

Lecture 22 - April 7, 2023

- Function Approximation in Reinforcement Learning

- Basics of Function Approximations



- Least square Policy Iteration (LSPI)

- Deep Q-Network (DQN) → Finite action

- Deep Q-Network (DQN)

- Double DQN

- Prioritized DQN

- Dueling DQN

- Deep Policy Gradients (DPG) → Large/continuous Action

Project 3 → Due April 14

HWS → Due April 18

TA's office hour:

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

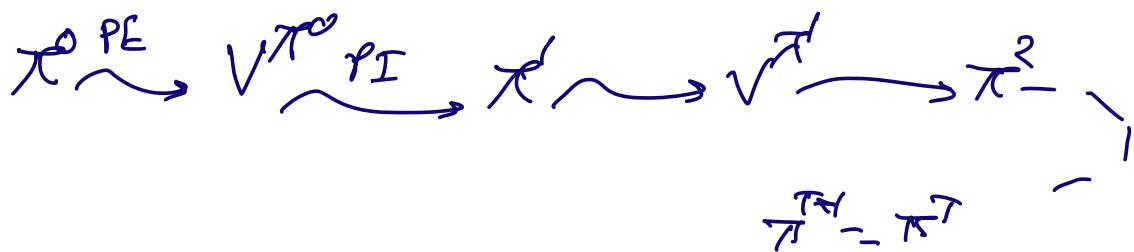
Fridays, 2pm-3pm (virtual)

Batch Learning

$$D = \{(s_0, a_0, r_0, s_1), \dots\}$$

\Downarrow
 π^*

LSPIT



Policy Evaluation:

$$V_{\pi}(s) = \sum_{s'} P(s'|s, \pi(s)) [R + \gamma V_{\pi}(s')]$$

Policy Improve-

$$\pi'(s) = \underset{a \in A}{\operatorname{argmax}} \frac{\sum_{s'} P(s'|s, a) [R + \gamma V_{\pi}(s')]}{Q_{\pi}(s, a)}$$

$$Q_{\pi}(s, a) = \sum_{s'} P(s'|s, a) [R + \gamma V_{\pi}(s')]$$

$$= \sum_{s'} P(s'|s, a) [R + \gamma Q_{\pi}(s', \pi(s'))]$$

$$Q^{\pi} = R + \gamma M Q^{\pi}$$

$\downarrow |S| \times |A|$ $|S| \times |A|$

$$\begin{bmatrix} Q^T(s_1, a_1) & \dots & Q^T(s_n, a^{(A)}) \\ Q^T(s'^1, a') & \dots & Q^T(s'^n, a^{(A)}) \end{bmatrix}$$

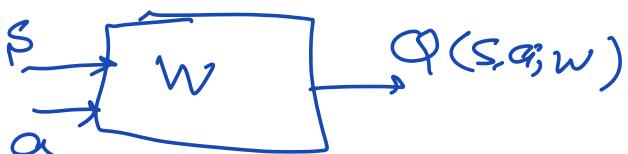
LSPI

$$Q(s, a; w) = \sum_{j=1}^n \phi_j(s, a) w_j = \underbrace{\Phi^T(s, a)}_{\text{Basis function}} w$$

$$\Phi_{(s, a)} = \begin{bmatrix} 1 \\ s \\ a \\ s \otimes a \\ s^2 \otimes a \\ \exp(sa) \end{bmatrix}_{N \times 1}$$

$$Q(s, a; w) = \Phi_{(s, a)}^T w$$

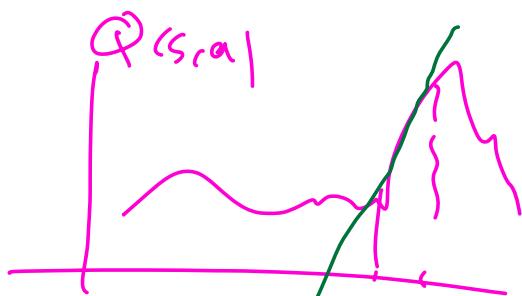
$$Q(s, a; w) = [1 \ s \ a \ s a \ s^2 a \ \exp(sa)] \begin{bmatrix} w_1 \\ \vdots \\ w_6 \end{bmatrix} = w_1 + w_2 s + w_3 a + w_4 s a + w_5 s^2 a + w_6 \exp(sa)$$



