

*(for Policy Optimizer)*

# **Approximate Value Iteration and Variants**

Chen-Yu Wei

# Value Iteration

For  $k = 1, 2, \dots$

$$\forall s, a, \quad Q_k(s, a) \leftarrow \underbrace{R(s, a)}_{\text{unknown}} + \gamma \sum_{s'} \underbrace{P(s'|s, a)}_{\text{unknown}} \max_{a'} Q_{k-1}(s', a')$$

**Idea:** In each iteration, use multiple samples to estimate the right-hand side.

# Value Iteration with Samples

$$Q_{\theta_k}(x, a) \approx R(x, a)$$

For  $k = 1, 2, \dots$

$$(x_i, a_i, r_i)$$

Obtain  $N$  samples  $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$  where  $\mathbb{E}[r_i] = R(s_i, a_i)$ ,  $s'_i \sim P(\cdot | s_i, a_i)$

Perform **regression** on  $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$  to find  $Q_k$  such that

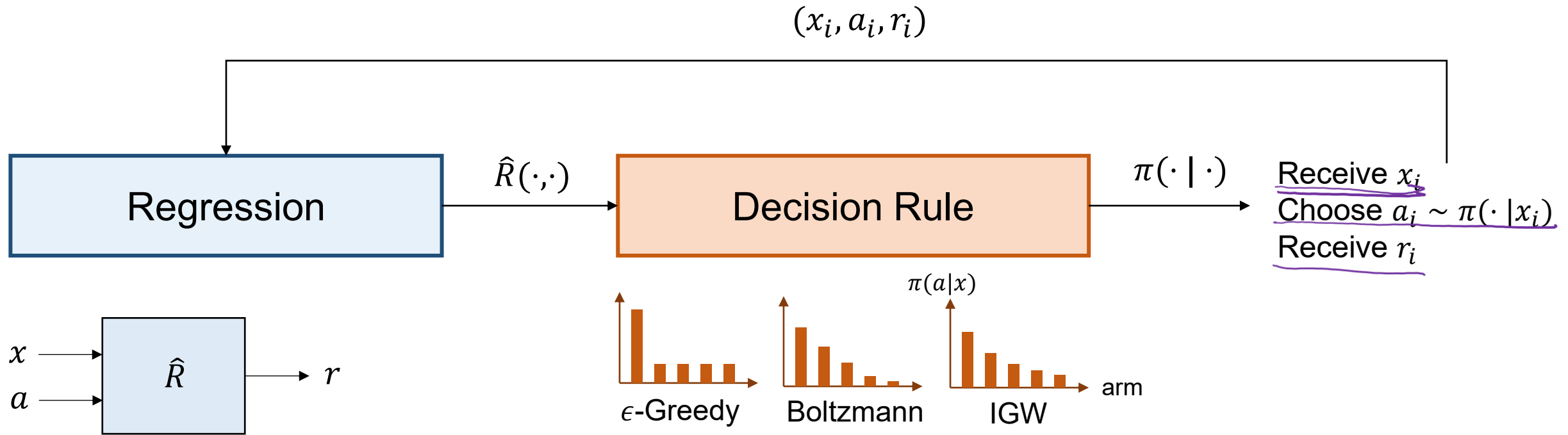
$$\forall s, a, \quad Q_k(s, a) \approx R(s, a) + \gamma \sum_{s'} P(s' | s, a) \max_{a'} Q_{k-1}(s', a')$$

Perform one iteration  
of Value Iteration

Parameterize  $Q_k = Q_{\theta_k}$

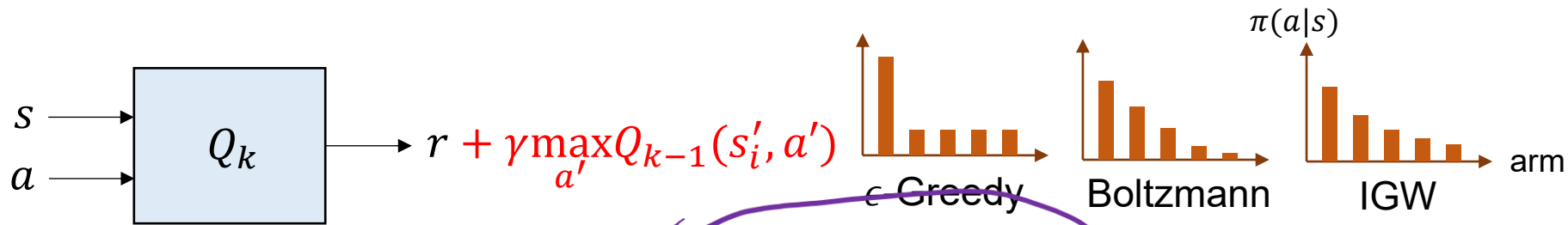
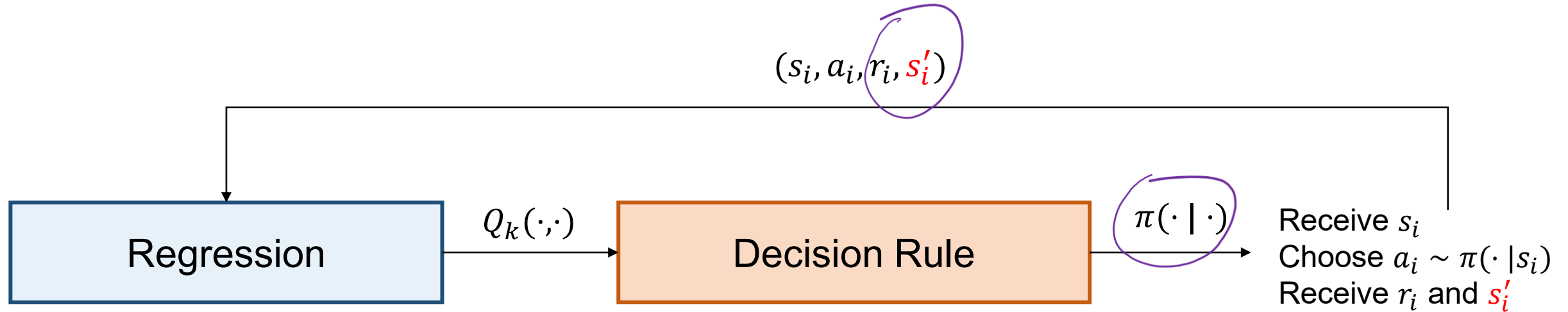
Find  $\theta_k = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N \left( \underbrace{Q_{\theta}(s_i, a_i)} - \underbrace{r_i - \gamma \max_{a'} Q_{\theta_{k-1}}(s'_i, a')}_{} \right)^2$

# Recall: Contextual Bandits with Regression



Train  $\hat{R}$  such that  $\hat{R}(x_i, a_i) \approx r_i = R(x_i, a_i) + \text{noise}$   
 $\hat{R}(x_i, a_i) \approx R(x_i, a_i)$

# Value Iteration with Regression



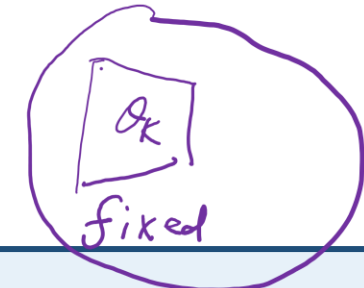
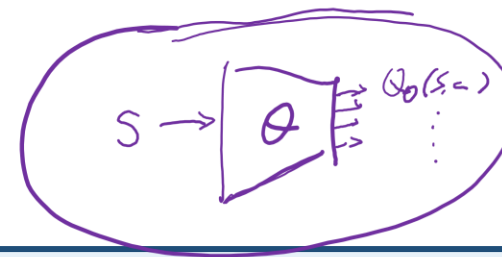
Train  $Q_k$  such that  $Q_k(s_i, a_i) \approx r_i + \gamma \max_{a'} Q_{k-1}(s'_i, a')$  =  $\square + \text{noise}$

This is just one iteration of Value Iteration

Handwritten notes in purple ink show the expectation form of the Bellman optimality equation:

$$Q_k(s_i, a_i) \approx \mathbb{E}_{s' \sim P(\cdot | s_i, a_i)} [r + \gamma \max_{a'} Q_{k-1}(s', a')]$$

