

Lecture 7: Temporal-Difference Methods

Dr. Wen Fuxi

Introduction

"If one had to identify one idea as central and novel to RL, it would undoubtedly be TD learning."



Richard Sutton

- ▶ This lecture introduces **Temporal-Difference (TD)** learning, one of the most well-known methods in Reinforcement Learning (RL).
- ▶ **Monte Carlo (MC)** learning is the first model-free method. TD learning is the second model-free method. TD has some advantages compared to MC.

2. TD learning of state values
3. TD learning of action values: Sarsa
4. TD learning of action values: n -step Sarsa
5. TD learning of optimal action values: Q-learning
6. A unified point of view
7. Summary

2. TD learning of state values
3. TD learning of action values: Sarsa
4. TD learning of action values: n -step Sarsa
5. TD learning of optimal action values: Q-learning
6. A unified point of view
7. Summary

TD learning of state values - Algorithm description

Problem statement:

- ▶ Given policy π , the aim is to estimate the state values $\{v_\pi(s)\}_{s \in \mathcal{S}}$ under π .
- ▶ Experience samples: $(s_0, r_1, s_1, \dots, s_t, r_{t+1}, s_{t+1}, \dots)$ or $\{(s_t, r_{t+1}, s_{t+1})\}_t$ generated by π .

Important notations:

$v_t(s_t)$ is the estimated state value of $v_\pi(s_t)$

Deriving the Bellman equation

Consider a random trajectory:

$$S_t \xrightarrow{A_t} R_{t+1}, S_{t+1} \xrightarrow{A_{t+1}} R_{t+2}, S_{t+2} \xrightarrow{A_{t+2}} R_{t+3}, \dots$$

The return G_t can be written as

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots, \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots), \\ &= R_{t+1} + \gamma G_{t+1}, \end{aligned}$$

Then, it follows from the definition of the state value that

$$\begin{aligned} v_\pi(s) &= \mathbb{E}[G_t | S_t = s] \\ &= \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \mathbb{E}[R_{t+1} | S_t = s] + \gamma \mathbb{E}[G_{t+1} | S_t = s] \end{aligned}$$

TD learning of state values - The idea of the algorithm

First, a new expression of the Bellman equation.

The definition of the state value of π is

$$v_\pi(s) = \mathbb{E}[R_{t+1} + \gamma G_{t+1} \mid S_t = s], \quad s \in \mathcal{S} \quad (1)$$

where G is the discounted return. Since

$$\mathbb{E}[G_{t+1} \mid S = s] = \sum_a \pi(a \mid s) \sum_{s'} p(s' \mid s, a) v_\pi(s') = \mathbb{E}[v_\pi(S_{t+1}) \mid S_t = s],$$

where $S_{t+1} = s'$ is the next state, we can rewrite (1) as

$$v_\pi(s) = \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s], \quad s \in \mathcal{S}. \quad (2)$$

Equation (2) is another expression of the Bellman equation. It is sometimes referred to as the Bellman expectation equation, an important tool for designing and analyzing TD algorithms.

TD Policy Evaluation

Bootstrapping technique in TD

$$v^\pi(s_t) \leftarrow (1 - \alpha)v^\pi(s_t) + \alpha \cdot G_t$$

Bootstrapping

Return : $G_t = \sum_{i=0}^{+\infty} \gamma^i r_{t+i}$

- Return G_t comes from interaction with the environment
- Recall the self-consistency condition

$$\mathbb{E}_\pi\{G_t|s\} = v^\pi(s) = \mathbb{E}_\pi\{r + \gamma v^\pi(s')\}$$

Sample return to replace the actual expected return
- Use estimate $V^\pi(s)$ to replace true value $v^\pi(s)$

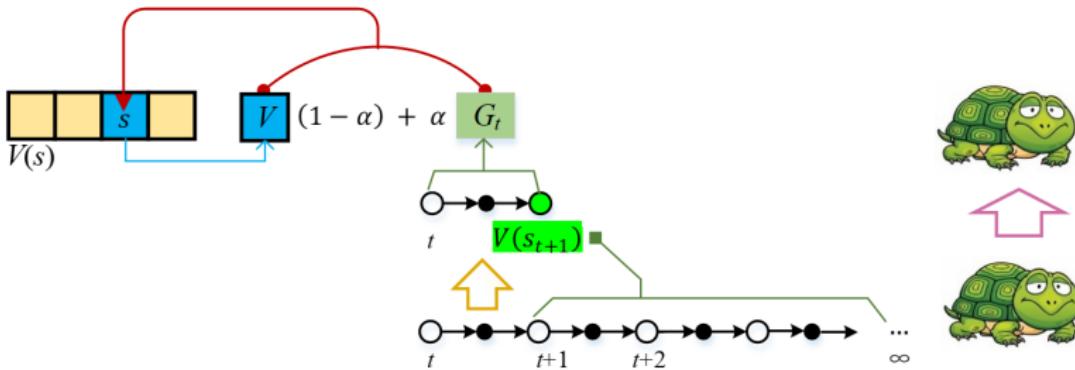
$$V^\pi(s_t) \leftarrow V^\pi(s_t) + \alpha \frac{(r_{t+1} + \gamma V^\pi(s_{t+1}) - V^\pi(s_t))}{G_t}$$

New estimate

Reward Prediction Error

TD Policy Evaluation

One-step TD for the State-value function



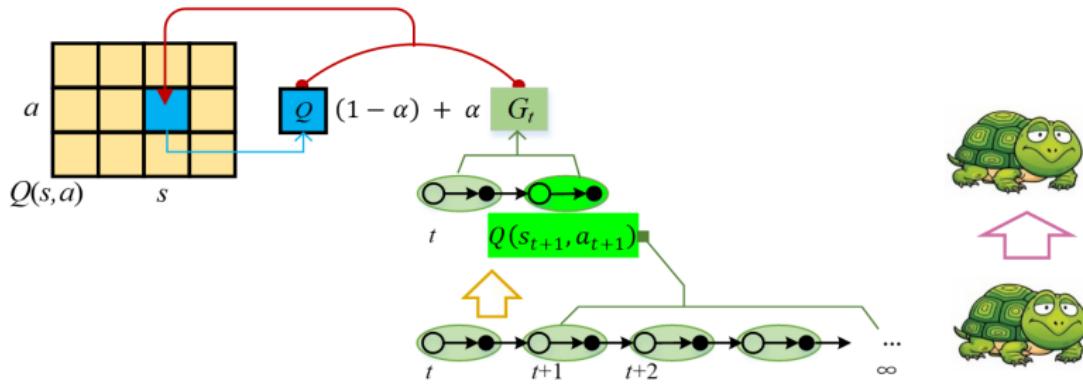
Use new sample $(s_t, a_t, r_{t+1}, s_{t+1})$ to update state-value $V(s)$

$$V^\pi(s_t) \leftarrow V^\pi(s_t) + \alpha(r_{t+1} + \gamma V^\pi(s_{t+1}) - V^\pi(s_t))$$

New Experience \leftarrow History estimate

TD Policy Evaluation

One-step TD for the Action-value function



Use new sample $(s_t, a_t, r_t, s_{t+1}, a_{t+1})$ to update state-value $Q(s, a)$

$$Q^\pi(s_t, a_t) \leftarrow \underline{Q^\pi(s_t, a_t)} + \alpha \underbrace{(r_{t+1} + \gamma Q^\pi(s_{t+1}, a_{t+1}) - Q^\pi(s_t, a_t))}_{\text{History estimate}}$$

New Experience ← → History estimate

TD learning of state values - Algorithm description

The TD learning algorithm is

$$v_{t+1}(s_t) = v_t(s_t) - \alpha_t(s_t) \left[v_t(s_t) - [r_{t+1} + \gamma v_t(s_{t+1})] \right], \quad (1)$$

$$v_{t+1}(s) = v_t(s), \quad \forall s \neq s_t, \quad (2)$$

where $t = 0, 1, 2, \dots$

Here, $v_t(s_t)$ is the estimated state value of $v_\pi(s_t)$; $\alpha_t(s_t)$ is the learning rate of s_t at time t .

- ▶ At time t , only the value of the visited state s_t is updated whereas the values of the unvisited states $s \neq s_t$ remain unchanged.
- ▶ The update in (2) will be omitted when the context is clear.

TD learning of state values - Algorithm properties

The TD algorithm can be annotated as

$$\underbrace{v_{t+1}(s_t)}_{\text{new estimate}} = \underbrace{v_t(s_t)}_{\text{current estimate}} - \alpha_t(s_t) \left[\overbrace{v_t(s_t) - [\underbrace{r_{t+1} + \gamma v_t(s_{t+1})}_{\text{TD target } \bar{v}_t}]^{\text{TD error } \delta_t}} \right] \quad (3)$$

Here,

$$\bar{v}_t \doteq r_{t+1} + \gamma v_t(s_{t+1})$$

is called the **TD target**.

$$\delta_t \doteq v_t(s_t) - [r_{t+1} + \gamma v_t(s_{t+1})] = v_t(s_t) - \bar{v}_t$$

is called the **TD error**.

Observation: The new estimate $v_{t+1}(s_t)$ is a combination of the current estimate $v_t(s_t)$ and the TD error.

TD learning of state values - Algorithm properties

First, why is \bar{v}_t called the TD target?

That is because the algorithm drives $v(s_t)$ towards \bar{v}_t .

To see that,

$$\begin{aligned} v_{t+1}(s_t) &= v_t(s_t) - \alpha_t(s_t) [v_t(s_t) - \bar{v}_t] \\ \implies v_{t+1}(s_t) - \bar{v}_t &= v_t(s_t) - \bar{v}_t - \alpha_t(s_t) [v_t(s_t) - \bar{v}_t] \\ \implies v_{t+1}(s_t) - \bar{v}_t &= [1 - \alpha_t(s_t)] [v_t(s_t) - \bar{v}_t] \\ \implies |v_{t+1}(s_t) - \bar{v}_t| &= |1 - \alpha_t(s_t)| |v_t(s_t) - \bar{v}_t| \end{aligned}$$

Since $\alpha_t(s_t)$ is a small positive number, we have

$$0 < 1 - \alpha_t(s_t) < 1$$

Therefore,

$$|v_{t+1}(s_t) - \bar{v}_t| \leq |v_t(s_t) - \bar{v}_t|$$

which means $v(s_t)$ is driven towards \bar{v}_t !

TD learning of state values - Algorithm properties

Second, what is the interpretation of the TD error?

$$\delta_t = v_t(s_t) - [r_{t+1} + \gamma v_t(s_{t+1})]$$

- It reflects the difference between v_t and v_π . To see that, denote

$$\delta_{\pi,t} \doteq v_\pi(s_t) - [r_{t+1} + \gamma v_\pi(s_{t+1})]$$

Note that

$$\mathbb{E}[\delta_{\pi,t} \mid S_t = s_t] = v_\pi(s_t) - \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s_t] = 0$$

- If $v_t(s_t) = v_\pi(s_t)$, then δ_t should be zero (in the expectation sense).
- Hence, if δ_t is not zero, then v_t is not equal to v_π .
- The TD error can be interpreted as innovation, which means new information obtained from the experience (s_t, r_{t+1}, s_{t+1}) .

TD learning of state values - Algorithm properties

Other properties

- ▶ The TD algorithm in (3) only estimates the state value of a given policy.
- ▶ It does not estimate the action values.
- ▶ It does not search for optimal policies.
- ▶ This algorithm will be extended to estimate action values and then search for optimal policies later in this lecture.
- ▶ The TD algorithm in (3) is fundamental for understanding more complex TD algorithms.

TD learning of state values - The idea of the algorithm

Q: What does this TD algorithm do mathematically?

A: It is a model-free algorithm for solving the Bellman equation of a given policy π .

- ▶ Chapter 2 introduces the model-based algorithm for solving the Bellman equation, which combines a closed-form solution with an iterative algorithm.

TD learning of state values - Algorithm convergence

Theorem (Convergence of TD Learning)

By the TD algorithm (1), $v_t(s)$ converges with probability 1 to $v_\pi(s)$ for all $s \in \mathcal{S}$ as $t \rightarrow \infty$ if $\sum_t \alpha_t(s) = \infty$ and $\sum_t \alpha_t^2(s) < \infty$ for all $s \in \mathcal{S}$.

- ▶ This theorem says the state value can be found by the TD algorithm for a given a policy π .
- ▶ $\sum_t \alpha_t(s) = \infty$ and $\sum_t \alpha_t^2(s) < \infty$ must be valid for all $s \in \mathcal{S}$.
- ▶ For condition $\sum_t \alpha_t(s) = \infty$: At time step t ,
 - ◊ If $s = s_t$, then $\alpha_t(s) > 0$;
 - ◊ If $s \neq s_t$, then $\alpha_t(s) = 0$.As a result, $\sum_t \alpha_t(s) = \infty$ requires every state must be visited an infinite (or sufficiently many) number of times.
- ▶ For condition $\sum_t \alpha_t^2(s) < \infty$: In practice, the learning rate α is often selected as a small constant. In this case, the condition that $\sum_t \alpha_t^2(s) < \infty$ is invalid anymore. When α is constant, it can still be shown that the algorithm converges in the sense of expectation sense.

TD learning of state values - Algorithm properties

While TD learning and MC learning are both model-free, what are the advantages and disadvantages of TD learning compared to MC learning?

Table: Comparison between TD learning and MC learning.

TD/Sarsa learning	MC learning
Online: TD learning is online. It can update the state/action values Immediately after receiving a reward.	Offline: MC learning is offline. It has to wait until an episode has been completely collected.
Continuing tasks: Since TD learning is online, it can handle Both episodic and continuing tasks.	Episodic tasks: Since MC learning is offline, it can only handle episodic tasks that have terminating states.

TD learning of state values - Algorithm properties

While TD learning and MC learning are both model-free, what are the advantages and disadvantages of TD learning compared to MC learning?

Table: Comparison between TD learning and MC learning (continued).

TD/Sarsa learning	MC learning
Bootstrapping: TD bootstraps because the update of a value relies on the previous estimate of this value. Hence, it requires initial guesses.	Non-bootstrapping: MC is not bootstrapping, because it can directly estimate state/action values without any initial guess.
Low estimation variance: TD has lower than MC because there are fewer random variables. For instance, Sarsa requires $R_{t+1}, S_{t+1}, A_{t+1}$.	High estimation variance: To estimate $q_\pi(s_t, a_t)$, we need samples of $R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$. Suppose the length of each episode is L . There are $ \mathcal{A} ^L$ possible episodes.

2. TD learning of state values

3. TD learning of action values: Sarsa

4. TD learning of action values: n -step Sarsa

5. TD learning of optimal action values: Q-learning

6. A unified point of view

7. Summary

TD learning of action values - Sarsa

- ▶ The TD algorithm introduced in the last section can only estimate state values.
- ▶ Next, we introduce Sarsa, an algorithm that can directly estimate action values.
- ▶ We will also see how to use Sarsa to find optimal policies.

Sarsa - Algorithm

First, our aim is to estimate the action values of a given policy π .

Suppose we have some experience $\{(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1})\}_t$.

We can use the following Sarsa algorithm to estimate the action values:

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t) [q_t(s_t, a_t) - [r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})]]$$

$$q_{t+1}(s, a) = q_t(s, a), \quad \forall (s, a) \neq (s_t, a_t)$$

where $t = 0, 1, 2, \dots$

- ▶ $q_t(s_t, a_t)$ is an estimate of $q_\pi(s_t, a_t)$;
- ▶ $\alpha_t(s_t, a_t)$ is the learning rate depending on s_t, a_t .

Sarsa - Algorithm

Why is this algorithm called Sarsa?

That is because each step of the algorithm involves $(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1})$. Sarsa is the abbreviation of state-action-reward-state-action.

What is the relationship between Sarsa and the previous TD learning algorithm?

We can obtain Sarsa by replacing the state value estimate $v(s)$ in the TD algorithm with the action value estimate $q(s, a)$. As a result, Sarsa is an action-value version of the TD algorithm.

What does the Sarsa algorithm do mathematically?

The expression of Sarsa suggests that it is a stochastic approximation algorithm solving the following equation:

$$q_\pi(s, a) = \mathbb{E} [R + \gamma q_\pi(S', A') | s, a], \quad \forall s, a.$$

This is another expression of the Bellman equation expressed in terms of action values.

Sarsa - Algorithm

Theorem (Convergence of Sarsa learning)

By the Sarsa algorithm, $q_t(s, a)$ converges with probability 1 to the action value $q_\pi(s, a)$ as $t \rightarrow \infty$ for all (s, a) if $\sum_t \alpha_t(s, a) = \infty$ and $\sum_t \alpha_t^2(s, a) < \infty$ for all (s, a) .

Remarks:

- ▶ This theorem says that the action values can be found by Sarsa for a given a policy π .

Sarsa - Implementation

The ultimate goal of RL is to find optimal policies. To do that, we can combine Sarsa with a policy improvement step. The combined algorithm is also called Sarsa.

Pseudocode: Policy searching by Sarsa

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 an 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 for (s_t, a_t) :

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t) [q_t(s_t, a_t) - (r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1}))]$$

Update policy for s_t :

$$\pi_{t+1}(a | s_t) = 1 - \frac{\epsilon}{\epsilon} \frac{\epsilon}{|\mathcal{A}(s_t)|} (|\mathcal{A}(s_t)| - 1) \text{ if } a = \arg \max_a q_{t+1}(s_t, a)$$

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

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

Sarsa - Implementation

Remarks about this algorithm:

- ▶ The policy of s_t is updated immediately after $q(s_t, a_t)$ is updated. This is based on the idea of generalized policy iteration.
- ▶ The policy is ϵ -greedy instead of greedy to balance exploitation and exploration well.

Be clear about the core idea and complication:

- ▶ The core idea is simple: to use an algorithm to solve the Bellman equation for a given policy.
- ▶ The complication emerges when we try to find optimal policies and work efficiently.

Sarsa - Examples

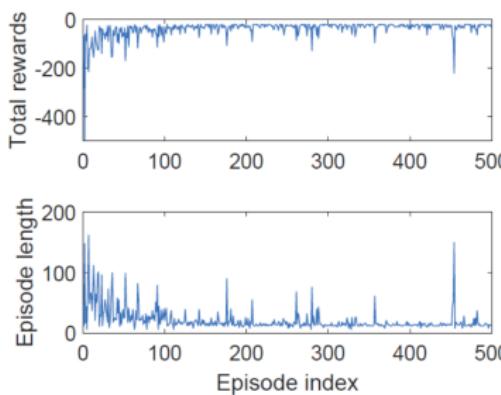
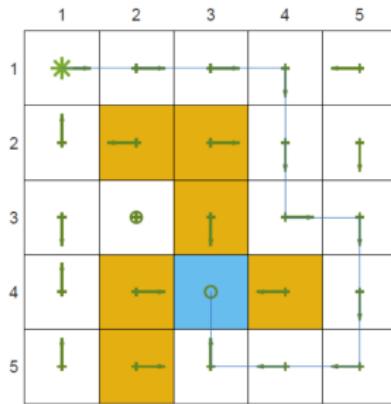
Task description:

- ▶ The task is to find a good path from a specific starting state to the target state.
- This task is different from all the previous tasks where we need to find out the optimal policy for every state!
- Each episode starts from the top-left state and end in the target state.
- In the future, pay attention to what the task is.
- ▶ $r_{\text{target}} = 0$, $r_{\text{forbidden}} = r_{\text{boundary}} = -10$, and $r_{\text{other}} = -1$. The learning rate is $\alpha = 0.1$ and the value of ϵ is 0.1.

Sarsa - Examples

Results:

- ▶ The left figures above show the final policy obtained by Sarsa.
Not all states have the optimal policy.
- ▶ The right figures show the total reward and length of every episode.
The metric of total reward per episode will be frequently used.



2. TD learning of state values

3. TD learning of action values: Sarsa

4. TD learning of action values: *n*-step Sarsa

5. TD learning of optimal action values: Q-learning

6. A unified point of view

7. Summary

TD learning of action values: *n*-step Sarsa

***n*-step Sarsa can unify Sarsa and Monte Carlo learning**

The definition of action value is

$$q_{\pi}(s, a) = \mathbb{E} [G_t \mid S_t = s, A_t = a].$$

The discounted return G_t can be written in different forms as

Sarsa $\leftarrow G_t^{(1)} = R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}),$

$$G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 q_{\pi}(S_{t+2}, A_{t+2}),$$

$$\vdots$$

n-step Sarsa $\leftarrow G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^n q_{\pi}(S_{t+n}, A_{t+n}),$

$$\vdots$$

MC $\leftarrow G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$

It should be noted that $G_t = G_t^{(1)} = G_t^{(2)} = G_t^{(n)} = G_t^{(\infty)}$, where the superscripts merely indicate the different decomposition structures of G_t .

TD learning of action values: *n*-step Sarsa

- ▶ Sarsa aims to solve

$$q_{\pi}(s, a) = \mathbb{E} \left[G_t^{(1)} \mid s, a \right] = \mathbb{E} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid s, a].$$

- ▶ MC learning aims to solve

$$q_{\pi}(s, a) = \mathbb{E} \left[G_t^{(\infty)} \mid s, a \right] = \mathbb{E} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid s, a].$$

- ▶ An intermediate algorithm called *n*-step Sarsa aims to solve

$$q_{\pi}(s, a) = \mathbb{E} \left[G_t^{(n)} \mid s, a \right] = \mathbb{E} [R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n q_{\pi}(S_{t+n}, A_{t+n}) \mid s, a].$$

- ▶ The algorithm of *n*-step Sarsa is

$$\begin{aligned} q_{t+1}(s_t, a_t) &= q_t(s_t, a_t) \\ &- \alpha_t(s_t, a_t) [q_t(s_t, a_t) - [r_{t+1} + \gamma r_{t+2} + \dots + \gamma^n q_t(s_{t+n}, a_{t+n})]]. \end{aligned}$$

- *n*-step Sarsa becomes the (one-step) Sarsa algorithm when $n = 1$.
- *n*-step Sarsa becomes the MC learning algorithm when $n = \infty$.

TD learning of action values: *n*-step Sarsa

- ▶ Data: *n*-step Sarsa needs $(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}, \dots, r_{t+n}, s_{t+n}, a_{t+n})$.
- ▶ Since $(r_{t+n}, s_{t+n}, a_{t+n})$ has not been collected at time t , we are not able to implement *n*-step Sarsa at step t . We need to wait until time $t + n$ to update the q-value of (s_t, a_t) :

$$q_{t+n}(s_t, a_t) = q_{t+n-1}(s_t, a_t)$$

$$- \alpha_{t+n-1}(s_t, a_t) [q_{t+n-1}(s_t, a_t) - [r_{t+1} + \gamma r_{t+2} + \dots + \gamma^n q_{t+n-1}(s_{t+n}, a_{t+n})]]$$

- ▶ Since *n*-step Sarsa includes Sarsa and MC learning as two extreme cases, its performance is a blend of Sarsa and MC learning:
 - If n is large, its performance is close to MC learning and hence has a large variance but a small bias.
 - If n is small, its performance is close to Sarsa and hence has a relatively large bias due to the initial guess and relatively low variance.
- ▶ Finally, *n*-step Sarsa is also for policy evaluation. It can be combined with the policy improvement step to search for optimal policies.

2. TD learning of state values

3. TD learning of action values: Sarsa

4. TD learning of action values: n -step Sarsa

5. TD learning of optimal action values: Q-learning

6. A unified point of view

7. Summary

TD learning of optimal action values: Q-learning

- ▶ Next, we introduce Q-learning, one of the most widely used RL algorithms.
- ▶ Sarsa can estimate the action values of a given policy. It must be combined with a policy improvement step to find optimal policies.
- ▶ Q-learning can directly estimate optimal action values and hence optimal policies.

Q-learning - Algorithm

The Q-learning algorithm is

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t) \left[q_t(s_t, a_t) - \left[r_{t+1} + \gamma \max_{a \in \mathcal{A}} q_t(s_{t+1}, a) \right] \right],$$
$$q_{t+1}(s, a) = q_t(s, a), \quad \forall (s, a) \neq (s_t, a_t),$$

Q-learning is very similar to SARSA. They are different only in terms of the TD target:

- ▶ The TD target in **Q-learning** is $r_{t+1} + \gamma \max_{a \in \mathcal{A}} q_t(s_{t+1}, a)$
- ▶ The TD target in **Sarsa** is $r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})$

Q-learning - Algorithm

What does Q-learning do mathematically?

It aims to solve

$$q(s, a) = \mathbb{E} \left[R_{t+1} + \gamma \max_a q(S_{t+1}, a) \mid S_t = s, A_t = a \right], \quad \forall s, a$$

This is the Bellman optimality equation expressed in terms of action values.

As a result, Q-learning can directly estimate the optimal action values, rather than the action values of a given policy.

Off-policy vs on-policy

Before further studying Q-learning, we first introduce two important concepts:

On-policy learning and off-policy learning.

There exist two policies in a TD learning task:

- ▶ The behavior policy is used to generate experience samples.
- ▶ The target policy is constantly updated toward an optimal policy.

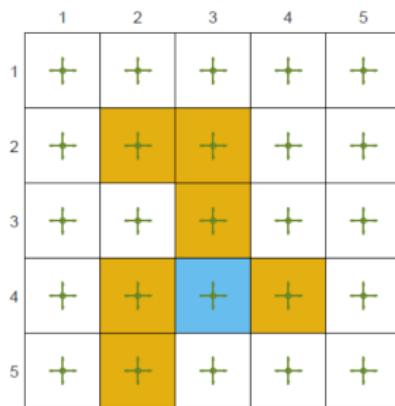
On-policy vs off-policy

- ▶ When the **behavior policy** is the same as the **target policy**, this kind of learning is called on-policy.
- ▶ When they are different, the learning is referred to as off-policy.

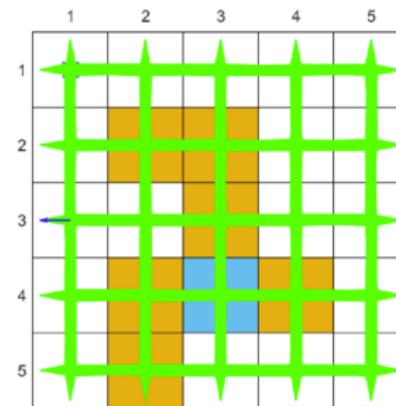
Off-policy vs on-policy

Advantages of off-policy learning:

- ▶ It can search for optimal policies based on the experience samples generated by any other policies.
- ▶ Example: The behavior policy is exploratory, so that we can generate episodes visiting every state-action pair sufficiently many times.



(a) Exploratory behavior policy



(b) Generated episode

Off-policy vs on-policy

How to judge if a TD algorithm is on-policy or off-policy?

- ▶ First, check what math problem the algorithm aims to solve.
- ▶ Second, check what experience samples the algorithm requires.

It deserves special attention because it may not be very clear to beginners.

Off-policy vs on-policy

- ▶ Sarsa aims to evaluate a given policy π by solving

$$q_{\pi}(s, a) = \mathbb{E} [R + \gamma q_{\pi}(S', A') \mid s, a], \quad \forall s, a.$$

where $R \sim p(R \mid s, a)$, $S' \sim p(S' \mid s, a)$, $A' \sim \pi(A' \mid S')$.

- ▶ MC aims to evaluate a given policy π by solving

$$q_{\pi}(s, a) = \mathbb{E} [R_{t+1} + \gamma R_{t+2} + \dots \mid S_t = s, A_t = a], \quad \forall s, a.$$

where the samples are generated by π .

Sarsa and MC are on-policy

- ▶ π is the behavior policy because we need the experience samples generated by π to estimate the action values of π .
- ▶ π is also the target policy because it is updated continuously to approach the optimal policy.

Off-policy vs on-policy

Q-learning is off-policy.

- ▶ First, Q-learning aims to solve the Bellman optimality equation

$$q(s, a) = \mathbb{E} \left[R_{t+1} + \gamma \max_a q(S_{t+1}, a) \mid S_t = s, A_t = a \right], \quad \forall s, a.$$

- ▶ Second, the algorithm is

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t) \left[q_t(s_t, a_t) - \left[r_{t+1} + \gamma \max_{a \in \mathcal{A}} q_t(s_{t+1}, a) \right] \right]$$

which requires $(s_t, a_t, r_{t+1}, s_{t+1})$.

- ▶ The behavior policy is the one for generating a_t in s_t . It can be any policy.

Q-learning - Implementation

Since Q-learning is off-policy, it can be implemented in an either off-policy or on-policy fashion.

Pseudocode: Policy searching by Q-learning (on-policy version)

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 q-value for (s_t, a_t) :

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t) \left[q_t(s_t, a_t) - \left(r_{t+1} + \gamma \max_a q_t(s_{t+1}, a) \right) \right]$$

Update policy for s_t :

$$\pi_{t+1}(a | s_t) = 1 - \frac{\epsilon}{|\mathcal{A}(s_t)|} (|\mathcal{A}(s_t)| - 1) \text{ if } a = \arg \max_a q_{t+1}(s_t, a)$$

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

Q-learning - Algorithm

Pseudocode: Optimal policy search by Q-learning (off-policy version)

Goal: Learn an optimal target policy π_T for all states from the experience samples generated by π_b .

For each episode $\{s_0, a_0, r_1, s_1, a_1, r_2, \dots\}$ generated by π_b , do

For each step $t = 0, 1, 2, \dots$ of the episode, do

Update q -value for (s_t, a_t) :

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t) \left[q(s_t, a_t) - \left(r_{t+1} + \gamma \max_a q_t(s_{t+1}, a) \right) \right]$$

Update target policy for s_t :

$$\pi_{T,t+1}(a | s_t) = 1 \text{ if } a = \arg \max_a q_{t+1}(s_t, a)$$

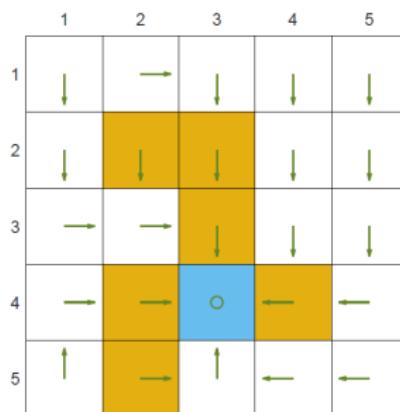
$$\pi_{T,t+1}(a | s_t) = 0 \text{ otherwise}$$

Q-learning - Examples

Task description:

- ▶ The task in these examples is to find an optimal policy for all the states.
- ▶ The reward setting is $r_{\text{boundary}} = r_{\text{forbidden}} = -1$, and $r_{\text{target}} = 1$. The discount rate is $\gamma = 0.9$. The learning rate is $\alpha = 0.1$.

Ground truth: an optimal policy and the corresponding optimal state values.



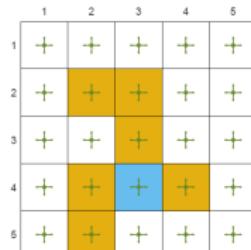
(a) Optimal policy

	1	2	3	4	5
1	5.8	5.6	6.2	6.5	5.8
2	6.5	7.2	8.0	7.2	6.5
3	7.2	8.0	10.0	8.0	7.2
4	8.0	10.0	10.0	10.0	8.0
5	7.2	9.0	10.0	9.0	8.1

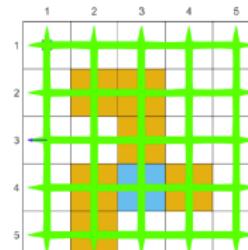
(b) Optimal state value

Q-learning - Examples

The behavior policy and the generated experience (10^5 steps):

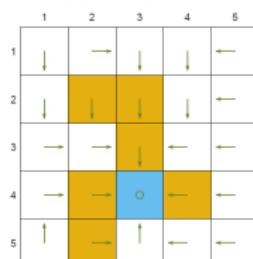


(a) Behavior policy

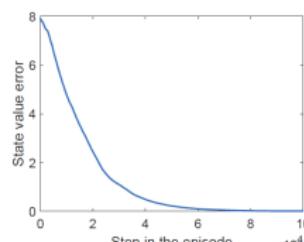


(b) Generated episode

The policy found by off-policy Q-learning:



(a) Estimated policy

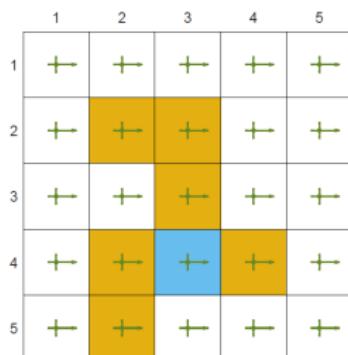


(b) State value error

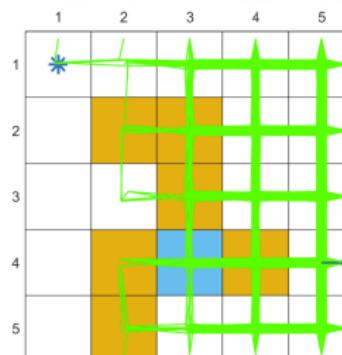
Q-learning - Examples

The importance of exploration: episodes of 10^5 steps

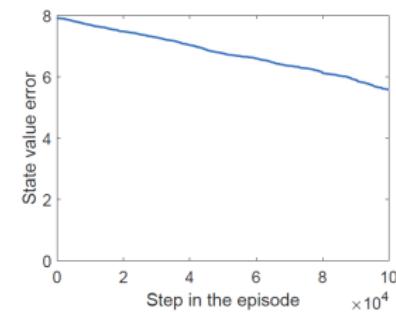
If the policy is not sufficiently exploratory, the samples are not good.



(a) Behavior policy $\epsilon = 0.5$

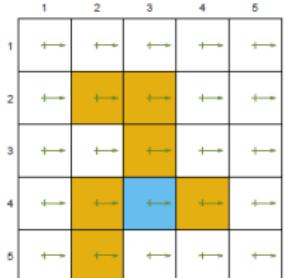


(b) Generated episode

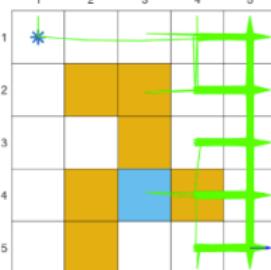


(c) Q-learning result

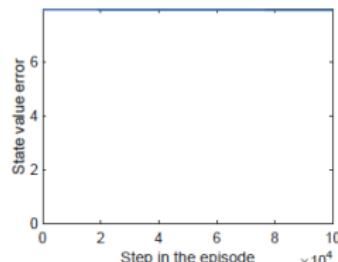
Q-learning - Examples



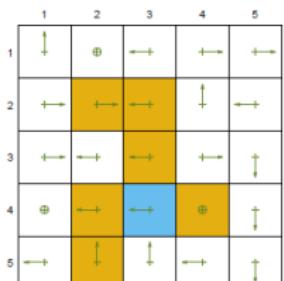
(a) Behavior policy
 $\epsilon = 0.1$



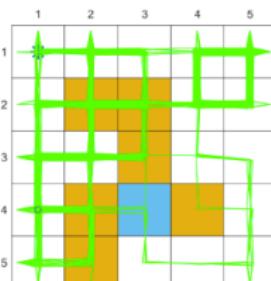
(b) Generated episode



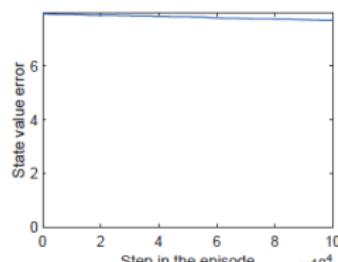
(c) Q-learning result



(a) Behavior policy
 $\epsilon = 0.1$



(b) Generated episode



(c) Q-learning result

2. TD learning of state values

3. TD learning of action values: Sarsa

4. TD learning of action values: n -step Sarsa

5. TD learning of optimal action values: Q-learning

6. A unified point of view

7. Summary

A unified point of view

All the algorithms we introduced in this lecture can be expressed in a unified expression:

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) - \alpha_t(s_t, a_t) [q_t(s_t, a_t) - \bar{q}_t]$$

where \bar{q}_t is the TD target. Different TD algorithms have different \bar{q}_t .

Algorithm	Expression of \bar{q}_t
Sarsa	$\bar{q}_t = r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})$
n -step Sarsa	$\bar{q}_t = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^n q_t(s_{t+n}, a_{t+n})$
Q-learning	$\bar{q}_t = r_{t+1} + \gamma \max_a q_t(s_{t+1}, a)$
Monte Carlo	$\bar{q}_t = r_{t+1} + \gamma r_{t+2} + \dots$

Remark: The MC method can also be expressed in this unified expression by setting $\alpha_t(s_t, a_t) = 1$. In particular, the expression is $q_{t+1}(s_t, a_t) = \bar{q}_t$.

A unified point of view

All the TD algorithms can be viewed as stochastic approximation algorithms solving the Bellman equation or the Bellman optimality equation:

Algorithm	Equation to solve
Sarsa	BE: $q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) S_t = s, A_t = a]$
n -step Sarsa	BE: $q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n q_{\pi}(s_{t+n}, a_{t+n}) S_t = s, A_t = a]$
Q-learning	BOE: $q(s, a) = \mathbb{E}[R_{t+1} + \gamma \max_a q(S_{t+1}, a) S_t = s, A_t = a]$
Monte Carlo	BE: $q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots S_t = s, A_t = a]$

2. TD learning of state values
3. TD learning of action values: Sarsa
4. TD learning of action values: n -step Sarsa
5. TD learning of optimal action values: Q-learning
6. A unified point of view
7. Summary

Summary

- ▶ Introduced various TD learning algorithms
- ▶ Their expressions, math interpretations, implementation, relationship, examples
- ▶ Unified point of view