$$\begin{aligned}
 Q^{\pi} &= R + \gamma M Q^{\pi} \\
 \Phi^T \downarrow \quad \Phi^T & \quad \Phi^T \\
 \Phi^T w^{\pi} &= R + \gamma M \Phi^T w^{\pi} \\
 \left[\begin{array}{c} \Phi^T(s', a') \dots \Phi^T(s', a'^{|A|}) \\ \vdots \end{array} \right] & \\
 \Phi^T (\Phi - \gamma M \Phi) w^{\pi} &= \Phi^T R
 \end{aligned}$$

$$w^{\pi} = \underbrace{\left[\Phi^T \mu (\Phi - \gamma M \Phi) \right]^A}_{A} \underbrace{\Phi^T \mu R}_{b}$$

State distribution under Policy π



$$w^{\pi} = A^{-1} b \rightarrow \begin{cases} A = \Phi^T \mu (\Phi - \gamma M \Phi) \\ b = \Phi^T \mu R \end{cases}$$

$D = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^L \leftarrow$ Batch Data available

$$\begin{aligned}\hat{A} &= \frac{1}{L} \sum_{i=1}^L \phi(s_i, a_i) [\phi(s_i, a_i) - \gamma \phi(s'_i, \pi(s'_i))] \\ \hat{b} &= \frac{1}{L} \sum_{i=1}^L \phi(s_i, a_i) r_{i+1}\end{aligned}$$

$$\omega^\pi = \hat{A}^{-1} \hat{b}$$

$$\begin{array}{c} D \\ \omega_0 \xrightarrow{\pi^0} \underbrace{A^{\pi^1} b}_{\omega_1} \xrightarrow{\pi^1} A^{\pi^2} b \xrightarrow{\pi^2} \omega_2 \end{array}$$

$$\pi^*(s) = \arg \max \mathbb{Q}^{\pi^0}(s, a)$$

Random ω^* Fixed $\hat{\Phi}(s, a)$
Policy Evaluation

$$\omega^* = \bar{\omega} \rightarrow \hat{Q}_{(s, a)}^T = \hat{\Phi}_{(s, a)}^T \bar{\omega}$$

$$\pi_s^* = \underset{a \in A}{\operatorname{argmax}} \hat{Q}_{(s, a)}^T = \underset{a \in A}{\operatorname{argmax}} \hat{\Phi}_{(s, a)}^T \bar{\omega}$$

Policy Improvement

$$A = \frac{1}{L} \sum_{c=1}^L \hat{\Phi}_{(s_c, a_c)} \left[\hat{\Phi}_{(s_c, a_c)} - \gamma \hat{\Phi}_{(s_{c+1}, \pi^*(s_{c+1}))}^T \right]$$

$$b = \frac{1}{L} \sum_{c=1}^L \hat{\Phi}_{(s_c, a_c)} r_{c+1}$$

$$\omega^* = A^{-1} b$$

Example of LSPI:

$$S \in [-2, 2]$$

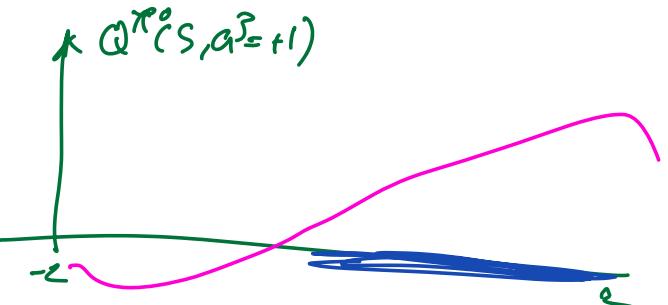
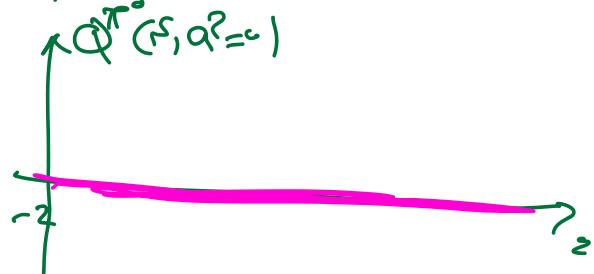
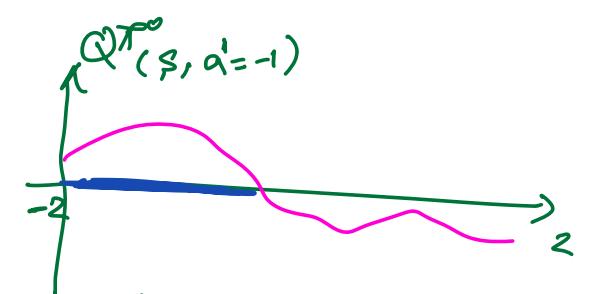
$$\Phi(s, a) = \begin{bmatrix} sa \\ |s+1|/a^2 \end{bmatrix}$$

