# Lecture 8: Value Function Methods

Shiyu Zhao

Department of Artificial Intelligence Westlake University



### 1 Motivating examples: from table to function

- 2 Algorithm for state value estimation
  - Objective function
  - Optimization algorithms
  - Selection of function approximators
  - Illustrative examples
  - Summary of the story
  - Theoretical analysis (optional)
- 3 Sarsa with function approximation
- 4 Q-learning with function approximation
- 5 Deep Q-learning

### 6 Summary

### 1 Motivating examples: from table to function

### 2 Algorithm for state value estimation

- Objective function
- Optimization algorithms
- Selection of function approximators
- Illustrative examples
- Summary of the story
- Theoretical analysis (optional)
- 3 Sarsa with function approximation
- 4 Q-learning with function approximation
- 5 Deep Q-learning

#### 6 Summary

#### So far in this book, state and action values are represented by tables.

#### For example, state value:

- Advantage: intuitive and easy to analyze
- Disadvantage: difficult to handle large or continuous state or action spaces. Two aspects: 1) storage; 2) generalization ability

So far in this book, state and action values are represented by tables.

• For example, state value:

State	$s_1$	$s_2$	 $s_n$
Value	$v_{\pi}(s_1)$	$v_{\pi}(s_2)$	 $v_{\pi}(s_n)$

- Advantage: intuitive and easy to analyze
- Disadvantage: difficult to handle large or continuous state or action spaces. Two aspects: 1) storage; 2) generalization ability

So far in this book, state and action values are represented by tables.

• For example, state value:

State	$s_1$	$s_2$	 $s_n$
Value	$v_{\pi}(s_1)$	$v_{\pi}(s_2)$	 $v_{\pi}(s_n)$

	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$
$s_1$	$q_{\pi}(s_1, a_1)$	$q_{\pi}(s_1, a_2)$	$q_{\pi}(s_1, a_3)$	$q_{\pi}(s_1, a_4)$	$q_{\pi}(s_1, a_5)$
:		•		•	
<i>s</i> 9	$q_{\pi}(s_9, a_1)$	$q_{\pi}(s_9, a_2)$	$q_{\pi}(s_9, a_3)$	$q_{\pi}(s_9, a_4)$	$q_{\pi}(s_9, a_5)$

- Advantage: intuitive and easy to analyze
- Disadvantage: difficult to handle large or continuous state or action spaces. Two aspects: 1) storage; 2) generalization ability

So far in this book, state and action values are represented by tables.

• For example, state value:

State	$s_1$	$s_2$	 $s_n$
Value	$v_{\pi}(s_1)$	$v_{\pi}(s_2)$	 $v_{\pi}(s_n)$

• For example, action value:

-	$a_1$	$a_2$	<i>a</i> 3	$a_4$	$a_5$
$s_1$	$q_{\pi}(s_1, a_1)$	$q_{\pi}(s_1, a_2)$	$q_{\pi}(s_1, a_3)$	$q_{\pi}(s_1, a_4)$	$q_{\pi}(s_1, a_5)$
÷				•	-
<i>s</i> 9	$q_{\pi}(s_9, a_1)$	$q_{\pi}(s_9, a_2)$	$q_{\pi}(s_9, a_3)$	$q_{\pi}(s_9, a_4)$	$q_{\pi}(s_9, a_5)$

#### • Advantage: intuitive and easy to analyze

 Disadvantage: difficult to handle large or continuous state or action spaces. Two aspects: 1) storage; 2) generalization ability

So far in this book, state and action values are represented by tables.

• For example, state value:

State	$s_1$	$s_2$	 $s_n$
Value	$v_{\pi}(s_1)$	$v_{\pi}(s_2)$	 $v_{\pi}(s_n)$

	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$
$s_1$	$q_{\pi}(s_1, a_1)$	$q_{\pi}(s_1, a_2)$	$q_{\pi}(s_1, a_3)$	$q_{\pi}(s_1, a_4)$	$q_{\pi}(s_1, a_5)$
:		•		•	
<i>s</i> 9	$q_{\pi}(s_9, a_1)$	$q_{\pi}(s_9, a_2)$	$q_{\pi}(s_9, a_3)$	$q_{\pi}(s_9, a_4)$	$q_{\pi}(s_9, a_5)$

- Advantage: intuitive and easy to analyze
- Disadvantage: difficult to handle large or continuous state or action spaces. Two aspects: 1) storage; 2) generalization ability

Consider an example:

- There are n states:  $s_1, \ldots, s_n$ .
- The state values are  $v_{\pi}(s_1), \ldots, v_{\pi}(s_n)$ , where  $\pi$  is a given policy.
- n is very large!
- We hope to use a simple curve to approximate these values.

For example, we can use a simple straight line to fit the dots.



Suppose the equation of the straight line is

$$\hat{v}(s,w) = as + b = \underbrace{[s,1]}_{\phi^T(s)} \underbrace{\begin{bmatrix} a \\ b \end{bmatrix}}_{w} = \phi^T(s)w$$

For example, we can use a simple straight line to fit the dots.



Suppose the equation of the straight line is

$$\hat{v}(s,w) = as + b = \underbrace{[s,1]}_{\phi^T(s)} \underbrace{\begin{bmatrix} a \\ b \end{bmatrix}}_{w} = \phi^T(s)w$$

For example, we can use a simple straight line to fit the dots.



Suppose the equation of the straight line is

$$\hat{v}(s,w) = as + b = \underbrace{[s,1]}_{\phi^{T}(s)} \underbrace{\begin{bmatrix} a \\ b \end{bmatrix}}_{w} = \phi^{T}(s)w$$

For example, we can use a simple straight line to fit the dots.



Suppose the equation of the straight line is

$$\hat{v}(s,w) = as + b = \underbrace{[s,1]}_{\phi^{T}(s)} \underbrace{\begin{bmatrix} a \\ b \end{bmatrix}}_{w} = \phi^{T}(s)w$$

### Difference 1: How to retrieve the value of a state

- When the values are represented by a table, we can directly read the value in the table.
- When the values are represented by a function, we need to input the state index *s* into the function and calculate the function value.

For example,  $s \to \phi(s) \to \phi^T(s)w = \hat{v}(s,w)$ 

### Difference 1: How to retrieve the value of a state

- When the values are represented by a table, we can directly read the value in the table.
- When the values are represented by a function, we need to input the state index *s* into the function and calculate the function value.

For example,  $s \to \phi(s) \to \phi^T(s)w = \hat{v}(s,w)$ 

### Difference 1: How to retrieve the value of a state

- When the values are represented by a table, we can directly read the value in the table.
- When the values are represented by a function, we need to input the state index *s* into the function and calculate the function value.

For example,  $s \to \phi(s) \to \phi^T(s) w = \hat{v}(s, w)$ 

### Difference 1: How to retrieve the value of a state

- When the values are represented by a table, we can directly read the value in the table.
- When the values are represented by a function, we need to input the state index *s* into the function and calculate the function value.



For example,  $s \to \phi(s) \to \phi^T(s) w = \hat{v}(s,w)$ 

### Difference 1: How to retrieve the value of a state

- When the values are represented by a table, we can directly read the value in the table.
- When the values are represented by a function, we need to input the state index *s* into the function and calculate the function value.



For example,  $s \rightarrow \phi(s) \rightarrow \phi^T(s) w = \hat{v}(s,w)$ 

### Difference 1: How to retrieve the value of a state

- When the values are represented by a table, we can directly read the value in the table.
- When the values are represented by a function, we need to input the state index *s* into the function and calculate the function value.



For example,  $s \to \phi(s) \to \phi^T(s) w = \hat{v}(s, w)$ 

- When the values are represented by a table, we can directly rewrite the value in the table.
- When the values are represented by a function, we must update w to change the values indirectly.
  - How to update w will be addressed in detail later.

- When the values are represented by a table, we can directly rewrite the value in the table.
- When the values are represented by a function, we must update w to change the values indirectly.
  - How to update w will be addressed in detail later.

- When the values are represented by a table, we can directly rewrite the value in the table.
- When the values are represented by a function, we must update w to change the values indirectly.
  - How to update w will be addressed in detail later.

- When the values are represented by a table, we can directly rewrite the value in the table.
- When the values are represented by a function, we must update w to change the values indirectly.
  - How to update  $\boldsymbol{w}$  will be addressed in detail later.

Difference between the tabular and function methods:





(b) Function method

**Benefit: generalization ability.** When we update  $\hat{v}(s, w)$  by changing w, the values of the neighboring states are also changed.

Difference between the tabular and function methods:





(b) Function method

**Benefit: generalization ability.** When we update  $\hat{v}(s, w)$  by changing w, the values of the neighboring states are also changed.

We can fit the points more precisely using high-order curves:

$$\hat{v}(s,w) = as^2 + bs + c = \underbrace{[s^2, s, 1]}_{\phi^T(s)} \underbrace{\begin{bmatrix} a \\ b \\ c \end{bmatrix}}_{w} = \phi^T(s)w.$$

- The dimensions of w and  $\phi(s)$  increase; the values may be fitted more accurately.
- Although  $\hat{v}(s, w)$  is nonlinear in s, it is linear in w. The nonlinearity is contained in  $\phi(s)$ .

We can fit the points more precisely using high-order curves:

$$\hat{v}(s,w) = as^2 + bs + c = \underbrace{[s^2, s, 1]}_{\phi^T(s)} \underbrace{\begin{bmatrix} a \\ b \\ c \end{bmatrix}}_{w} = \phi^T(s)w.$$

- The dimensions of w and  $\phi(s)$  increase; the values may be fitted more accurately.
- Although  $\hat{v}(s, w)$  is nonlinear in s, it is linear in w. The nonlinearity is contained in  $\phi(s)$ .

We can fit the points more precisely using high-order curves:

$$\hat{v}(s,w) = as^2 + bs + c = \underbrace{[s^2, s, 1]}_{\phi^T(s)} \underbrace{\begin{bmatrix} a \\ b \\ c \end{bmatrix}}_{w} = \phi^T(s)w.$$

- The dimensions of w and  $\phi(s)$  increase; the values may be fitted more accurately.
- Although  $\hat{v}(s, w)$  is nonlinear in s, it is linear in w. The nonlinearity is contained in  $\phi(s)$ .

We can fit the points more precisely using high-order curves:

$$\hat{v}(s,w) = as^2 + bs + c = \underbrace{[s^2, s, 1]}_{\phi^T(s)} \underbrace{\begin{bmatrix} a \\ b \\ c \end{bmatrix}}_{w} = \phi^T(s)w.$$

- The dimensions of w and  $\phi(s)$  increase; the values may be fitted more accurately.
- Although  $\hat{v}(s, w)$  is nonlinear in s, it is linear in w. The nonlinearity is contained in  $\phi(s)$ .

- Idea: Approximate the state and action values using parameterized functions: v̂(s, w) ≈ v<sub>π</sub>(s) where w ∈ ℝ<sup>m</sup> is the parameter vector.
- Key difference: How to retrieve and change the value of v(s)
- Advantages:
  - 1) Storage: The dimension of w may be much smaller than |S|.
  - 2) Generalization: When a state s is visited, the parameter w is updated so that the values of some other unvisited states can also be updated.

- Idea: Approximate the state and action values using parameterized functions: v̂(s, w) ≈ v<sub>π</sub>(s) where w ∈ ℝ<sup>m</sup> is the parameter vector.
- Key difference: How to retrieve and change the value of  $\boldsymbol{v}(\boldsymbol{s})$
- Advantages:
  - 1) Storage: The dimension of w may be much smaller than |S|.
  - 2) Generalization: When a state s is visited, the parameter w is updated so that the values of some other unvisited states can also be updated.

- Idea: Approximate the state and action values using parameterized functions: v̂(s, w) ≈ v<sub>π</sub>(s) where w ∈ ℝ<sup>m</sup> is the parameter vector.
- Key difference: How to retrieve and change the value of v(s)
- Advantages:
  - 1) Storage: The dimension of w may be much smaller than |S|.
  - 2) Generalization: When a state s is visited, the parameter w is updated so that the values of some other unvisited states can also be updated.

- Idea: Approximate the state and action values using parameterized functions: v̂(s, w) ≈ v<sub>π</sub>(s) where w ∈ ℝ<sup>m</sup> is the parameter vector.
- Key difference: How to retrieve and change the value of v(s)
- Advantages:
  - 1) Storage: The dimension of w may be much smaller than |S|.
  - 2) Generalization: When a state s is visited, the parameter w is updated so that the values of some other unvisited states can also be updated.

#### 1 Motivating examples: from table to function

### 2 Algorithm for state value estimation

- Objective function
- Optimization algorithms
- Selection of function approximators
- Illustrative examples
- Summary of the story
- Theoretical analysis (optional)
- 3 Sarsa with function approximation
- 4 Q-learning with function approximation
- 5 Deep Q-learning

#### 6 Summary

#### 1 Motivating examples: from table to function

### 2 Algorithm for state value estimation

- Objective function
- Optimization algorithms
- Selection of function approximators
- Illustrative examples
- Summary of the story
- Theoretical analysis (optional)
- 3 Sarsa with function approximation
- 4 Q-learning with function approximation
- 5 Deep Q-learning

#### 6 Summary
- Let  $v_{\pi}(s)$  and  $\hat{v}(s, w)$  be the true state value and the estimated state value, respectively.
- Our goal is to find an optimal w so that  $\hat{v}(s, w)$  can best approximate  $v_{\pi}(s)$  for every s.
- This is a policy evaluation problem. Later we will extend to policy improvement.

- The first step is to define an objective function.
- The second step is to derive algorithms for optimizing the objective function.

- Let  $v_{\pi}(s)$  and  $\hat{v}(s, w)$  be the true state value and the estimated state value, respectively.
- Our goal is to find an optimal w so that  $\hat{v}(s, w)$  can best approximate  $v_{\pi}(s)$  for every s.
- This is a policy evaluation problem. Later we will extend to policy improvement.

- The first step is to define an objective function.
- The second step is to derive algorithms for optimizing the objective function.

- Let  $v_{\pi}(s)$  and  $\hat{v}(s, w)$  be the true state value and the estimated state value, respectively.
- Our goal is to find an optimal w so that  $\hat{v}(s, w)$  can best approximate  $v_{\pi}(s)$  for every s.
- This is a policy evaluation problem. Later we will extend to policy improvement.

- The first step is to define an objective function.
- The second step is to derive algorithms for optimizing the objective function.

- Let  $v_{\pi}(s)$  and  $\hat{v}(s, w)$  be the true state value and the estimated state value, respectively.
- Our goal is to find an optimal w so that  $\hat{v}(s, w)$  can best approximate  $v_{\pi}(s)$  for every s.
- This is a policy evaluation problem. Later we will extend to policy improvement.

- The first step is to define an objective function.
- The second step is to derive algorithms for optimizing the objective function.

- Let  $v_{\pi}(s)$  and  $\hat{v}(s, w)$  be the true state value and the estimated state value, respectively.
- Our goal is to find an optimal w so that  $\hat{v}(s, w)$  can best approximate  $v_{\pi}(s)$  for every s.
- This is a policy evaluation problem. Later we will extend to policy improvement.

- The first step is to define an objective function.
- The second step is to derive algorithms for optimizing the objective function.

- Let  $v_{\pi}(s)$  and  $\hat{v}(s, w)$  be the true state value and the estimated state value, respectively.
- Our goal is to find an optimal w so that  $\hat{v}(s, w)$  can best approximate  $v_{\pi}(s)$  for every s.
- This is a policy evaluation problem. Later we will extend to policy improvement.

- The first step is to define an objective function.
- The second step is to derive algorithms for optimizing the objective function.

- Let  $v_{\pi}(s)$  and  $\hat{v}(s, w)$  be the true state value and the estimated state value, respectively.
- Our goal is to find an optimal w so that  $\hat{v}(s, w)$  can best approximate  $v_{\pi}(s)$  for every s.
- This is a policy evaluation problem. Later we will extend to policy improvement.

- The first step is to define an objective function.
- The second step is to derive algorithms for optimizing the objective function.

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2].$$

- Our goal is to find the best w that can minimize J(w).
- The expectation is with respect to the random variable  $S \in S$ . What is the probability distribution of S?
  - This is new. We have not discussed the probability distribution of states so far.
  - There are several ways to define the probability distribution of  $\boldsymbol{S}.$

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2].$$

- Our goal is to find the best w that can minimize J(w).
- The expectation is with respect to the random variable  $S \in S$ . What is the probability distribution of S?
  - This is new. We have not discussed the probability distribution of states so far.
  - There are several ways to define the probability distribution of  $\boldsymbol{S}.$

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2].$$

- Our goal is to find the best w that can minimize J(w).
- The expectation is with respect to the random variable  $S \in S$ . What is the probability distribution of S?
  - This is new. We have not discussed the probability distribution of states so far.
  - There are several ways to define the probability distribution of  $\boldsymbol{S}.$

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2].$$

- Our goal is to find the best w that can minimize J(w).
- The expectation is with respect to the random variable  $S \in S$ . What is the probability distribution of S?
  - This is new. We have not discussed the probability distribution of states so far.
  - There are several ways to define the probability distribution of  $\boldsymbol{S}.$

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2].$$

- Our goal is to find the best w that can minimize J(w).
- The expectation is with respect to the random variable  $S \in S$ . What is the probability distribution of S?
  - This is new. We have not discussed the probability distribution of states so far.
  - There are several ways to define the probability distribution of  $\boldsymbol{S}.$

- That is to treat all the states to be equally important by setting the probability of each state as 1/|S|.
- In this case, the objective function becomes

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^{2}] = \frac{1}{|S|} \sum_{s \in S} (v_{\pi}(s) - \hat{v}(s, w))^{2}.$$

- Drawback:
  - The states may not be equally important. For example, some states may be rarely visited by a policy. Hence, this way does not consider the real dynamics of the Markov process under the given policy.

- That is to treat all the states to be equally important by setting the probability of each state as  $1/|\mathcal{S}|$ .
- In this case, the objective function becomes

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^{2}] = \frac{1}{|S|} \sum_{s \in S} (v_{\pi}(s) - \hat{v}(s, w))^{2}.$$

- Drawback:
  - The states may not be equally important. For example, some states may be rarely visited by a policy. Hence, this way does not consider the real dynamics of the Markov process under the given policy.

- That is to treat all the states to be equally important by setting the probability of each state as 1/|S|.
- In this case, the objective function becomes

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^{2}] = \frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} (v_{\pi}(s) - \hat{v}(s, w))^{2}.$$

- Drawback:
  - The states may not be equally important. For example, some states may be rarely visited by a policy. Hence, this way does not consider the real dynamics of the Markov process under the given policy.

- That is to treat all the states to be equally important by setting the probability of each state as 1/|S|.
- In this case, the objective function becomes

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^{2}] = \frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} (v_{\pi}(s) - \hat{v}(s, w))^{2}.$$

- Drawback:
  - The states may not be equally important. For example, some states may be rarely visited by a policy. Hence, this way does not consider the real dynamics of the Markov process under the given policy.

- Stationary distribution is an important concept that will be frequently used in this course. It describes the long-run behavior of a Markov process.
- Let  $\{d_{\pi}(s)\}_{s \in S}$  denote the stationary distribution of the Markov process under policy  $\pi$ . By definition,  $d_{\pi}(s) \ge 0$  and  $\sum_{s \in S} d_{\pi}(s) = 1$ .
- The objective function can be rewritten as

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2] = \sum_{s \in S} d_{\pi}(s)(v_{\pi}(s) - \hat{v}(s, w))^2.$$

This function is a weighted squared error.

- Stationary distribution is an important concept that will be frequently used in this course. It describes the long-run behavior of a Markov process.
- Let {d<sub>π</sub>(s)}<sub>s∈S</sub> denote the stationary distribution of the Markov process under policy π. By definition, d<sub>π</sub>(s) ≥ 0 and Σ<sub>s∈S</sub> d<sub>π</sub>(s) = 1.
- The objective function can be rewritten as

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2] = \sum_{s \in S} d_{\pi}(s)(v_{\pi}(s) - \hat{v}(s, w))^2.$$

This function is a weighted squared error.

- Stationary distribution is an important concept that will be frequently used in this course. It describes the long-run behavior of a Markov process.
- Let {d<sub>π</sub>(s)}<sub>s∈S</sub> denote the stationary distribution of the Markov process under policy π. By definition, d<sub>π</sub>(s) ≥ 0 and ∑<sub>s∈S</sub> d<sub>π</sub>(s) = 1.
- The objective function can be rewritten as

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2] = \sum_{s \in S} d_{\pi}(s)(v_{\pi}(s) - \hat{v}(s, w))^2.$$

This function is a weighted squared error.

- Stationary distribution is an important concept that will be frequently used in this course. It describes the long-run behavior of a Markov process.
- Let {d<sub>π</sub>(s)}<sub>s∈S</sub> denote the stationary distribution of the Markov process under policy π. By definition, d<sub>π</sub>(s) ≥ 0 and ∑<sub>s∈S</sub> d<sub>π</sub>(s) = 1.
- The objective function can be rewritten as

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^{2}] = \sum_{s \in S} d_{\pi}(s)(v_{\pi}(s) - \hat{v}(s, w))^{2}.$$

This function is a weighted squared error.

- Stationary distribution is an important concept that will be frequently used in this course. It describes the long-run behavior of a Markov process.
- Let {d<sub>π</sub>(s)}<sub>s∈S</sub> denote the stationary distribution of the Markov process under policy π. By definition, d<sub>π</sub>(s) ≥ 0 and ∑<sub>s∈S</sub> d<sub>π</sub>(s) = 1.
- The objective function can be rewritten as

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^{2}] = \sum_{s \in S} d_{\pi}(s)(v_{\pi}(s) - \hat{v}(s, w))^{2}.$$

This function is a weighted squared error.

- Distribution: Distribution of the state
- Stationary: Long-run behavior
- *Summary*: after the agent runs a long time following a policy, the probability that the agent is at any state can be described by this distribution.

- Stationary distribution is also called steady-state distribution, or limiting distribution.
- It is critical to understand the value function method.
- It is also important for the policy gradient method in the next lecture.

- Distribution: Distribution of the state
- Stationary: Long-run behavior
- *Summary*: after the agent runs a long time following a policy, the probability that the agent is at any state can be described by this distribution.

- Stationary distribution is also called steady-state distribution, or limiting distribution.
- It is critical to understand the value function method.
- It is also important for the policy gradient method in the next lecture.

- Distribution: Distribution of the state
- Stationary: Long-run behavior
- *Summary*: after the agent runs a long time following a policy, the probability that the agent is at any state can be described by this distribution.

- Stationary distribution is also called steady-state distribution, or limiting distribution.
- It is critical to understand the value function method.
- It is also important for the policy gradient method in the next lecture.

- Distribution: Distribution of the state
- Stationary: Long-run behavior
- *Summary*: after the agent runs a long time following a policy, the probability that the agent is at any state can be described by this distribution.

- Stationary distribution is also called steady-state distribution, or limiting distribution.
- It is critical to understand the value function method.
- It is also important for the policy gradient method in the next lecture.