# Value Iteration with Samples



For  $k = 1, 2, \dots$

For  $i = 1, 2, \dots, N$ :

Choose action  $a_i \sim \text{EG}(Q_{\theta_k}(s_i, \cdot))$  // or BE or IGW

Receive reward  $r_i \sim R(s_i, a_i)$  and  $s'_i \sim P(\cdot | s_i, a_i)$

$s_{i+1} = s'_i$  if episode continues,  $s_{i+1} \sim \rho$  if episode ends

$\theta \leftarrow \theta_k$

For  $m = 1, 2, \dots, M$ :

Randomly pick an  $i$  (or a mini-batch) from  $\{1, 2, \dots, N\}$

$\theta \leftarrow \theta - \alpha \nabla_\theta \left( Q_\theta(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta_k}(s'_i, a') \right)^2$

$\theta_{k+1} \leftarrow \theta$

$(s_i, a_i, r_i, s'_i)$

main network  
online

target network

Data collection

$$\sum_{i=1}^N \left( Q_\theta(s_i, a_i) - Q_{\theta_k}(s_i, a_i) \right)^2$$

Perform one iteration  
of Value Iteration

Target network

**2<sup>nd</sup> for-loop:** trying to find  $\theta_{k+1} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N \left( Q_\theta(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta_k}(s'_i, a') \right)^2$

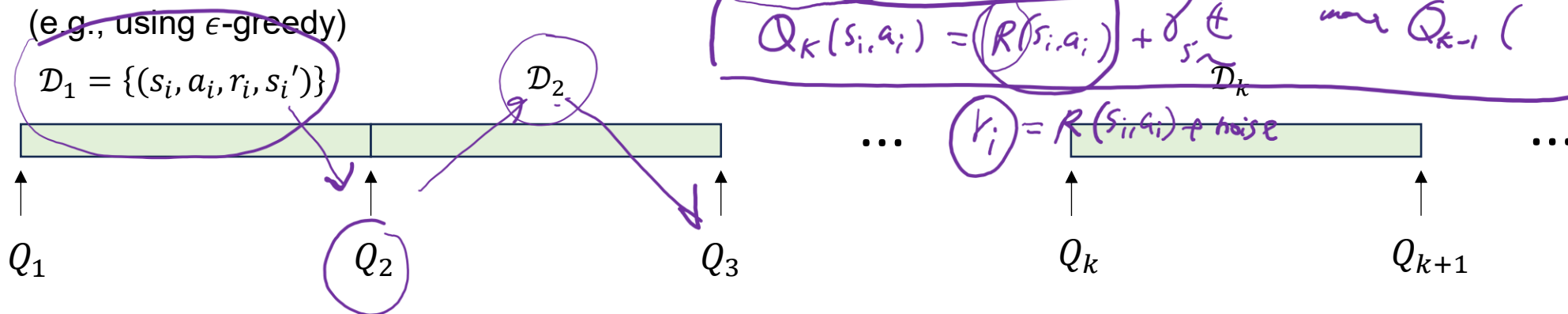
# Reusing Samples

$$\forall s, a: Q_k(s, a) = \underbrace{R(s, a)}_{\text{supp.}(s, a)} + \gamma \mathbb{E}_{a'} \max_{a'} Q_{k-1}(s', a')$$

for any  $S_i, a_i$  collected in previous iters.

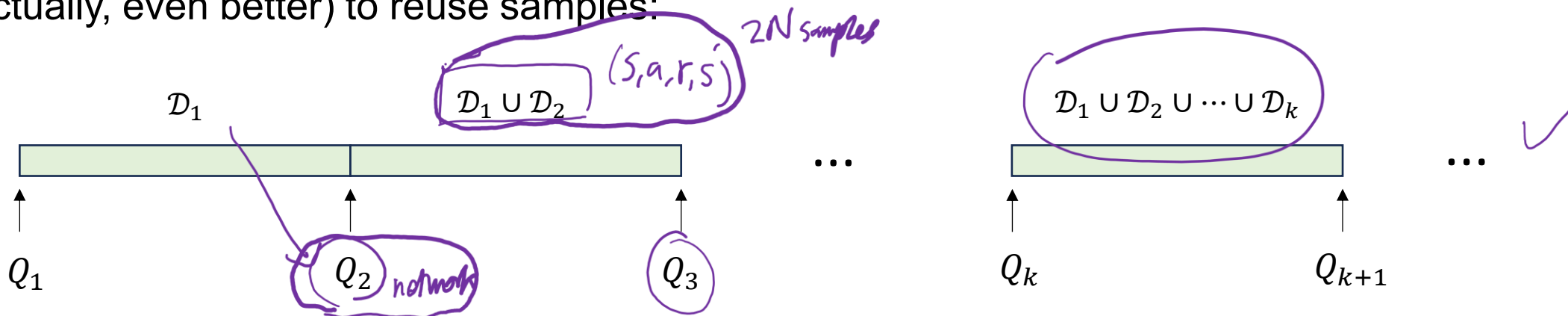
$$Q_k(s_i, a_i) = (R(s_i, a_i) + \gamma \sum_{s' \in \mathcal{D}_k} Q_{k-1}(s', a_i))$$

$$y_i = R(s_i, g_i) + \text{noise}$$



The algorithm in the previous slide only use  $\mathcal{D}_k$  to train  $\theta_{k+1}$ .

However, as the reward function  $R$  and transition  $P$  remains unchanged, it is valid (actually, even better) to reuse samples:



# Benefits of Reusing Samples

- Improving data efficiency
  - Every sample is used multiple times in training – just like we usually go through multiple epochs for supervised learning tasks.
- The ~~buffer~~ will consist of a wider range of state-actions
  - It allows <sup>Set of samples</sup> better approximation of

$$\forall s, a, \quad Q_{k+1}(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_k(s', a')$$



# Value Iteration with Reused Samples (= Deep Q-Learning or DQN)

Initialize  $\mathcal{B} = \{\}$   $\leftarrow$  Replay buffer

For  $k = 1, 2, \dots$

For  $i = 1, 2, \dots, N$ :

Choose action  $a_i \sim \text{EG}(Q_{\theta_k}(s_i, \cdot))$  // or BE or IGW

Receive reward  $r_i \sim R(s_i, a_i)$  and  $s'_i \sim P(\cdot | s_i, a_i)$

$s_{i+1} = s'_i$  if episode continues,  $s_{i+1} \sim \rho$  if episode ends

Push  $(s_i, a_i, r_i, s'_i)$  to  $\mathcal{B}$

$\theta \leftarrow \theta_k$

For  $m = 1, 2, \dots, M$ :

Randomly pick an  $i$  (or a mini-batch) from  $\mathcal{B}$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \left( \underbrace{Q_{\theta}(s_i, a_i)} - r_i - \gamma \max_{a'} \underbrace{Q_{\theta_k}(s'_i, a')} \right)^2$$

$\theta_{k+1} \leftarrow \theta$

HW3 task

Data collection

Perform one iteration  
of Value Iteration

Target network

# Another Popular Implementation

HW3 task

Initialize  $\mathcal{B} = \{\}$   $\leftarrow$  Replay buffer

For  $k = 1, 2, \dots$

For  $i = 1, 2, \dots, N$ :

Choose action  $a_i \sim \text{EG}(Q_\theta(s_i, \cdot))$

Receive reward  $r_i \sim R(s_i, a_i)$  and  $s'_i \sim P(\cdot | s_i, a_i)$

$s_{i+1} = s'_i$  if episode continues,  $s_{i+1} \sim \rho$  if episode ends

Push  $(s_i, a_i, r_i, s'_i)$  to  $\mathcal{B}$

For  $m = 1, 2, \dots, M$ :

Randomly pick an  $i$  (or a mini-batch) from  $\mathcal{B}$

$$\theta \leftarrow \theta - \nabla_\theta \left( Q_\theta(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\bar{\theta}}(s'_i, a') \right)^2$$

$$\bar{\theta} \leftarrow (1 - \tau)\bar{\theta} + \tau\theta$$

$\bar{\theta}$   
 $\uparrow$   
Target network

$$\tau \approx 0.01$$

better approximate VI

# When Does DQN Succeed?

DQN tries to approximate **Value Iteration** by solving

$$\theta_{k+1} = \underset{\theta}{\operatorname{argmin}} \sum_{i \in \mathcal{B}} \left( Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta_k}(s'_i, a') \right)^2 \quad (1)$$

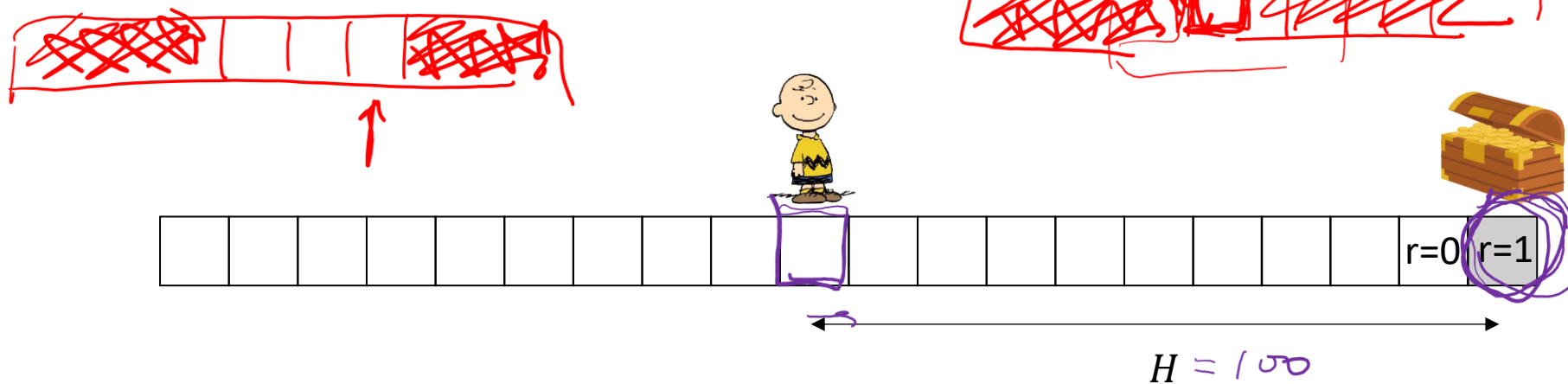
The true Value Iteration:

$$\forall s, a, \quad Q_{k+1}(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_k(s', a') \quad (2)$$

Under what conditions can (1) well approximate (2)?

- $\mathcal{B}$  should contain a wide range of state-action pairs (a challenge of **exploration**) ✓
- $Q_{\theta_{k+1}}(s, a)$  should recover  $R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta_k}(s', a')$  well for all state-actions (a challenge of **function approximation**, or **generalization**)

# 1. Exploration in MDPs (Not Easy)



Environment:

- Fixed-horizon MDP with episode length  $H$
- Initial state at 0
- A single rewarding state at state  $H$
- Actions: Go LEFT or RIGHT

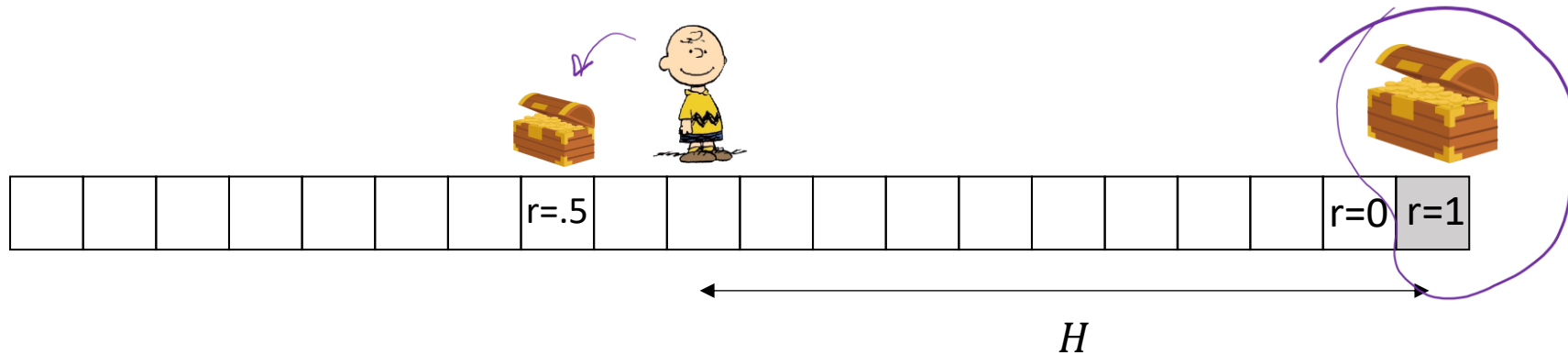
$$\frac{1}{2^H}$$

Suppose we perform DQN with  $\epsilon$ -greedy with random initialization

$\Rightarrow$  On average, we need  $2^H$  episodes to see the reward

(before that, we won't make any meaningful update and will just do random walk around state 0)

# 1. Exploration in MDPs (Not Easy)



Environment:

- Fixed-horizon MDP with episode length  $H$
- Initial state at 0
- A single rewarding state at state  $H$
- Actions: Go LEFT or RIGHT

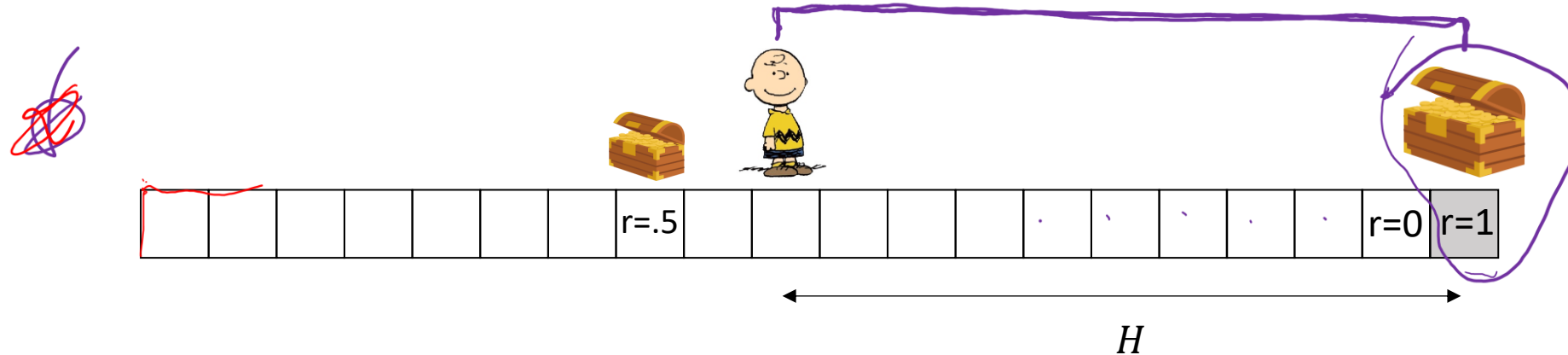
Suppose we perform DQN with  $\epsilon$ -greedy with random initialization

⇒ On average, we need  $2^H$  episodes to see the reward

(before that, we won't make any meaningful update and will just do random walk around state 0)

# 1. Exploration in MDPs (Not Easy)

$$\text{loss} = \text{dist}(\text{Learner}, \text{goal})$$



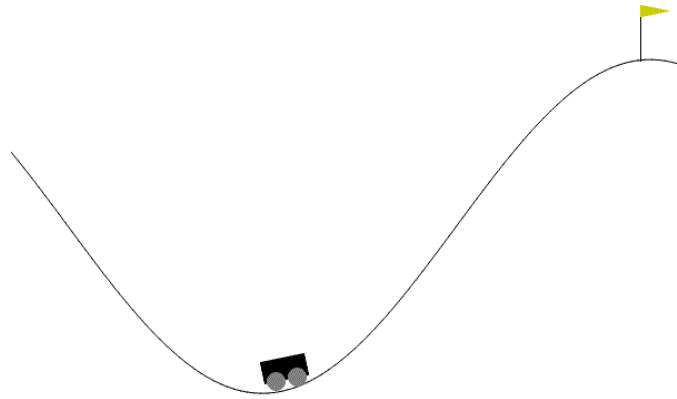
Key issue:

- The  $\epsilon$ -greedy strategy (or BE, IGW) performs **action-space** exploration but not **state-space** exploration.
- This problem becomes more severe when the reward signal is **sparse** and the horizon length is **long**.
- To solve this, we usually require the **exploration bonus** (like UCB, TS), or a better **reward design**. (We will discuss them much later in the course)

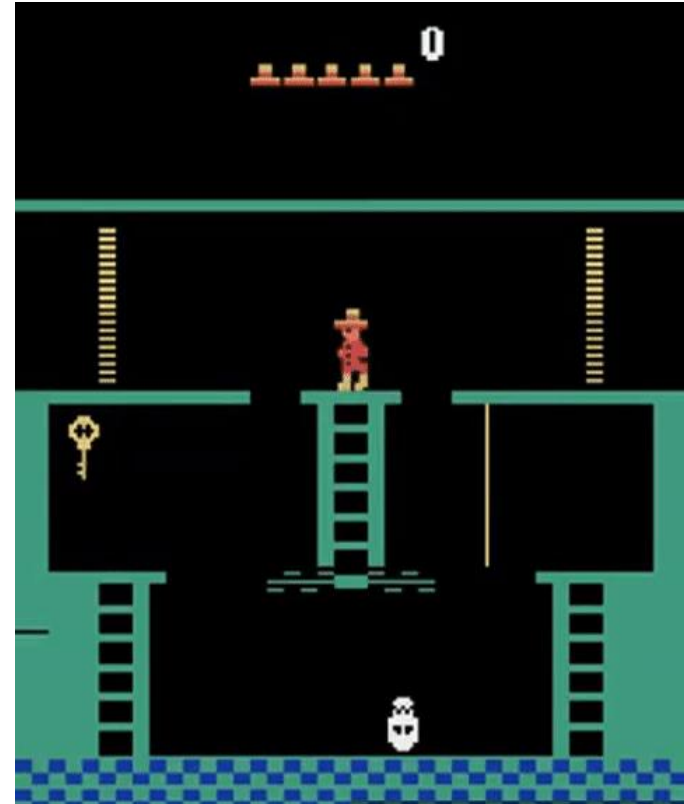
At this point (for the discussion of DQN), we pretend that EG, BE, or IGW will lead to sufficient exploration over the **state space**.

# 1. Exploration in MDPs (Not Easy)

Classic sparse-reward environments:



Mountain Car



Montezuma's Revenge

## 2. Function Approximation

To make DQN well approximate VI, we need

$$\forall s, a \quad Q_{\theta_{k+1}}(s, a) \approx R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta_k}(s', a')$$

### ( $\epsilon$ -approximate) Bellman Completeness

an assumption both on the MDP and the function expressiveness

$$\forall \theta', \exists \theta \quad \forall s, a, \quad \left| Q_{\theta}(s, a) - \left( R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta'}(s', a') \right) \right| \leq \epsilon$$

This allows us to quantify the regression error in each iteration.



## 2. Function Approximation

In HW1 you have shown

$\epsilon$ -Greedy ensures

$$\text{Regret} \lesssim \epsilon T + \sqrt{\frac{AT \cdot \text{Err}}{\epsilon}}$$

Regression error

$$\text{Err} = \sum_{t=1}^T \left( \hat{R}_t(x_t, a_t) - R(x_t, a_t) \right)^2$$

In value-based contextual bandits, the requirement / assumption for function approximation is

$$\exists \theta \quad \forall x, a \quad R_{\theta}(x, a) \approx R(x, a)$$

In value-based MDPs, the requirement / assumption for function approximation is

$$\forall \theta', \exists \theta \quad \forall s, a \quad Q_{\theta}(s, a) \approx R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta'}(s', a')$$

# Analysis of DQN assuming sufficient exploration and Bellman Completeness

Recall the analysis for the exact Value Iteration:

1. Value Iteration will terminate.

$$|Q_k(s, a) - Q_{k-1}(s, a)| \leq \epsilon \quad \forall s, a$$

2. When it terminates, it holds that

$$|Q_k(s, a) - Q^*(s, a)| \leq \frac{\epsilon}{1 - \gamma} \quad \forall s, a$$

3. When it terminates, it holds that

$$V^*(s) - V^{\hat{\pi}}(s) \leq \frac{2\epsilon}{(1 - \gamma)^2} \quad \forall s$$

where  $\hat{\pi}(s) = \operatorname{argmax}_a Q_k(s, a)$

$$\begin{aligned} & \max_{s,a} |Q_k(s, a) - Q_{k-1}(s, a)| \\ & \leq \gamma \max_{s,a} |Q_{k-1}(s, a) - Q_{k-2}(s, a)| \end{aligned}$$

$$\text{ValueError} \leq \frac{1}{1 - \gamma} \text{BellmanError}$$

$$\text{Suboptimality} \leq \frac{1}{1 - \gamma} \text{ValueError}$$

# DQN can be offline

information

$$R(s,a) \quad P(s'|s,a)$$

Let  $\mathcal{B}$  consists of  $(s, a, r, s')$  tuples collected by a mixture of arbitrary policies.

Data collection

For  $k = 1, 2, \dots$

$$\theta \leftarrow \theta_k$$

For  $m = 1, 2, \dots, M$ :

Randomly pick an  $i$  (or a mini-batch) from  $\mathcal{B}$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \left( Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta_k}(s'_i, a') \right)^2$$

$$\theta_{k+1} \leftarrow \theta$$

Perform Value Iteration

Again, its success relies on 1)  $\mathcal{B}$  contains data with sufficiently wide range of state-actions, 2) Bellman completeness.

The same theoretical analysis applies.

# **Handling the Non-Ideal Case**

# When DQN cannot well-approximate VI

In practice,

- We may not be able to collect sufficiently wide range of state-actions
- Bellman completeness may not hold

In either case, we may not have

$$\forall s, a \quad Q_{\theta_{k+1}}(s, a) \approx R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta_k}(s', a')$$

This makes our previous analysis based on VI fails.

# When DQN cannot well-approximate VI

In this case,  $Q_{\theta_k}(s, a)$  tends to **overestimate**  $Q^*(s, a)$ , and the greedy policy  $\hat{\pi}(s) = \operatorname{argmax}_a Q_{\theta_k}(s, a)$  could be very bad.

In iteration  $k$ , we try to approximate

$$Q_{\theta_k}(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q_{\theta_{k-1}}(s', a') + \boxed{\text{err}_k(s, a)}$$

In the next iteration, we try to app

$$Q_{\theta_{k+1}}(s, a) \approx R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} \left( \underbrace{Q_{\theta_k}(s', a')}_{\substack{\text{overestimation} \\ \text{error}}} \right) + \boxed{\phantom{\text{err}_k(s, a')}} + \text{err}_k(s, a')$$

# When DQN cannot well-approximate VI

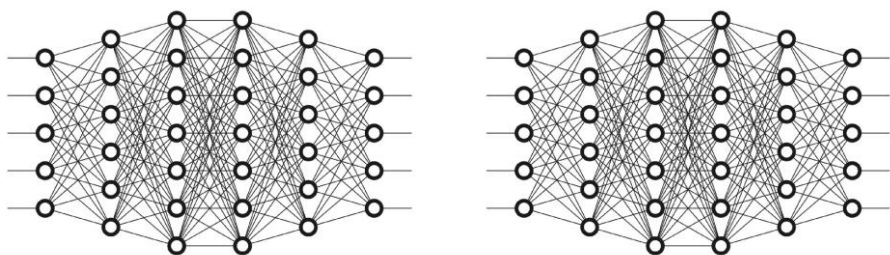
Such “seeking the error” behavior is due to “bootstrapping”

- An issue only in MDP but not in bandits

To prevent overestimation, two strategies are

- Double Q-learning: decorrelating the choice of argmax action and the error of the value function
- Conservative Q-learning: being conservative

