

# Lecture 5: Value-based RL

Last time we saw:

Actor-critic methods.

$$D_{i+1} = D_i + \gamma$$

$$\eta \mathbb{E}_{\tau \sim P_{\pi}(\tau)} \left[ \sum_{i=0}^T \nabla_{\theta} \log \pi_{\theta}(a_i | s_i) \cdot \underbrace{A^{\pi}(s_i, a_i)}_{\text{advantage fn.}} \right]$$

$a \sim \pi_{\theta}(a|s)$

The advantage is defined as:

$$A^{\pi}(s_i, a_i) = \underbrace{r(s_i, a_i) + \gamma V^{\pi}(s_{i+1}) - V^{\pi}(s_i)}$$

this is the Q-value  
of  $(s_i, a_i)$

$$\Rightarrow A^{\pi}(s_i, a_i) = \underbrace{Q^{\pi}(s_i, a_i)}_{\text{how good action } a_i \text{ is}} - \underbrace{V^{\pi}(s_i)}_{\text{avg. value of an action}}$$

## Back to the RL algorithm skeleton

- ① Collect data.
  - ② Estimate return (advantage) by evaluating the policy
  - ③ Improve the policy
- how do we do this?

$A^\pi$  using rollouts from  $\pi$

→  $R(s_i, a_i) - V_\phi(s_i)$  high variance

## Value-based methods

"off-policy"

- ② → estimate  $A^\pi$  as perfectly as possible
- ③ → update policy to  $\arg\max$


$A = \{a_1, \dots, a_N\}$

$$\pi^*(s) \leftarrow \max_a A^\pi(s, a)$$

# Skeleton of a value-based RL algo

① Collect data in an env.  
using current policy  $\pi$   
(+ some exploration) "noise"

② Train  $A^\pi(s, a) \Rightarrow ??$   $\underline{\underline{A^\pi_\phi(s, a)}}$

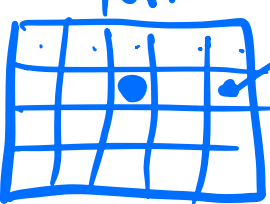
③ update  $\pi$  as  $\pi' := \mathbb{1}_{\{a = \arg\max_a A^\pi(s, a)\}}$   
 $a_1 \dots a_n$   
  
 $A^\pi_\phi(s, a_1) \dots a_n$

Theoretical algorithm:

"Dynamic programming"

$$\Rightarrow Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(s'|s, a)} [V(s')]$$

How can we implement this in a concrete algorithm?

$|S|$    $Q(s, a)$   
 $|S| = \text{finite}, |A| = \text{finite}$   
 $\leftarrow r + \gamma$

## Value iteration

while converged:

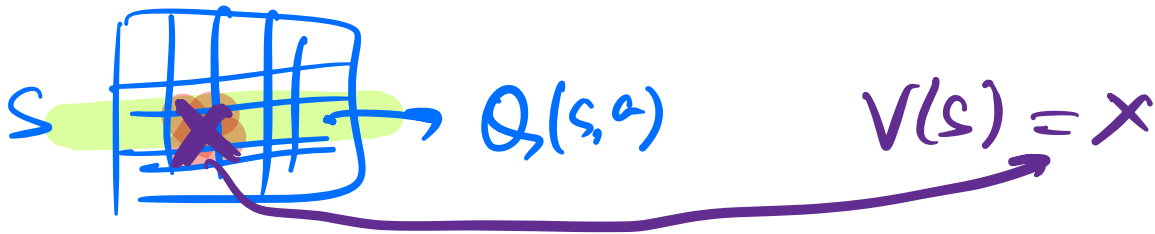
①

$$\underline{Q(s,a)} \leftarrow r(s,a) + \gamma E_{s'} [V(s')]$$

②

$$\underline{V(s)} \leftarrow \max_a Q(s,a)$$

A what does this mean?



## Fitted Q-iteration

Do this using neural networks!

$$Q(s,a) \leftarrow r(s,a) + \gamma E_{s'} [V(s')]$$

Train  $Q_0$  to satisfy:

$$Q_0(s,a) \leftarrow r(s,a) + \gamma E_{s'} [\max_{a'} Q_0(s',a')]$$

How do I write this down as a loss fn??

Draw inspiration from TD learning:  
gradient descent

using  $Q_0$

$$Q_1 \leftarrow \min_{\theta} \mathbb{E}_{\substack{s, a \sim \text{somewhere} \\ s' \sim p(\cdot | s, a)}} \left[ \left( Q_{\theta}(s, a) - \text{RHS} \right)^2 \right]$$

anything

target value

$$Q_{\theta}(s, a)$$

$$\text{RHS} \Rightarrow y(s, a) = r(s, a) + \gamma \max_{a'} Q(s', a')$$

$$(s, a, r, s') \quad \sum_{(s, a, r, s')} (Q(s, a) - y)^2$$

This algorithm can use data from any policy!!

$$\mathbb{E}_{\substack{s, a \sim \text{somewhere} \\ s'}} \left[ \text{squared error} \right]$$

Important: This algorithm is not doing gradient descent!

$$\boxed{\text{FQI} \neq \min_{\theta} \mathcal{L}(\theta)}$$

## References:

Martin Riedmüller. "Neural FQI" paper, 2009.

Fu\*, Kumar\* et al. "Diagnosing Bottlenecks in Deep Q-learning Algos" ICML 2019.

---

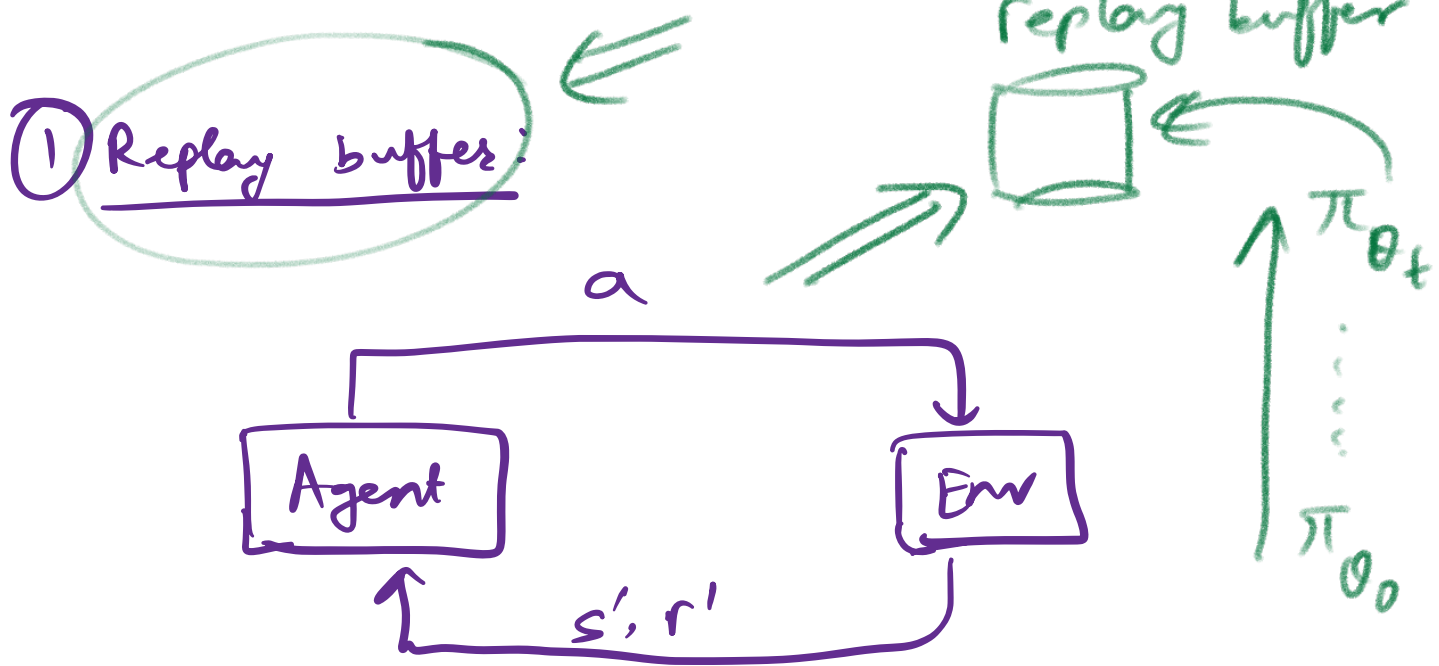
## Deep Q-Networks (DQN)

FQI w/ gradients.

$$\text{TD-Error}(\theta) = \mathbb{E}_{s,a,s' \sim \mathcal{D}} \left[ \left( Q_{\theta}(s,a) - y(s,a) \right)^2 \right]$$

$$\theta_{t+1} \leftarrow \theta_t - \alpha \nabla_{\theta} \text{TD-Error}(\theta) \Big|_{\theta=\theta_t}$$

Problem with doing this on on-policy data??



store  $(s, a, r, s')$  in a dataset of states!

- 1) collect data using  $\pi_t$
  - 2) update  $\theta_{\theta_t}$  using  $(s, a, r, s') \in \mathcal{D}$
- $a \sim \pi_{\theta_0} / \pi_{\theta_1} / \dots / \pi_{\theta_t}$

② Target networks:

- do one gradient step on  $(s, a, r, s')$
- update target network.

Soft targets:

$$\min_{\theta_t} \frac{(\theta_{\theta_t}(s, a) - y_{\theta_{t-1}}(s, a))^2}{\quad}$$

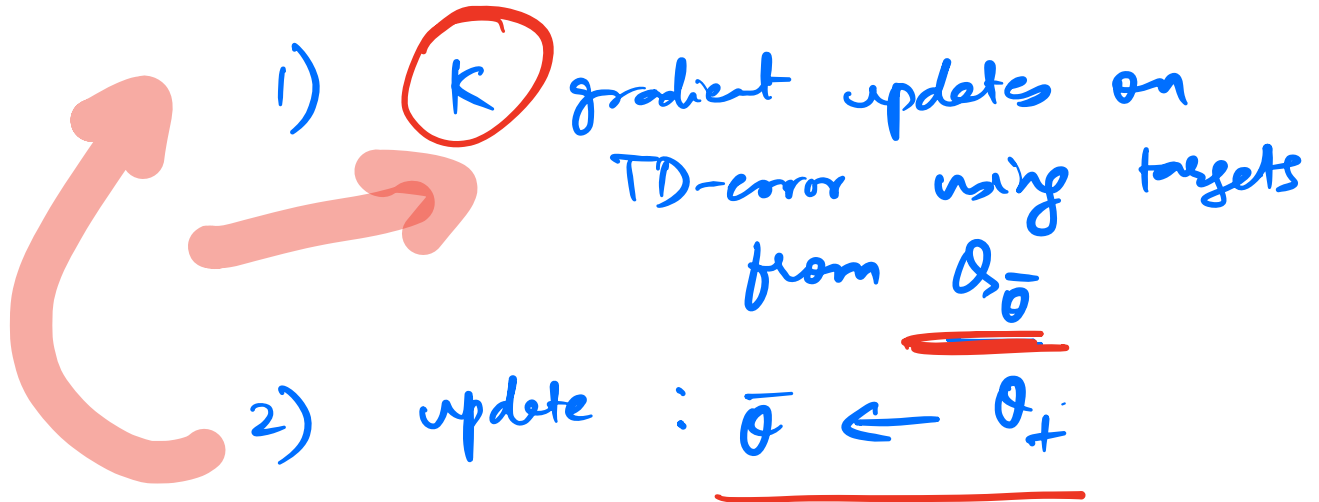
few gradient updates

& release tracks

Hard targets:

$Q_0$

"target" :  $Q_{\bar{\theta}}$



$K=1 \rightarrow$  instabilities / correlations b/w target &  $Q$ -function

Challenges w/ DQNs:

- ① Best way to sample from replay buffers
- ② Noise & stochasticity  
"the overestimation problem"