# Objective function - Stationary distribution

Illustrative example:

- Given a policy shown in the figure.
- Let n<sub>π</sub>(s) denote the number of times that s has been visited in a very long episode generated by π.
- Then,  $d_{\pi}(s)$  can be approximated by

$$d_{\pi}(s) \approx \frac{n_{\pi}(s)}{\sum_{s' \in \mathcal{S}} n_{\pi}(s')}$$



# Objective function - Stationary distribution

Illustrative example:

- Given a policy shown in the figure.
- Let n<sub>π</sub>(s) denote the number of times that s has been visited in a very long episode generated by π.
- Then,  $d_{\pi}(s)$  can be approximated by

$$d_{\pi}(s) \approx \frac{n_{\pi}(s)}{\sum_{s' \in \mathcal{S}} n_{\pi}(s')}$$



# Objective function - Stationary distribution

Illustrative example:

- Given a policy shown in the figure.
- Let n<sub>π</sub>(s) denote the number of times that s has been visited in a very long episode generated by π.
- Then,  $d_{\pi}(s)$  can be approximated by

$$d_{\pi}(s) \approx \frac{n_{\pi}(s)}{\sum_{s' \in \mathcal{S}} n_{\pi}(s')}$$



Figure: Long-run behavior of an  $\epsilon$ -greedy policy with  $\epsilon = 0.5$ .

$$d_{\pi}^{T} = d_{\pi}^{T} P_{\pi}$$

For this example, we have  $P_{\pi}$  as

$$P_{\pi} = \begin{bmatrix} 0.3 & 0.1 & 0.6 & 0 \\ 0.1 & 0.3 & 0 & 0.6 \\ 0.1 & 0 & 0.3 & 0.6 \\ 0 & 0.1 & 0.1 & 0.8 \end{bmatrix}.$$

It can be calculated that the left eigenvector for the eigenvalue of one is

$$d_{\pi} = \begin{bmatrix} 0.0345, 0.1084, 0.1330, 0.7241 \end{bmatrix}^{2}$$

$$d_{\pi}^{T} = d_{\pi}^{T} P_{\pi}$$

For this example, we have  $P_{\pi}$  as

$$P_{\pi} = \begin{bmatrix} 0.3 & 0.1 & 0.6 & 0 \\ 0.1 & 0.3 & 0 & 0.6 \\ 0.1 & 0 & 0.3 & 0.6 \\ 0 & 0.1 & 0.1 & 0.8 \end{bmatrix}$$

It can be calculated that the left eigenvector for the eigenvalue of one is

$$d_{\pi} = \begin{bmatrix} 0.0345, 0.1084, 0.1330, 0.7241 \end{bmatrix}$$

$$d_{\pi}^{T} = d_{\pi}^{T} P_{\pi}$$

For this example, we have  $P_{\pi}$  as

$$P_{\pi} = \begin{bmatrix} 0.3 & 0.1 & 0.6 & 0 \\ 0.1 & 0.3 & 0 & 0.6 \\ 0.1 & 0 & 0.3 & 0.6 \\ 0 & 0.1 & 0.1 & 0.8 \end{bmatrix}$$

It can be calculated that the left eigenvector for the eigenvalue of one is

$$d_{\pi} = \left[0.0345, 0.1084, 0.1330, 0.7241\right]^{T}$$

$$d_{\pi}^{T} = d_{\pi}^{T} P_{\pi}$$

For this example, we have  $P_{\pi}$  as

$$P_{\pi} = \begin{bmatrix} 0.3 & 0.1 & 0.6 & 0 \\ 0.1 & 0.3 & 0 & 0.6 \\ 0.1 & 0 & 0.3 & 0.6 \\ 0 & 0.1 & 0.1 & 0.8 \end{bmatrix}$$

It can be calculated that the left eigenvector for the eigenvalue of one is

$$d_{\pi} = \left[0.0345, 0.1084, 0.1330, 0.7241\right]^{T}$$

# Outline

### 1 Motivating examples: from table to function

# 2 Algorithm for state value estimation

Objective function

## Optimization algorithms

- Selection of function approximators
- Illustrative examples
- Summary of the story
- Theoretical analysis (optional)
- 3 Sarsa with function approximation
- 4 Q-learning with function approximation
- 5 Deep Q-learning

#### 6 Summary

While we have the objective function, the next step is to optimize it.

• To minimize the objective function J(w), we can use the gradient-descent algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w J(w_k)$$

The true gradient is

$$\nabla_w J(w) = \nabla_w \mathbb{E}[(v_\pi(S) - \hat{v}(S, w))^2]$$

While we have the objective function, the next step is to optimize it.

• To minimize the objective function J(w), we can use the gradient-descent algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w J(w_k)$$

The true gradient is

$$\nabla_{w} J(w) = \nabla_{w} \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^{2}]$$

While we have the objective function, the next step is to optimize it.

• To minimize the objective function J(w), we can use the gradient-descent algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w J(w_k)$$

The true gradient is

$$\nabla_w J(w) = \nabla_w \mathbb{E}[(v_\pi(S) - \hat{v}(S, w))^2]$$
• To minimize the objective function J(w), we can use the gradient-descent algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w J(w_k)$$

The true gradient is

$$\nabla_w J(w) = \nabla_w \mathbb{E}[(v_\pi(S) - \hat{v}(S, w))^2]$$
$$= \mathbb{E}[\nabla_w (v_\pi(S) - \hat{v}(S, w))^2]$$

• To minimize the objective function J(w), we can use the gradient-descent algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w J(w_k)$$

The true gradient is

$$\nabla_w J(w) = \nabla_w \mathbb{E}[(v_\pi(S) - \hat{v}(S, w))^2]$$
  
=  $\mathbb{E}[\nabla_w (v_\pi(S) - \hat{v}(S, w))^2]$   
=  $2\mathbb{E}[(v_\pi(S) - \hat{v}(S, w))(-\nabla_w \hat{v}(S, w))]$ 

• To minimize the objective function J(w), we can use the gradient-descent algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w J(w_k)$$

The true gradient is

$$\nabla_{w}J(w) = \nabla_{w}\mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^{2}]$$
  
=  $\mathbb{E}[\nabla_{w}(v_{\pi}(S) - \hat{v}(S, w))^{2}]$   
=  $2\mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))(-\nabla_{w}\hat{v}(S, w))]$   
=  $-2\mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))\nabla_{w}\hat{v}(S, w)]$ 

• To minimize the objective function J(w), we can use the gradient-descent algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w J(w_k)$$

The true gradient is

$$\nabla_{w}J(w) = \nabla_{w}\mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^{2}]$$
  
=  $\mathbb{E}[\nabla_{w}(v_{\pi}(S) - \hat{v}(S, w))^{2}]$   
=  $2\mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))(-\nabla_{w}\hat{v}(S, w))]$   
=  $-2\mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))\nabla_{w}\hat{v}(S, w)]$ 

The true gradient above involves the calculation of an expectation.

$$w_{k+1} = w_k + \alpha_k \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))\nabla_w \hat{v}(S, w)]$$

$$\downarrow$$

$$w_{t+1} = w_t + \alpha_t (v_{\pi}(s_t) - \hat{v}(s_t, w_t))\nabla_w \hat{v}(s_t, w_t)$$

- The samples are expected to satisfy the stationary distribution. In practice, they may not satisfy.
- This algorithm is not implementable because it requires the true state value  $v_{\pi}$ , which is the unknown to be estimated.
- We can replace  $v_{\pi}(s_t)$  with an approximation so that the algorithm is implementable.

$$w_{k+1} = w_k + \alpha_k \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))\nabla_w \hat{v}(S, w)]$$

$$\downarrow$$

$$w_{t+1} = w_t + \alpha_t (v_{\pi}(s_t) - \hat{v}(s_t, w_t))\nabla_w \hat{v}(s_t, w_t)$$

- The samples are expected to satisfy the stationary distribution. In practice, they may not satisfy.
- This algorithm is not implementable because it requires the true state value  $v_{\pi}$ , which is the unknown to be estimated.
- We can replace  $v_{\pi}(s_t)$  with an approximation so that the algorithm is implementable.

$$w_{k+1} = w_k + \alpha_k \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))\nabla_w \hat{v}(S, w)]$$

$$\downarrow$$

$$w_{t+1} = w_t + \alpha_t (v_{\pi}(s_t) - \hat{v}(s_t, w_t))\nabla_w \hat{v}(s_t, w_t)$$

- The samples are expected to satisfy the stationary distribution. In practice, they may not satisfy.
- This algorithm is not implementable because it requires the true state value  $v_{\pi}$ , which is the unknown to be estimated.
- We can replace  $v_{\pi}(s_t)$  with an approximation so that the algorithm is implementable.

$$w_{k+1} = w_k + \alpha_k \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))\nabla_w \hat{v}(S, w)]$$

$$\downarrow$$

$$w_{t+1} = w_t + \alpha_t (v_{\pi}(s_t) - \hat{v}(s_t, w_t))\nabla_w \hat{v}(s_t, w_t)$$

- The samples are expected to satisfy the stationary distribution. In practice, they may not satisfy.
- This algorithm is not implementable because it requires the true state value  $v_{\pi}$ , which is the unknown to be estimated.
- We can replace  $v_{\pi}(s_t)$  with an approximation so that the algorithm is implementable.

• First, Monte Carlo learning with function approximation Let  $g_t$  be the discounted return starting from  $s_t$  in the episode. Then,  $g_t$  can be used to approximate  $v_{\pi}(s_t)$ . The algorithm becomes

$$w_{t+1} = w_t + \alpha_t (g_t - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t).$$

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t).$$

## • First, Monte Carlo learning with function approximation

Let  $g_t$  be the discounted return starting from  $s_t$  in the episode. Then,  $g_t$  can be used to approximate  $v_{\pi}(s_t)$ . The algorithm becomes

$$w_{t+1} = w_t + \alpha_t (g_t - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t).$$

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t).$$

• First, Monte Carlo learning with function approximation Let  $g_t$  be the discounted return starting from  $s_t$  in the episode. Then,  $g_t$  can be used to approximate  $v_{\pi}(s_t)$ . The algorithm becomes

$$w_{t+1} = w_t + \alpha_t (g_t - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t).$$

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t).$$

• First, Monte Carlo learning with function approximation Let  $g_t$  be the discounted return starting from  $s_t$  in the episode. Then,  $g_t$  can be used to approximate  $v_{\pi}(s_t)$ . The algorithm becomes

$$w_{t+1} = w_t + \alpha_t (g_t - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t).$$

### • Second, TD learning with function approximation

By the spirit of TD learning,  $r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t)$  can be viewed as an approximation of  $v_{\pi}(s_t)$ . Then, the algorithm becomes

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t).$$

• First, Monte Carlo learning with function approximation Let  $g_t$  be the discounted return starting from  $s_t$  in the episode. Then,  $g_t$  can be used to approximate  $v_{\pi}(s_t)$ . The algorithm becomes

$$w_{t+1} = w_t + \alpha_t (g_t - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t).$$

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t).$$

#### Pseudocode: TD learning of state values with function approximation

**Initialization:** A function  $\hat{v}(s, w)$  that is differentiable in w. Initial parameter  $w_0$ . **Goal:** Learn the true state values of a given policy  $\pi$ .

For each episode  $\{(s_t, r_{t+1}, s_{t+1})\}_t$  generated by  $\pi$ , do For each sample  $(s_t, r_{t+1}, s_{t+1})$ , do In the general case,  $w_{t+1} = w_t + \alpha_t \left[r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t)\right] \nabla_w \hat{v}(s_t, w_t)$ 

It can only estimate the state values of a given policy, but it is important to understand other algorithms introduced later.

#### Pseudocode: TD learning of state values with function approximation

**Initialization:** A function  $\hat{v}(s, w)$  that is differentiable in w. Initial parameter  $w_0$ . **Goal:** Learn the true state values of a given policy  $\pi$ .

```
For each episode \{(s_t, r_{t+1}, s_{t+1})\}_t generated by \pi, do
For each sample (s_t, r_{t+1}, s_{t+1}), do
In the general case,
w_{t+1} = w_t + \alpha_t [r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t)] \nabla_w \hat{v}(s_t, w_t)
In the linear case,
w_{t+1} = w_t + \alpha_t [r_{t+1} + \gamma \phi^T(s_{t+1})w_t - \phi^T(s_t)w_t] \phi(s_t)
```

It can only estimate the state values of a given policy, but it is important to understand other algorithms introduced later.

#### Pseudocode: TD learning of state values with function approximation

**Initialization:** A function  $\hat{v}(s, w)$  that is differentiable in w. Initial parameter  $w_0$ . **Goal:** Learn the true state values of a given policy  $\pi$ .

```
For each episode \{(s_t, r_{t+1}, s_{t+1})\}_t generated by \pi, do

For each sample (s_t, r_{t+1}, s_{t+1}), do

In the general case,

w_{t+1} = w_t + \alpha_t [r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t)] \nabla_w \hat{v}(s_t, w_t)

In the linear case,

w_{t+1} = w_t + \alpha_t [r_{t+1} + \gamma \phi^T(s_{t+1})w_t - \phi^T(s_t)w_t] \phi(s_t)
```

It can only estimate the state values of a given policy, but it is important to understand other algorithms introduced later.

# Outline

### 1 Motivating examples: from table to function

## 2 Algorithm for state value estimation

- Objective function
- Optimization algorithms
- Selection of function approximators
- Illustrative examples
- Summary of the story
- Theoretical analysis (optional)
- 3 Sarsa with function approximation
- 4 Q-learning with function approximation
- 5 Deep Q-learning

#### 6 Summary

• The first approach, which was widely used before, is to use a linear function

$$\hat{v}(s,w) = \phi^T(s)w$$

- The second approach, which is **widely used nowadays**, is to use a neural network as a nonlinear function approximator.
  - For example, the input is s, the output is  $\hat{v}(s,w)$ , and the parameter is w.

• The first approach, which was widely used before, is to use a linear function

$$\hat{v}(s,w) = \phi^T(s)w$$

- The second approach, which is **widely used nowadays**, is to use a neural network as a nonlinear function approximator.
  - For example, the input is s, the output is  $\hat{v}(s,w)$ , and the parameter is w.

• The first approach, which was widely used before, is to use a linear function

$$\hat{v}(s,w) = \phi^T(s)w$$

- The second approach, which is **widely used nowadays**, is to use a neural network as a nonlinear function approximator.
  - For example, the input is s, the output is  $\hat{v}(s, w)$ , and the parameter is w.

