

Lecture 24 - April 14, 2023

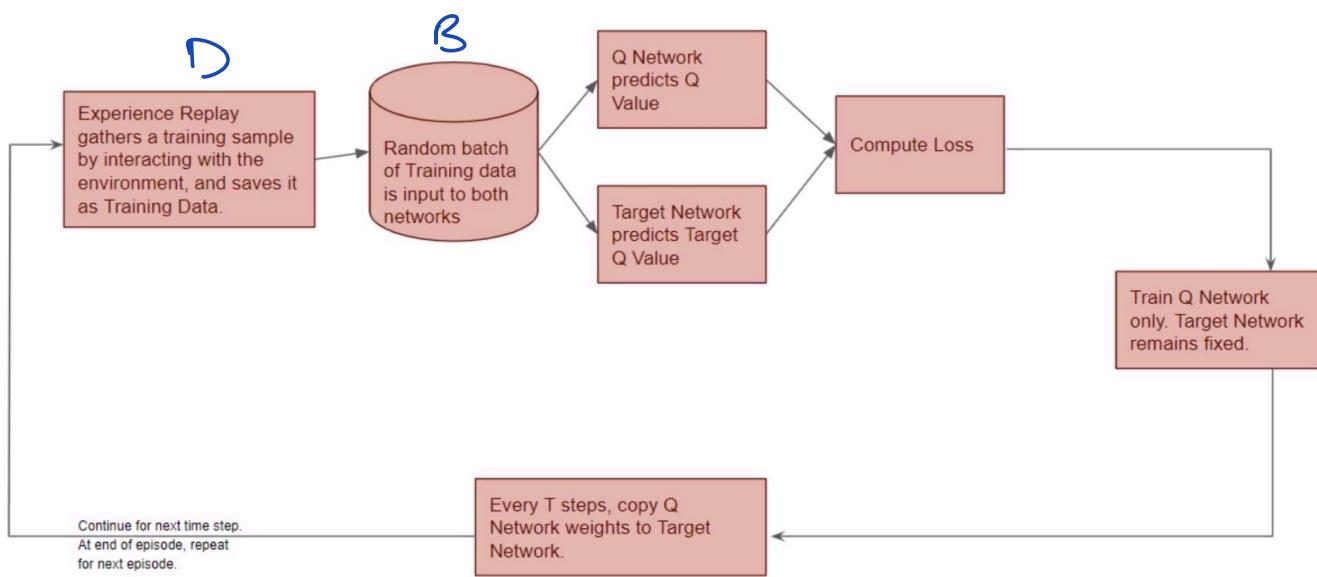
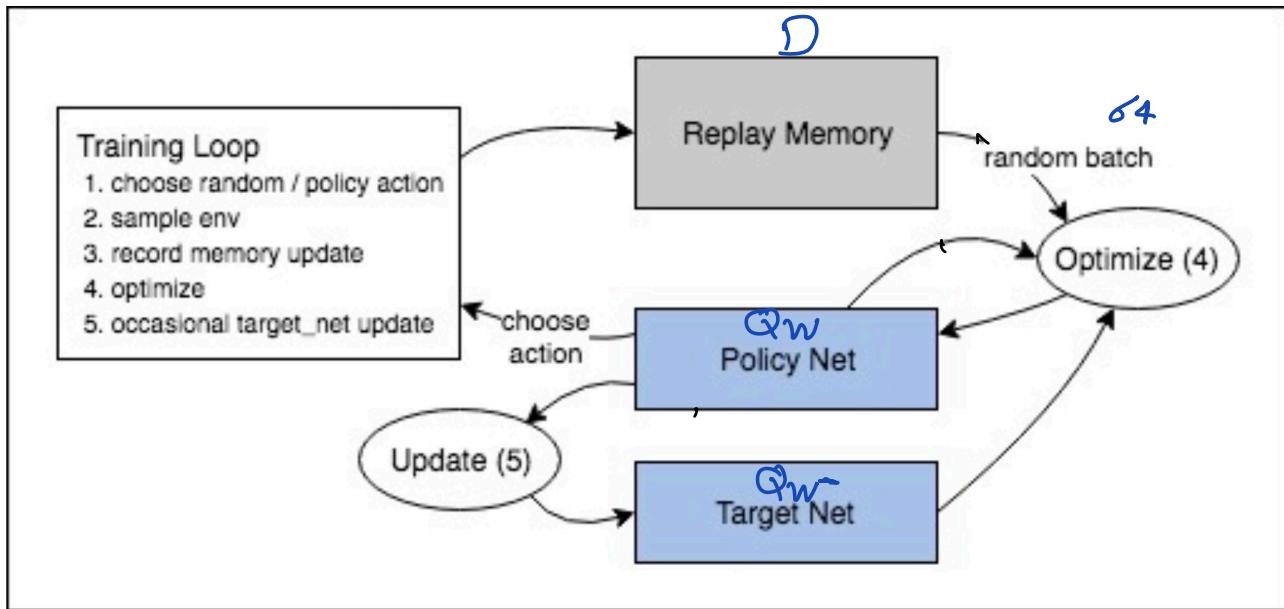
- Deep Q-Network (DQN) → Finite action
 - Deep Q-Network (DQN)
 - Double DQN
 - Prioritized DQN
 - Dueling DQN
 - Noisy-Net DQN
- Deep Policy Gradients (DPG) → Large/continuous Action
 - REINFORCE
 - REINFORCE with Baseline
 - Advantage Actor Critic (A2C)
 - Deep Deterministic Policy Gradient (DDPG)

Project 3 → Due April 17

HWS → Due April 18

TA's office hour:

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



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., ϵ -greedy)

 Take action a , observe r, s'

$$D \leftarrow D \cup (s, a, s', r)$$

$$s \leftarrow s'$$

 If $\text{mod}(\text{step}, \text{trainfreq}) == 0$: 64

 sample batch B from D

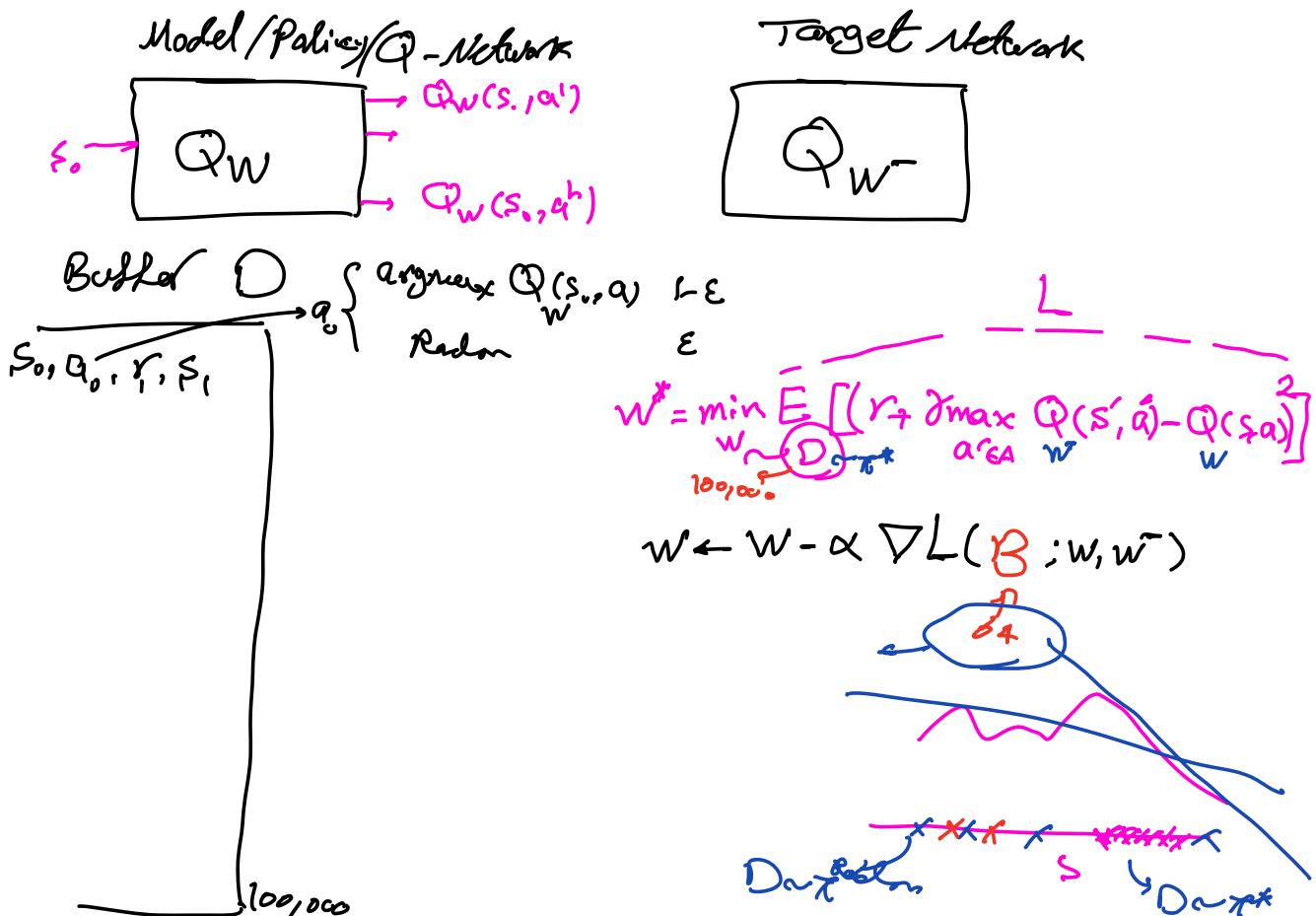
$$w \leftarrow w - \alpha \nabla_w L(B; w, w^-) \rightarrow Q_w$$

 if $\text{mod}(\text{step}, \text{copyfreq}) == 0$:

$$w^- \leftarrow w \quad \text{spanning circle}$$

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

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



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., ϵ -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

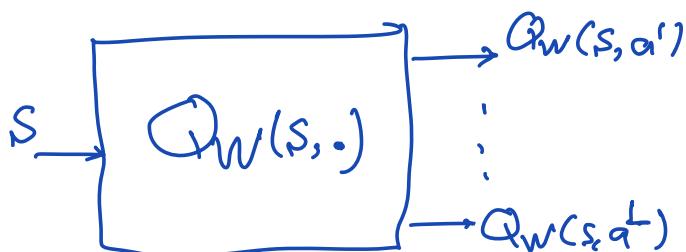
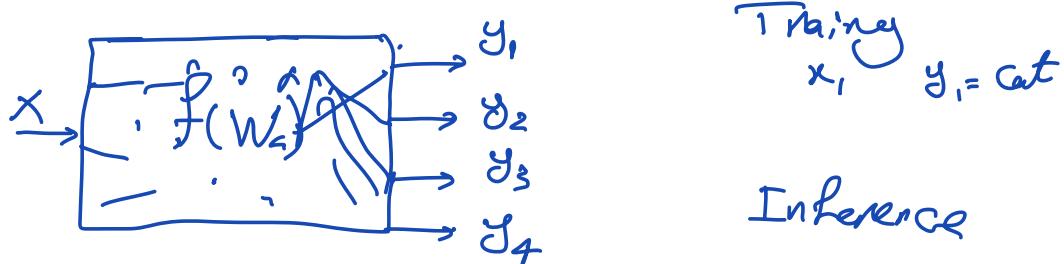
$$w \leftarrow w - \alpha \nabla_w L(B; w, w^-)$$

 if $\text{mod}(\text{step}, \text{copyfreq}) == 0$:

$$w^- \leftarrow w$$

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

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



$D \sim \mathcal{B}$ $s_0, a_0, r_1, s_1, \dots$

$$W \leftarrow W + \alpha \frac{1}{|B|} \sum_{S, a, r, S' \in B} (r + \gamma \max_{a'} \tilde{Q}(S', a') - \tilde{Q}(S, a)) \nabla_{W,W} Q(S, a)$$

$s_0, q_0, r_0, \varepsilon_0$

$$\rightarrow S_i = \left| r_i + \max_{a'} Q_W(s_i, a') - Q_W(s_i, a_i) \right|$$