$$A = \{-1, 0, 1\}$$

$$D = \{(s_0=0, a_0=-1, r_1=1, s_1=1), (s_1=1, a_1=1, r_1=0.9, s_2=-1)\}$$

$$\omega_s^{\text{Random}} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow Q^{\pi^*}(s, a) = \Phi^T(s, a) \omega_s = [sa \ |s+1|/a^2] \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$= sa + |s+1|/a^2$$



$$\Phi(s, a) = \begin{bmatrix} sa \\ |s+1|/a^2 \end{bmatrix}$$

$$D = \{(s_0=0, a_0=-1, r_1=1, s_1=1), (s_1=1, a_1=1, r_1=0, s_2=-1)\}$$

$$\hat{A} = \frac{1}{L} \sum_{c=0}^L \phi^T(s_c, a_c) (\phi(s_c, a_c) - \gamma \phi(s_{c+1}, \pi(s_{c+1})))^T$$

$$= \frac{1}{2} \left[\underbrace{\phi^T(s_0, a_0)}_{\begin{matrix} s_0 \\ a_0 \\ 1 \end{matrix}} \left(\phi(s_0, a_0) - \gamma \phi(s_1, \pi(s_1)) \right)^T \right]$$

$$\phi(s, a) = \begin{bmatrix} s \\ a \end{bmatrix} = \begin{bmatrix} x-1 \\ (-1)^2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \phi(s_1, \pi(s_1)) = \begin{bmatrix} s_1 \\ \pi(s_1) \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1^2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$+ \phi^T(s_1, a_1) \left(\phi(s_1, a_1) - \gamma \phi(s_2, \pi(s_2)) \right)^T$$

$$\phi(s_1, a_1) = \begin{bmatrix} s_1 \\ a_1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1^2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \phi(s_2, \pi(s_2)) = \begin{bmatrix} s_2 \\ \pi(s_2) \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1^2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$= \begin{bmatrix} -0.05 & 1 \\ -0.35 & 1.6 \end{bmatrix} = \begin{bmatrix} -1x-1 \\ 0x(-1)^2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\pi(s) = \arg \max Q^\pi(s, a) = \arg \max_{a \in A} \phi^T(s, a) w_0$$

I only need these for s'_c

$$\pi'(s_1) = \pi'(1) = \arg \max_{a \in \{-1, 0, 1\}} \phi^T(s_1, a) \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$= \underset{a \in \{-1, 0, 1\}}{\operatorname{argmax}} \left[S_1 a - |S_1 + 1| a^2 \right] \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$= \underset{a \in \{-1, 0, 1\}}{\operatorname{argmax}} \frac{S_1 a - |S_1 + 1| a^2}{a + 2a^2}$$

$$= \underset{a \in \{-1, 0, 1\}}{\operatorname{argmax}} \left\{ \begin{array}{l} \overbrace{-1 + 2(-1)^2}^{a = -1} \quad \overbrace{0 + 2 \cdot 0^2}^{a = 0} \quad \overbrace{1 + 2(1)^2}^{a = 1} \end{array} \right\} = 1$$

$$\pi(S_2) = \underset{a \in \{-1, 0, 1\}}{\operatorname{argmax}} \phi^T(S_2, a) \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$= \underset{a \in \{-1, 0, 1\}}{\operatorname{argmax}} S_2 a + |S_2 + 1| a^2 = \underset{a \in \{-1, 0, 1\}}{\operatorname{argmax}} -a + 0 \cdot a^2$$

$$= \underset{a \in \{-1, 0, 1\}}{\operatorname{argmax}} \left\{ \begin{array}{l} \overbrace{-1}^{a = -1} \quad \overbrace{0}^{a = 0} \quad \overbrace{-1}^{a = 1} \end{array} \right\} = -1$$

$$b = \frac{1}{L} \sum_{c=0}^L \phi(S_i, a_i) r_{c+1}$$

$$= \frac{1}{2} \left(\underbrace{\phi(S_0, a_0)}_{\begin{bmatrix} 1 \\ 1 \end{bmatrix}} \frac{r_1}{1} + \phi(S_1, a_1) \underbrace{\begin{bmatrix} 1 \\ 2 \end{bmatrix} r_2}_{0.9} \right) = \begin{bmatrix} 0.45 \\ 1.4 \end{bmatrix}$$

$$\omega^l = \hat{A}^{-1} \hat{b} = \begin{bmatrix} -0.05 & 1 \\ 0.35 & 1.6 \end{bmatrix}^{-1} \begin{bmatrix} 0.43 \\ 1.4 \end{bmatrix} = \begin{bmatrix} -1.58 \\ 0.53 \end{bmatrix}$$

$$Q(s, a) = \phi^T(s, a) \omega^l = \begin{bmatrix} s^0 & 1 \\ |s+1| & a^2 \end{bmatrix}^T \begin{bmatrix} -1.58 \\ 0.53 \end{bmatrix}$$

$$= -1.58 s^0 + 0.53 |s+1| a^2$$

$$D = \{(s_0=0, a_0=-1, r_1=1, s_1=1), (s_1=1, a_1=1, r_1=0.8, s_2=-1)\}$$

$$\hat{A} = \frac{1}{L} \sum_{i=0}^{L-1} \phi^T(s_i, a_i) (\phi(s_i, a_i) - \gamma \phi(s_{i+1}, \pi(s_{i+1})))^T$$

$$= \frac{1}{2} \left[\underbrace{\phi^T(s_0, a_0)}_{\begin{bmatrix} 1 \\ 1 \end{bmatrix}} (\phi(s_0, a_0) - \gamma \phi(s_1, \pi(s_1)))^T \right]$$

$$\phi(s_0, a_0) = \begin{bmatrix} s_0 & a_0 \\ |s_0+1| & a_0^2 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \phi(s_1, \pi(s_1)) = \begin{bmatrix} s_1 & \pi^2(s_1) \\ |s_1+1| & \pi^2(s_1) \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 2 & 3 \end{bmatrix}$$

$$+ \phi^T(s_1, a_1) (\phi(s_1, a_1) - \gamma \phi(s_2, \pi(s_2)))^T$$

$$\phi(s_1, a_1) = \begin{bmatrix} s_1 & a_1 \\ |s_1+1| & a_1^2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 2 & 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \phi(s_2, \pi(s_2)) = \begin{bmatrix} s_2 & \pi^2(s_2) \\ |s_2+1| & \pi^2(s_2) \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ 3 & 3 \end{bmatrix}$$

$$= \begin{bmatrix} -0.05 & 1 \\ -0.35 & 1.6 \end{bmatrix}, \quad (\text{P}_2 + 1) \left(\mathbb{E}(S_2) \right)^T = \begin{bmatrix} -1x+1 \\ 0x(+1)^2 \end{bmatrix} - \begin{bmatrix} -1 \\ 0 \end{bmatrix}$$

$$\pi^2(s) = \arg\max_a Q^{\pi}(s, a) = \arg\max_{a \in A} \phi(s, a)^T w,$$

$$= \arg\max_{a \in A} \{-1.58 s a + 0.53 / (s+1) a^2\}$$

$$\pi^2(s_1 = 1) = \arg\max_a \{-1.58 a + 2 \times 0.53 a^2\} = -1$$

$$\pi^2(s_2 = -1) = \arg\max_a \{+1.58 a\} = 1$$

$$A = \begin{bmatrix} 0.95 & 1 \\ 2.35 & 1.6 \end{bmatrix}$$

b = same

$$\hat{w} = \hat{A}^{-1} \hat{b} = \begin{bmatrix} 0.82 \\ -0.33 \end{bmatrix}$$

```

LSPI ( $D, k, \phi, \gamma, \epsilon, \pi_0$ ) // Learns a policy from samples

//  $D$  : Source of samples  $(s, a, r, s')$ 
//  $k$  : Number of basis functions
//  $\phi$  : Basis functions
//  $\gamma$  : Discount factor
//  $\epsilon$  : Stopping criterion
//  $\pi_0$  : Initial policy, given as  $w_0$  (default:  $w_0 = 0$ )

 $\pi' \leftarrow \pi_0$  //  $w' \leftarrow w_0$ 

repeat
     $\pi \leftarrow \pi'$  //  $w \leftarrow w'$ 
     $\pi' \leftarrow \text{LSTDQ} (D, k, \phi, \gamma, \pi)$  //  $w' \leftarrow \text{LSTDQ} (D, k, \phi, \gamma, w)$ 
until ( $\pi \approx \pi'$ ) // until ( $\|w - w'\| < \epsilon$ )

return  $\pi$  // return  $w$ 

```

LSTDQ (D, k, ϕ, γ, π) // Learns \hat{Q}^π from samples

```

//  $D$  : Source of samples  $(s, a, r, s')$ 
//  $k$  : Number of basis functions
//  $\phi$  : Basis functions
//  $\gamma$  : Discount factor
//  $\pi$  : Policy whose value function is sought

 $\tilde{\mathbf{A}} \leftarrow \mathbf{0}$  //  $(k \times k)$  matrix
 $\tilde{b} \leftarrow \mathbf{0}$  //  $(k \times 1)$  vector

for each  $(s, a, r, s') \in D$ 
     $\tilde{\mathbf{A}} \leftarrow \tilde{\mathbf{A}} + \phi(s, a) \left( \phi(s, a) - \gamma \phi(s', \pi(s')) \right)^\top$ 
     $\tilde{b} \leftarrow \tilde{b} + \phi(s, a)r$ 
     $\overbrace{\arg\max_{a \in A} \underbrace{\hat{Q}(s', a)}_{\phi(s', a) w}}$ 

 $\tilde{w}^\pi \leftarrow \tilde{\mathbf{A}}^{-1} \tilde{b}$ 

return  $\tilde{w}^\pi$ 

```

Figure 5: The LSTDQ algorithm.

Deep Reinforcement Learning

Tableau Q-Learning

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):

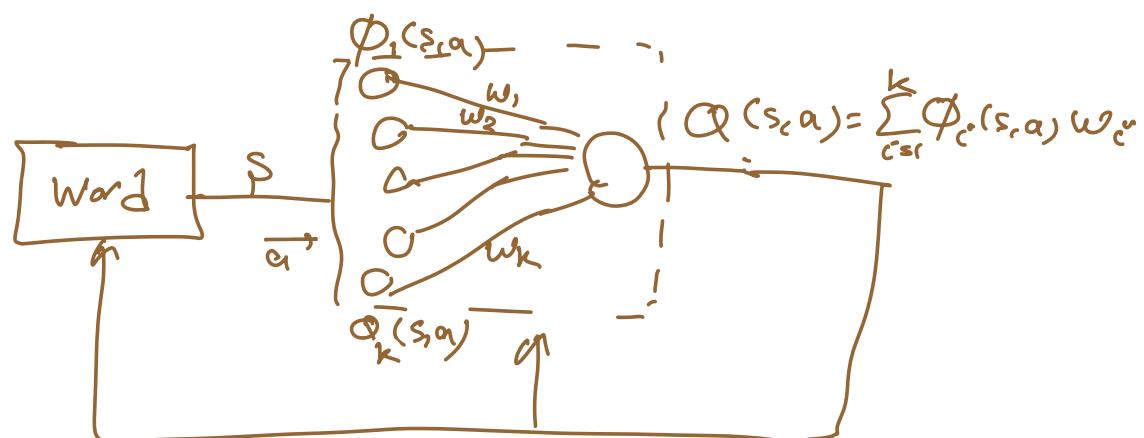
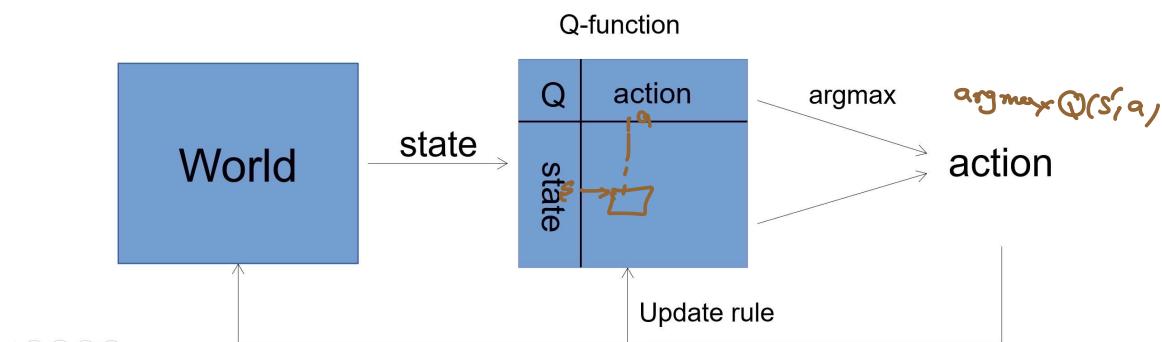
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)

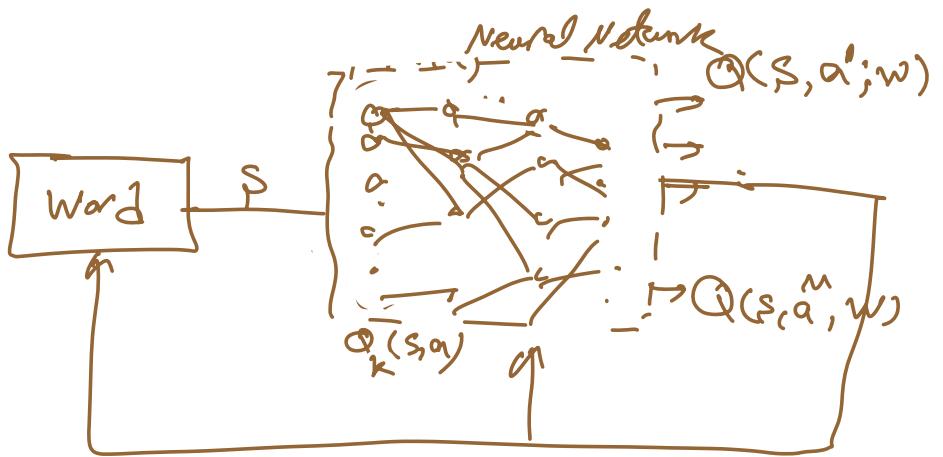
 Take action A , observe R, S'

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

$$S \leftarrow S'$$

 until S is terminal





Q-learning

$$Q(s, a) = Q(s, a) + \alpha \left[\underbrace{r + \gamma \max_{a'} Q(s', a') - Q(s, a)}_{\text{Q-learning Error}} \right]$$

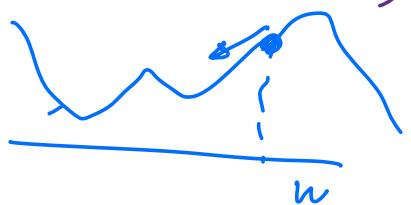
① Model $Q_w(s, a)$

② Loss function

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

, decrease

Stoch Gradient Descent



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

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

$$w \leftarrow \bar{w} + \alpha (r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a)) \nabla_w Q_w(s, a) \Big|_{w=\bar{w}}$$

Bare bone DQN (Does 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. e-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:

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

$$s, a, s', r \quad y = r + \gamma \max_{a'} Q_w(s', a')$$

}

}