• The first approach, which was widely used before, is to use a linear function

$$\hat{v}(s,w) = \phi^T(s)w$$

- The second approach, which is **widely used nowadays**, is to use a neural network as a nonlinear function approximator.
  - For example, the input is s, the output is  $\hat{v}(s,w)$ , and the parameter is w.

• The first approach, which was widely used before, is to use a linear function

$$\hat{v}(s,w) = \phi^T(s)w$$

- The second approach, which is **widely used nowadays**, is to use a neural network as a nonlinear function approximator.
  - For example, the input is s, the output is  $\hat{v}(s, w)$ , and the parameter is w.

• The first approach, which was widely used before, is to use a linear function

$$\hat{v}(s,w) = \phi^T(s)w$$

- The second approach, which is **widely used nowadays**, is to use a neural network as a nonlinear function approximator.
  - For example, the input is s, the output is  $\hat{v}(s, w)$ , and the parameter is w.

In the linear case where  $\hat{v}(s,w)=\phi^T(s)w$  , we have

 $\nabla_w \hat{v}(s, w) = \phi(s).$ 

Substituting the gradient into the TD algorithm

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t)$$

yields

$$w_{t+1} = w_t + \alpha_t \big[ r_{t+1} + \gamma \phi^T(s_{t+1}) w_t - \phi^T(s_t) w_t \big] \phi(s_t),$$

which is the algorithm of TD learning with linear function approximation. It is called TD-Linear in our course.

In the linear case where  $\hat{v}(s,w)=\phi^T(s)w$  , we have

$$\nabla_w \hat{v}(s, w) = \phi(s).$$

Substituting the gradient into the TD algorithm

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t)$$

yields

$$w_{t+1} = w_t + \alpha_t \big[ r_{t+1} + \gamma \phi^T(s_{t+1}) w_t - \phi^T(s_t) w_t \big] \phi(s_t),$$

which is the algorithm of TD learning with linear function approximation. It is called TD-Linear in our course.

In the linear case where  $\hat{v}(s,w)=\phi^{T}(s)w$  , we have

$$\nabla_w \hat{v}(s, w) = \phi(s).$$

Substituting the gradient into the TD algorithm

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t)$$

#### yields

$$w_{t+1} = w_t + \alpha_t \big[ r_{t+1} + \gamma \phi^T(s_{t+1}) w_t - \phi^T(s_t) w_t \big] \phi(s_t),$$

which is the algorithm of TD learning with linear function approximation. It is called TD-Linear in our course. In the linear case where  $\hat{v}(s,w)=\phi^T(s)w$  , we have

$$\nabla_w \hat{v}(s, w) = \phi(s).$$

Substituting the gradient into the TD algorithm

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t)$$

yields

$$w_{t+1} = w_t + \alpha_t \big[ r_{t+1} + \gamma \phi^T(s_{t+1}) w_t - \phi^T(s_t) w_t \big] \phi(s_t),$$

which is the algorithm of TD learning with linear function approximation. It is called TD-Linear in our course. In the linear case where  $\hat{v}(s,w)=\phi^{T}(s)w$  , we have

$$\nabla_w \hat{v}(s, w) = \phi(s).$$

Substituting the gradient into the TD algorithm

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t)$$

yields

$$w_{t+1} = w_t + \alpha_t \big[ r_{t+1} + \gamma \phi^T(s_{t+1}) w_t - \phi^T(s_t) w_t \big] \phi(s_t),$$

which is the algorithm of TD learning with linear function approximation. It is called TD-Linear in our course.

- Disadvantages of linear function methods:
  - Difficult to select appropriate feature vectors.
- Advantages of linear function methods:
  - The theoretical properties of the TD algorithm in the linear case can be much better understood than in the nonlinear case.
  - Linear function approximation is still powerful in the sense that the tabular representation is a special case of linear function representation.

- Disadvantages of linear function methods:
  - Difficult to select appropriate feature vectors.
- Advantages of linear function methods:
  - The theoretical properties of the TD algorithm in the linear case can be much better understood than in the nonlinear case.
  - Linear function approximation is still powerful in the sense that the tabular representation is a special case of linear function representation.

- Disadvantages of linear function methods:
  - Difficult to select appropriate feature vectors.
- Advantages of linear function methods:
  - The theoretical properties of the TD algorithm in the linear case can be much better understood than in the nonlinear case.
  - Linear function approximation is still powerful in the sense that the tabular representation is a special case of linear function representation.

# We next show that tabular representation is a special case of linear function representation. Hence, the tabular and function representations are unified!

• Consider a special feature vector for state s:

 $\phi(s) = e_s \in \mathbb{R}^{|\mathcal{S}|},$ 

where  $e_s$  is a vector with the sth entry as 1 and the others as 0. • In this case,

$$\hat{v}(s,w) = \phi^T(s)w = e_s^T w = w(s),$$

where w(s) is the sth entry of w.

We next show that tabular representation is a special case of linear function representation. Hence, the tabular and function representations are unified!

• Consider a special feature vector for state s:

$$\phi(s) = e_s \in \mathbb{R}^{|\mathcal{S}|},$$

where  $e_s$  is a vector with the *s*th entry as 1 and the others as 0.

• In this case,

$$\hat{v}(s,w) = \phi^T(s)w = e_s^T w = w(s),$$

where w(s) is the sth entry of w.

We next show that tabular representation is a special case of linear function representation. Hence, the tabular and function representations are unified!

• Consider a special feature vector for state s:

$$\phi(s) = e_s \in \mathbb{R}^{|\mathcal{S}|},$$

where  $e_s$  is a vector with the *s*th entry as 1 and the others as 0.

• In this case,

$$\hat{v}(s,w) = \phi^T(s)w = e_s^T w = w(s),$$

where w(s) is the *s*th entry of w.

#### Recall that the TD-Linear algorithm is

$$w_{t+1} = w_t + \alpha_t \big[ r_{t+1} + \gamma \phi^T(s_{t+1}) w_t - \phi^T(s_t) w_t \big] \phi(s_t),$$

• When  $\phi(s_t) = e_s$ , the above algorithm becomes

$$w_{t+1} = w_t + \alpha_t \left( r_{t+1} + \gamma w_t(s_{t+1}) - w_t(s_t) \right) e_{s_t}.$$

This is a vector equation that merely updates the  $s_t$ th entry of  $w_t$ .

• Multiplying  $e_{s_t}^T$  on both sides of the equation gives

$$w_{t+1}(s_t) = w_t(s_t) + \alpha_t \left( r_{t+1} + \gamma w_t(s_{t+1}) - w_t(s_t) \right),$$

which is exactly the tabular TD algorithm (which is called TD-Table here). Summary: TD-Linear becomes TD-Table if we select a special feature vector. Recall that the TD-Linear algorithm is

$$w_{t+1} = w_t + \alpha_t \big[ r_{t+1} + \gamma \phi^T(s_{t+1}) w_t - \phi^T(s_t) w_t \big] \phi(s_t),$$

• When  $\phi(s_t) = e_s$ , the above algorithm becomes

$$w_{t+1} = w_t + \alpha_t \left( r_{t+1} + \gamma w_t(s_{t+1}) - w_t(s_t) \right) e_{s_t}.$$

This is a vector equation that merely updates the  $s_t$ th entry of  $w_t$ .

• Multiplying  $e_{s_t}^T$  on both sides of the equation gives

 $w_{t+1}(s_t) = w_t(s_t) + \alpha_t \left( r_{t+1} + \gamma w_t(s_{t+1}) - w_t(s_t) \right),$ 

which is exactly the tabular TD algorithm (which is called TD-Table here). Summary: TD-Linear becomes TD-Table if we select a special feature vector.
Recall that the TD-Linear algorithm is

$$w_{t+1} = w_t + \alpha_t \big[ r_{t+1} + \gamma \phi^T(s_{t+1}) w_t - \phi^T(s_t) w_t \big] \phi(s_t),$$

• When  $\phi(s_t) = e_s$ , the above algorithm becomes

$$w_{t+1} = w_t + \alpha_t \left( r_{t+1} + \gamma w_t(s_{t+1}) - w_t(s_t) \right) e_{s_t}.$$

This is a vector equation that merely updates the  $s_t$ th entry of  $w_t$ .

• Multiplying  $e_{s_t}^T$  on both sides of the equation gives

$$w_{t+1}(s_t) = w_t(s_t) + \alpha_t \left( r_{t+1} + \gamma w_t(s_{t+1}) - w_t(s_t) \right),$$

which is exactly the tabular TD algorithm (which is called TD-Table here). Summary: TD-Linear becomes TD-Table if we select a special feature vector. Recall that the TD-Linear algorithm is

$$w_{t+1} = w_t + \alpha_t \big[ r_{t+1} + \gamma \phi^T(s_{t+1}) w_t - \phi^T(s_t) w_t \big] \phi(s_t),$$

• When  $\phi(s_t) = e_s$ , the above algorithm becomes

$$w_{t+1} = w_t + \alpha_t \left( r_{t+1} + \gamma w_t(s_{t+1}) - w_t(s_t) \right) e_{s_t}.$$

This is a vector equation that merely updates the  $s_t$ th entry of  $w_t$ .

• Multiplying  $e_{s_t}^T$  on both sides of the equation gives

$$w_{t+1}(s_t) = w_t(s_t) + \alpha_t \left( r_{t+1} + \gamma w_t(s_{t+1}) - w_t(s_t) \right),$$

which is exactly the tabular TD algorithm (which is called TD-Table here). Summary: TD-Linear becomes TD-Table if we select a special feature vector.

# Outline

#### 1 Motivating examples: from table to function

# 2 Algorithm for state value estimation

- Objective function
- Optimization algorithms
- Selection of function approximators

#### Illustrative examples

- Summary of the story
- Theoretical analysis (optional)
- 3 Sarsa with function approximation
- 4 Q-learning with function approximation
- 5 Deep Q-learning

#### 6 Summary

	1	2	3	4	5
1	+	+	+	+	+
2	+	+	+	+	+
3	+	+	+	+	+
4	+	+	+	+	+
5	+	+	+	+	+

- Given a policy:  $\pi(a|s) = 0.2$  for any s, a
- Our aim is to estimate the state values of this policy (policy evaluation problem).
- There are 25 state values in total. We next show that we can use less than 25 parameters to approximate 25 state values.
- Set  $r_{\text{forbidden}} = r_{\text{boundary}} = -1$ ,  $r_{\text{target}} = 1$ , and  $\gamma = 0.9$ .

	1	2	3	4	5
1	+	+	+	+	+
2	+	+	+	+	+
3	+	+	+	+	$^+$
4	+	+	+	+	+
5	+	+	+	+	+

- Given a policy:  $\pi(a|s) = 0.2$  for any s, a
- Our aim is to estimate the state values of this policy (policy evaluation problem).
- There are 25 state values in total. We next show that we can use less than 25 parameters to approximate 25 state values.
- Set  $r_{\text{forbidden}} = r_{\text{boundary}} = -1$ ,  $r_{\text{target}} = 1$ , and  $\gamma = 0.9$ .

	1	2	3	4	5
1	+	+	+	+	+
2	+	+	+	+	+
3	+	+	+	+	$^+$
4	+	+	+	+	+
5	+	+	+	+	+

- Given a policy:  $\pi(a|s) = 0.2$  for any s, a
- Our aim is to estimate the state values of this policy (policy evaluation problem).
- There are 25 state values in total. We next show that we can use less than 25 parameters to approximate 25 state values.
- Set  $r_{\text{forbidden}} = r_{\text{boundary}} = -1$ ,  $r_{\text{target}} = 1$ , and  $\gamma = 0.9$ .

	1	2	3	4	5
1	+	+	+	+	+
2	+	+	+	+	+
3	+	+	+	+	$^+$
4	+	+	+	+	+
5	+	+	+	+	+

- Given a policy:  $\pi(a|s) = 0.2$  for any s, a
- Our aim is to estimate the state values of this policy (policy evaluation problem).
- There are 25 state values in total. We next show that we can use less than 25 parameters to approximate 25 state values.
- Set  $r_{\text{forbidden}} = r_{\text{boundary}} = -1$ ,  $r_{\text{target}} = 1$ , and  $\gamma = 0.9$ .

