



# ST455: Reinforcement Learning

## Lecture 5: TD Learning with Function Approximation

Chengchun Shi

# Lecture Outline

---

1. Introduction to Value Function Approximation
2. Gradient Descent-based Methods
3. Fitted Q-Iteration

# Lecture Outline

---

**1. Introduction to Value Function Approximation**

2. Gradient Descent-based Methods

3. Fitted Q-Iteration

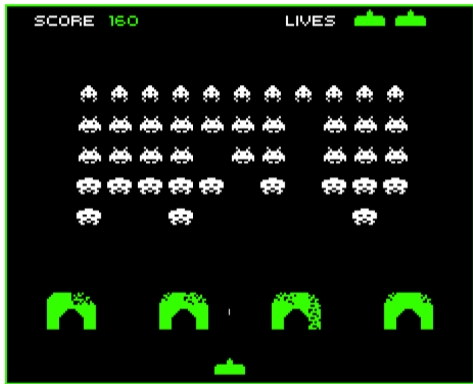
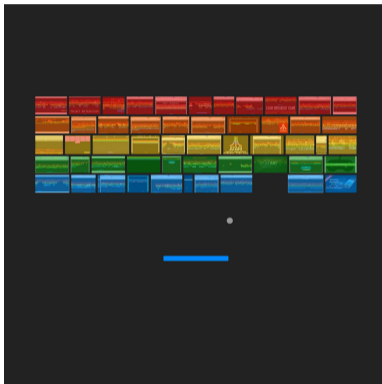
# Limitations of Tabular Methods

---

- So far, we studied reinforcement learning methods using a **tabular representation**
  - Focus on **finite** MDPs
  - Value function represented by a **table**
  - Each state  $s$  has an entry for value  $V(s)$
  - Each state-action pair  $(a, s)$  has an entry for value  $Q(s, a)$
- **Limitations** of tabular methods
  - Cannot handle **large-scale** RL problems or **continuous** state space
  - **Scalability**: computation time and storage needed to maintain estimates
  - **Slow learning**: learning the value of each state individually

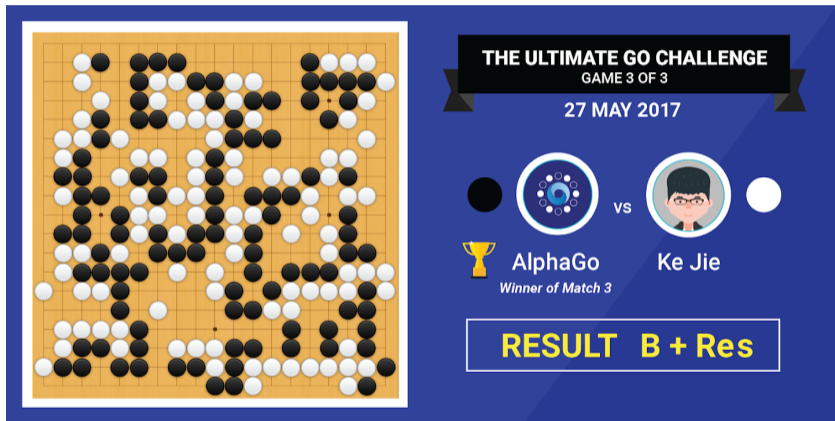
# Large Scale RL Problems (Examples)

---



- Image-valued observations (e.g.,  $210 \times 160$  pixel image frames, 129 colours)

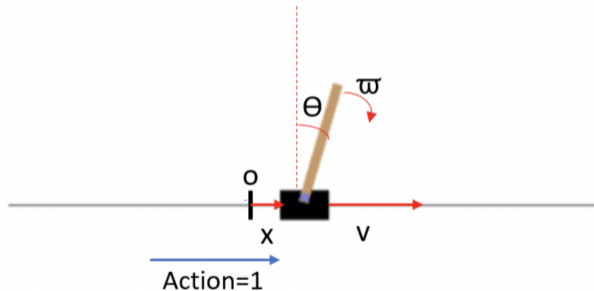
# Large Scale RL Problems (Examples)



- $19 \times 19 = 361$  Go board, each location (empty, black or white)  $\rightarrow 3^{361} \approx 10^{170}$  states

# Continuous State Space (Examples)

frame: 53, Obs: (0.018, 0.669, 0.286, 0.618)  
Action: 1.0, Cumulative Reward: 47.0, Done: 1



- $S_t$ :  $x$  (Position);  $v$  (velocity);  $\theta$  (Angle);  $\omega$  (Angular velocity)
- All components are **continuous**

# Function Approximation

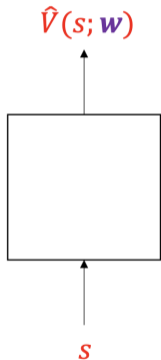
---

- Estimate a value function using a **parametric** approximator function
- $\hat{V}(\mathbf{s}; \boldsymbol{\omega})$  as an approximator for the value function  $V^\pi(\mathbf{s})$
- $\hat{Q}(\mathbf{s}, \mathbf{a}; \boldsymbol{\omega})$  as an approximator for the value function  $Q^\pi(\mathbf{s}, \mathbf{a})$
- Dimension of  $\boldsymbol{\omega}$  much smaller than the state space size. Represents a **tradeoff**:
  - **Bias** (approximation error) usually decreases with dimension of  $\boldsymbol{\omega}$
  - **Variance** (estimation error) usually increase with dimension of  $\boldsymbol{\omega}$
- Update parameter  $\boldsymbol{\omega}$  to find a good approximation using a learning method
  - Eg., MC or TD methods
- Function approximation studied in supervised learning
  - Integrate known methods in reinforcement learning



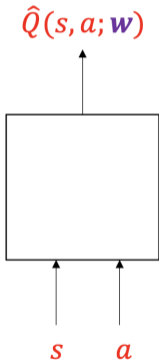
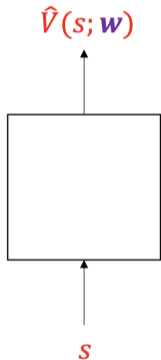
# Types of Value Function Approximator

---



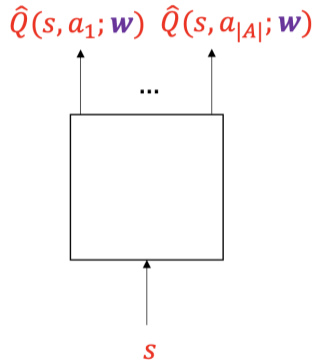
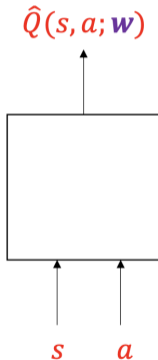
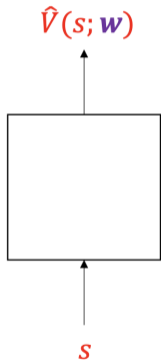
# Types of Value Function Approximator (Cont'd)

---



# Types of Value Function Approximator (Cont'd)

---



# Types of Value Function Approximator (Cont'd)

---

- **Linear combinations of features**
- **Neural networks**
- Decision tree, random forest, boosting
- Nearest neighbor, kernel methods
- ...

# Linear Methods

---

- **Table lookup** is a special case of linear value function approximation in finite MDP

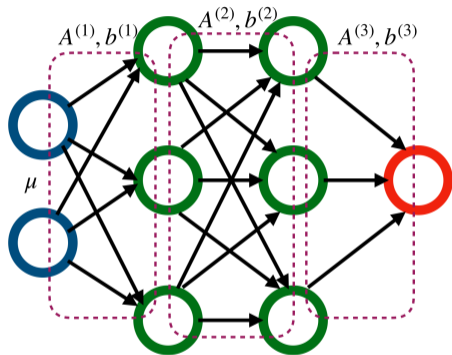
$$\phi(\mathbf{S}) = \begin{bmatrix} \mathbb{I}(\mathbf{S} = \mathbf{s}_1) \\ \mathbb{I}(\mathbf{S} = \mathbf{s}_2) \\ \vdots \\ \mathbb{I}(\mathbf{S} = \mathbf{s}_n) \end{bmatrix}$$

- Parameter vector  $\omega$  gives value of each individual state

$$\hat{V}(\mathbf{s}; \omega) = (\omega_1, \omega_2, \dots, \omega_n)^\top \phi(\mathbf{S})$$

- **Equivalent** to tabular methods
- Other popular choices for  $\phi$ : RBFsampler, splines, polynomials, etc.

# Neural Networks



$$\mu^{(1)} = \sigma(A^{(1)}\mu + b^{(1)}) \in \mathbb{R}^3$$

$$\mu^{(2)} = \sigma(A^{(2)}\mu^{(1)} + b^{(2)}) \in \mathbb{R}^3$$

$$y = A^{(3)}\mu^{(2)} + b^{(3)} \in \mathbb{R}$$

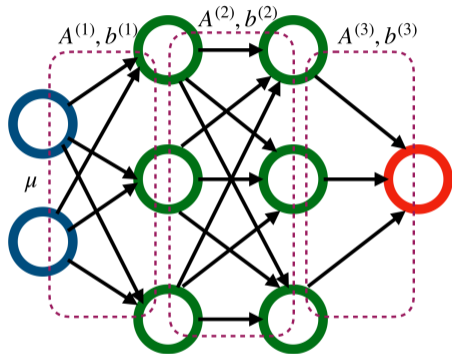
**Figure:** Illustration of fully-connected neural networks with two hidden layers and three hidden nodes per hidden layer. Here  $\mu$  is the 2-dimensional input,  $A^{(\ell)}$  and  $b^{(\ell)}$  denote the corresponding parameters to produce the linear transformation for the  $(\ell - 1)$ th layer, and  $\sigma$  denotes the element-wise nonlinear transformation function (e.g., sigmoid or ReLU)

# Linear v.s. Neural Networks

---

- Neural networks has **universal approximation property** [Barron, 1993]
- Able to approximate both **smooth** and **nonsmooth** (e.g., step function) functions [Imaizumi and Fukumizu, 2019]
- Difficult to **optimize** (using back propagation)
- Linear methods are **computationally efficient** to implement
- Requires **feature engineering** to have good approximation property

## Linear v.s. Neural Networks (Cont'd)



$$\mu^{(1)} = \sigma(A^{(1)}\mu + b^{(1)}) \in \mathbb{R}^3$$

$$\mu^{(2)} = \sigma(A^{(2)}\mu^{(1)} + b^{(2)}) \in \mathbb{R}^3$$

$$y = A^{(3)}\mu^{(2)} + b^{(3)} \in \mathbb{R}$$

- When  $A^{(1)}$ ,  $b^{(1)}$ ,  $A^{(2)}$  and  $b^{(2)}$  are fixed, can treat  $\sigma(A^{(2)}\mu^{(1)} + b^{(2)})$  as features and employ linear method to estimate  $A^{(3)}$  and  $b^{(3)}$
- Neural networks are **adaptive**: the features involve parameters that are adaptively constructed based on the data



# Lecture Outline

---

1. Introduction to Value Function Approximation

**2. Gradient Descent-based Methods**

3. Fitted Q-Iteration

# Function Approximation in RL

---

- Consider the **policy evaluation** problem:  $s \rightarrow V^\pi(s)$
- RL Examples:
  - MC:  $S_t \rightarrow G_t$
  - TD(0):  $S_t \rightarrow R_t + \gamma V(S_{t+1})$
  - TD( $\lambda$ ):  $S_t \rightarrow G_t^\lambda$
- Like supervised learning: **feature vector**  $\rightarrow$  **response**
- Unique characteristics of RL: online learning, nonstationary target functions (value function changes with policy)

# Parameter Estimation

---

- **Goal:** find a parameter  $\omega$  that minimizes a given **error** function  $J : \omega \rightarrow \mathbb{R}$
- **Mean squared error** (common supervised learning objective):

$$J(\omega) = \frac{1}{2} \mathbb{E}_{\mathbf{s} \sim \mu} [\widehat{\mathbf{V}}(\mathbf{s}; \omega) - \mathbf{V}^{\pi}(\mathbf{s})]^2$$

where  $\mu$  is a distribution on the state space  
(specifies how the error is distributed over different states)

- Common choice for  $\mu$ : equal to the **on-policy** distribution
  - Distribution of states encountered under policy  $\pi$
  - Minimize the error that occur while following the policy

# Gradient-Descent Methods

---

- Assume  $\hat{V}(\mathbf{s}; \omega)$  is differentiable with respect to  $\omega$  for each  $\mathbf{s}$
- Consider first a simple case where the training examples are  $\mathbf{S}_t \rightarrow \mathbf{V}^\pi(\mathbf{S}_t)$ 
  - Input examples give the exact value of the state value
- Gradient-Descent Algorithm:  $\omega_{t+1} = \omega_t - \alpha_t \nabla_\omega J(\omega_t)$

$$\omega_{t+1} = \omega_t - \alpha_t \mathbb{E}_{\mathbf{s} \sim \mu} \left[ \left( \mathbf{V}^\pi(\mathbf{s}) - \hat{V}(\mathbf{s}; \omega_t) \right) \nabla_\omega \hat{V}(\mathbf{s}; \omega_t) \right]$$

- Stochastic Gradient-Descent Algorithm:  $\omega_{t+1} = \omega_t - \alpha_t \nabla_\omega J(\omega_t)$

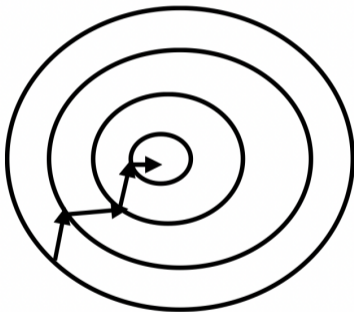
$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( \mathbf{V}^\pi(\mathbf{S}_t) - \hat{V}(\mathbf{S}_t; \omega_t) \right) \nabla_\omega \hat{V}(\mathbf{S}_t; \omega_t) \right]$$

where  $\mathbf{S}_t$  is distributed according to  $\mu$

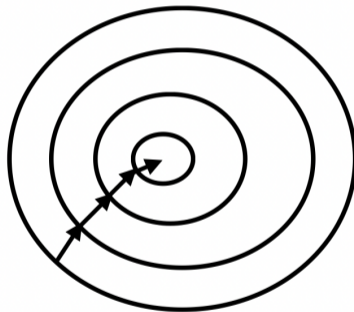
# Gradient-Descent Methods (Cont'd)

---

**Stochastic Gradient Descent**



**Gradient Descent**



- Each point represents a parameter
- Circle represents parameters with the same loss function

# Gradient-Descent Methods (Cont'd)

---

- Assume training examples  $\mathbf{S}_t \rightarrow \nu_t$  where  $\nu_t$  is a target, some approx of  $\mathbf{V}^\pi(\mathbf{S}_t)$
- Stochastic Gradient-Descent Algorithm:  $\omega_{t+1} = \omega_t - \alpha_t \nabla_\omega J(\omega)$

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( \nu_t - \hat{\mathbf{V}}(\mathbf{S}_t; \omega_t) \right) \nabla_\omega \hat{\mathbf{V}}(\mathbf{S}_t; \omega_t) \right]$$

where  $\mathbf{S}_t$  is distributed according to  $\mu$

- Some sufficient conditions for convergence to local minimum
  - Standard assumptions on step size:  $\sum \alpha_t = \infty$  &  $\sum \alpha_t^2 < \infty$  [Robbins and Monro, 1951]
  - $\nu_t$  is unbiased to  $\mathbf{V}^\pi(\mathbf{S}_t)$

# Monte-Carlo with Value Function Approximation

---

- Return  $G_t$  is an **unbiased**, noisy sample of true value  $V^\pi(S_t)$
- Can therefore apply **supervised learning** to “training data”

$$\langle S_1, G_1 \rangle, \langle S_2, G_2 \rangle, \dots, \langle S_T, G_T \rangle$$

- Applying **gradient-descent methods**

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( G_t - \hat{V}(S_t; \omega_t) \right) \nabla_{\omega} \hat{V}(S_t; \omega_t) \right]$$

- Monte Carlo evaluation **converges** to a local minimum

# TD Learning with Value Function Approximation

---

- The TD-target  $R_t + \gamma \hat{V}(\mathbf{S}_{t+1}; \omega)$  is a **biased** sample of true value  $V^\pi(\mathbf{S}_t)$
- Can still apply **supervised learning** to “training data”

$$\langle \mathbf{S}_1, R_1 + \gamma \hat{V}(\mathbf{S}_2; \omega) \rangle, \langle \mathbf{S}_2, R_2 + \gamma \hat{V}(\mathbf{S}_3; \omega) \rangle, \dots, \langle \mathbf{S}_T, R_T + \gamma \hat{V}(\mathbf{S}_{T+1}; \omega) \rangle$$

- Applying **gradient-descent methods**

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( R_t + \gamma \hat{V}(\mathbf{S}_{t+1}; \omega_t) - \hat{V}(\mathbf{S}_t; \omega_t) \right) \nabla_\omega \hat{V}(\mathbf{S}_t; \omega_t) \right]$$

- Linear TD(0) **converges** to a global minimum [Tsitsiklis and Van Roy, 1997]



# TD( $\lambda$ ) with Value Function Approximation

---

- The  $\lambda$ -return  $G_t^\lambda$  is a **biased** sample of true value  $V^\pi(S_t)$
- Can again apply **supervised learning** to “training data”

$$\langle S_1, G_1^\lambda \rangle, \langle S_2, G_2^\lambda \rangle, \dots, \langle S_T, G_T^\lambda \rangle$$

- Applying **gradient-descent methods**

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( G_t^\lambda - \hat{V}(S_t; \omega_t) \right) \nabla_{\omega} \hat{V}(S_t; \omega_t) \right]$$

- Linear TD( $\lambda$ ) **converges** to a global minimum [Tsitsiklis and Van Roy, 1997]

# Linear Function Approximation

---

- Linear features:  $\phi(\mathbf{s})$  (e.g., polynomials, trigonometric polynomials, B-splines)
- MC update rule:

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( G_t - \phi^\top(\mathbf{S}_t)\omega_t \right) \phi(\mathbf{S}_t) \right]$$

- TD(0) update rule:

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( R_t + \gamma \phi^\top(\mathbf{S}_{t+1})\omega_t - \phi^\top(\mathbf{S}_t)\omega_t \right) \phi(\mathbf{S}_t) \right]$$

- TD( $\lambda$ ) update rule:

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( G_t^\lambda - \phi^\top(\mathbf{S}_t)\omega_t \right) \phi(\mathbf{S}_t) \right]$$

# Nonlinear Function Approximation

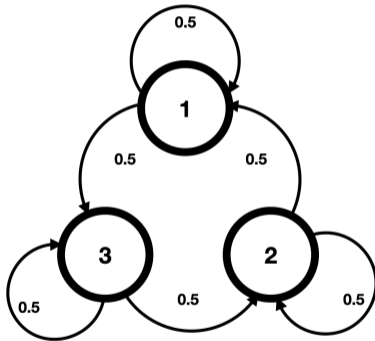
---

- Unlike linear methods, gradient-based TD learning algorithm with nonlinear approximation may diverge
- Proof by an example [Tsitsiklis and Van Roy, 1997]

# Bad Example

---

- Markov chain with state space  $\{1, 2, 3\}$  and transition matrix



- All instantaneous rewards equal to zero
- The value function is zero for any state and policy

## Bad Example (Cont'd)

---

- Nonlinear approximator  $\widehat{\mathbf{V}}(\boldsymbol{\omega}) = (\widehat{\mathbf{V}}(\mathbf{1}; \boldsymbol{\omega}), \widehat{\mathbf{V}}(\mathbf{2}; \boldsymbol{\omega}), \widehat{\mathbf{V}}(\mathbf{3}; \boldsymbol{\omega}))^\top$  for some scalar  $\boldsymbol{\omega}$
- Arbitrary initial value  $\sum_{\mathbf{s}} \widehat{\mathbf{V}}(\mathbf{s}; \boldsymbol{\omega}) = \mathbf{0}$
- Use an ordinary differential equation (ODE) model for parametrization

$$\frac{d\widehat{\mathbf{V}}(\boldsymbol{\omega})}{d\boldsymbol{\omega}} = (\mathbf{Q} + \varepsilon \mathbf{I}) \widehat{\mathbf{V}}(\boldsymbol{\omega})$$

for some small constant  $\varepsilon > \mathbf{0}$  and

$$\mathbf{Q} = \begin{pmatrix} 1 & 1/2 & 3/2 \\ 3/2 & 1 & 1/2 \\ 1/2 & 3/2 & 1 \end{pmatrix}$$

- If the RHS of ODE does not involve  $\widehat{\mathbf{V}}$ , reduces to the linear model

## Bad Example (Cont'd)

---

- Recall that the value function equals zero
- The mean squared error objective function

$$J(\omega) = \sum_{s=1}^3 \left[ \widehat{V}(s; \omega) - V^{\pi}(s) \right]^2 = \sum_{s=1}^3 \widehat{V}^2(s; \omega)$$

- It can be shown that under TD(0) update [Tsitsiklis and Van Roy, 1997]
  - $\omega_t$  increases with  $t$
  - $\sum_{s=1}^3 \widehat{V}^2(s; \omega)$  increases with  $\omega$
  - $J(\omega_t)$  diverges with  $t$

# Convergence of Prediction Algorithms

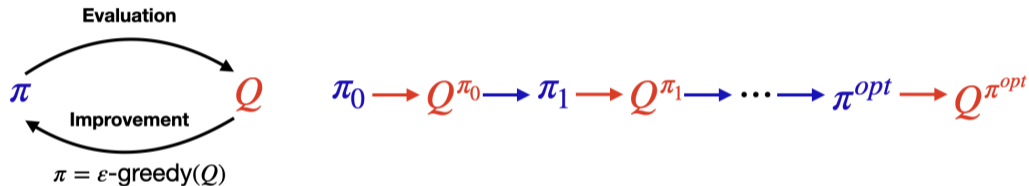
---

Algorithm	Tabular	Linear	Non-linear
MC	✓	✓	✓
TD( $\mathbf{0}$ )	✓	✓	✗
TD( $\lambda$ )	✓	✓	✗

Source: Silver, UCL RL course,  
<https://www.davidsilver.uk/wp-content/uploads/2020/03/FA.pdf>

# Control with Gradient Descent-based Methods

---



- **Policy evaluation**: approximate action-state value function  $Q^\pi = \hat{Q}(\cdot, \cdot; \omega)$
- **Policy improvement**:  $\varepsilon$ -greedy policy improvement



# Action-State Value Function Approximation

---

- Approximate action-state value function

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \hat{Q}(\mathbf{s}, \mathbf{a}; \omega)$$

- Minimise mean-squared error

$$J(\omega) = \frac{1}{2} \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mu} [\hat{Q}(\mathbf{s}, \mathbf{a}; \omega) - Q^\pi(\mathbf{s}, \mathbf{a})]^2$$

- Use stochastic gradient descent

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( \mathbf{q}_t - \hat{Q}(\mathbf{S}_t, \mathbf{A}_t; \omega_t) \right) \nabla_{\omega} \hat{Q}(\mathbf{S}_t, \mathbf{A}_t; \omega_t) \right]$$

for some target  $\mathbf{q}_t$ , some approx of  $Q^\pi(\mathbf{S}_t, \mathbf{A}_t)$

# Value Function Approximation (Cont'd)

---

- For MC, the target is the return  $G_t$

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( G_t - \hat{Q}(S_t, A_t; \omega_t) \right) \nabla_{\omega} \hat{Q}(S_t, A_t; \omega_t) \right]$$

- For TD(0) (SARSA), the target is  $R_t + \gamma \hat{Q}(S_{t+1}, A_{t+1}; \omega)$

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( R_t + \gamma \hat{Q}(S_{t+1}, A_{t+1}; \omega_t) - \hat{Q}(S_t, A_t; \omega_t) \right) \nabla_{\omega} \hat{Q}(S_t, A_t; \omega_t) \right]$$

- For TD( $\lambda$ ) (SARSA( $\lambda$ )), the target is  $Q_t^\lambda$

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( Q_t^\lambda - \hat{Q}(S_t, A_t; \omega_t) \right) \nabla_{\omega} \hat{Q}(S_t, A_t; \omega_t) \right]$$

# Linear Function Approximation (Cont'd)

---

- For MC, the target is the return  $G_t$

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( G_t - \phi^\top(\mathbf{S}_t, \mathbf{A}_t) \omega_t \right) \phi(\mathbf{S}_t, \mathbf{A}_t) \right]$$

- For TD(0) (SARSA), the target is  $R_t + \gamma \phi^\top(\mathbf{S}_{t+1}, \mathbf{A}_{t+1}) \omega$

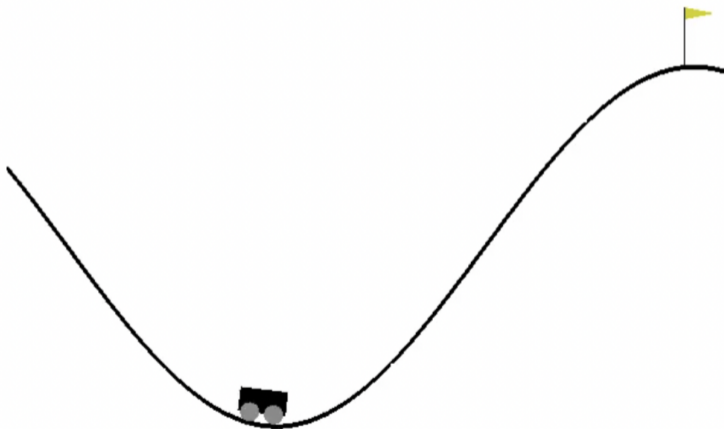
$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( R_t + \gamma \phi^\top(\mathbf{S}_{t+1}, \mathbf{A}_{t+1}) \omega_t - \phi^\top(\mathbf{S}_t, \mathbf{A}_t) \omega_t \right) \phi(\mathbf{S}_t, \mathbf{A}_t) \right]$$

- For TD( $\lambda$ ) (SARSA( $\lambda$ )), the target is  $Q_t^\lambda$

$$\omega_{t+1} = \omega_t - \alpha_t \left[ \left( Q_t^\lambda - \phi^\top(\mathbf{S}_t, \mathbf{A}_t) \omega_t \right) \phi(\mathbf{S}_t, \mathbf{A}_t) \right]$$

# The Mountain Car Example

---



# Convergence of Control Algorithms

---

Algorithm	Tabular	Linear	Non-linear
MC	✓	✓	✓
SARSA	✓	✓	✗
Q-Learning	✓	✗	✗

Source: Silver, UCL RL course,  
<https://www.davidsilver.uk/wp-content/uploads/2020/03/FA.pdf>

# Lecture Outline

---

1. Introduction to Value Function Approximation
2. Gradient Descent-based Methods
- 3. Fitted Q-Iteration**

# Limitations of Gradient-based Control Methods

---

- MC control allows both linear and nonlinear approximation, but is **inefficient**
  - $G_t$  suffers from **large variance**, so is the estimated Q-function
- SARSA is efficient, but cannot allow **nonlinear** approximation
  - When using linear approximation, the estimator can suffer from **large bias**

