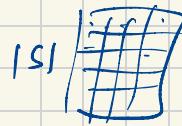


# Lecture - 6 : Value-based RL

Last time :

## ① Value-iteration:

$$\pi(s) = \underset{a}{\operatorname{argmax}} Q(s, a)$$

$$Q(s, a) \leftarrow r(s, a) + \gamma V(s')$$


$$V(s) \leftarrow \max_{a'} Q(s, a')$$

## ② Fitted Q-iteration: (neural fQI)

{ for  $i \geq 0, \dots$  :

$$Q_{i+1} \leftarrow \underset{\theta}{\operatorname{argmin}} \mathbb{E}_{s, a, s'} \left[ (Q_\theta(s, a) - y)^2 \right]$$

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

$Q_i$  used here

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

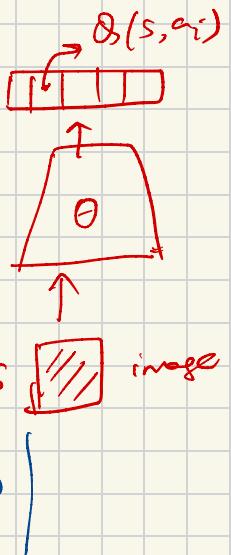
## ③ DQN:

- 1) replay buffer  $\rightarrow$
- 2) target networks



history of  
 $(s, a, s')$

## Outline of DQN algorithm:



### Exploration:

$$\pi_\theta(s)$$

$$\underset{a}{\operatorname{argmax}} \ Q_\theta(s, a)$$

"exp-greedy"  $\rightarrow$  prob.  $\epsilon$

$$a \sim \text{unif}(a)$$

### Update:



replay  
buffer

$$(s, a, r, s')$$

### Target networks

$$\textcircled{1} \quad \theta_0 = \theta_0^{\text{tgt}}$$

$$\textcircled{2} \quad \theta_0 \rightarrow \text{grad. descent}$$

$$\min_{\theta} E_{s, a, r, s' \sim \mathcal{D}} \left[ (Q_\theta - (r + \gamma^{\max_a} Q_\theta(s, a)))^2 \right]$$

$\theta_{i+1}$

main

$\theta_i$

target  
network

### update target net.

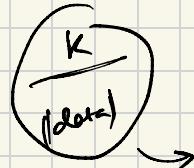
$$\theta_K^{\text{tgt}} = \theta_K$$

" $\theta_i$ "

(long)

exploring  
updates

(few)



can't be  
too large

### UTD rate:

updates - to - data

$\textcircled{1}$

Hard target update:

vs soft

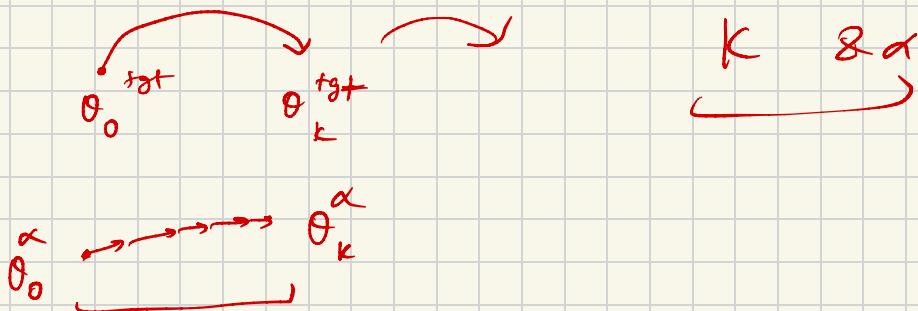
$$\alpha = 0.005$$

Soft target updates:  $\{ \theta_{i+1}^{\text{tgt}} \leftarrow (1-\alpha) \theta_i^{\text{tgt}} + \alpha \theta_i \}$   
"soft"

$$\theta' \leftarrow (1-\alpha) \theta + \alpha \theta_{\text{new}}$$

hard  $\{ \theta_k^{\text{tgt}} = \theta_k, \quad \theta_{2k}^{\text{tgt}} = \theta_{2k} \dots \dots \}$

Why prefer one over the other?



