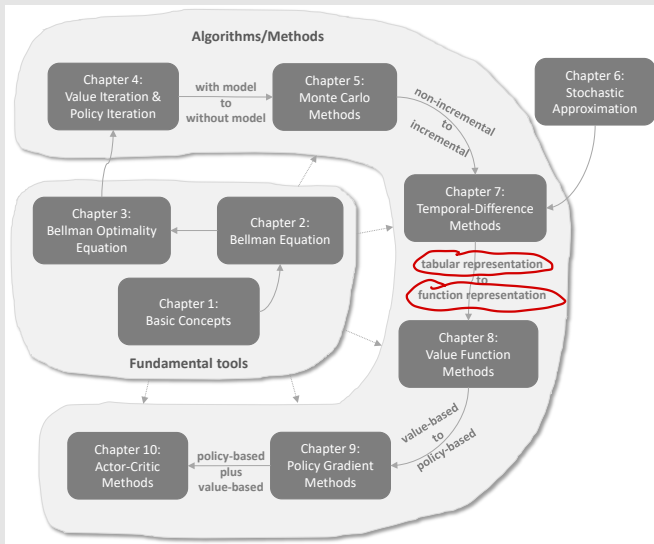


Lecture 8: Value Function Methods

Shiyu Zhao

Department of Artificial Intelligence
Westlake University

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 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

Motivating examples: from table to function

就是我在代码中实现的 $Q\{(s,a):0 \text{ for } s \text{ in states for } a \text{ in actions}\}$

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

- For example, state value:

State	s_1	s_2	\dots	s_n
Value	$v_\pi(s_1)$	$v_\pi(s_2)$	\dots	$v_\pi(s_n)$

- For example, action value:

	a_1	a_2	a_3	a_4	a_5
s_1	$q_\pi(s_1, a_1)$	$q_\pi(s_1, a_2)$	$q_\pi(s_1, a_3)$	$q_\pi(s_1, a_4)$	$q_\pi(s_1, a_5)$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
s_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 所有的 state-action pair
- Disadvantage:** difficult to handle large or continuous state or action spaces.

Two aspects: 1) storage; 2) generalization ability

例如在计算 action value 时, 需要访问几乎

Consider an example:

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

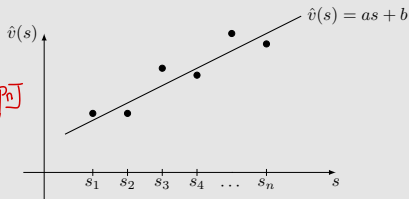
Motivating examples: from table to function

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

table 需要储存大量的值

straight line 只需要储存

2个值 $\phi^T(s)$ 和 w 即可



Suppose the equation of the straight line is

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

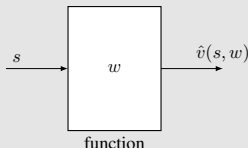
w is the **parameter vector**; $\phi(s)$ the **feature vector** of s ; $\hat{v}(s, w)$ is **linear** in w .

Motivating examples: from table to function

Difference between the tabular and function methods:

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

- Benefit: storage. We do not need to store $|\mathcal{S}|$ state values. We only need to store a lower-dimensional w .

Difference between the tabular and function methods:

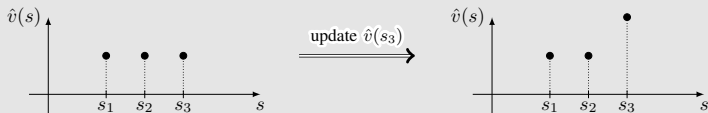
Difference 2: How to update the value of a state

- 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.

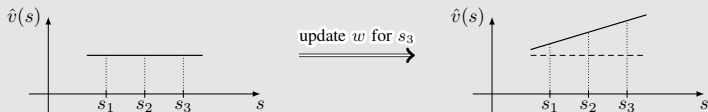
Motivating examples: from table to function

Difference between the tabular and function methods:

Difference 2: How to update the value of a state



(a) Tabular method



(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.

Motivating examples: from table to function

线条拟合: 节省空间, 但并不准确。用更高阶的线条去拟合会更准确。

The benefits are **not free**. It comes with a **cost**: the state values can not be represented accurately. This is why this method is called **approximation**.

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

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

In this case,

- 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 $\phi(s)$.
⇧ 对 w 是线性的

Quick summary:

- **Idea:** Approximate the state and action values using **parameterized functions**: $\hat{v}(s, w) \approx v_{\pi}(s)$ where $w \in \mathbb{R}^m$ 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 $|\mathcal{S}|$.
 - 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

找到一个 optimal w , 使 $\hat{v}(s)$ 接近 $v(s)$

Introduce in a more formal way:

- 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.

To find the optimal w , we need **two steps**.

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

Objective function

目标函数 (MSE)

即使 state space 是连续的, 我们在训练时并不会均匀地访问所有 states。

The objective function is

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

- Our goal is to find the best w that can minimize $J(w)$.
- The expectation is with respect to the random variable $S \in \mathcal{S}$.

比如在 Grid World 中, 靠近 goal 的区域经常访问

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 S .

平均分布，每个 state 是平等重要的

The first way is to use a uniform distribution.

- 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

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

- **Drawback:**

根据 状态分布 $d_{\pi}(s)$ 加权期望

- **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.

Objective function

策略 π 下, 访问 s 的概率

稳态分布: 长期运行, $d_\pi(s)$ 不再改变

The second way is to use the stationary distribution.

- 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 $\{d_\pi(s)\}_{s \in \mathcal{S}}$ denote the stationary distribution of the Markov process under policy π . By definition, $d_\pi(s) \geq 0$ and $\sum_{s \in \mathcal{S}} d_\pi(s) = 1$.
- The objective function can be rewritten as

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

This function is a weighted squared error.

- Since more frequently visited states have higher values of $d_\pi(s)$, their weights in the objective function are also higher than those rarely visited states.

More explanation about stationary distribution:

- *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.

Remarks:

- 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 $n_\pi(s)$ 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

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

s 被访问的次数
episode 的长度

一个具有探索性的 policy

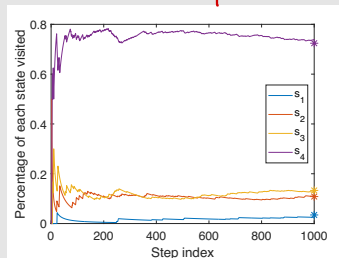
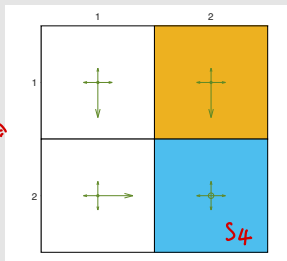


Figure: Long-run behavior of an ϵ -greedy policy with $\epsilon = 0.5$.

Objective function - Stationary distribution

贝尔曼公式: $V_\pi = r_\pi + \gamma \underline{P_\pi} V_\pi$

The converged values can be predicted because they are the entries of d_π :

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

For this example, we have P_π as

状态转移矩阵 $\rightarrow P_\pi = \begin{bmatrix} 0.3 & 0.1 & 0.6 & 0 \\ 0.1 & 0.3 & 0 & 0.6 \\ 0.1 & 0 & 0.3 & 0.6 \\ 0 & 0.1 & 0.1 & 0.8 \end{bmatrix}$.

每个元素代表 $S \rightarrow S'$ 的转移概率

eg: $S_2 \rightarrow S_1$ 的概率为 0.1

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

$$d_\pi = [0.0345, 0.1084, 0.1330, 0.7241]^T$$

A comprehensive introduction can be found in my book.

前面2页讲了求 stationary distribution of S 有2种方法:

① 通过长期运行求: $d_{\pi}(s) = \frac{s \text{ 的访问次数}}{\text{episode 长度}}$

② 通过状态转移矩阵 P_{π} :

给定 3 个状态的转移矩阵 (在策略 π 下):

$$P = \begin{bmatrix} 0.7 & 0.2 & 0.1 \\ 0.1 & 0.6 & 0.3 \\ 0.2 & 0.3 & 0.5 \end{bmatrix}$$

设稳态分布 $\pi = [x, y, z]$, 满足

$$\pi = \pi P, \quad x + y + z = 1.$$

把 $\pi = \pi P$ 展开:

$$\begin{cases} x = 0.7x + 0.1y + 0.2z \\ y = 0.2x + 0.6y + 0.3z \\ z = 0.1x + 0.3y + 0.5z \end{cases} \iff \begin{cases} 0.3x - 0.1y - 0.2z = 0 \\ -0.2x + 0.4y - 0.3z = 0 \\ x + y + z = 1 \end{cases}$$

化简 (乘以10) 更好看:

$$\begin{cases} 3x - y - 2z = 0 \Rightarrow y = 3x - 2z \\ -2x + 4y - 3z = 0 \Rightarrow y = \frac{1}{2}x + \frac{3}{4}z \\ x + y + z = 1 \end{cases}$$

两式求 $x : z$:

$$3x - 2z = \frac{1}{2}x + \frac{3}{4}z \Rightarrow 2.5x = 2.75z \Rightarrow x = \frac{11}{10}z.$$

$$\text{再代回 } y = 3x - 2z = 3 \cdot \frac{11}{10}z - 2z = \frac{13}{10}z.$$

比值为 $x : y : z = 11 : 13 : 10$ 。归一化 (和为 34):

$$\pi = \left(\frac{11}{34}, \frac{13}{34}, \frac{10}{34} \right) \approx (0.3235, 0.3824, 0.2941).$$

- 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

Optimization algorithms

求 $J(w) = E[(v_\pi(S) - \hat{v}(S, w))^2]$ 的最小值 \longrightarrow 求 $\nabla J(w^*) = 0 \longrightarrow$ 用梯度下降求 w^*

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

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

The true gradient above involves the calculation of an expectation.

→ 用一个样本近似真实梯度 .

We can use the **stochastic gradient** to replace the true gradient:

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

\Downarrow

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

where s_t is a sample of S . Here, $2\alpha_t$ is merged to α_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_π , which is the unknown to be estimated.
- We can **replace** $v_\pi(s_t)$ with an approximation so that the algorithm is implementable.

不能用, 因为 $V_\pi(s_t)$ 不知道
这是 true state value .

In particular,

- 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

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

$$\downarrow R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

- 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

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

$$V_\pi(s_t), \text{ i.e. TD target}$$

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 π .

For each episode $\{(s_t, r_{t+1}, s_{t+1})\}_t$ generated by π , 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**,

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

$$\exists w^* \rightarrow \nabla J(w^*) = 0 \rightarrow \min J(w) \rightarrow V_{\pi}(s) \approx \hat{V}(s, w)$$

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

- 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

Selection of function approximators

目的 $\hat{V}(s, w) \rightarrow V_{\pi}(s)$, 但如何选择 $\hat{V}(s, w)$?

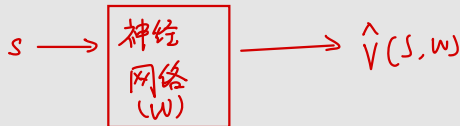
An important question that has not been answered: **How to select the function $\hat{v}(s, w)$?**

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

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

Here, $\phi(s)$ is the feature vector, which can be a polynomial basis, Fourier basis, ... (see my book for details). We have seen in the motivating example and will see again in the illustrative examples later.

- 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 .



Linear function approximation

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

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

yields

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

which is the algorithm of TD learning with linear function approximation.

It is called **TD-Linear** in our course.

↓ (s, w) 是 linear function: 难选 $\phi(s)$

- **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 :

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

$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix}$$

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

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

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

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

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

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

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.

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

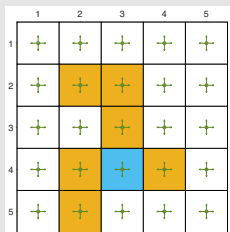
4 Q-learning with function approximation

5 Deep Q-learning

6 Summary

Illustrative examples

Consider a 5x5 grid-world example:



- 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 $\gamma = 0.9$.

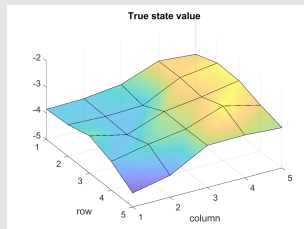
Illustrative examples

Ground truth: 用 模型参数 + 状态方程, 求出 state values 的真实值.

- The true state values and the 3D visualization

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

	1	2	3	4	5
1	-3.8	-3.8	-3.6	-3.1	-3.2
2	-3.8	-3.8	-3.8	-3.1	-2.9
3	-3.6	-3.9	-3.4	-3.2	-2.9
4	-3.9	-3.6	-3.4	-2.9	-3.2
5	-4.5	-4.2	-3.4	-3.4	-3.5

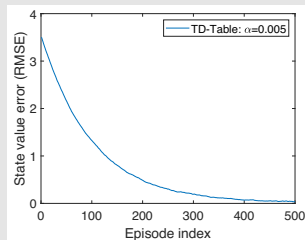
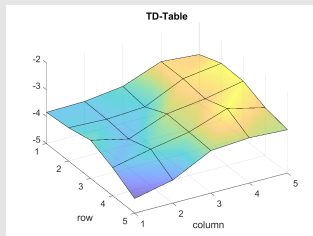


Experience samples:

- 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):



TD-Linear:

- How to apply the TD-Linear algorithm?

- Feature vector selection:

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

最后得到的 $\hat{v}(s, w)$ 是平面

- In this case, the approximated state value is

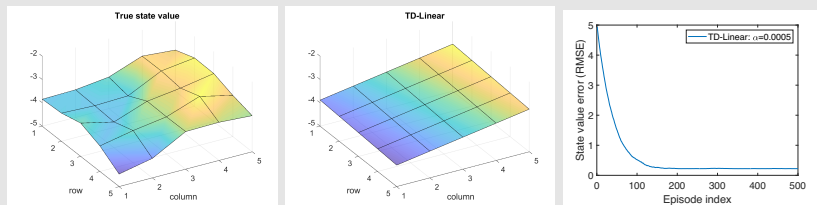
$$\hat{v}(s, w) = \phi^T(s)w = [1, x, y] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = w_1 + w_2x + w_3y.$$

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

Illustrative examples

TD-Linear:

- 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!

To enhance the approximation ability, we can use **high-order feature vectors** and hence **more parameters**.

- For example, we can consider

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

In this case,

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

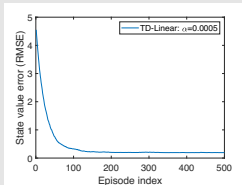
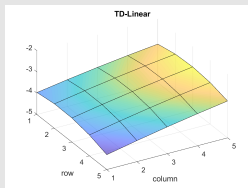
which corresponds to a quadratic surface.

- We can further increase the dimension of the feature vector:

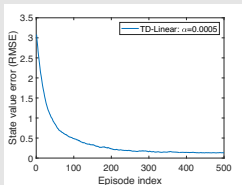
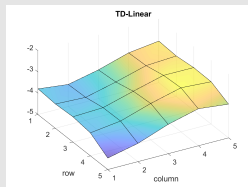
$$\phi(s) = [1, x, y, x^2, y^2, xy, x^3, y^3, x^2y, xy^2]^T \in \mathbb{R}^{10}. \quad \text{3-order}$$

Illustrative examples

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



The above figure: $\phi(s) \in \mathbb{R}^6$



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

More examples and features are given in the book.

- 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

Summary of the story

Up to now, we finished the story of TD learning with value function approximation.

- 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

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

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

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

Although this story is very helpful to understand the basic idea, it is not mathematically rigorous.

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

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

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**

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

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

- **Objective function 3: Projected Bellman error**

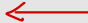
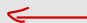
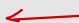

$$J_{PBE}(w) = \|\hat{v}(w) - \underline{MT_\pi(\hat{v}(w))}\|_D^2,$$

where M is a projection matrix.

M是个投影矩阵，强行让 $\hat{v}(w)$ 和 $T_\pi(\hat{v}(w))$ 的维数相同

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

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

- 1 Motivating examples: from table to function
- 2 Algorithm for state value estimation  state value (with policy π)
 - Objective function
 - Optimization algorithms
 - Selection of function approximators
 - Illustrative examples
 - Summary of the story
 - Theoretical analysis (optional)
- 3 Sarsa with function approximation  action value
- 4 Q-learning with function approximation  optimal action value
- 5 Deep Q-learning 
- 6 Summary

Sarsa with function approximation

So far, we merely considered **state value estimation**. That is

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

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

The **Sarsa algorithm** with value function approximation is

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

This is the same as the algorithm we introduced previously in this lecture **更新Q表格** except that \hat{v} is replaced by \hat{q} .

Tabular Sarsa: $Q_t = Q_t + \alpha [r_{t+1} + \gamma Q_{t+1} - Q_t]$ **更新Q表格**

函数近似 Sarsa: $w_{t+1} = w_t + \alpha [r_{t+1} + \gamma \hat{Q}_{t+1} - \hat{Q}_t] \nabla_w \hat{Q}_t$ **更新w**

Sarsa with function approximation

给出 π \longrightarrow action value \longrightarrow optimal q \longrightarrow update π

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 π_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 ($t = 0, 1, 2, \dots$) 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):

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

Update policy:

$$\pi_{t+1}(a|s_t) = 1 - \frac{\epsilon}{|\mathcal{A}(s_t)|} (|\mathcal{A}(s_t)| - 1) \text{ if } a = \arg \max_{a \in \mathcal{A}(s_t)} \hat{q}(s_t, a, w_{t+1})$$

$$\pi_{t+1}(a|s_t) = \frac{\epsilon}{|\mathcal{A}(s_t)|} \text{ otherwise}$$

$$s_t \leftarrow s_{t+1}, a_t \leftarrow a_{t+1}$$

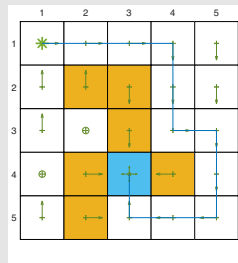
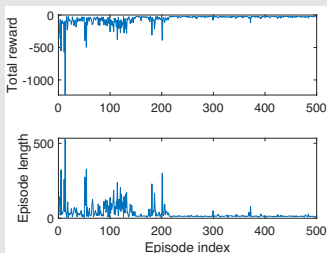
由 $\hat{q}(s_t, a, w_{t+1})$
算出新的 q 值

不直接更新 Q 表格
而是更新 w .

Sarsa with function approximation

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.

- 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

Q-learning with function approximation

Tabular Q-Learning: $Q_{t+1} = Q_t + d_t [(r_{t+1} + \gamma \max_a Q_{t+1}) - Q_t]$

$$\pi(S_{t+1}) = \arg \max_a Q_{t+1}$$

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

The q-value update rule is

TD target
↙

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

which is the same as Sarsa except that $\hat{q}(s_{t+1}, a_{t+1}, w_t)$ is replaced by $\max_{a \in \mathcal{A}(s_{t+1})} \hat{q}(s_{t+1}, a, w_t)$.

Q-learning with function approximation

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

Initialization: Initial parameter w_0 . Initial policy π_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 ($t = 0, 1, 2, \dots$) 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):

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

 Update policy:

$$\pi_{t+1}(a|s_t) = 1 - \frac{\epsilon}{|\mathcal{A}(s_t)|} (|\mathcal{A}(s_t)| - 1) \text{ if } a = \arg \max_{a \in \mathcal{A}(s_t)} \hat{q}(s_t, a, w_{t+1})$$
$$\pi_{t+1}(a|s_t) = \frac{\epsilon}{|\mathcal{A}(s_t)|} \text{ otherwise}$$

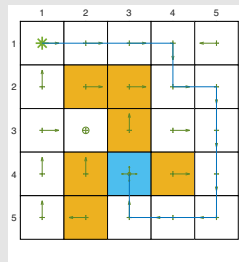
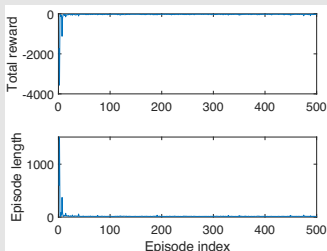
ϵ -greedy

Q-learning with function approximation

linear, 不是 DQN

Illustrative example:

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



- 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

Deep Q-learning or deep Q-network (DQN):

- 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:

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

because of the way of training a network.

上-页用的是 linear ϕ^T

可以. 但不推荐直接用这个算法



计算 $\nabla_w \hat{q}$ 十分困难 (反向传播)

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

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

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

贝尔曼最优方程

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

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

$\nabla \hat{q}$, 比较简单

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

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

- 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 solve this problem, we can assume that w in y is fixed (at least for a while) when we calculate the gradient.

为了求 ∇y , 我们假设 w 在 y 中是个常数 (一小段时间)

To do that, we can introduce two networks.

- One is a **main network** representing $\hat{q}(s, a, w)$ ← - 一直在更新
- The other is a **target network** $\hat{q}(s, a, w_T)$. ← 间断式更新, 更新时从 main 复制: $w_T = w$

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.

When w_T is fixed, the gradient of J can be easily obtained as (有个-2被省略了)

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

- 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.

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).
→ 在这个 batch 中, w_T 不变
- For every (s, a, r, s') , we can calculate the desired output as

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

Therefore, we obtain a mini-batch of data:

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

求 w^* 使下面 $J(w)$ 最小

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

Technique 2: Experience replay

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. *不是按顺序,而是 uniform 的随机拿*.
- The draw of samples, or called **experience replay**, should follow a **uniform distribution**.

Question: Why is experience replay necessary in deep Q-learning? Why does the replay must follow a uniform distribution?

Answer: The answers lie in the objective function.

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

- $R \sim p(R|S, A)$, $S' \sim p(S'|S, A)$: 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.

experience replay 是因为训练 samples 没有 uniformly collected, 所以随机打乱拿 samples.

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.

Revisit the tabular case: (S, a) 的 distribution .

- 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 π_b .

Store the experience samples generated by π_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? off-policy \rightarrow 全部计算完之后, 再用 optimal q 更新
- The network input and output are different from the DQN paper.

Illustrative example:

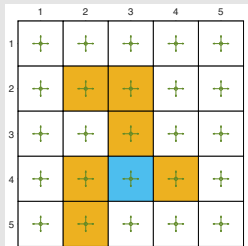
- We need to learn optimal action values for every state-action pair.
- Once the optimal action values are obtained, the optimal greedy policy can be obtained immediately.

Setup:

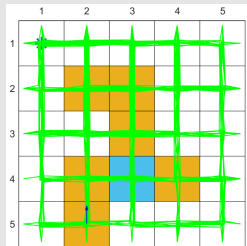
- One single episode is used to train the network.
- This episode is generated by an exploratory behavior policy shown in Fig. (a).
- The episode only has 1,000 steps! The tabular Q-learning requires 100,000 steps.
- A shallow neural network with one single hidden layer is used as a nonlinear approximator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

See details in the book.

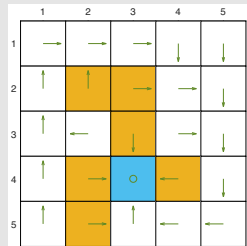
Deep Q-learning



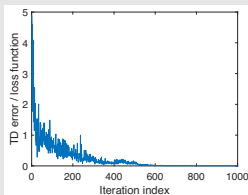
The behavior policy.



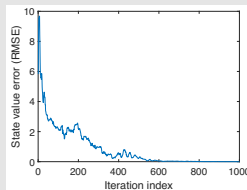
An episode of 1,000 steps.



The obtained policy.

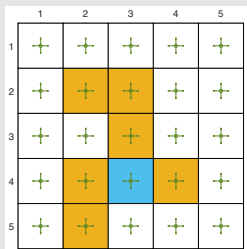


The TD error converges to zero.

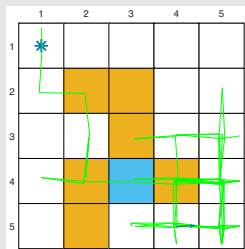


The state estimation error converges to zero.

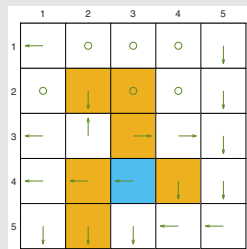
What if we only use a single episode of 100 steps? Insufficient data



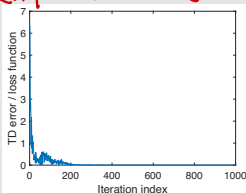
The behavior policy.



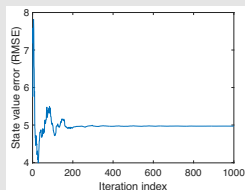
An episode of 100 steps.



The final policy.



The TD error converges to zero.



The state error does not converge to zero.

TD error 只是神经网络对目前为止的 Batch 的拟合，不代表所有的。

- 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.