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

B
→

$$\delta_b = |r_b + \max_{a'} Q_w(s_1, a') - Q_w(s_0, a_0)|$$

Prioritized DQN:

$$P(i) = \frac{\delta_i + c}{\sum_j \delta_j + c}$$

Ranked \rightarrow

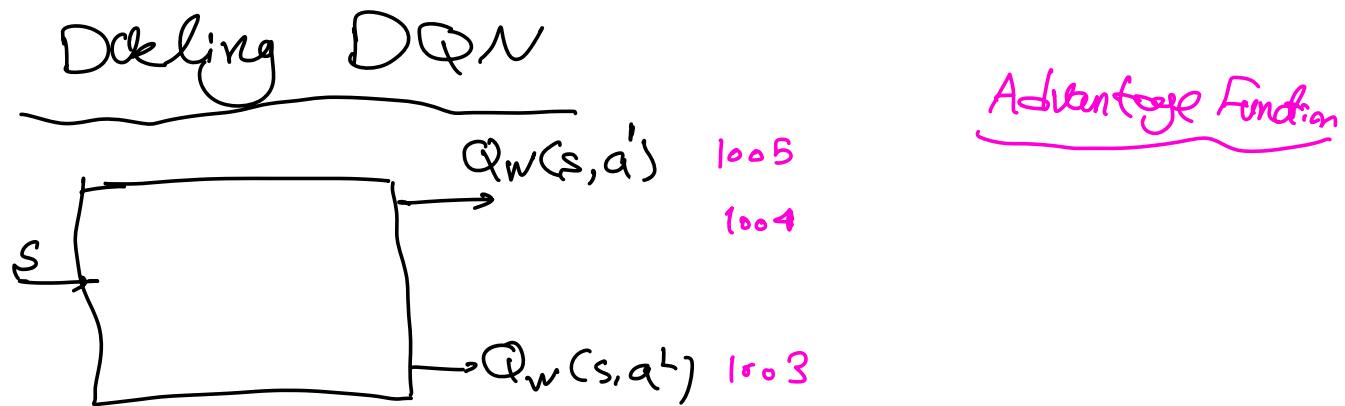
$$P(c) \propto \frac{1}{\text{Rank}(c)}$$

$$P_i = \frac{P(i)}{\sum_{j=1}^{|D|} P(j)} \rightarrow \text{level of priority}$$

$$\text{Random } \mathcal{B} \rightarrow W \leftarrow W + \alpha \frac{1}{|\mathcal{B}|} \sum_{s, a, r, s' \in \mathcal{B}} (r + \gamma \max_{a'} Q(s', a') - Q(s, a)) \nabla_w Q(s, a)$$

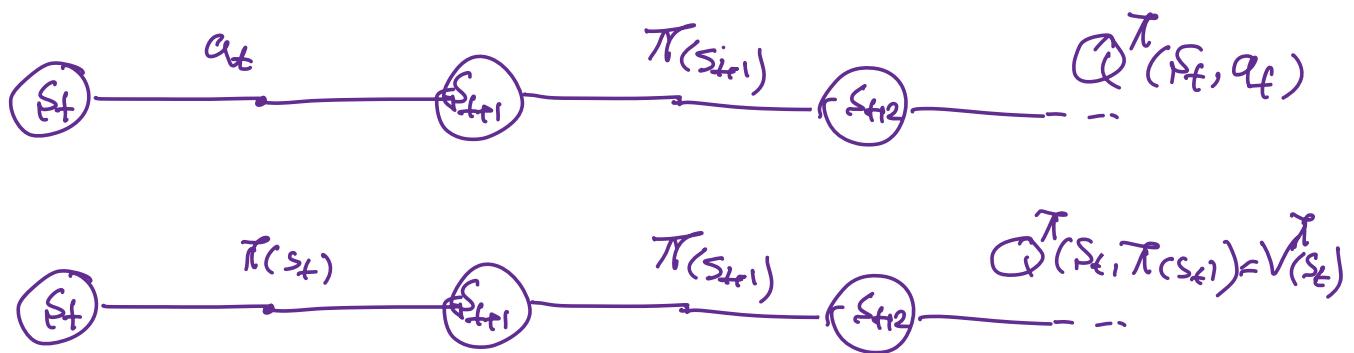
$$\text{Priority} \Rightarrow W \leftarrow W + \alpha \frac{1}{|\mathcal{B}|} \sum_{\substack{s, a, r, s' \in \mathcal{B} \\ P_{s \rightarrow s', r}}} (r + \gamma \max_{a'} Q(s', a') - Q(s, a)) \nabla_w Q(s, a)$$

- Periodically Updating \mathcal{S}_c
- only for minibatch selected
- stochastic update $\rightarrow 1000$ according to \mathcal{S}_c



$$Q(s, a) = V(s) + A(s, a)$$

\leftarrow advantage function



$$V^\pi(15) = 99$$

$$Q^\pi(15, U) = -(-1 + 99) = 98$$

$$Q^\pi(15, R) = -1 + 99 = 98$$

$$Q^\pi(15, D) = -1 + 99 = 98$$

9 ↓	10 ←	11 ←	12 ↓
8 ↓		14 ↓	13 ↓
7 ←		16	15 ←
6 ↑	5 ←		
4 ↑	3 ↑	2 ←	1 ←

■ Wall ■ Bump ■ Goal

$$\left\{ \begin{array}{l} A^\pi(15, U) = \frac{Q^\pi(15, U)}{98} - \frac{V^\pi(15)}{99} = -2 \\ A^\pi(15, R) = -1 \\ A^\pi(15, D) = -1 \\ A^\pi(15, L) = Q^\pi(15, L) - V^\pi(15) = 0 \end{array} \right.$$

$$\underset{a \in A}{\operatorname{Argmax}} Q^\pi(s, a) = \underset{a \in A}{\operatorname{argmax}} A^\pi(s, a) + V^\pi(s)$$

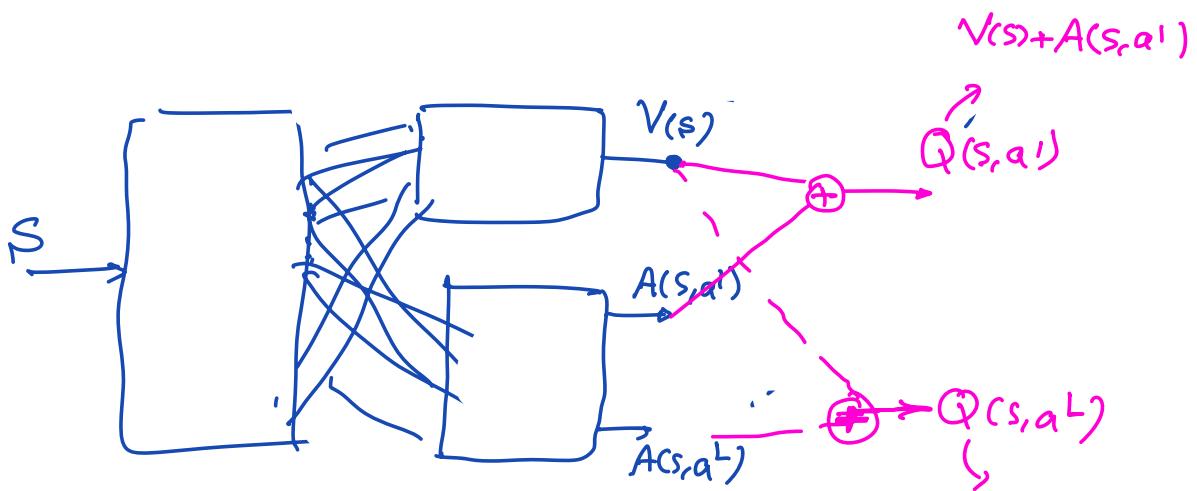
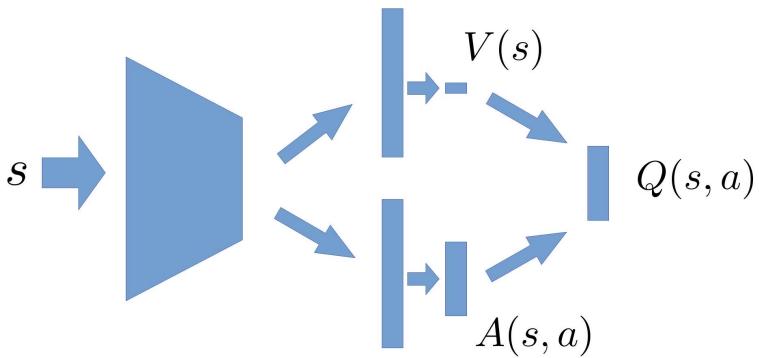
$$V^\pi(11) = -100$$

$$Q^\pi(11, R) = 98$$

$$Q^\pi(11, L) = -100$$

$$Q^\pi(11, U) = 88$$

$$\left\{ \begin{array}{l} A^\pi(11, R) = \frac{Q^\pi(11, R)}{98} - \frac{V^\pi(11)}{-100} \\ \Rightarrow A^\pi_{11,R} = 198 \end{array} \right.$$



$$\left. \begin{array}{l} V(s) = 6 \\ A(s, a^1) = 0 \\ A(s, a^2) = 1 \end{array} \right\} \rightarrow \begin{array}{l} V(s) + A(s, a^1) \\ Q(s, a^1) = 6 \\ Q(s, a^2) = 7 \end{array}$$

possible solutions

$$V(s) = 1$$

$$A(s, a^1) = 5$$

$$A(s, a^2) = 6$$

$$V(s) = 5$$

$$A(s, a^1) = 1$$

$$A(s, a^2) = 2$$

$$\max_{a \in A} Q(s, a) = V(s)$$

$\xrightarrow{a^{\max}}$

$$A(s, a^{\max}) = V(s) - Q(s, a^{\max}) = 0$$

$$Q_W(s, a) = V_W(s) + A_W(s, a) - \max_{a' \in A} A_W(s, a')$$

$$V(s) = 6$$

$$A(s, a^1) = 0 \xrightarrow[\text{NN}]{\text{Prediction}}$$

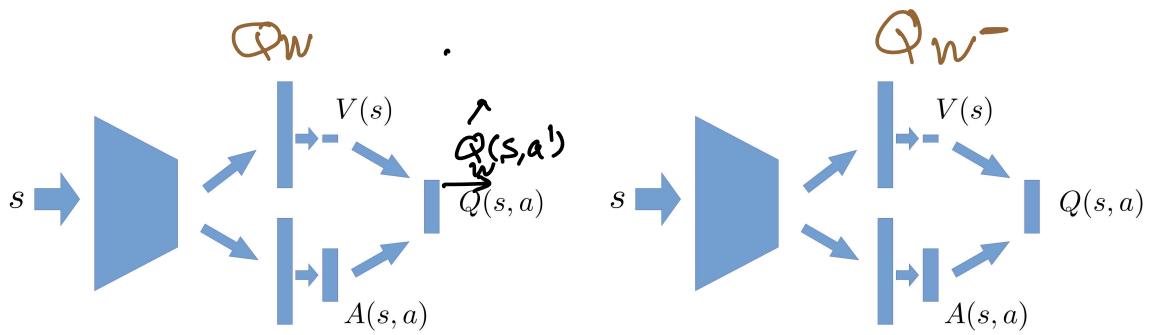
$$Q(s, a^1) = 6 + 0 = 6$$

$$Q(s, a^2) = 6 + 1 = 7$$

$$\max_{a \in A} Q(s, a) = V(s)$$

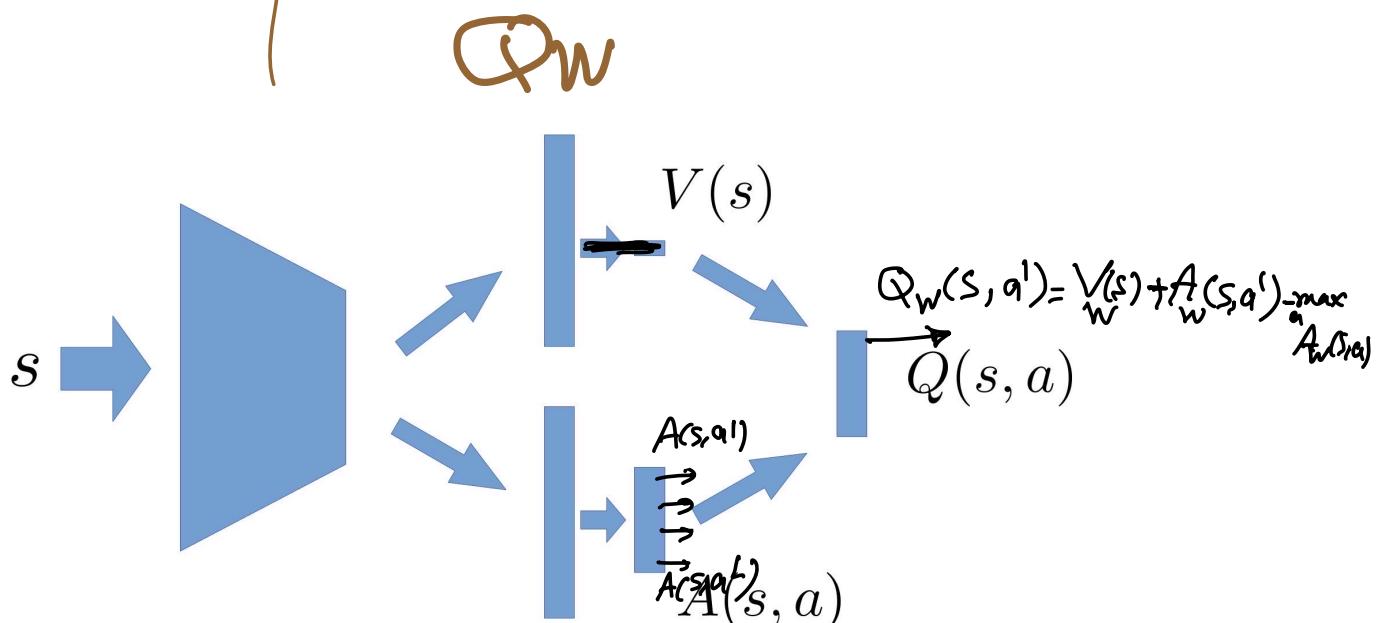
$$Q_W(s, a^1) = \underbrace{V_W(s)}_6 + \underbrace{A_W(s, a^1)}_0 - \underbrace{\max_{a' \in A} A_W(s, a')}_1 = 5$$

$$Q_W(s, a^2) = \underbrace{V_W(s)}_6 + \underbrace{A_W(s, a^2)}_1 - \underbrace{\max_{a' \in A} A_W(s, a')}_1 = 6$$



$$L(B; w, \bar{w}) = \frac{1}{B} \sum_{s, a, r, s' \in B} (y_i - Q_w(s, a))^2$$

$\overbrace{\hspace{10em}}$



100 5

1004

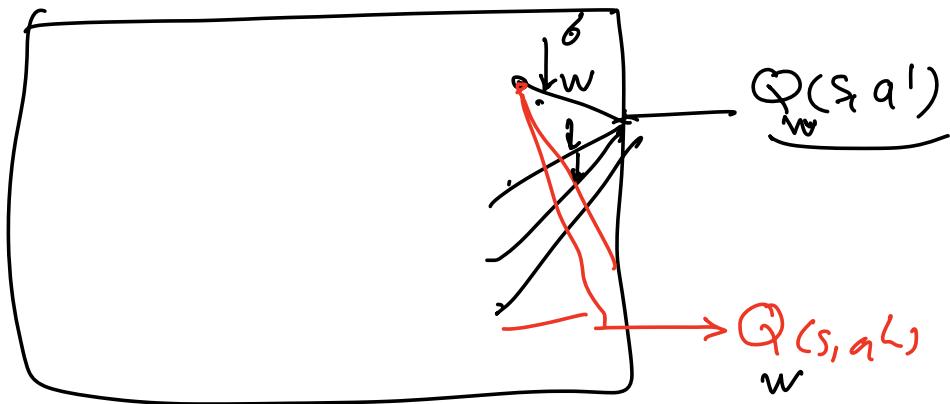
1002

$$Q_W(s, a) = V_W(s) + A_W(s, a) - \frac{1}{|A|} \sum_{a' \in A} A_W(s, a')$$

$$\begin{aligned} V(s) &= 6 \\ A(s, a^1) &= 0 \\ A(s, a^2) &= 1 \end{aligned} \quad \begin{aligned} Q_W(s, a^1) &= \underbrace{V_W(s)}_6 + \underbrace{A(s, a^1)}_0 - \frac{1}{2} \left(\underbrace{A(s, a^1)}_0 + \underbrace{A(s, a^2)}_1 \right) \\ &= 5.5 \end{aligned}$$

$$\begin{aligned} Q_W(s, a^2) &= \underbrace{V_W(s)}_6 + \underbrace{A(s, a^2)}_1 - \frac{1}{2} \left(\underbrace{A(s, a^1)}_0 + \underbrace{A(s, a^2)}_1 \right) \\ &= 6.5 \end{aligned}$$

Noisy-Net DQN



$$y = \mu_w + \sigma_w \odot \epsilon$$

\uparrow

$$\mathcal{N}(0, 1)$$

Instead of argmax \Rightarrow Random layer at the end of NN

$$s_0, a \xrightarrow{\arg\max Q_w(s, a)}$$