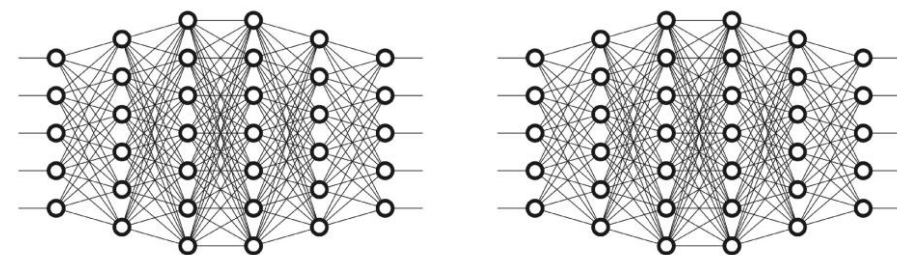
# Double DQN (v1)



$\theta_1$

$\bar{\theta}_1$

$$\text{loss} = \left( Q_{\theta_1}(s, a) - r - \gamma \max_{a'} Q_{\bar{\theta}_1}(s', a') \right)^2$$



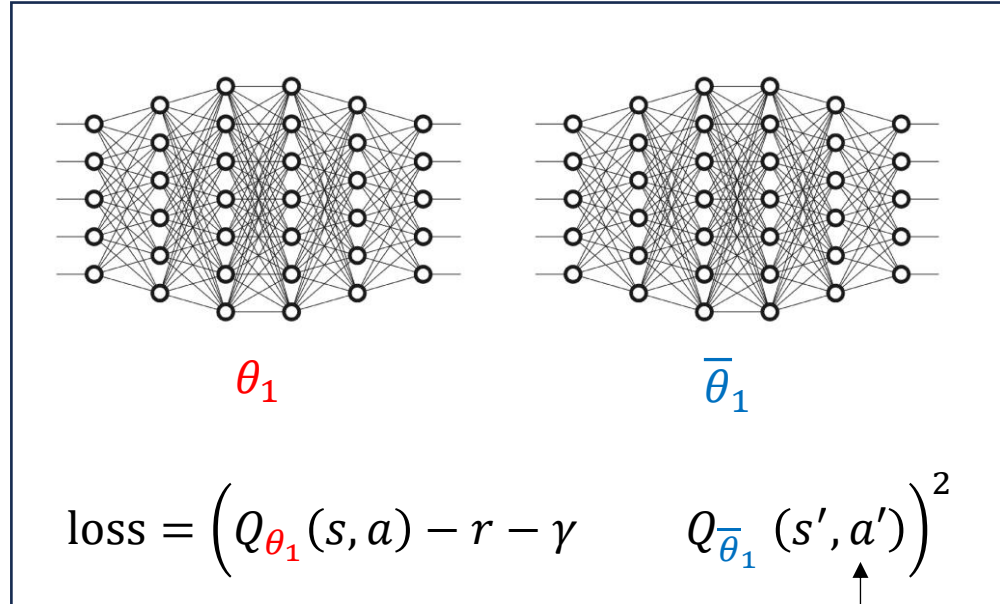
$\theta_2$

$\bar{\theta}_2$

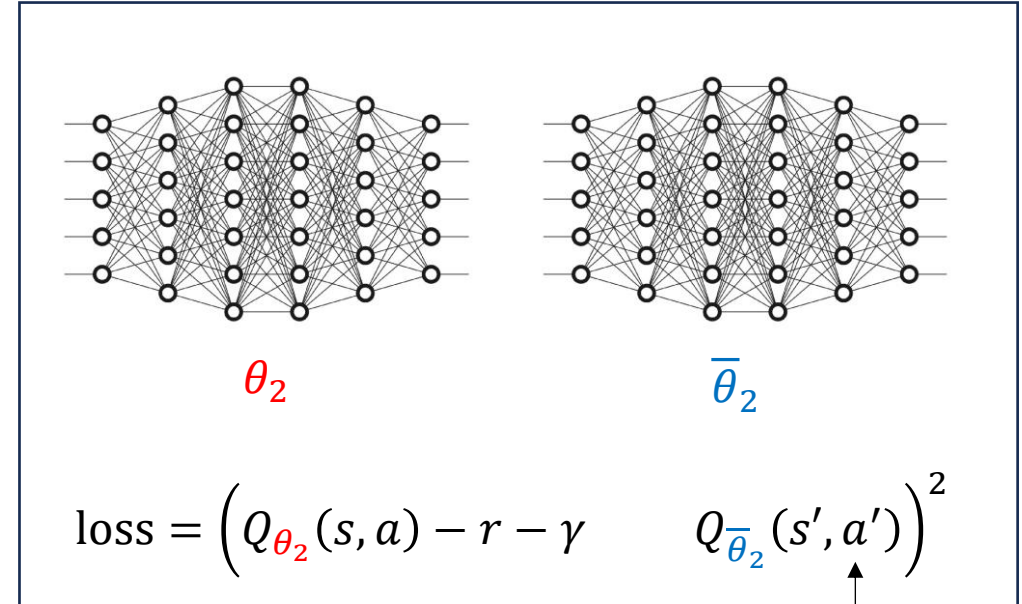
$$\text{loss} = \left( Q_{\theta_2}(s, a) - r - \gamma \max_{a'} Q_{\bar{\theta}_2}(s', a') \right)^2$$



# Double DQN (v1)

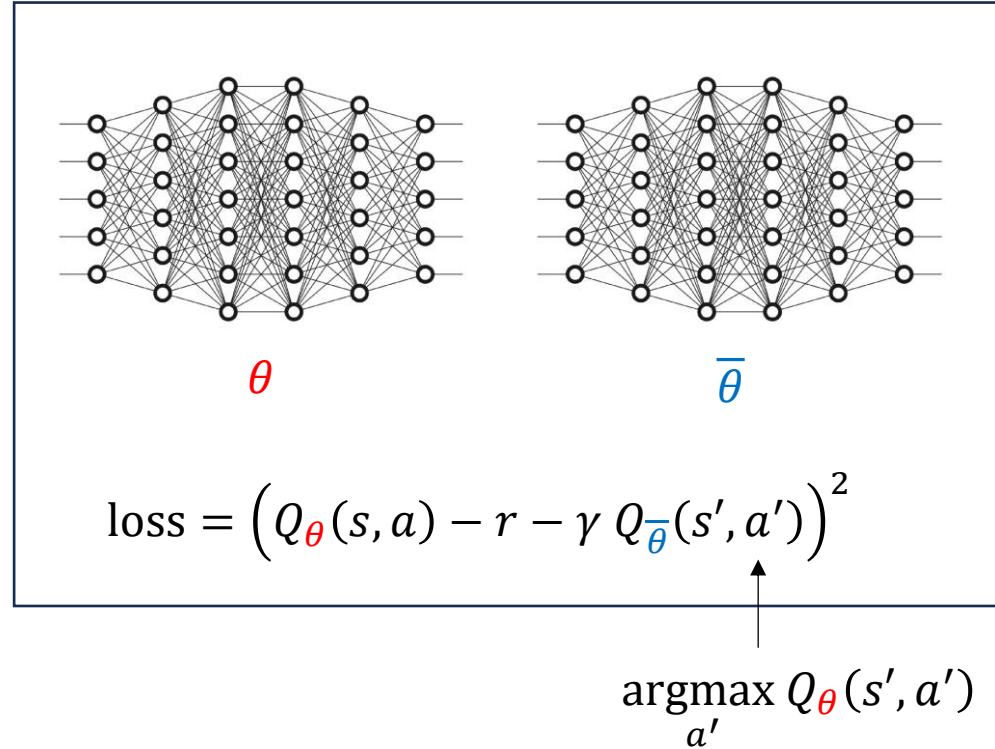


$$\underbrace{\text{argmax}_{a'} Q_{\bar{\theta}_2}(s', a')}_{a'}$$



$$\underbrace{\text{argmax}_{a'} Q_{\bar{\theta}_1}(s', a')}_{a'}$$

# Double DQN (v2)



Hado van Hasselt, Arthur Guez, David Silver. Deep Reinforcement Learning with Double Q-learning. 2015.