# Efficient Control with Function Approximation

---

- **Batch** (offline) setting with pre-collected data  $\{S_t, A_t, R_t, S_{t+1}\}_t$
- **Main idea** of Q-learning: learn the optimal Q-function  $Q^{\pi^{\text{opt}}}$  based on the Bellman **optimality** equation and derive the optimal policy

$$\pi^{\text{opt}}(s) = \arg \max_a Q^{\pi^{\text{opt}}}(s, a)$$

- Bellman optimality equation

$$Q^{\pi^{\text{opt}}}(S_t, A_t) = \mathbb{E} \left[ R_t + \gamma \max_a Q^{\pi^{\text{opt}}}(S_{t+1}, a) \mid S_t, A_t \right]$$

- Supervised learning is sample efficient in **batch** settings
- Use supervised learning to learn  $Q^{\pi^{\text{opt}}}$  by solving Bellman optimality equation



# Challenge

---

- Bellman optimality equation

$$Q^{\pi^{\text{opt}}}(S_t, A_t) = \mathbb{E} \left[ R_t + \gamma \max_a Q^{\pi^{\text{opt}}}(S_{t+1}, a) \mid S_t, A_t \right]$$

- Both LHS and RHS involve  $Q^{\pi^{\text{opt}}}$
- A naive approach: minimize the **mean squared Bellman error**

$$\sum_t \left[ R_t + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]^2$$

This would yield a **biased** estimator!

# The Bellman Error is Not Learnable

---

- For a given random variable  $Z$ ,  $\mathbb{E}Z^2 = (\mathbb{E}Z)^2 + \text{Var}(Z)$
- The mean squared Bellman error can be decomposed into **squared bias** + **variance**

$$\begin{aligned} & \mathbb{E} \left[ \left\{ R_t + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right\}^2 \middle| A_t, S_t \right] \\ &= \left[ \mathbb{E} \left\{ R_t + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \middle| A_t, S_t \right\} \right]^2 \\ &+ \text{Var} \left[ R_t + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \middle| A_t, S_t \right] \end{aligned}$$

- The second line is **zero** when  $Q = Q^{\pi^{\text{opt}}}$
- The third line is **nonzero** for any  $Q$  and is a function of  $Q$  as well
- There is no guarantee  $Q^{\pi^{\text{opt}}}$  is the minimizer

# Fitted Q-Iteration [Riedmiller, 2005]

---

- Bellman optimality equation

$$Q^{\pi^{\text{opt}}}(S_t, A_t) = \mathbb{E} \left[ R_t + \gamma \max_a Q^{\pi^{\text{opt}}}(S_{t+1}, a) \mid S_t, A_t \right]$$

Both LHS and RHS involve  $Q^{\pi^{\text{opt}}}$

- **Main idea:** Fix  $Q^{\pi^{\text{opt}}}$  on the RHS
- **Repeat** the following
  1. Compute  $\hat{Q}$  as the argmin of

$$\arg \min_Q \sum_t \left[ R_t + \gamma \max_a \tilde{Q}(S_{t+1}, a) - Q(S_t, A_t) \right]^2$$

2. Set  $\tilde{Q} = \hat{Q}$

## Fitted Q-Iteration (Cont'd)

---

- During each iteration, consider the objective function

$$\begin{aligned} & \mathbb{E} \left[ R_t + \gamma \max_a \tilde{Q}(S_{t+1}, a) - Q(S_t, A_t) \mid A_t, S_t \right]^2 \\ &= \left[ \mathbb{E} \left\{ R_t + \gamma \max_a \tilde{Q}(S_{t+1}, a) - Q(S_t, A_t) \mid A_t, S_t \right\} \right]^2 \\ &+ \text{Var} \left[ R_t + \gamma \max_a \tilde{Q}(S_{t+1}, a) - Q(S_t, A_t) \mid A_t, S_t \right] \end{aligned}$$

- When  $\tilde{Q}$  is close to  $Q^{\pi^{\text{opt}}}$ , the second line is **small** when  $Q = Q^{\pi^{\text{opt}}}$
- The third line is the same for any  $Q$ , since  $\tilde{Q}$  is fixed

# Fitted Q-Iteration: Algorithm

---

- **Initialization:**  $\hat{Q}, \tilde{Q}$  arbitrary,  $k = 0$

- **While ( $k < K$ ) Repeat**

Generated data  $\{(S_t, A_t, R_t, S_{t+1})\}$  using policy derived from  $\hat{Q}$  (e.g.,  $\epsilon$ -greedy)

Compute  $\hat{Q}$  as the argmin of

$$\arg \min_{\hat{Q}} \sum_t \left[ R_t + \gamma \max_a \tilde{Q}(S_{t+1}, a) - \hat{Q}(S_t, A_t) \right]^2$$

Set  $\tilde{Q} = \hat{Q}$

# Advantages of Fitted Q-Iteration

---

- **Flexibility:** any supervised learning method (e.g., deep learning, boosting, random forest) is applicable to learn the Q-function during each iteration.
  - Gradient Descent-based methods require the Q-function model to be a smooth function of the model parameters
- **Efficiency:** borrows the strength of supervised learning for sample-efficient estimation. Allows high-dimensional state information.