	1	2	3	4	5
1	+	+	+	+	+
2	+	+	+	+	+
3	+	+	+	+	+
4	+	+	+	+	+
5	+	+	+	+	+

- Given a policy:  $\pi(a|s) = 0.2$  for any s, a
- Our aim is to estimate the state values of this policy (policy evaluation problem).
- There are 25 state values in total. We next show that we can use less than 25 parameters to approximate 25 state values.
- Set  $r_{\text{forbidden}} = r_{\text{boundary}} = -1$ ,  $r_{\text{target}} = 1$ , and  $\gamma = 0.9$ .

# Ground truth:

• The true state values and the 3D visualization



Experience samples:

- 500 episodes were generated following the given policy.
- Each episode has 500 steps and starts from a randomly selected state-action pair following a uniform distribution.

## Ground truth:

• The true state values and the 3D visualization



Experience samples:

- 500 episodes were generated following the given policy.
- Each episode has 500 steps and starts from a randomly selected state-action pair following a uniform distribution.

# Ground truth:

• The true state values and the 3D visualization



Experience samples:

- 500 episodes were generated following the given policy.
- Each episode has 500 steps and starts from a randomly selected state-action pair following a uniform distribution.

# TD-Table:

• For comparison, the results by the tabular TD algorithm (called TD-Table here):



- How to apply the TD-Linear algorithm?
  - Feature vector selection:

$$\phi(s) = \begin{bmatrix} 1\\ x\\ y \end{bmatrix} \in \mathbb{R}^3.$$

- In this case, the approximated state value is

$$\hat{v}(s,w) = \phi^T(s)w = [1, x, y] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = w_1 + w_2 x + w_3 y.$$

Remark:  $\phi(s)$  can also be defined as  $\phi(s) = [x,y,1]^T$ , where the order of the elements does not matter.

- How to apply the TD-Linear algorithm?
  - Feature vector selection:

$$\phi(s) = \begin{bmatrix} 1\\ x\\ y \end{bmatrix} \in \mathbb{R}^3.$$

- In this case, the approximated state value is

$$\hat{v}(s,w) = \phi^T(s)w = [1, x, y] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = w_1 + w_2 x + w_3 y.$$

Remark:  $\phi(s)$  can also be defined as  $\phi(s) = [x,y,1]^T$ , where the order of the elements does not matter.

- How to apply the TD-Linear algorithm?
  - Feature vector selection:

$$\phi(s) = \begin{bmatrix} 1\\ x\\ y \end{bmatrix} \in \mathbb{R}^3.$$

- In this case, the approximated state value is

$$\hat{v}(s,w) = \phi^T(s)w = [1, x, y] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = w_1 + w_2 x + w_3 y.$$

Remark:  $\phi(s)$  can also be defined as  $\phi(s) = [x,y,1]^T$ , where the order of the elements does not matter.

- How to apply the TD-Linear algorithm?
  - Feature vector selection:

$$\phi(s) = \begin{bmatrix} 1\\ x\\ y \end{bmatrix} \in \mathbb{R}^3.$$

- In this case, the approximated state value is

$$\hat{v}(s,w) = \phi^T(s)w = [1, x, y] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = w_1 + w_2 x + w_3 y.$$

Remark:  $\phi(s)$  can also be defined as  $\phi(s) = [x, y, 1]^T$ , where the order of the elements does not matter.

• Results by the TD-Linear algorithm:



- The trend is right, but there are errors due to limited approximation ability!
- We are trying to use a plane to approximate a non-plane surface!

• Results by the TD-Linear algorithm:



- The trend is right, but there are errors due to limited approximation ability!
- We are trying to use a plane to approximate a non-plane surface!

• For example, we can consider

$$\phi(s) = [1, x, y, x^2, y^2, xy]^T \in \mathbb{R}^6.$$

In this case,

$$\hat{v}(s,w) = \phi^{T}(s)w = w_1 + w_2x + w_3y + w_4x^2 + w_5y^2 + w_6xy$$

which corresponds to a quadratic surface.

$$\phi(s) = [1, x, y, x^2, y^2, xy, x^3, y^3, x^2y, xy^2]^T \in \mathbb{R}^{10}.$$

• For example, we can consider

$$\phi(s) = [1, x, y, x^2, y^2, xy]^T \in \mathbb{R}^6.$$

In this case,

$$\hat{v}(s,w) = \phi^{T}(s)w = w_{1} + w_{2}x + w_{3}y + w_{4}x^{2} + w_{5}y^{2} + w_{6}xy$$

which corresponds to a quadratic surface.

$$\phi(s) = [1, x, y, x^2, y^2, xy, x^3, y^3, x^2y, xy^2]^T \in \mathbb{R}^{10}.$$

• For example, we can consider

$$\phi(s) = [1, x, y, x^2, y^2, xy]^T \in \mathbb{R}^6.$$

In this case,

$$\hat{v}(s,w) = \phi^{T}(s)w = w_1 + w_2x + w_3y + w_4x^2 + w_5y^2 + w_6xy$$

which corresponds to a quadratic surface.

$$\phi(s) = [1, x, y, x^2, y^2, xy, x^3, y^3, x^2y, xy^2]^T \in \mathbb{R}^{10}.$$

• For example, we can consider

$$\phi(s) = [1, x, y, x^2, y^2, xy]^T \in \mathbb{R}^6.$$

In this case,

$$\hat{v}(s,w) = \phi^T(s)w = w_1 + w_2x + w_3y + w_4x^2 + w_5y^2 + w_6xy$$

which corresponds to a quadratic surface.

$$\phi(s) = [1, x, y, x^2, y^2, xy, x^3, y^3, x^2y, xy^2]^T \in \mathbb{R}^{10}$$

Results by the TD-Linear algorithm with higher-order feature vectors:



More examples and features are given in the book.

Results by the TD-Linear algorithm with higher-order feature vectors:



The above figure:  $\phi(s) \in \mathbb{R}^{10}$ 

More examples and features are given in the book.

Results by the TD-Linear algorithm with higher-order feature vectors:



The above figure:  $\phi(s) \in \mathbb{R}^{10}$ 

More examples and features are given in the book.

# Outline

#### 1 Motivating examples: from table to function

## 2 Algorithm for state value estimation

- Objective function
- Optimization algorithms
- Selection of function approximators
- Illustrative examples

#### Summary of the story

- Theoretical analysis (optional)
- 3 Sarsa with function approximation
- 4 Q-learning with function approximation
- 5 Deep Q-learning

#### 6 Summary

1) This story started from the objective function:

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2]$$

The objective function suggests that it is a policy evaluation problem.

2) The gradient-descent algorithm is

$$w_{t+1} = w_t + \alpha_t (v_\pi(s_t) - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t)$$

 The true value function, which is unknown, in the algorithm is replaced by an approximation, leading to the algorithm:

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t)$$

1) This story started from the objective function:

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2]$$

The objective function suggests that it is a policy evaluation problem.

2) The gradient-descent algorithm is

$$w_{t+1} = w_t + \alpha_t (v_\pi(s_t) - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t)$$

 The true value function, which is unknown, in the algorithm is replaced by an approximation, leading to the algorithm:

 $w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t)$ 

1) This story started from the objective function:

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2]$$

The objective function suggests that it is a policy evaluation problem.

2) The gradient-descent algorithm is

$$w_{t+1} = w_t + \alpha_t (v_\pi(s_t) - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t)$$

 The true value function, which is unknown, in the algorithm is replaced by an approximation, leading to the algorithm:

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t)$$

1) This story started from the objective function:

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2]$$

The objective function suggests that it is a policy evaluation problem.

2) The gradient-descent algorithm is

$$w_{t+1} = w_t + \alpha_t (v_\pi(s_t) - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t)$$

 The true value function, which is unknown, in the algorithm is replaced by an approximation, leading to the algorithm:

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t)$$

1) This story started from the objective function:

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2]$$

The objective function suggests that it is a policy evaluation problem.

2) The gradient-descent algorithm is

$$w_{t+1} = w_t + \alpha_t (v_\pi(s_t) - \hat{v}(s_t, w_t)) \nabla_w \hat{v}(s_t, w_t)$$

 The true value function, which is unknown, in the algorithm is replaced by an approximation, leading to the algorithm:

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t)$$

# Outline

#### 1 Motivating examples: from table to function

# 2 Algorithm for state value estimation

- Objective function
- Optimization algorithms
- Selection of function approximators
- Illustrative examples
- Summary of the story
- Theoretical analysis (optional)
- 3 Sarsa with function approximation
- 4 Q-learning with function approximation
- 5 Deep Q-learning

#### 6 Summary

• The algorithm

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \right] \nabla_w \hat{v}(s_t, w_t)$$

does not minimize the following objective function:

$$J(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2]$$

# Theoretical analysis (optional)

#### Different objective functions:

• Objective function 1: True value error

$$J_E(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2] = \|\hat{v}(w) - v_{\pi}\|_D^2$$

• Objective function 2: Bellman error

$$J_{BE}(w) = \|\hat{v}(w) - (r_{\pi} + \gamma P_{\pi} \hat{v}(w))\|_{D}^{2} \doteq \|\hat{v}(w) - T_{\pi}(\hat{v}(w))\|_{D}^{2},$$

where  $T_{\pi}(x) \doteq r_{\pi} + \gamma P_{\pi} x$ 

• Objective function 3: Projected Bellman error

$$J_{PBE}(w) = \|\hat{v}(w) - MT_{\pi}(\hat{v}(w))\|_{D}^{2},$$

where M is a projection matrix.

- The TD-Linear algorithm minimizes the projected Bellman error. More details are omitted here. Interested readers can check my book. Different objective functions:

• Objective function 1: True value error

$$J_E(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2] = \|\hat{v}(w) - v_{\pi}\|_D^2$$

• Objective function 2: Bellman error

$$J_{BE}(w) = \|\hat{v}(w) - (r_{\pi} + \gamma P_{\pi} \hat{v}(w))\|_{D}^{2} \doteq \|\hat{v}(w) - T_{\pi}(\hat{v}(w))\|_{D}^{2},$$

where  $T_{\pi}(x) \doteq r_{\pi} + \gamma P_{\pi} x$ 

• Objective function 3: Projected Bellman error

$$J_{PBE}(w) = \|\hat{v}(w) - MT_{\pi}(\hat{v}(w))\|_{D}^{2},$$

where M is a projection matrix.

- The TD-Linear algorithm minimizes the projected Bellman error. Nore details are omitted here. Interested readers can check my book. Different objective functions:

• Objective function 1: True value error

$$J_E(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2] = \|\hat{v}(w) - v_{\pi}\|_D^2$$

• Objective function 2: Bellman error

$$J_{BE}(w) = \|\hat{v}(w) - (r_{\pi} + \gamma P_{\pi}\hat{v}(w))\|_{D}^{2} \doteq \|\hat{v}(w) - T_{\pi}(\hat{v}(w))\|_{D}^{2},$$

where  $T_{\pi}(x) \doteq r_{\pi} + \gamma P_{\pi} x$ 

• Objective function 3: Projected Bellman error

$$J_{PBE}(w) = \|\hat{v}(w) - MT_{\pi}(\hat{v}(w))\|_{D}^{2},$$

where M is a projection matrix.

- The TD-Linear algorithm minimizes the projected Bellman error.
Different objective functions:

• Objective function 1: True value error

$$J_E(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2] = \|\hat{v}(w) - v_{\pi}\|_D^2$$

• Objective function 2: Bellman error

$$J_{BE}(w) = \|\hat{v}(w) - (r_{\pi} + \gamma P_{\pi}\hat{v}(w))\|_{D}^{2} \doteq \|\hat{v}(w) - T_{\pi}(\hat{v}(w))\|_{D}^{2},$$

where  $T_{\pi}(x) \doteq r_{\pi} + \gamma P_{\pi} x$ 

• Objective function 3: Projected Bellman error

$$J_{PBE}(w) = \|\hat{v}(w) - MT_{\pi}(\hat{v}(w))\|_{D}^{2},$$

#### where M is a projection matrix.

- The TD-Linear algorithm minimizes the projected Bellman error. Aore details are omitted here. Interested readers can check my book. Different objective functions:

• Objective function 1: True value error

$$J_E(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2] = \|\hat{v}(w) - v_{\pi}\|_D^2$$

• Objective function 2: Bellman error

$$J_{BE}(w) = \|\hat{v}(w) - (r_{\pi} + \gamma P_{\pi}\hat{v}(w))\|_{D}^{2} \doteq \|\hat{v}(w) - T_{\pi}(\hat{v}(w))\|_{D}^{2},$$

where  $T_{\pi}(x) \doteq r_{\pi} + \gamma P_{\pi} x$ 

• Objective function 3: Projected Bellman error

$$J_{PBE}(w) = \|\hat{v}(w) - MT_{\pi}(\hat{v}(w))\|_{D}^{2},$$

where M is a projection matrix.

- The TD-Linear algorithm minimizes the projected Bellman error.

More details are omitted here. Interested readers can check my book.

Different objective functions:

• Objective function 1: True value error

$$J_E(w) = \mathbb{E}[(v_{\pi}(S) - \hat{v}(S, w))^2] = \|\hat{v}(w) - v_{\pi}\|_D^2$$

• Objective function 2: Bellman error

$$J_{BE}(w) = \|\hat{v}(w) - (r_{\pi} + \gamma P_{\pi}\hat{v}(w))\|_{D}^{2} \doteq \|\hat{v}(w) - T_{\pi}(\hat{v}(w))\|_{D}^{2},$$

where  $T_{\pi}(x) \doteq r_{\pi} + \gamma P_{\pi} x$ 

• Objective function 3: Projected Bellman error

$$J_{PBE}(w) = \|\hat{v}(w) - MT_{\pi}(\hat{v}(w))\|_{D}^{2},$$

where M is a projection matrix.

- The TD-Linear algorithm minimizes the projected Bellman error.

More details are omitted here. Interested readers can check my book.

## Outline

#### 1 Motivating examples: from table to function

## 2 Algorithm for state value estimation

- Objective function
- Optimization algorithms
- Selection of function approximators
- Illustrative examples
- Summary of the story
- Theoretical analysis (optional)

#### 3 Sarsa with function approximation

- 4 Q-learning with function approximation
- 5 Deep Q-learning

#### 6 Summary

 $\hat{v}(s) \approx v_{\pi}(s), \quad s \in \mathcal{S}$ 

To search for optimal policies, we need to estimate action values.

The Sarsa algorithm with value function approximation is

$$w_{t+1} = w_t + \alpha_t \Big[ r_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t).$$

$$\hat{v}(s) \approx v_{\pi}(s), \quad s \in \mathcal{S}$$

#### To search for optimal policies, we need to estimate action values.

The Sarsa algorithm with value function approximation is

$$w_{t+1} = w_t + \alpha_t \Big[ r_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t).$$

$$\hat{v}(s) \approx v_{\pi}(s), \quad s \in \mathcal{S}$$

To search for optimal policies, we need to estimate action values.

The Sarsa algorithm with value function approximation is

$$w_{t+1} = w_t + \alpha_t \Big[ r_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t).$$

$$\hat{v}(s) \approx v_{\pi}(s), \quad s \in \mathcal{S}$$

To search for optimal policies, we need to estimate action values.

The Sarsa algorithm with value function approximation is

$$w_{t+1} = w_t + \alpha_t \Big[ r_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t).$$

## Sarsa with function approximation

To search for optimal policies, we can combine policy evaluation and policy improvement.

To search for optimal policies, we can combine policy evaluation and policy improvement.

Pseudocode: Sarsa with function approximation

Initialization: Initial parameter  $w_0$ . Initial policy  $\pi_0$ .  $\alpha_t = \alpha > 0$  for all  $t. \ \epsilon \in (0, 1)$ . Goal: Learn an optimal policy to lead the agent to the target state from an initial state  $s_0$ . For each episode, do Generate  $a_0$  at  $s_0$  following  $\pi_0(s_0)$ If  $s_t$  (t = 0, 1, 2, ...) is not the target state, do Collect the experience sample  $(r_{t+1}, s_{t+1}, a_{t+1})$  given  $(s_t, a_t)$ : generate  $r_{t+1}, s_{t+1}$  by interacting with the environment; generate  $a_{t+1}$  following  $\pi_t(s_{t+1})$ . Update q-value (update parameter):  $w_{t+1} = w_t + \alpha_t \Big[ r_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t)$ Update policy:  $\pi_{t+1}(a|s_t) = 1 - \frac{\varepsilon}{|\mathcal{A}(s_t)|} (|\mathcal{A}(s_t)| - 1)$  if  $a = \arg \max_{a \in \mathcal{A}(s_t)} \hat{q}(s_t, a, w_{t+1})$   $\pi_{t+1}(a|s_t) = \frac{\varepsilon}{|\mathcal{A}(s_t)|}$  otherwise  $s_t \leftarrow s_{t+1}, a_t \leftarrow a_{t+1}$  Illustrative example:

- Sarsa with linear function approximation:  $\hat{q}(s,a,w) = \phi^T(s,a) w$
- $\gamma = 0.9$ ,  $\epsilon = 0.1$ ,  $r_{\text{boundary}} = r_{\text{forbidden}} = -10$ ,  $r_{\text{target}} = 1$ ,  $\alpha = 0.001$ .



For details, please see the book.

## Outline

#### 1 Motivating examples: from table to function

## 2 Algorithm for state value estimation

- Objective function
- Optimization algorithms
- Selection of function approximators
- Illustrative examples
- Summary of the story
- Theoretical analysis (optional)
- **3** Sarsa with function approximation
- 4 Q-learning with function approximation
- 5 Deep Q-learning

#### 6 Summary

# Similar to Sarsa, tabular Q-learning can also be extended to the case of value function approximation.

The q-value update rule is  $w_{t+1} = w_t + \alpha_t \Big[ r_{t+1} + \gamma \max_{a \in \mathcal{A}(s_{t+1})} \hat{q}(s_{t+1}, a, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t),$ which is the same as Sarsa except that  $\hat{q}(s_{t+1}, a_{t+1}, w_t)$  is replaced by  $\max_{a \in \mathcal{A}(s_{t+1})} \hat{q}(s_{t+1}, a, w_t).$  Similar to Sarsa, tabular Q-learning can also be extended to the case of value function approximation.

The q-value update rule is 
$$\begin{split} w_{t+1} &= w_t + \alpha_t \Big[ r_{t+1} + \gamma \max_{a \in \mathcal{A}(s_{t+1})} \hat{q}(s_{t+1}, a, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t), \\ \text{which is the same as Sarsa except that } \hat{q}(s_{t+1}, a_{t+1}, w_t) \text{ is replaced by} \\ \max_{a \in \mathcal{A}(s_{t+1})} \hat{q}(s_{t+1}, a, w_t). \end{split}$$

#### Pseudocode: Q-learning with function approximation (on-policy version)

Initialization: Initial parameter  $w_0$ . Initial policy  $\pi_0$ .  $\alpha_t = \alpha > 0$  for all t.  $\epsilon \in (0, 1)$ . Goal: Learn an optimal path to lead the agent to the target state from an initial state  $s_0$ .

For each episode, do

If 
$$s_t$$
  $(t = 0, 1, 2, ...)$  is not the target state, do  
Collect the experience sample  $(a_t, r_{t+1}, s_{t+1})$  given  $s_t$ : generate  $a_t$  following  $\pi_t(s_t)$ ; generate  $r_{t+1}, s_{t+1}$  by interacting with the environment.  
Update value (update parameter):  
 $w_{t+1} = w_t + \alpha_t \Big[ r_{t+1} + \gamma \max_{a \in \mathcal{A}(s_{t+1})} \hat{q}(s_{t+1}, a, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t)$   
Update policy:  
 $\pi_{t+1}(a|s_t) = 1 - \frac{\varepsilon}{|\mathcal{A}(s_t)|} (|\mathcal{A}(s_t)| - 1)$  if  $a = \arg \max_{a \in \mathcal{A}(s_t)} \hat{q}(s_t, a, w_{t+1})$   
 $\pi_{t+1}(a|s_t) = \frac{\varepsilon}{|\mathcal{A}(s_t)|}$  otherwise

Illustrative example:

- Q-learning with linear function approximation:  $\hat{q}(s, a, w) = \phi^T(s, a)w$
- $\gamma = 0.9$ ,  $\epsilon = 0.1$ ,  $r_{\text{boundary}} = r_{\text{forbidden}} = -10$ ,  $r_{\text{target}} = 1$ ,  $\alpha = 0.001$ .



## Outline

#### 1 Motivating examples: from table to function

## 2 Algorithm for state value estimation

- Objective function
- Optimization algorithms
- Selection of function approximators
- Illustrative examples
- Summary of the story
- Theoretical analysis (optional)
- 3 Sarsa with function approximation
- 4 Q-learning with function approximation

## 5 Deep Q-learning

#### 6 Summary

- One of the earliest and most successful algorithms that introduce deep neural networks into RL.
- The role of neural networks is to be a nonlinear function approximator.
- Different from the following algorithm:

 $w_{t+1} = w_t + \alpha_t \Big[ r_{t+1} + \gamma \max_{a \in \mathcal{A}(s_{t+1})} \hat{q}(s_{t+1}, a, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t)$ 

because of the way of training a network.

- One of the earliest and most successful algorithms that introduce deep neural networks into RL.
- The role of neural networks is to be a nonlinear function approximator.
- Different from the following algorithm:

 $w_{t+1} = w_t + \alpha_t \Big[ r_{t+1} + \gamma \max_{a \in \mathcal{A}(s_{t+1})} \hat{q}(s_{t+1}, a, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t)$ 

because of the way of training a network.

- One of the earliest and most successful algorithms that introduce deep neural networks into RL.
- The role of neural networks is to be a nonlinear function approximator.
- Different from the following algorithm:

 $w_{t+1} = w_t + \alpha_t \Big[ r_{t+1} + \gamma \max_{a \in \mathcal{A}(s_{t+1})} \hat{q}(s_{t+1}, a, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t)$ because of the way of training a network.

- One of the earliest and most successful algorithms that introduce deep neural networks into RL.
- The role of neural networks is to be a nonlinear function approximator.
- Different from the following algorithm:

 $w_{t+1} = w_t + \alpha_t \Big[ r_{t+1} + \gamma \max_{a \in \mathcal{A}(s_{t+1})} \hat{q}(s_{t+1}, a, w_t) - \hat{q}(s_t, a_t, w_t) \Big] \nabla_w \hat{q}(s_t, a_t, w_t)$ 

because of the way of training a network.

Deep Q-learning aims to minimize the objective function/loss function:

$$w_{t+1} = w_t + \alpha_t \left[ r_{t+1} + \gamma \max_{a \in \mathcal{A}(s_{t+1})} \hat{q}(s_{t+1}, a, w_t) - \hat{q}(s_t, a_t, w_t) \right] \nabla_w \hat{q}(s_t, a_t, w_t)$$

$$\Downarrow$$

$$J(w) = \mathbb{E} \left[ \left( R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w) \right)^2 \right]$$

where (S, A, R, S') are random variables.

#### How to minimize the objective function? Gradient-descent!

- How to calculate the gradient of the objective function? Tricky!
- That is because, in this objective function

$$J(w) = \mathbb{E}\left[ \left( R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w) \right)^2 \right],$$

the parameter w not only appears in  $\hat{q}(S, A, w)$  but also in

$$y \doteq R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w)$$

• Since the optimal *a* depends on *w*,

$$\nabla_w y \neq \gamma \max_{a \in \mathcal{A}(S')} \nabla_w \hat{q}(S', a, w)$$

How to minimize the objective function? Gradient-descent!

- How to calculate the gradient of the objective function? Tricky!
- That is because, in this objective function

$$J(w) = \mathbb{E}\left[ \left( R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w) \right)^2 \right],$$

the parameter w not only appears in  $\hat{q}(S, A, w)$  but also in

$$y \doteq R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w)$$

• Since the optimal *a* depends on *w*,

$$\nabla_w y \neq \gamma \max_{a \in \mathcal{A}(S')} \nabla_w \hat{q}(S', a, w)$$

How to minimize the objective function? Gradient-descent!

- How to calculate the gradient of the objective function? Tricky!
- That is because, in this objective function

$$J(w) = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w)\right)^2\right],$$

the parameter w not only appears in  $\hat{q}(S,A,w)$  but also in

$$y \doteq R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w)$$

Since the optimal a depends on w,

$$\nabla_w y \neq \gamma \max_{a \in \mathcal{A}(S')} \nabla_w \hat{q}(S', a, w)$$

How to minimize the objective function? Gradient-descent!

- How to calculate the gradient of the objective function? Tricky!
- That is because, in this objective function

$$J(w) = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w)\right)^2\right],$$

the parameter w not only appears in  $\hat{q}(S, A, w)$  but also in

$$y \doteq R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w)$$

• Since the optimal *a* depends on *w*,

$$\nabla_w y \neq \gamma \max_{a \in \mathcal{A}(S')} \nabla_w \hat{q}(S', a, w)$$

How to minimize the objective function? Gradient-descent!