# Conservative Q-learning (CQL)

$\mathcal{B}$  = offline dataset

(offline RL)

$$\theta_{k+1} = \operatorname{argmin}_{\theta} \sum_{i \in \mathcal{B}} \left( Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta_k}(s'_i, a') \right)^2$$

$$+ \alpha \sum_{i \in \mathcal{B}} \left( \log \left( \sum_a \exp(Q_{\theta}(s_i, a)) \right) - Q_{\theta}(s_i, a_i) \right)$$

$$= \operatorname{argmin}_{\theta} \sum_{i \in \mathcal{B}} \left( Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta_k}(s'_i, a') \right)^2$$

$$+ \alpha \sum_{i \in \mathcal{B}} \left( \max_{\mu} \sum_a \mu(a|s_i) Q_{\theta}(s_i, a) - Q_{\theta}(s_i, a_i) - \text{KL}(\mu(\cdot | s_i), \text{Unif}) \right)$$

$\max_a Q_{\theta}(s_i, a)$

# Comparison

- Double-Q: make the  $\operatorname{argmax}_a Q_\theta(s, a)$  choice decoupled from  $\theta$
- Conservative-Q: mitigate the overestimation of  $\max_a Q_\theta(s_i, a)$  over  $Q_\theta(s_i, a_i)$

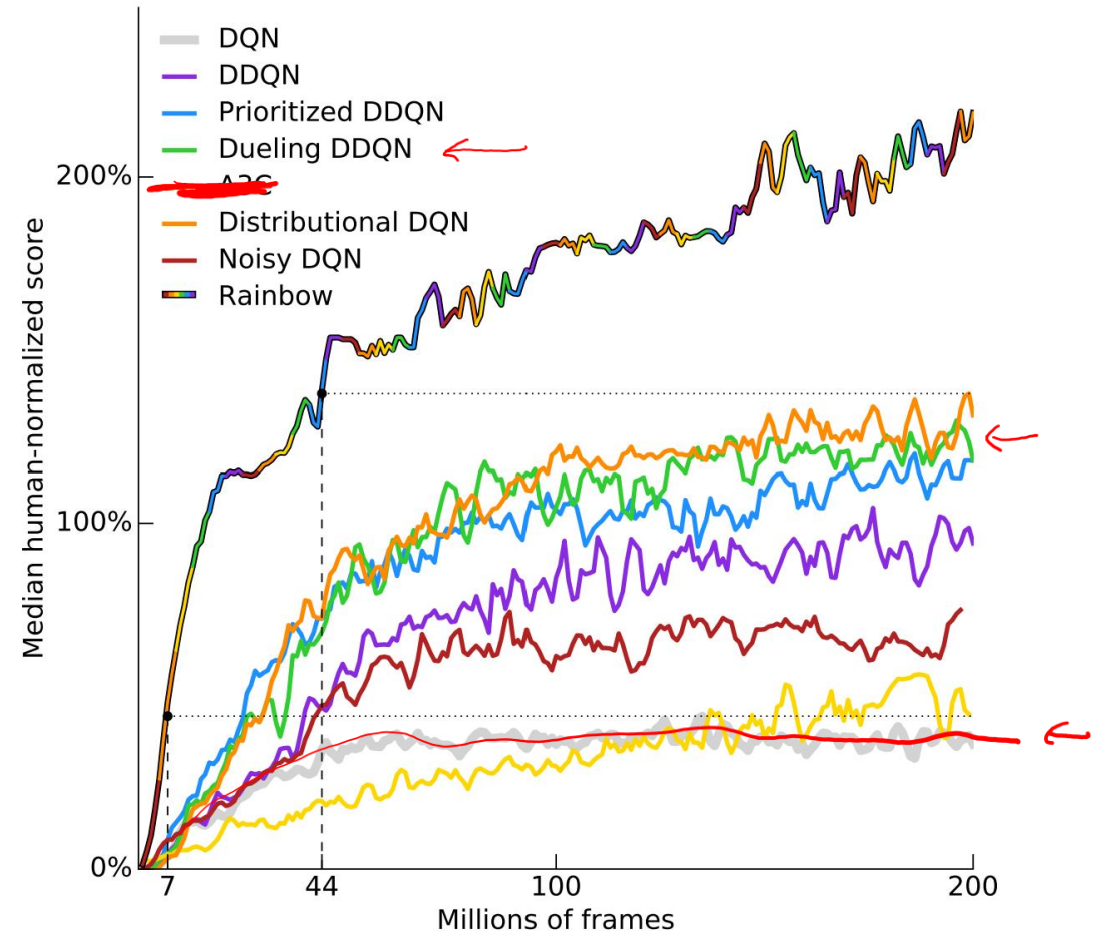
# Summary for DQN

- Motivation: approximating Value Iteration using **samples** and **function approximation**
- Standard elements: target network, replay buffer
- Work as desired when both of the following conditions hold:
  - The learner is able to obtain exploratory data (online or offline)
  - Neural network is sufficiently expressive: Bellman completeness
- When the conditions above do not hold
  - Tends to overestimate  $Q$  values and suggest arbitrary actions
- Solutions
  - Double Q-learning
  - Conservative Q-learning

# Improvements on DQN

- Dueling DDQN
- Prioritized replay
- Distributional DQN
- ...

Rainbow: Combining Improvements in Deep Reinforcement Learning. 2018.



## **Other Variants that Fail**

# An Unstable Variant

DQN without target network

For  $k = 1, 2, \dots$

Randomly pick an  $i$  (or a mini-batch) from  $\mathcal{B}$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \left( Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\bar{\theta}}(s'_i, a') \right)^2$$

$$\bar{\theta} \leftarrow \theta$$

$$\tau = 1$$

$$M = 1$$

cf. DQN with target network

For  $k = 1, 2, \dots$

Randomly pick an  $i$  (or a mini-batch) from  $\mathcal{B}$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \left( Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\bar{\theta}}(s'_i, a') \right)^2$$

$$\bar{\theta} \leftarrow (1 - \tau) \bar{\theta} + \tau \theta$$

$$\tau = 0.01$$

For  $k = 1, 2, \dots$

$$\theta \leftarrow \bar{\theta}$$

For  $m = 1, \dots, M$ :

Randomly pick an  $i$  (or a mini-batch) from  $\mathcal{B}$

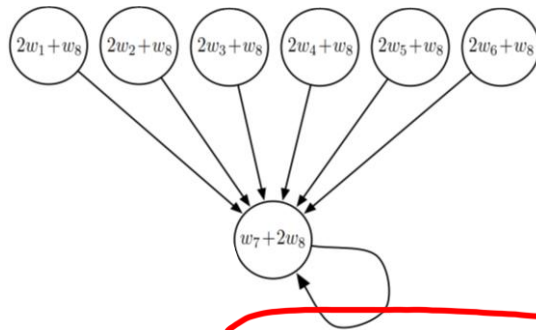
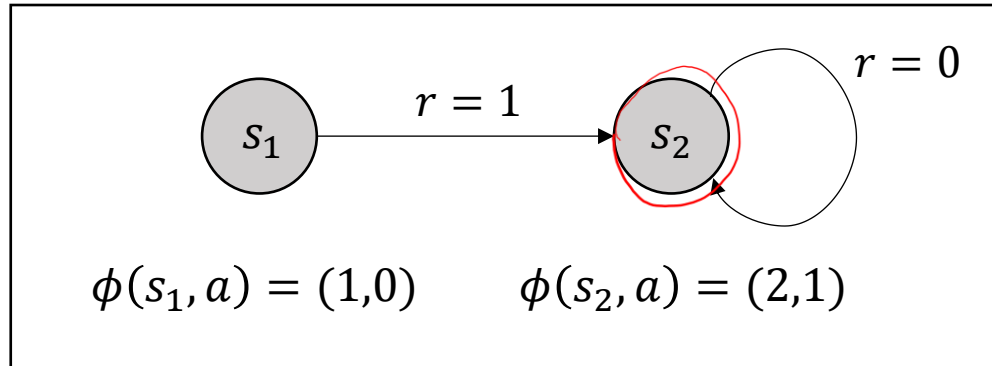
$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \left( Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\bar{\theta}}(s'_i, a') \right)^2$$

$$\bar{\theta} \leftarrow \theta$$



# An Unstable Variant

**Diverges** even when exploration assumption and Bellman completeness hold



Simplified from the “Baird’s counterexample”  
(see Sutton and Barto Section 11.2)

$$Q_{\theta}(s, a) = \phi(s, a)^T \theta$$

$$\rightarrow Q^*(s, a) = \phi(s, a)^T \theta^*$$

$$Q^*(s_2, a) = 0 = (2, 1) \cdot (\theta_1^*, \theta_2^*) = 2\theta_1^* + \theta_2^*$$

$$Q^*(s_1, a) = 1 = (1, 0) \cdot (\theta_1^*, \theta_2^*) = \theta_1^*$$

$$\Rightarrow \begin{cases} \theta_1^* = 1 \\ \theta_2^* = -2 \end{cases}$$

# The Effect of Target Network

Let  $KN = 100000$

For  $k = 1, 2, \dots, K$

$$\theta_k \leftarrow \theta$$

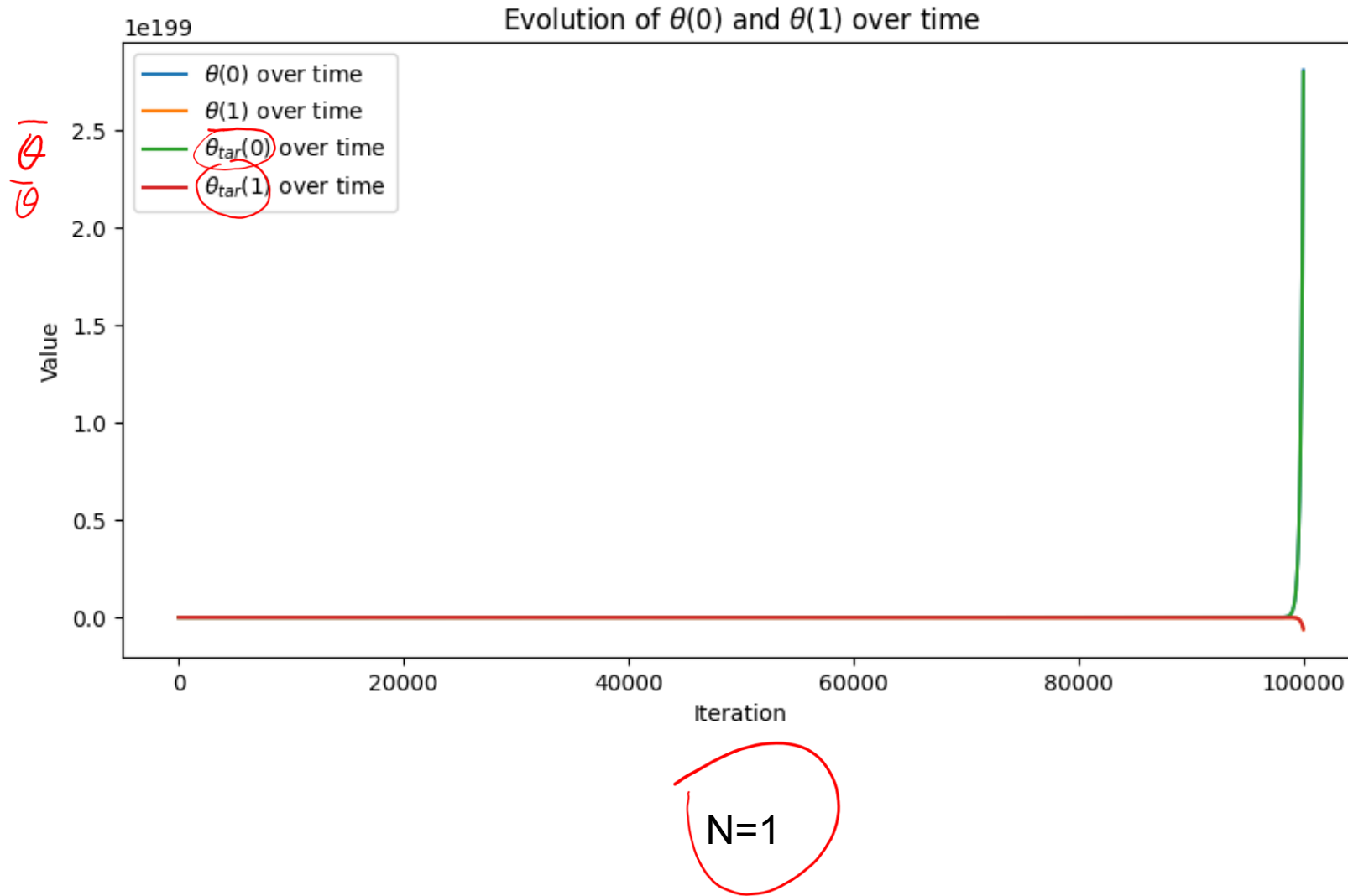
For  $i = 1, \dots, N$ :

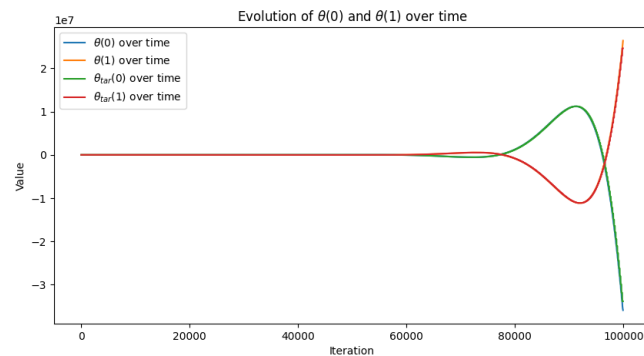
Sample  $(s, a, r, s') \sim \text{Uniform} \{(s_1, a, 1, s_2), (s_2, a, 0, s_2)\}$

$$\theta \leftarrow \theta - \alpha \left( \underline{\phi(s, a)^\top \theta} - r - \gamma \phi(s', a)^\top \theta_k \right) \phi(s, a)$$

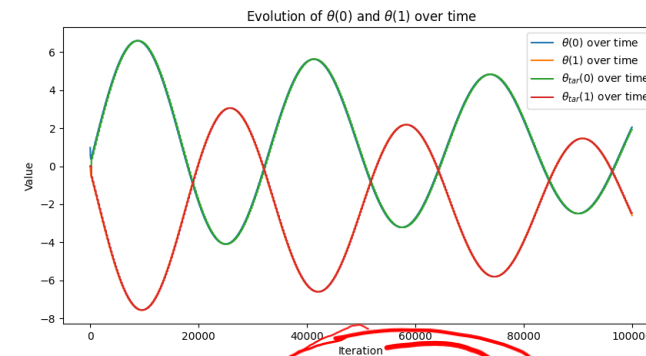
$$\theta_{k+1} \leftarrow \theta$$

# The Effect of Target Network

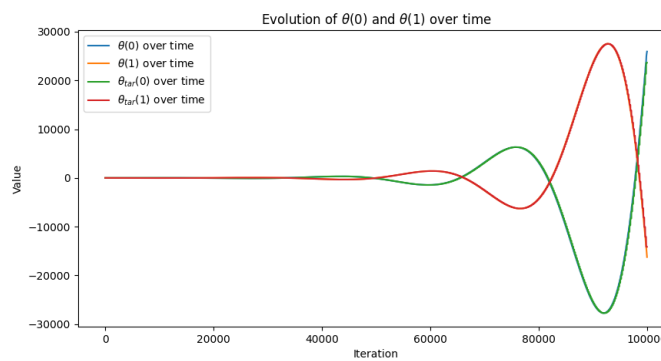




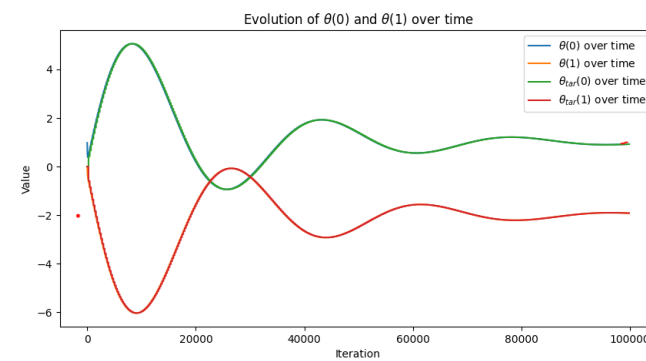
N=150



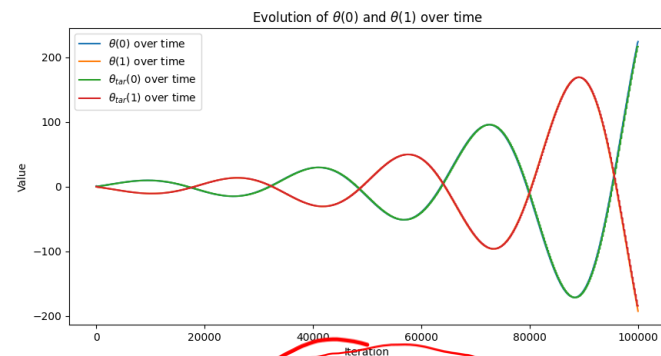
N=210



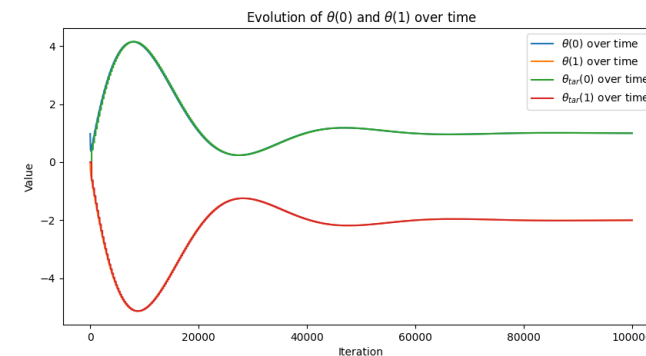
N=170



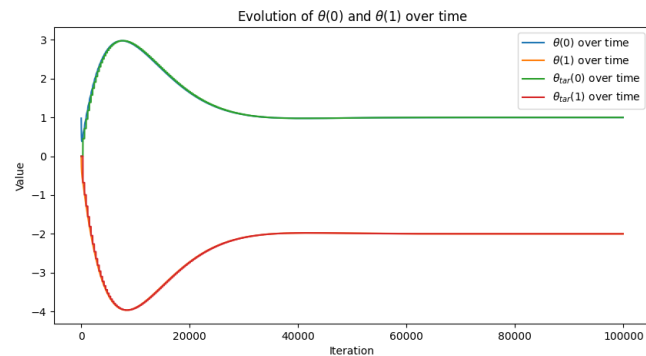
N=230



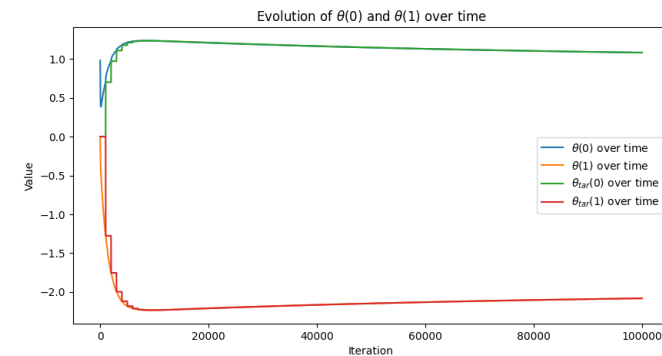
N=190



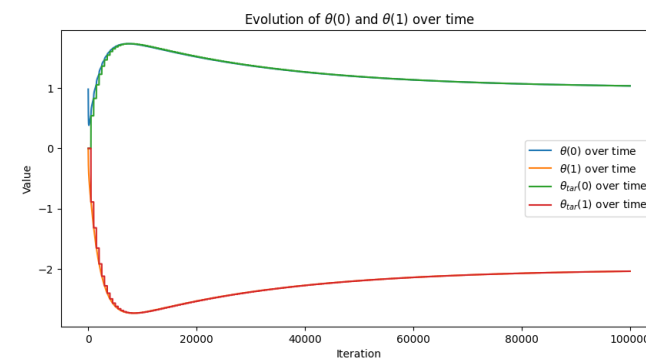
N=250



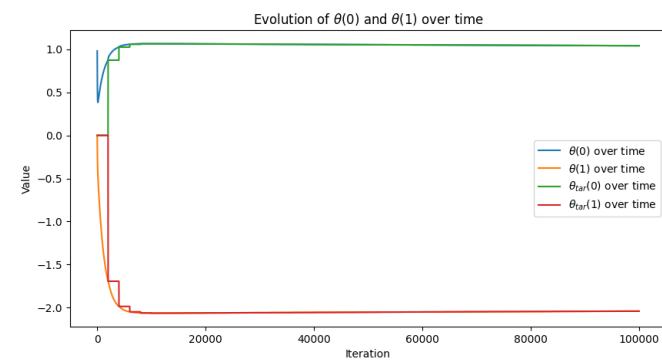
N=300



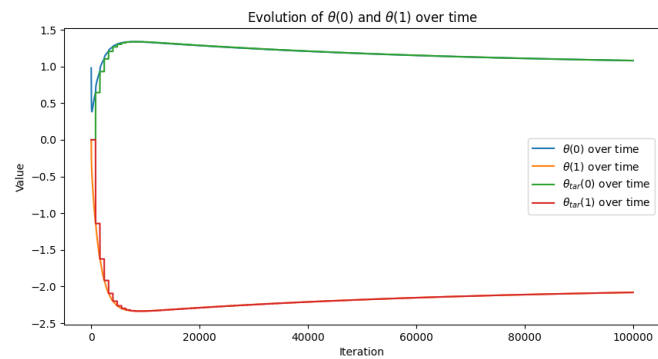
N=1000



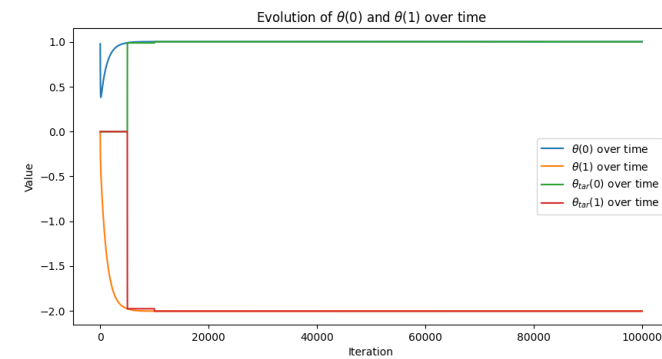
N=500



N=2000



N=800



N=5000

# A Biased Variant

Residual Gradient (DQN without **stop gradient**)

For  $k = 1, 2, \dots$

Randomly pick an  $i$  (or a mini-batch) from  $\mathcal{B}$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \left( Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta}(s'_i, a') \right)^2$$

*cf.* standard DQN

For  $k = 1, 2, \dots$

Randomly pick an  $i$  (or a mini-batch) from  $\mathcal{B}$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \left( Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\bar{\theta}}(s'_i, a') \right)^2$$

$$\bar{\theta} \leftarrow (1 - \tau) \bar{\theta} + \tau \theta$$

For  $k = 1, 2, \dots$

$$\theta \leftarrow \bar{\theta}$$

For  $m = 1, \dots, M$ :

Randomly pick an  $i$  (or a mini-batch) from  $\mathcal{B}$

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \left( Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\bar{\theta}}(s'_i, a') \right)^2$$

$$\bar{\theta} \leftarrow \theta$$

# A Biased Variant

This variant will converge (as it is similar to standard SGD), but the solution it converges to could be undesirable.

The underlying loss function of the Residual Gradient algorithm is

$$\sum_{i \in \mathcal{B}} \left( Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta}(s'_i, a') \right)^2$$

Value Iteration

$$Q_{\theta}(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \max_{a'} Q_{\theta}(s', a') \right]$$

$$\Rightarrow \text{Desired loss} = \left( Q_{\theta}(s, a) - R(s, a) - \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \max_{a'} Q_{\theta}(s', a') \right] \right)^2$$

$$\left( Q_{\theta}(s, a) - R(s, a) - \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \max_{a'} Q_{\theta}(s', a') \right] \right)^2 + \left( \gamma \max_{a'} Q_{\theta}(s', a') - \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ \max_{a'} Q_{\theta}(s', a') \right] \right)^2$$

# Variants that Fail

- Both variants, while look somewhat reasonable, deviate from the idea of Value Iteration.



# Roadmap