# Challenges with DQN:

"Double Q-learning"

## ① Overestimation in Q-values

error accumulation

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

↑  
tgt  
θ

not trained on all  $s'$

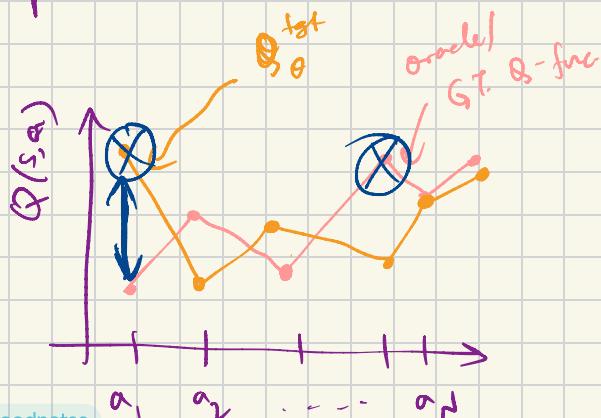
not trained on all  $a'$

Compounding error

what is this term supposed to be?

optimal value from  $s'$  ideal  $\boxed{V^*(s')}$  but  
 $Q$  is  $Q_\theta \neq Q^*$

$\max_{a'} Q_\theta(s', a')$  → erroneous & we amplify the tve errors



$Q_0^{tgt}$  → erroneous  
 $\Downarrow$   
 $Q_0$  → erroneous  
 $\Downarrow$   
 $Q_1$  → erroneous  
 $\Downarrow$   
 $Q_K^{tgt}$  → erroneous

## Double DQN

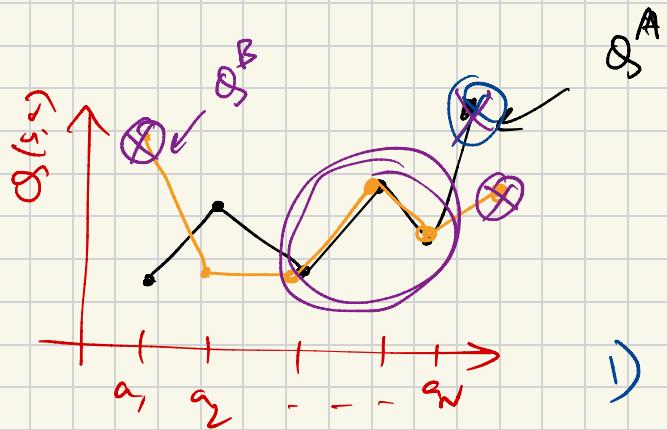
How can we fix this problem?

Idea: use argmax actions from one Q-function  
 & compute Q-values under the other

$$\frac{(Q_{\theta}^A - y^A)^2}{Q_{\theta}^A} + \frac{(Q_{\theta}^B - y^B)^2}{Q_{\theta}^B}$$

$Q_{\theta}^{\text{tgt}, A}$        $Q_{\theta}^{\text{tgt}, B}$

initialized differently.



$$\max_a Q_t(s, a)$$

$$1) \arg\max_a Q_t(s, a) := a^*$$

$$2) Q_t^{\text{tgt}, B}(s, a^*)$$

$$y = r(s, a) + \gamma \max_{a'} Q_t^{\text{tgt}}(s', a')$$

$$y^A = r(s, a) + \gamma Q_t^{\text{tgt}, A}(s', a'^*)$$

Another way to fix this problem.

Idea: reduce error compounding by using N-step returns.

$$\left( Q_0(s, a) - y(s, a) \right)^2$$

$$E_{s, a, s', r \sim \pi} [ \dots ]$$

$$\underline{y(s, a)} = \underline{r(s)} + \gamma \max_{\substack{\text{tgt} \\ a'}} Q(s', a') \quad \downarrow \quad \begin{array}{l} \text{arg max}_{a'} \\ \text{if seen in the replay buffer} \end{array}$$

$$\underline{Q(s, a)} = \underline{r(s, a)} + \gamma \underline{r(s', a')} + \gamma^2 \max_{a''} Q^{\text{tgt}}(s'', a'')$$

$$\underline{Q(s, a)}$$

"N"

$$r(s, a) + \gamma r(s_{01}, a_{01}) + \gamma^2 r(s_{02}, a_{02}) + \dots + \gamma^N \max_{a_{N+1}} Q^{\text{tgt}}(s_{N+1}, a_{N+1})$$

## Practical details for Q-learning

- ① Coverage really helps.
- ② High UTD can destabilize training!
- ③ Improved loss functions:

Huber loss:

$$L(x) = \begin{cases} \frac{x^2}{2}, & |x| < \delta \\ \delta|x| - \frac{\delta^2}{2}, & |x| \geq \delta \end{cases}$$

Cross-entropy loss:

## Advanced content:

### Residual gradient vs TD-learning

$$\min_{\theta} \mathbb{E} \left[ \| Q_{\theta}(s, a) - (r + \gamma Q_{\theta}(s', a')) \|_2^2 \right]$$

$$\Rightarrow \min_{\theta} \mathbb{E}_{s, a, s'} \left[ (Q_{\theta}(s, a) - \gamma Q_{\theta}(s', a') - r(s, a))^2 \right]$$

What's the difference?

At convergence, in theory:

In practice: