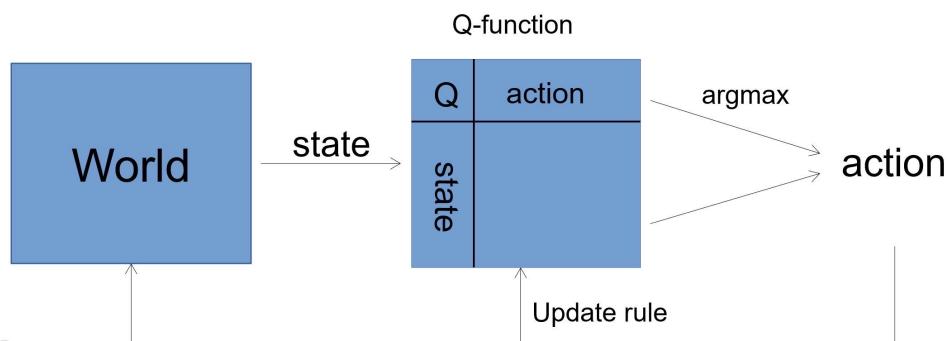
Fridays, 2pm-3pm (virtual)

## Q-learning = Tabular

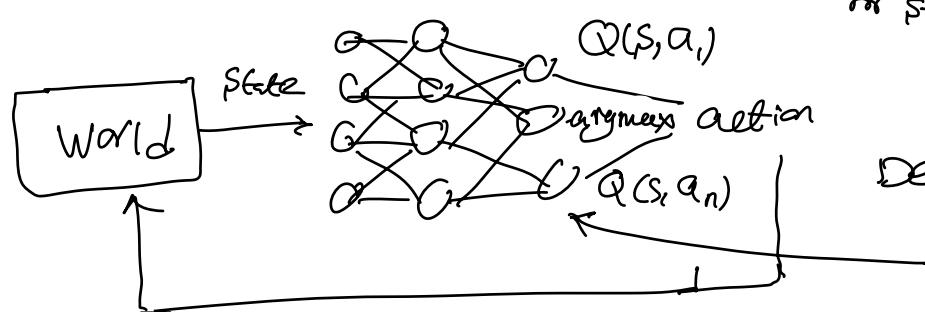
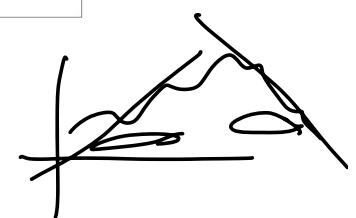
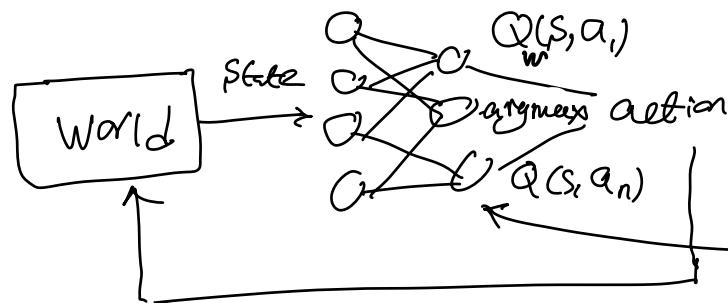
Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$   
 Repeat (for each episode):  
 Initialize  $S$   
 Repeat (for each step of episode):  
 Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)  
 Take action  $A$ , observe  $R, S'$   

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$
  

$$S \leftarrow S'$$
  
 until  $S$  is terminal



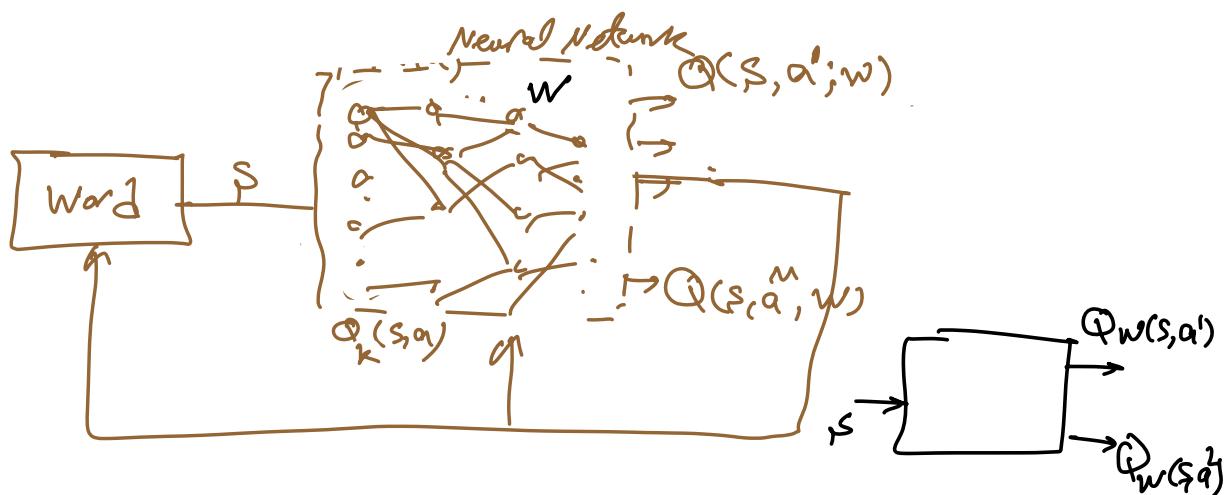
## Q-Learning with Linear Function Approximation



hand crafted  
basis function  
for state-action  
pair

Deep Q-network

State  $\leftarrow$  big & complex  $\leftarrow$  raw data (Images)



Q-Learning Update:

$$Q(s, a) = Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

*Q-learning Target*

*Q-learning error*

① model  $Q_w(s, a)$  with NN with weights  $w$ .