- How to calculate the gradient of the objective function? Tricky!
- That is because, in this objective function

$$J(w) = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w)\right)^2\right],$$

the parameter w not only appears in  $\hat{q}(S, A, w)$  but also in

$$y \doteq R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w)$$

• Since the optimal *a* depends on *w*,

$$\nabla_w y \neq \gamma \max_{a \in \mathcal{A}(S')} \nabla_w \hat{q}(S', a, w)$$

#### To do that, we can introduce two networks.

- One is a main network representing  $\hat{q}(s, a, w)$
- The other is a target network  $\hat{q}(s, a, w_T)$ .

The objective function in this case degenerates to

$$J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w)\right)^2\right],\$$

where  $w_T$  is the target network parameter.

$$\nabla_{w}J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w)\right) \nabla_{w} \hat{q}(S, A, w)\right].$$

- The basic idea of deep Q-learning is to use the gradient-descent algorithm to minimize the objective function.
- However, such an optimization process evolves some important techniques that deserve special attention.

To do that, we can introduce two networks.

- One is a main network representing  $\hat{q}(s,a,w)$
- The other is a target network  $\hat{q}(s, a, w_T)$ .

The objective function in this case degenerates to

$$J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w)\right)^2\right],\$$

where  $w_T$  is the target network parameter.

$$\nabla_{w}J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w)\right) \nabla_{w} \hat{q}(S, A, w)\right].$$

- The basic idea of deep Q-learning is to use the gradient-descent algorithm to minimize the objective function.
- However, such an optimization process evolves some important techniques that deserve special attention.

To do that, we can introduce two networks.

- One is a main network representing  $\hat{q}(s,a,w)$
- The other is a target network  $\hat{q}(s, a, w_T)$ .

The objective function in this case degenerates to

$$J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w)\right)^2\right],\$$

#### where $w_T$ is the target network parameter.

$$\nabla_{w}J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w)\right) \nabla_{w} \hat{q}(S, A, w)\right].$$

- The basic idea of deep Q-learning is to use the gradient-descent algorithm to minimize the objective function.
- However, such an optimization process evolves some important techniques that deserve special attention.

To do that, we can introduce two networks.

- One is a main network representing  $\hat{q}(s,a,w)$
- The other is a target network  $\hat{q}(s, a, w_T)$ .

The objective function in this case degenerates to

$$J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w)\right)^2\right],\$$

where  $w_T$  is the target network parameter.

$$\nabla_{w}J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_{T}) - \hat{q}(S, A, w)\right) \nabla_{w} \hat{q}(S, A, w)\right].$$

- The basic idea of deep Q-learning is to use the gradient-descent algorithm to minimize the objective function.
- However, such an optimization process evolves some important techniques that deserve special attention.

To do that, we can introduce two networks.

- One is a main network representing  $\hat{q}(s,a,w)$
- The other is a target network  $\hat{q}(s, a, w_T)$ .

The objective function in this case degenerates to

$$J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w)\right)^2\right],\$$

where  $w_T$  is the target network parameter.

$$\nabla_{w}J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_{T}) - \hat{q}(S, A, w)\right) \nabla_{w} \hat{q}(S, A, w)\right].$$

- The basic idea of deep Q-learning is to use the gradient-descent algorithm to minimize the objective function.
- However, such an optimization process evolves some important techniques that deserve special attention.

To do that, we can introduce two networks.

- One is a main network representing  $\hat{q}(s,a,w)$
- The other is a target network  $\hat{q}(s, a, w_T)$ .

The objective function in this case degenerates to

$$J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w)\right)^2\right],\$$

where  $w_T$  is the target network parameter.

$$\nabla_{w}J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w_{T}) - \hat{q}(S, A, w)\right) \nabla_{w} \hat{q}(S, A, w)\right].$$

- The basic idea of deep Q-learning is to use the gradient-descent algorithm to minimize the objective function.
- However, such an optimization process evolves some important techniques that deserve special attention.

## Deep Q-learning - Two networks

## **Technique 1:** Two networks, a main network and a target network. Why is it used?

- The mathematical reason has been explained when we calculate the gradient. Implementation details:
- Let w and  $w_T$  denote the parameters of the main and target networks, respectively. They are set to be the same initially.
- In every iteration, we draw a mini-batch of samples  $\{(s, a, r, s')\}$  from the replay buffer (will be explained later).
- For every (s, a, r, s'), we can calculate the desired output as

$$y_T \doteq r + \gamma \max_{a \in \mathcal{A}(s')} \hat{q}(s', a, w_T)$$

Therefore, we obtain a mini-batch of data:

$$\{(s, a, y_T)\}$$

• Use  $\{(s, a, y_T)\}$  to train the network so as to minimize  $(y_T - \hat{q}(s, a, w))^2$ .

# Technique 1: Two networks, a main network and a target network. Why is it used?

• The mathematical reason has been explained when we calculate the gradient.

#### Implementation details:

- Let w and w<sub>T</sub> denote the parameters of the main and target networks, respectively. They are set to be the same initially.
- In every iteration, we draw a mini-batch of samples  $\{(s, a, r, s')\}$  from the replay buffer (will be explained later).
- For every (s, a, r, s'), we can calculate the desired output as

$$y_T \doteq r + \gamma \max_{a \in \mathcal{A}(s')} \hat{q}(s', a, w_T)$$

Therefore, we obtain a mini-batch of data:

$$\{(s, a, y_T)\}$$

• Use  $\{(s, a, y_T)\}$  to train the network so as to minimize  $(y_T - \hat{q}(s, a, w))^2$ .

Technique 1: Two networks, a main network and a target network. Why is it used?

• The mathematical reason has been explained when we calculate the gradient.

## Implementation details:

- Let w and  $w_T$  denote the parameters of the main and target networks, respectively. They are set to be the same initially.
- In every iteration, we draw a mini-batch of samples  $\{(s, a, r, s')\}$  from the replay buffer (will be explained later).
- For every (s, a, r, s'), we can calculate the desired output as

$$y_T \doteq r + \gamma \max_{a \in \mathcal{A}(s')} \hat{q}(s', a, w_T)$$

Therefore, we obtain a mini-batch of data:

$$\{(s, a, y_T)\}$$

• Use  $\{(s, a, y_T)\}$  to train the network so as to minimize  $(y_T - \hat{q}(s, a, w))^2$ .
Technique 1: Two networks, a main network and a target network. Why is it used?

• The mathematical reason has been explained when we calculate the gradient.

#### Implementation details:

- Let w and  $w_T$  denote the parameters of the main and target networks, respectively. They are set to be the same initially.
- In every iteration, we draw a mini-batch of samples  $\{(s, a, r, s')\}$  from the replay buffer (will be explained later).
- For every (s, a, r, s'), we can calculate the desired output as

$$y_T \doteq r + \gamma \max_{a \in \mathcal{A}(s')} \hat{q}(s', a, w_T)$$

Therefore, we obtain a mini-batch of data:

$$\{(s, a, y_T)\}$$

• Use  $\{(s, a, y_T)\}$  to train the network so as to minimize  $(y_T - \hat{q}(s, a, w))^2$ .

Technique 1: Two networks, a main network and a target network. Why is it used?

• The mathematical reason has been explained when we calculate the gradient.

#### Implementation details:

- Let w and  $w_T$  denote the parameters of the main and target networks, respectively. They are set to be the same initially.
- In every iteration, we draw a mini-batch of samples  $\{(s, a, r, s')\}$  from the replay buffer (will be explained later).
- $\bullet\,$  For every (s,a,r,s'), we can calculate the desired output as

$$y_T \doteq r + \gamma \max_{a \in \mathcal{A}(s')} \hat{q}(s', a, w_T)$$

Therefore, we obtain a mini-batch of data:

$$\{(s, a, y_T)\}$$

• Use  $\{(s, a, y_T)\}$  to train the network so as to minimize  $(y_T - \hat{q}(s, a, w))^2$ .

Technique 1: Two networks, a main network and a target network. Why is it used?

• The mathematical reason has been explained when we calculate the gradient.

#### Implementation details:

- Let w and  $w_T$  denote the parameters of the main and target networks, respectively. They are set to be the same initially.
- In every iteration, we draw a mini-batch of samples  $\{(s, a, r, s')\}$  from the replay buffer (will be explained later).
- $\bullet\,$  For every (s,a,r,s'), we can calculate the desired output as

$$y_T \doteq r + \gamma \max_{a \in \mathcal{A}(s')} \hat{q}(s', a, w_T)$$

Therefore, we obtain a mini-batch of data:

$$\{(s, a, y_T)\}$$

• Use  $\{(s, a, y_T)\}$  to train the network so as to minimize  $(y_T - \hat{q}(s, a, w))^2$ .

**Question:** What is experience replay? **Answer:** 

- After we have collected some experience samples, we do NOT use these samples in the order they were collected.
- Instead, we store them in a set, called replay buffer  $\mathcal{B} \doteq \{(s, a, r, s')\}$
- Every time we train the neural network, we can draw a mini-batch of random samples from the replay buffer.
- The draw of samples, or called experience replay, should follow a uniform distribution.

# **Technique 2: Experience replay Question:** What is experience replay?

- After we have collected some experience samples, we do NOT use these samples in the order they were collected.
- Instead, we store them in a set, called replay buffer  $\mathcal{B} \doteq \{(s, a, r, s')\}$
- Every time we train the neural network, we can draw a mini-batch of random samples from the replay buffer.
- The draw of samples, or called experience replay, should follow a uniform distribution.

Question: What is experience replay?

- After we have collected some experience samples, we do NOT use these samples in the order they were collected.
- Instead, we store them in a set, called replay buffer  $\mathcal{B} \doteq \{(s, a, r, s')\}$
- Every time we train the neural network, we can draw a mini-batch of random samples from the replay buffer.
- The draw of samples, or called experience replay, should follow a uniform distribution.

Question: What is experience replay?

- After we have collected some experience samples, we do NOT use these samples in the order they were collected.
- Instead, we store them in a set, called replay buffer  $\mathcal{B} \doteq \{(s, a, r, s')\}$
- Every time we train the neural network, we can draw a mini-batch of random samples from the replay buffer.
- The draw of samples, or called experience replay, should follow a uniform distribution.

Question: What is experience replay?

- After we have collected some experience samples, we do NOT use these samples in the order they were collected.
- Instead, we store them in a set, called replay buffer  $\mathcal{B} \doteq \{(s, a, r, s')\}$
- Every time we train the neural network, we can draw a mini-batch of random samples from the replay buffer.
- The draw of samples, or called experience replay, should follow a uniform distribution.

Question: What is experience replay?

- After we have collected some experience samples, we do NOT use these samples in the order they were collected.
- Instead, we store them in a set, called replay buffer  $\mathcal{B} \doteq \{(s, a, r, s')\}$
- Every time we train the neural network, we can draw a mini-batch of random samples from the replay buffer.
- The draw of samples, or called experience replay, should follow a uniform distribution.

$$J = \mathbb{E}\left[\left(R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w)\right)^2\right]$$

- $R \sim p(R|S, A), S' \sim p(S'|S, A)$ : R and S are determined by the system model.
- $(S,A) \sim d \!\!: (S,A)$  is an index and treated as a single random variable
- The distribution of the state-action pair (S, A) is assumed to be uniform.
  - Why uniform distribution? Because no prior knowledge.
  - Can we use stationary distribution like before? No, since no policy is given.

$$J = \mathbb{E}\left[ \left( R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w) \right)^2 \right]$$

- $R \sim p(R|S, A), S' \sim p(S'|S, A)$ : R and S are determined by the system model.
- $(S,A) \sim d \!\!: (S,A)$  is an index and treated as a single random variable
- The distribution of the state-action pair (S, A) is assumed to be uniform.
  - Why uniform distribution? Because no prior knowledge.
  - Can we use stationary distribution like before? No, since no policy is given.

$$J = \mathbb{E}\left[ \left( R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w) \right)^2 \right]$$

- $R \sim p(R|S, A), S' \sim p(S'|S, A)$ : R and S are determined by the system model.
- $(S,A) \sim d$ : (S,A) is an index and treated as a single random variable
- The distribution of the state-action pair (S, A) is assumed to be uniform.
  - Why uniform distribution? Because no prior knowledge.
  - Can we use stationary distribution like before? No, since no policy is given.

$$J = \mathbb{E}\left[ \left( R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w) \right)^2 \right]$$

- $R \sim p(R|S, A), S' \sim p(S'|S, A)$ : R and S are determined by the system model.
- $(S, A) \sim d$ : (S, A) is an index and treated as a single random variable
- The distribution of the state-action pair (S, A) is assumed to be uniform.
  - Why uniform distribution? Because no prior knowledge.
  - Can we use stationary distribution like before? No, since no policy is given.

$$J = \mathbb{E}\left[ \left( R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w) \right)^2 \right]$$

- $R \sim p(R|S, A), S' \sim p(S'|S, A)$ : R and S are determined by the system model.
- $(S, A) \sim d$ : (S, A) is an index and treated as a single random variable
- The distribution of the state-action pair (S, A) is assumed to be uniform.
  - Why uniform distribution? Because no prior knowledge.
  - Can we use stationary distribution like before? No, since no policy is given.

$$J = \mathbb{E}\left[ \left( R + \gamma \max_{a \in \mathcal{A}(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w) \right)^2 \right]$$

- $R \sim p(R|S, A), S' \sim p(S'|S, A)$ : R and S are determined by the system model.
- $(S, A) \sim d$ : (S, A) is an index and treated as a single random variable
- The distribution of the state-action pair (S, A) is assumed to be uniform.
  - Why uniform distribution? Because no prior knowledge.
  - Can we use stationary distribution like before? No, since no policy is given.