# Theoretical Analysis of Fitted Q-Iteration

---

- Let  $\hat{Q}_k$  denote the Q-estimator during the  $k$ th iteration
- Error decomposition: **bias due to initialization** + **stochastic estimation error**
- The initialization bias  $\rightarrow \mathbf{0}$  as  $k \rightarrow \infty$
- The estimation error  $\rightarrow \mathbf{0}$  when supervised learning provides a **consistent** estimator at each iteration

# Theoretical Analysis of Fitted Q-Iteration (Cont'd)

---

- At the  $k$ th iteration,

$$\hat{Q}_k = \arg \min_Q \sum_t \left[ R_t + \gamma \max_a \hat{Q}_{k-1}(S_{t+1}, a) - Q(S_t, A_t) \right]^2$$

- Supervised learning target:

$$Q_k(s, a) = \mathbb{E} \left[ R_t + \gamma \max_a \hat{Q}_{k-1}(S_{t+1}, a) \mid S_t = s, A_t = a \right]$$



# Theoretical Analysis of Fitted Q-Iteration (Cont'd)

- A key inequality

$$\begin{aligned} \sup_{\mathbf{s}, \mathbf{a}} |\widehat{Q}_k(\mathbf{s}, \mathbf{a}) - Q^{\pi^{\text{opt}}}(\mathbf{s}, \mathbf{a})| &\leq \sup_{\mathbf{s}, \mathbf{a}} |\widehat{Q}_k(\mathbf{s}, \mathbf{a}) - Q_k(\mathbf{s}, \mathbf{a})| \\ &\quad + \gamma \sup_{\mathbf{s}, \mathbf{a}} |\widehat{Q}_{k-1}(\mathbf{s}, \mathbf{a}) - Q^{\pi^{\text{opt}}}(\mathbf{s}, \mathbf{a})| \end{aligned}$$

- Iteratively applying the inequality

$$\begin{aligned} \sup_{\mathbf{s}, \mathbf{a}} |\widehat{Q}_k(\mathbf{s}, \mathbf{a}) - Q^{\pi^{\text{opt}}}(\mathbf{s}, \mathbf{a})| &\leq \gamma^k \sup_{\mathbf{s}, \mathbf{a}} \underbrace{|\widehat{Q}_0(\mathbf{s}, \mathbf{a}) - Q^{\pi^{\text{opt}}}(\mathbf{s}, \mathbf{a})|}_{\text{Initialization Bias}} \\ &\quad + \sup_{\mathbf{s}, \mathbf{a}} \max_{j=\{1, \dots, k\}} \underbrace{|\widehat{Q}_j(\mathbf{s}, \mathbf{a}) - Q_j(\mathbf{s}, \mathbf{a})|}_{\text{Estimation Error}} \end{aligned}$$

# Summary

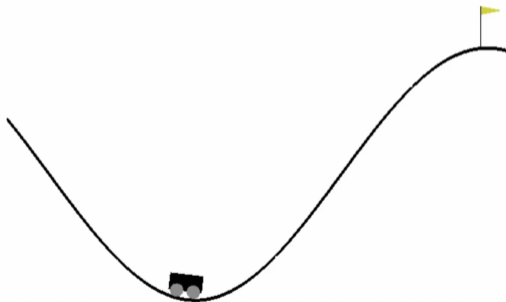
---

- Linear function approximation
- Gradient-based methods
- Gradient-based MC, TD, SARSA
- Neural networks
- Stochastic gradient-based methods
- Fitted Q-iteration

# Seminar Exercises

---

- Solution to HW4 (Deadline, Wed 12:00 PM)
- The mountain car example: gradient-based methods



# References I

---

- Andrew R Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information theory*, 39(3):930–945, 1993.
- Masaaki Imaizumi and Kenji Fukumizu. Deep neural networks learn non-smooth functions effectively. In *The 22nd international conference on artificial intelligence and statistics*, pages 869–878. PMLR, 2019.
- Martin Riedmiller. Neural fitted q iteration—first experiences with a data efficient neural reinforcement learning method. In *European conference on machine learning*, pages 317–328. Springer, 2005.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- John N Tsitsiklis and Benjamin Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE transactions on automatic control*, 42(5):674–690, 1997.

# Questions