② Define Loss function

$$L(s, a, s'; w) = \frac{1}{2} (r + \gamma \max_a Q_w(s', a) - Q_w(s, a))^2$$

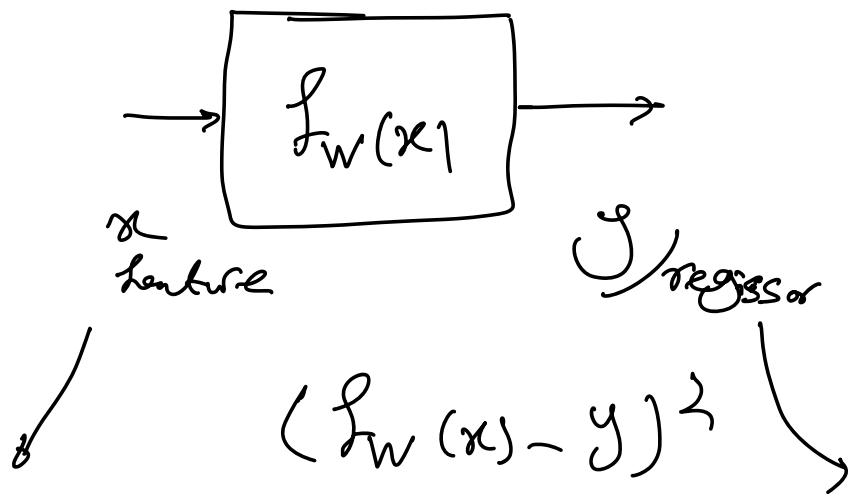
*y*

*Q is parameterized by weights w*

$$\nabla_w L(s, a, s'; w) \approx - (r + \gamma \max_a Q_w(s', a) - Q_w(s, a)) \nabla_w Q_w(s, a)$$

$$W \leftarrow W - \alpha \nabla_w L(s, a, s'; w)$$

$s, a, r, s'$



$$s_0, a_0 \quad Q_W(s_0, a_0) \quad y_0 = r + \max_{a \in A} Q_W(s', a)$$

$s_1, a_1$

$y_1$

$s_2, a_2$

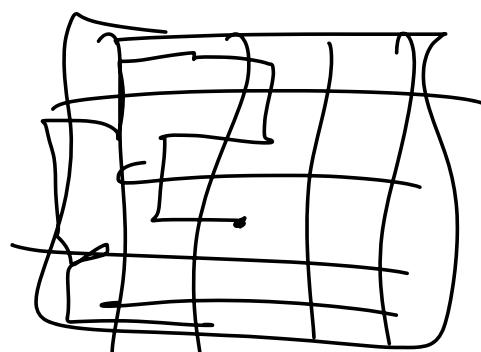
$y_2$

1

1

1

r



## Bone bones DQN (Do Not Converge)

Initialize  $Q(s, a; w)$  with random weights

Repeat (for each episode):

    Initialize  $s$

    Repeat (for each step of the episode):

        Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g.  $\epsilon$ -greedy)

        Take action  $a$ , observe  $r, s'$

$$w \leftarrow w - \alpha \nabla_w L(s, a, s'; w)$$

$$s \leftarrow s'$$

    Until  $s$  is terminal

Where:

$$\nabla_w L(s, a, s'; w) \approx - \left( r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a) \right) \nabla_w Q_w(s, a)$$

Deep Learning (and most Supervised learning) typically assume independent, identically distributed (iid) training data.

How about in the above algorithm   
 General Deep RL

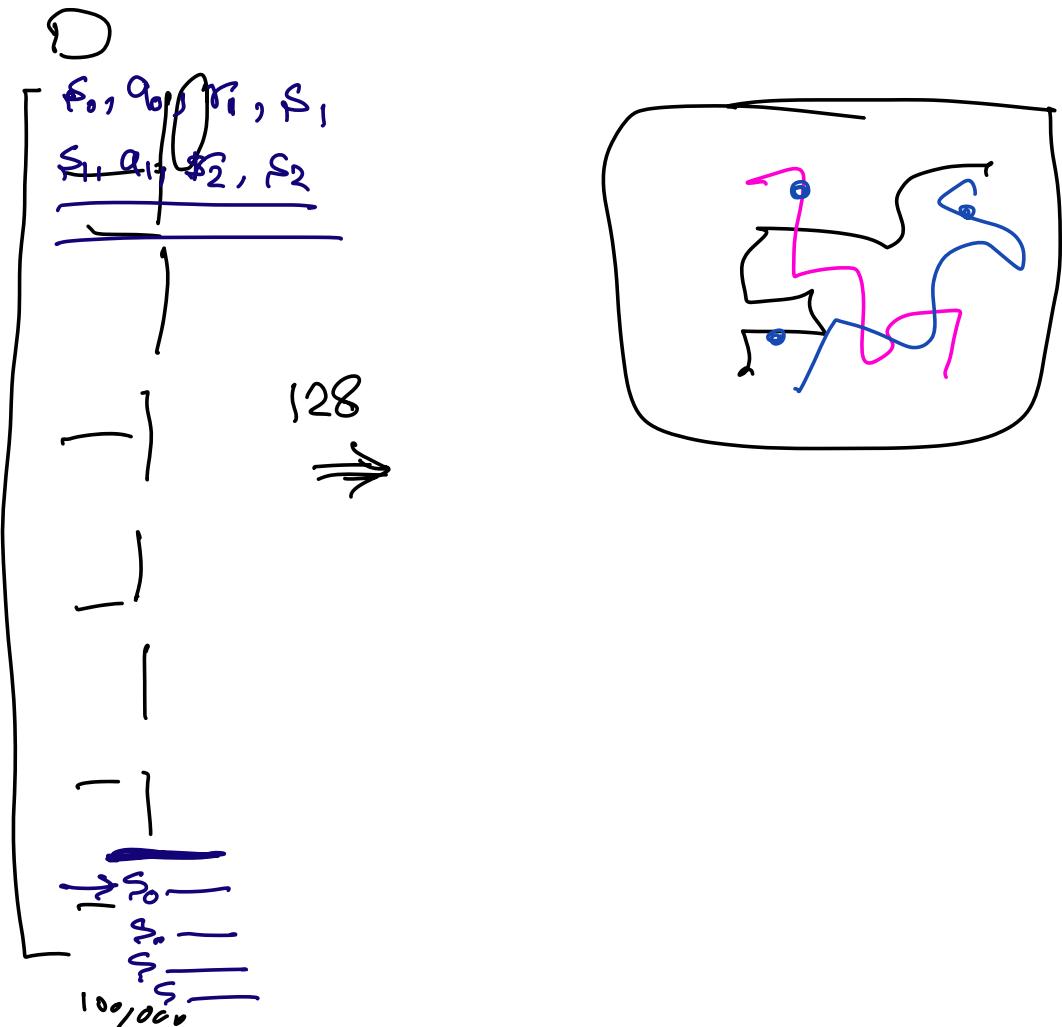
⇒ Solution: buffer experiences and then  
    Reply randomly them during training

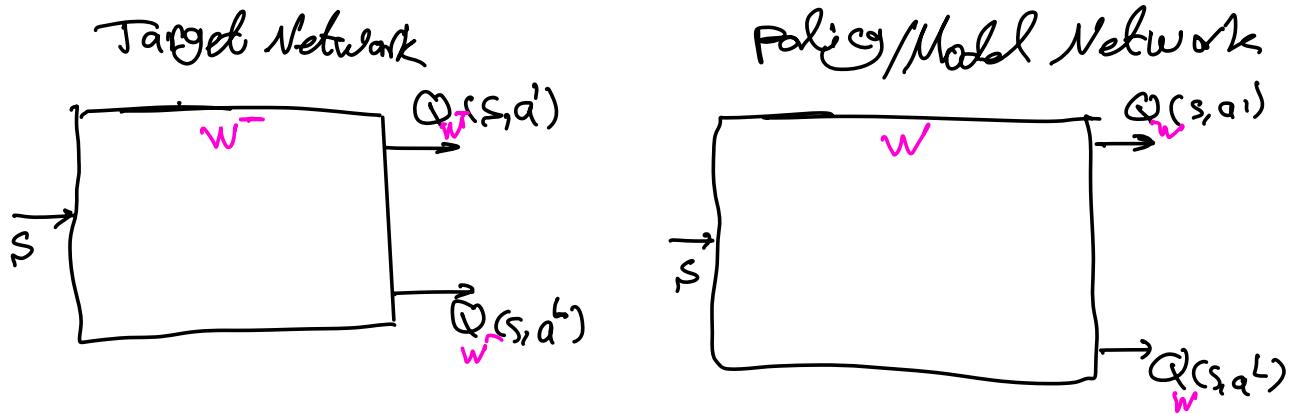
# LETTER

doi:10.1038/nature14236

## Human-level control through deep reinforcement learning

Volodymyr Mnih<sup>1\*</sup>, Koray Kavukcuoglu<sup>1\*</sup>, David Silver<sup>1\*</sup>, Andrei A. Rusu<sup>1</sup>, Joel Veness<sup>1</sup>, Marc G. Bellemare<sup>1</sup>, Alex Graves<sup>1</sup>, Martin Riedmiller<sup>1</sup>, Andreas K. Fidjeland<sup>1</sup>, Georg Ostrovski<sup>1</sup>, Stig Petersen<sup>1</sup>, Charles Beattie<sup>1</sup>, Amir Sadik<sup>1</sup>, Ioannis Antonoglou<sup>1</sup>, Helen King<sup>1</sup>, Dharshan Kumaran<sup>1</sup>, Daan Wierstra<sup>1</sup>, Shane Legg<sup>1</sup> & Demis Hassabis<sup>1</sup>



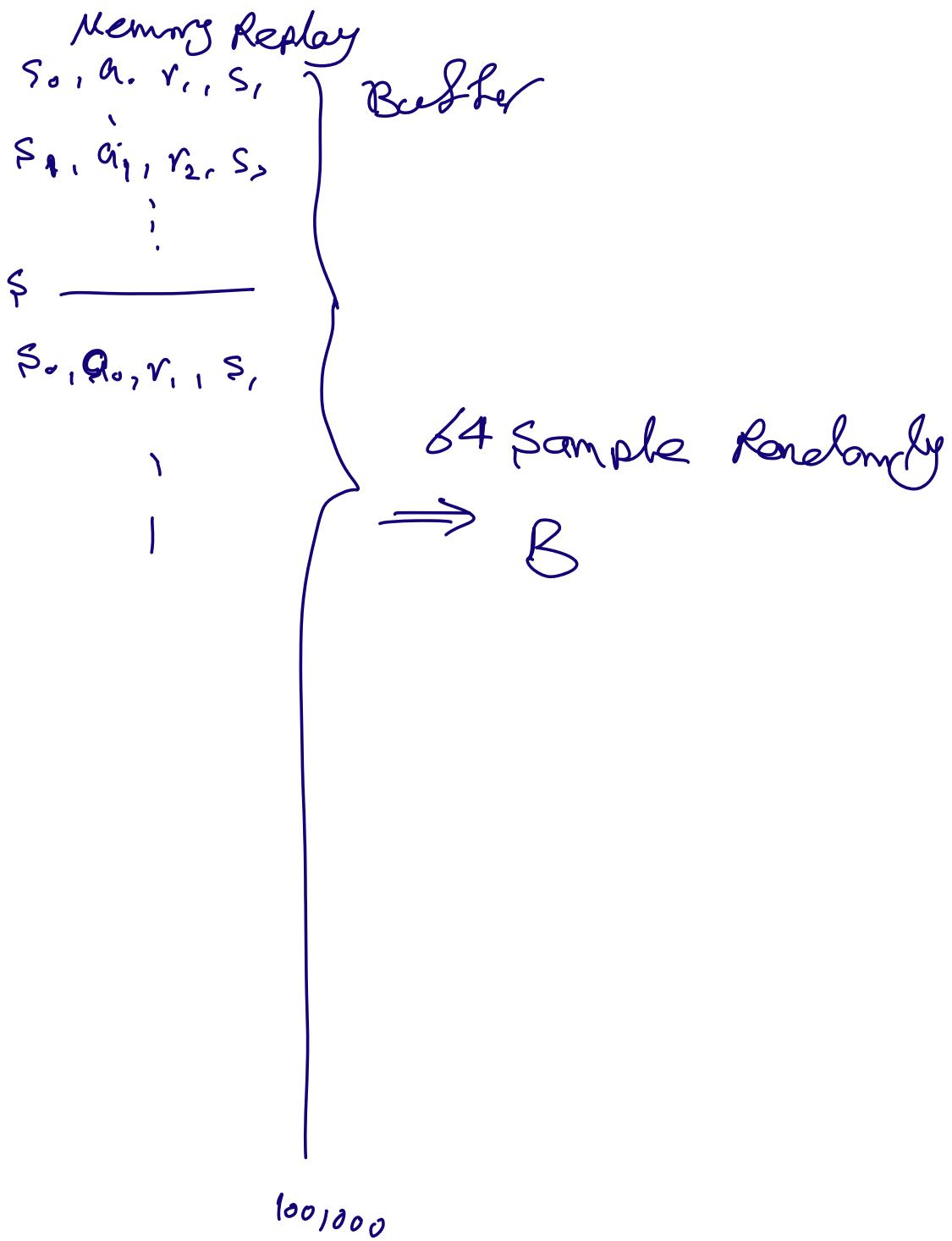


The same Structure

$$a_0 \sim \pi(s_0) = \begin{cases} \text{argmax } Q_W(s_0, a) & 1 - \epsilon \\ \text{Random} & \epsilon \end{cases}$$

Diagram illustrating the selection of action  $a_0$  from state  $s_0$ :

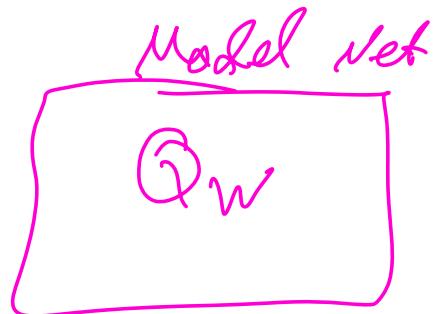
- An arrow points from  $s_0$  to  $a_0$ .
- An arrow points from  $a_0$  to  $r_i$ .
- An arrow points from  $r_i$  to  $s_1$ .



100,000

$$L(B; W, \tilde{W}) = \frac{1}{B} \sum_i^B (r_{i+1} + \gamma \max_{a \in \tilde{W}} Q_{\tilde{W}}(s_{i+1}, a) - Q_W(s_i, a))^2$$

$$W \leftarrow W + \alpha \frac{1}{B} \sum_{i=1}^{100,000} (y_i - Q_W(s_i, a_i)) \nabla Q_W(s_i, a_i)$$



$$r + \gamma \max_{a'} Q_{W^-}(s', a') = Q_W(s, a)$$

B

$s_1, a_1$

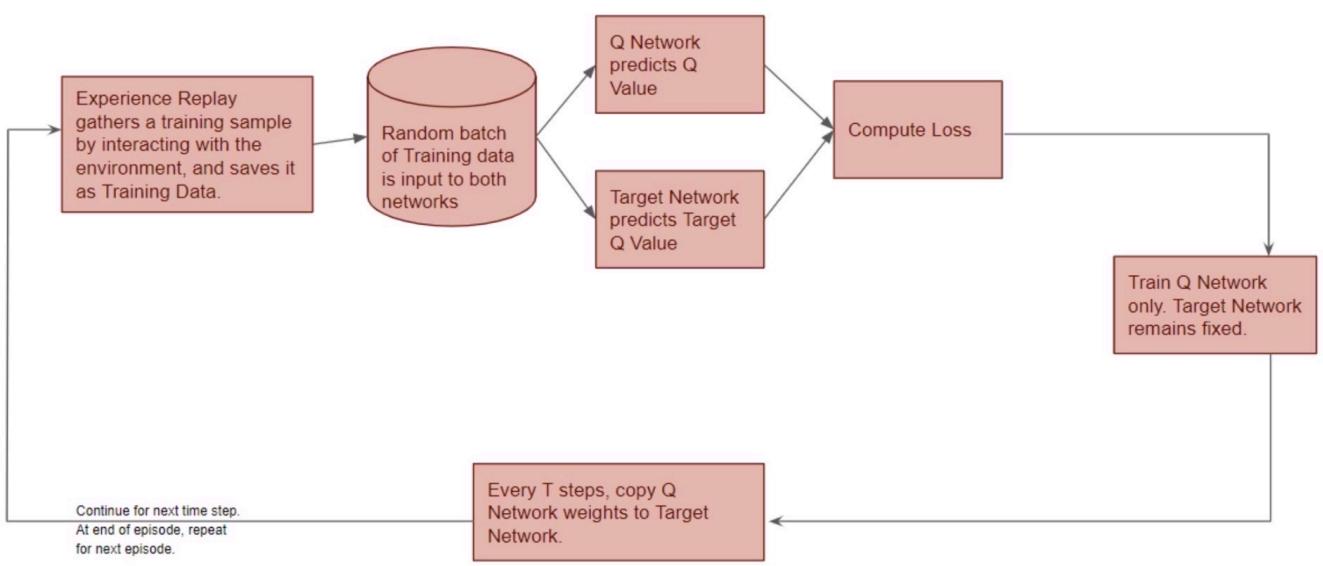
$$y_1 = r_1 + \gamma \max Q_W(s'_1, a')$$

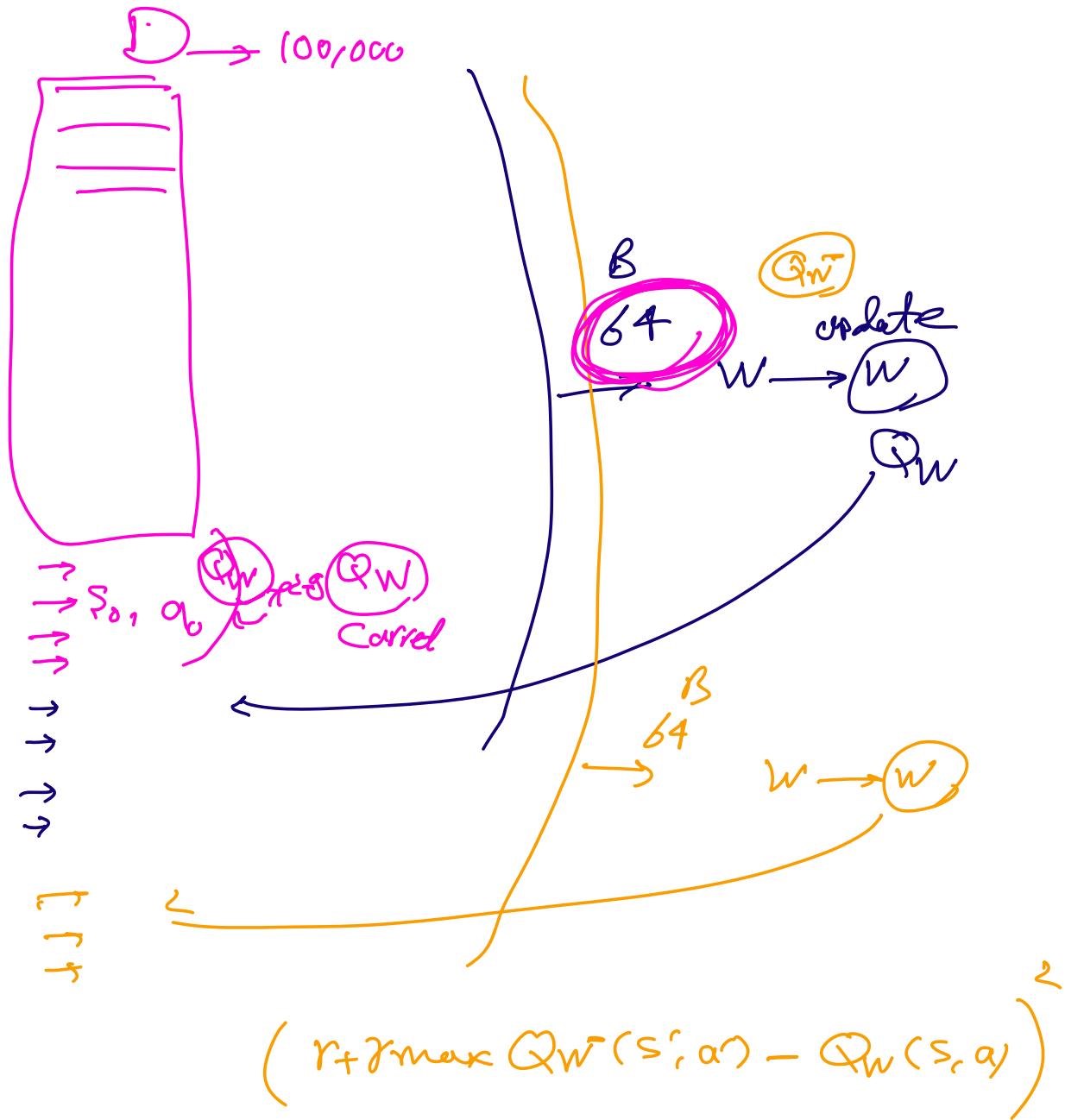
$\rightarrow$

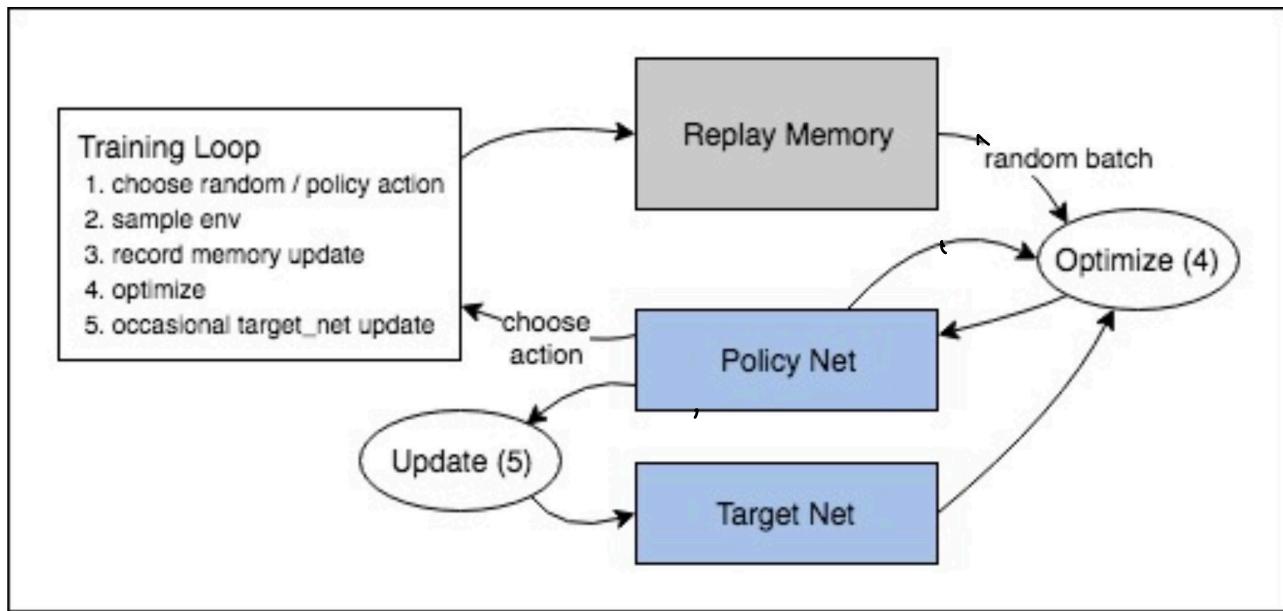
$\rightarrow$

$\rightarrow$

$\rightarrow$







# DQN

Initialize  $Q_w, Q_{w^-}$  with random weights  
 $D \leftarrow \emptyset$

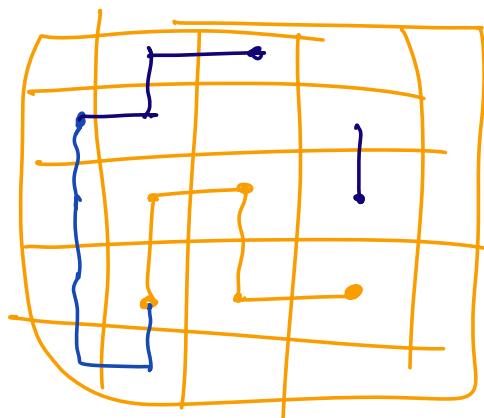
Repeat (for each episode):

- Initialize  $s$
- Repeat (for each step of the episode):
  - Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g., e-greedy)
  - Take action  $a$ , observe  $r, s'$
  - $D \leftarrow D \cup (s, a, s', r)$
  - $s \leftarrow s'$
  - If  $\text{mod}(\text{step}, \text{trainfreq}) == 0$ :  $\rightarrow 64 \text{ or } 128$ 
    - sample batch  $B$  from  $D$
    - $w \leftarrow w - \alpha \nabla_w L(B; w, w^-)$   $\rightarrow$  Default 500 step
  - if  $\text{mod}(\text{step}, \text{copyfreq}) == 0$ :
  - $w^- \leftarrow w$

Where:  $\nabla_w L(B; w, w^-) \approx -\frac{1}{|B|} \sum_{(s, a, s', r) \in B} (\text{target}(s'; w^-) - Q_w(s, a)) \nabla_w Q_w(s, a)$

$$\text{target}(s'; w^-) = r + \gamma \max_{a'} Q_{w^-}(s', a')$$

$D \sim \text{Experience } \times^*$   $Q_w = Q^*$



Two Tricks

Start  $\epsilon$  at 1 and reduce it to 0 over episode

Soft update of  $Q_W$

$$W' \leftarrow W\gamma + W(1-\gamma) \quad \gamma = 0.001$$

# DQN

Initialize  $Q_w, Q_{w^-}$  with random weights  
 $D \leftarrow \emptyset$   
 Repeat (for each episode):  
     Initialize  $s$   
     Repeat (for each step of the episode):  
         Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)  
         Take action  $a$ , observe  $r, s'$   
          $D \leftarrow D \cup (s, a, s', r)$   
          $s \leftarrow s'$   
         If  $\text{mod}(\text{step}, \text{trainfreq}) == 0$ :  
             sample batch  $B$  from  $D$   
              $w \leftarrow w - \alpha \nabla_w L(B; w, w^-)$   
             if  $\text{mod}(\text{step}, \text{copyfreq}) == 0$ :  
                  $w^- \leftarrow w$

Where:  $\nabla_w L(B; w, w^-) \approx -\frac{1}{|B|} \sum_{(s, a, s', r) \in B} (\text{target}(s'; w^-) - Q_w(s, a)) \nabla_w Q_w(s, a)$

$$\text{target}(s'; w^-) = r + \gamma \max_{a'} Q_{w^-}(s', a')$$

|                                  |     |              |
|----------------------------------|-----|--------------|
| $Q_1$ : Large Buffer             | vs. | Small Buffer |
| X Old Data                       |     | ✓ newer Data |
| ✓ less likely correlated samples |     | X Correlated |

$Q_2$ : Large Batch Size ( $B$ )

|           |             |
|-----------|-------------|
| Slower    | Faster      |
| Stability | Instability |

**Extended Data Table 1 | List of hyperparameters and their values**

| Hyperparameter                  | Value   | Description  |
|---------------------------------|---------|--|
| minibatch size                  | 32      | Number of training cases over which each stochastic gradient descent (SGD) update is computed.   |
| replay memory size              | 1000000 | SGD updates are sampled from this number of most recent frames.  |
| agent history length            | 4       | The number of most recent frames experienced by the agent that are given as input to the Q network.  |
| target network update frequency | 10000   | The frequency (measured in the number of parameter updates) with which the target network is updated (this corresponds to the parameter $C$ from Algorithm 1).                   |
| discount factor                 | 0.99    | Discount factor gamma used in the Q-learning update.   |
| action repeat                   | 4       | Repeat each action selected by the agent this many times. Using a value of 4 results in the agent seeing only every 4th input frame.   |
| update frequency                | 4       | The number of actions selected by the agent between successive SGD updates. Using a value of 4 results in the agent selecting 4 actions between each pair of successive updates. |
| learning rate                   | 0.00025 | The learning rate used by RMSProp.   |
| gradient momentum               | 0.95    | Gradient momentum used by RMSProp.   |
| squared gradient momentum       | 0.95    | Squared gradient (denominator) momentum used by RMSProp.   |
| min squared gradient            | 0.01    | Constant added to the squared gradient in the denominator of the RMSProp update.   |
| initial exploration             | 1       | Initial value of $\epsilon$ in $\epsilon$ -greedy exploration.   |
| final exploration               | 0.1     | Final value of $\epsilon$ in $\epsilon$ -greedy exploration.   |
| final exploration frame         | 1000000 | The number of frames over which the initial value of $\epsilon$ is linearly annealed to its final value.   |
| replay start size               | 50000   | A uniform random policy is run for this number of frames before learning starts and the resulting experience is used to populate the replay memory.                              |
| no-op max                       | 30      | Maximum number of "do nothing" actions to be performed by the agent at the start of an episode.  |

The values of all the hyperparameters were selected by performing an informal search on the games Pong, Breakout, Seaquest, Space Invaders and Beam Rider. We did not perform a systematic grid search owing to the high computational cost, although it is conceivable that even better results could be obtained by systematically tuning the hyperparameter values.

## Double DQN

### Double Q-Learning → Tabular

Initialize  $Q_1(s, a)$  and  $Q_2(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily

Initialize  $Q_1(\text{terminal-state}, \cdot) = Q_2(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

    Initialize  $S$

    Repeat (for each step of episode):

        Choose  $A$  from  $S$  using policy derived from  $Q_1$  and  $Q_2$  (e.g.,  $\varepsilon$ -greedy in  $Q_1 + Q_2$ )

        Take action  $A$ , observe  $R, S'$

        With 0.5 probability:

$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha \left( R + \gamma Q_2(S', \arg \max_a Q_1(S', a)) - Q_1(S, A) \right)$$

        else:

$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha \left( R + \gamma Q_1(S', \arg \max_a Q_2(S', a)) - Q_2(S, A) \right)$$

$$S \leftarrow S'$$

    until  $S$  is terminal

## Double DQN

$$\rightarrow \text{DQN} \quad w \leftarrow w - \alpha \nabla_w L(B; w, \bar{w})$$

$$\nabla_w L(B; w, \bar{w}) = -\frac{1}{B} \sum_{s, a, r, s' \in B} (\text{target}(s', \bar{w}) - Q_w(s, a)) \nabla_w Q_w(s, a)$$

$$\text{target}(s', \bar{w}) = r + \gamma \max_{a'} Q_{\bar{w}}(s', a')$$

$\downarrow \arg \max_{a'} Q_{\bar{w}}(s', a')$

Double DQN:

$$\nabla_w L(B; w, \bar{w}) = -\frac{1}{B} \sum_{s, a, r, s' \in B} (\text{target}(s', a'; \bar{w}, w) - Q_w(s, a)) \nabla_w Q_w(s, a)$$

$$\text{target}(s', a'; \bar{w}, w) = r + \gamma \underbrace{Q}_{\bar{w}}(s', \arg \max_{a'} Q_{\bar{w}}(s', a'))$$

# Double DQN

Initialize  $Q_w, Q_{w^-}$  with random weights  
 $D \leftarrow \emptyset$   
Repeat (for each episode):  
    Initialize  $s$   
    Repeat (for each step of the episode):  
        Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)  
        Take action  $a$ , observe  $r, s'$   
         $D \leftarrow D \cup (s, a, s', r)$   
         $s \leftarrow s'$   
        If  $\text{mod}(step, trainfreq) == 0$ :  
            sample batch  $B$  from  $D$   
             $w \leftarrow w - \alpha \nabla_w L(B; w, w^-)$   
            if  $\text{mod}(step, copyfreq) == 0$ :  
                 $w^- \leftarrow w$

Where:  $\nabla_w L(B; w, w^-) \approx -\frac{1}{|B|} \sum_{(s, a, s', r) \in B} (target(s', a'; w, w^-) - Q_w(s, a)) \nabla_w Q_w(s, a)$

$$target(s', a'; w, w^-) = r + \gamma Q_{w^-}(s', \arg \max_{a'} Q_w(s', a'))$$

## Prioritized DQN

### D Buffer

$s_0, a_0, r_1, s_1$

$s_1, a_1, r_2, s_2$

:

:

:

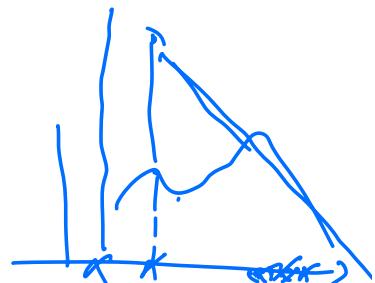
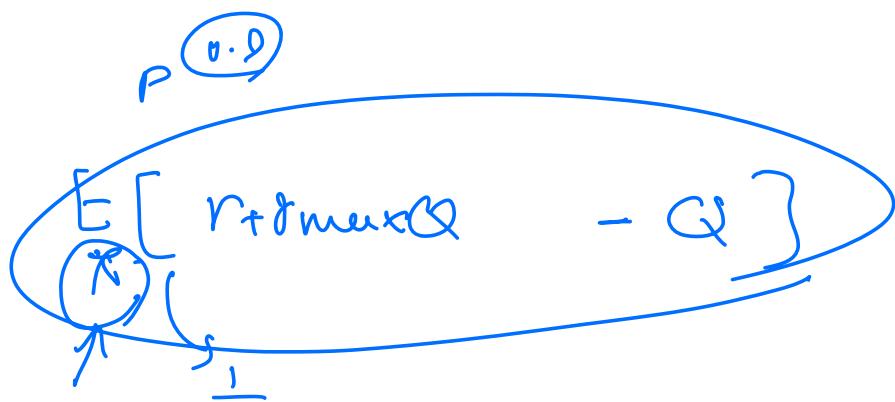
$s_D, a_D, r_{D+1}, s_{D+1}$

$$\rightarrow \delta_1 = |r_1 + \gamma \max_{a'} Q_W(s_1, a') - Q_W(s_0, a_0)|$$

$$\rightarrow \delta_2 = |r_2 + \gamma \max_{a' \in A} Q_W(s_2, a') - Q_W(s_1, a_1)|$$

$$\delta_D = |r_{D+1} + \gamma \max_{a' \in A} Q_W(s_{D+1}, a') - Q_W(s_D, a_D)|$$

$$W \leftarrow W + \frac{1}{B} \sum_B (r_t + \gamma \max_{a'} Q_W(s', a') - Q_W(s, a)) \nabla Q_W(s, a)$$



## D Buffer

$s_0, a_0, r_1, s_1$

$s_1, a_1, r_2, s_2$

$\vdots$

$s_D, a_D, r_{D+1}, s_{D+1}$

$$\rightarrow \delta_1 = |r_1 + \gamma \max_{a' \in A} Q_W(s_1, a') - Q_W(s_0, a_0)|$$

$$\rightarrow \delta_2 = |r_2 + \gamma \max_{a' \in A} Q_W(s_2, a') - Q_W(s_1, a_1)|$$

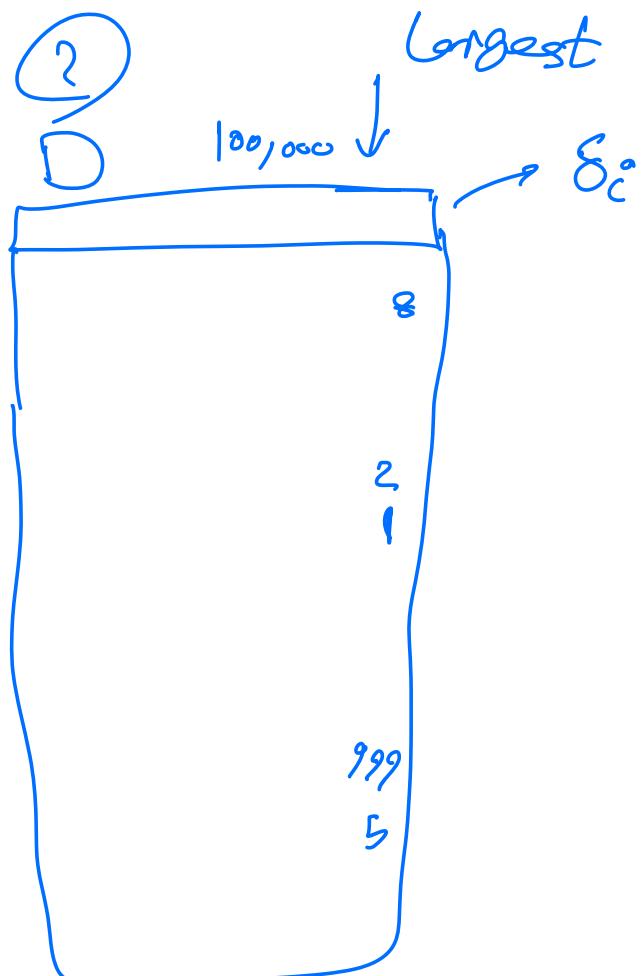
$$\delta_D = |r_{D+1} + \gamma \max_{a' \in A} Q_W(s_{D+1}, a') - Q_W(s_D, a_D)|$$

$$P \propto [ \delta_1, \dots, \delta_{D+1} ]$$

$$P(i) = \frac{\delta_i + c}{\sum_{j=1}^{D+1} \delta_j + c}$$

$0.9$   
 $0.01$   
 $0.01$        $\Rightarrow$   
 $\downarrow$   
 $\{$

## Ranked-Based



$$P(i) \propto \frac{1}{\text{Rank}(i)}$$

③

$P_c^i = \frac{P(i)^{\alpha}}{\sum_{j=1}^D P(j)^{\alpha}}$

Rank

$P(i)$

$\sum_{j=1}^D P(j)^{\alpha}$

$\delta_i^c$

$\alpha$ : level of prioritization

$\alpha = 0 \Rightarrow \text{EQN}$

$\alpha = 1$  : Full Probabil.

$$\min_w \mathbb{E}_{s,a,s',r \in D} \left[ (r + \gamma \max_{a'} Q_w^-(s', a') - Q_w(s, a))^2 \right]$$

$$\nabla L(B; w, \bar{w}) = -\frac{1}{|B|} \sum_{s,a,r,s' \in B} \frac{1}{P_{s,a,r,s'}} (r + \gamma \max_{\bar{w}} \hat{Q}_{\bar{w}}(s', a) - Q_w(s, a))$$

$\nabla_w Q_w$