## Answer (continued):

- However, the samples are not uniformly collected because they are generated consequently by certain policies.
- To break the correlation between consequent samples, we can use the experience replay technique by uniformly drawing samples from the replay buffer.
- This is the mathematical reason why experience replay is necessary and why the experience replay must be uniform.

#### Answer (continued):

- However, the samples are not uniformly collected because they are generated consequently by certain policies.
- To break the correlation between consequent samples, we can use the experience replay technique by uniformly drawing samples from the replay buffer.
- This is the mathematical reason why experience replay is necessary and why the experience replay must be uniform.

#### Answer (continued):

- However, the samples are not uniformly collected because they are generated consequently by certain policies.
- To break the correlation between consequent samples, we can use the experience replay technique by uniformly drawing samples from the replay buffer.
- This is the mathematical reason why experience replay is necessary and why the experience replay must be uniform.

- Question: Why does not tabular Q-learning require experience replay?
  - Answer: Because it does not require any distribution of S or A.
- Question: Why does Deep Q-learning involve distributions?
  - Answer: Because we need to define a *scalar* objective function  $J(w) = \mathbb{E}[*]$ , where  $\mathbb{E}$  is for all (S, A).
  - The tabular case aims to solve a set of equations for all (s, a) (Bellman optimality equation), whereas the deep case aims to optimize a scalar objective function.
- Question: Can we use experience replay in tabular Q-learning?
  - Answer: Yes, we can. And more sample efficient (why?)

- Question: Why does not tabular Q-learning require experience replay?
  - Answer: Because it does not require any distribution of S or A.
- Question: Why does Deep Q-learning involve distributions?
  - Answer: Because we need to define a *scalar* objective function  $J(w) = \mathbb{E}[*]$ , where  $\mathbb{E}$  is for all (S, A).
  - The tabular case aims to solve a set of equations for all (s, a) (Bellman optimality equation), whereas the deep case aims to optimize a scalar objective function.
- Question: Can we use experience replay in tabular Q-learning?
  - Answer: Yes, we can. And more sample efficient (why?)

- Question: Why does not tabular Q-learning require experience replay?
  - Answer: Because it does not require any distribution of S or A.
- Question: Why does Deep Q-learning involve distributions?
  - Answer: Because we need to define a *scalar* objective function  $J(w) = \mathbb{E}[*]$ , where  $\mathbb{E}$  is for all (S, A).
  - The tabular case aims to solve a set of equations for all (s, a) (Bellman optimality equation), whereas the deep case aims to optimize a scalar objective function.
- Question: Can we use experience replay in tabular Q-learning?
  - Answer: Yes, we can. And more sample efficient (why?)

- Question: Why does not tabular Q-learning require experience replay?
  - Answer: Because it does not require any distribution of S or A.
- Question: Why does Deep Q-learning involve distributions?
  - Answer: Because we need to define a *scalar* objective function  $J(w) = \mathbb{E}[*]$ , where  $\mathbb{E}$  is for all (S, A).
  - The tabular case aims to solve a set of equations for all (s, a) (Bellman optimality equation), whereas the deep case aims to optimize a scalar objective function.
- Question: Can we use experience replay in tabular Q-learning?
  - Answer: Yes, we can. And more sample efficient (why?)

- Question: Why does not tabular Q-learning require experience replay?
  - Answer: Because it does not require any distribution of S or A.
- Question: Why does Deep Q-learning involve distributions?
  - Answer: Because we need to define a *scalar* objective function  $J(w) = \mathbb{E}[*]$ , where  $\mathbb{E}$  is for all (S, A).
  - The tabular case aims to solve a set of equations for all (s, a) (Bellman optimality equation), whereas the deep case aims to optimize a scalar objective function.
- Question: Can we use experience replay in tabular Q-learning?
  - Answer: Yes, we can. And more sample efficient (why?)

- Question: Why does not tabular Q-learning require experience replay?
  - Answer: Because it does not require any distribution of S or A.
- Question: Why does Deep Q-learning involve distributions?
  - Answer: Because we need to define a *scalar* objective function  $J(w) = \mathbb{E}[*]$ , where  $\mathbb{E}$  is for all (S, A).
  - The tabular case aims to solve a set of equations for all (s, a) (Bellman optimality equation), whereas the deep case aims to optimize a scalar objective function.
- Question: Can we use experience replay in tabular Q-learning?
  - Answer: Yes, we can. And more sample efficient (why?)

- Question: Why does not tabular Q-learning require experience replay?
  - Answer: Because it does not require any distribution of S or A.
- Question: Why does Deep Q-learning involve distributions?
  - Answer: Because we need to define a *scalar* objective function  $J(w) = \mathbb{E}[*]$ , where  $\mathbb{E}$  is for all (S, A).
  - The tabular case aims to solve a set of equations for all (s, a) (Bellman optimality equation), whereas the deep case aims to optimize a scalar objective function.
- Question: Can we use experience replay in tabular Q-learning?
  - Answer: Yes, we can. And more sample efficient (why?)

```
Pseudocode: Deep Q-learning (off-policy version)

Initialization: A main network and a target network with the same initial parameter.

Goal: Learn an optimal target network to approximate the optimal action values from the

experience samples generated by a given behavior policy \pi_b.

Store the experience samples generated by \pi_b in a replay buffer \mathcal{B} = \{(s, a, r, s')\}

For each iteration, do

Uniformly draw a mini-batch of samples from \mathcal{B}

For each sample (s, a, r, s'), calculate the target value as y_T = r + \gamma \max_{a \in \mathcal{A}(s')} \hat{q}(s', a, w_T), where w_T is the parameter of the target network

Update the main network to minimize (y_T - \hat{q}(s, a, w))^2 using the mini-batch of

samples

Set w_T = w every C iterations
```

Remarks:

- Why no policy update?
- The network input and output are different from the DQN paper.

```
Pseudocode: Deep Q-learning (off-policy version)

Initialization: A main network and a target network with the same initial parameter.

Goal: Learn an optimal target network to approximate the optimal action values from the

experience samples generated by a given behavior policy \pi_b.

Store the experience samples generated by \pi_b in a replay buffer \mathcal{B} = \{(s, a, r, s')\}

For each iteration, do

Uniformly draw a mini-batch of samples from \mathcal{B}

For each sample (s, a, r, s'), calculate the target value as y_T = r + \gamma \max_{a \in \mathcal{A}(s')} \hat{q}(s', a, w_T), where w_T is the parameter of the target network

Update the main network to minimize (y_T - \hat{q}(s, a, w))^2 using the mini-batch of

samples

Set w_T = w every C iterations
```

Remarks:

- Why no policy update?
- The network input and output are different from the DQN paper.

```
Pseudocode: Deep Q-learning (off-policy version)

Initialization: A main network and a target network with the same initial parameter.

Goal: Learn an optimal target network to approximate the optimal action values from the

experience samples generated by a given behavior policy \pi_b.

Store the experience samples generated by \pi_b in a replay buffer \mathcal{B} = \{(s, a, r, s')\}

For each iteration, do

Uniformly draw a mini-batch of samples from \mathcal{B}

For each sample (s, a, r, s'), calculate the target value as y_T = r + \gamma \max_{a \in \mathcal{A}(s')} \hat{q}(s', a, w_T), where w_T is the parameter of the target network

Update the min network to minimize (y_T - \hat{q}(s, a, w))^2 using the mini-batch of

samples

Set w_T = w every C iterations
```

Remarks:

- Why no policy update?
- The network input and output are different from the DQN paper.

Illustrative example:

- We need to learn optimal action values for every state-action pair.
- Once the optimal action values are obtained, the optimal greedy policy can be obtained immediately.

- One single episode is used to train the network.
- This episode is generated by an exploratory behavior policy shown in Fig. (a).
- The episode only has 1,000 steps! The tabular Q-learning requires 100,000 steps.
- A shallow neural network with one single hidden layer is used as a nonlinear approximator of  $\hat{q}(s, a, w)$ . The hidden layer has 100 neurons.

- One single episode is used to train the network.
- This episode is generated by an exploratory behavior policy shown in Fig. (a).
- The episode only has 1,000 steps! The tabular Q-learning requires 100,000 steps.
- A shallow neural network with one single hidden layer is used as a nonlinear approximator of  $\hat{q}(s, a, w)$ . The hidden layer has 100 neurons.

- One single episode is used to train the network.
- This episode is generated by an exploratory behavior policy shown in Fig. (a).
- The episode only has 1,000 steps! The tabular Q-learning requires 100,000 steps.
- A shallow neural network with one single hidden layer is used as a nonlinear approximator of  $\hat{q}(s,a,w)$ . The hidden layer has 100 neurons.

- One single episode is used to train the network.
- This episode is generated by an exploratory behavior policy shown in Fig. (a).
- The episode only has 1,000 steps! The tabular Q-learning requires 100,000 steps.
- A shallow neural network with one single hidden layer is used as a nonlinear approximator of  $\hat{q}(s, a, w)$ . The hidden layer has 100 neurons.

- One single episode is used to train the network.
- This episode is generated by an exploratory behavior policy shown in Fig. (a).
- The episode only has 1,000 steps! The tabular Q-learning requires 100,000 steps.
- A shallow neural network with one single hidden layer is used as a nonlinear approximator of  $\hat{q}(s, a, w)$ . The hidden layer has 100 neurons.

- One single episode is used to train the network.
- This episode is generated by an exploratory behavior policy shown in Fig. (a).
- The episode only has 1,000 steps! The tabular Q-learning requires 100,000 steps.
- A shallow neural network with one single hidden layer is used as a nonlinear approximator of  $\hat{q}(s, a, w)$ . The hidden layer has 100 neurons.

# Deep Q-learning

	1	2	3	4	5
1	+	+	+	+	+
2	+	+	+	+	+
3	+	+	+	+	+
4	+	+	+	+	+
5	+	+	+	+	+

The behavior policy.

An episode of 1,000 steps.

The obtained policy.

The TD error converges to zero.

The state estimation error converges to zero.


The behavior policy.



An episode of 1,000 steps.

The obtained policy.



The state estimation error converges to zero.



The behavior policy.



An episode of 1,000 steps.



The obtained policy.

The TD error converges to zero.

The state estimation error converges to zero.



The behavior policy.



An episode of 1,000 steps.



The obtained policy.



The TD error converges to zero.

The state estimation error converges to zero.



The TD error converges to zero.



What if we only use a single episode of 100 steps? Insufficient data



The behavior policy.

An episode of 100 steps.

The final policy.

The TD error converges to zero.

What if we only use a single episode of 100 steps? Insufficient data



The behavior policy.



An episode of 100 steps.

The final policy.

#### What if we only use a single episode of 100 steps? Insufficient data



The behavior policy.



An episode of 100 steps.



The final policy.

#### What if we only use a single episode of 100 steps? Insufficient data



The behavior policy.



An episode of 100 steps.



The final policy.



The TD error converges to zero.

#### What if we only use a single episode of 100 steps? Insufficient data 2



1 2 3 4 5

3 4 5



The behavior policy.

An episode of 100 steps.

The final policy.



The TD error converges to zero.



The state error does not converge to zero.

## Outline

#### 1 Motivating examples: from table to function

#### 2 Algorithm for state value estimation

- Objective function
- Optimization algorithms
- Selection of function approximators
- Illustrative examples
- Summary of the story
- Theoretical analysis (optional)
- **3** Sarsa with function approximation
- 4 Q-learning with function approximation
- 5 Deep Q-learning

#### 6 Summary

This lecture introduces the method of value function approximation.

- First, understand the basic idea.
- Second, understand the basic algorithms.