# Lecture 7: Temporal-Difference Learning

Shiyu Zhao

Department of Artificial Intelligence Westlake University

## Outline



- This lecture introduces temporal-difference (TD) learning, which is 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.
- We will see how the stochastic approximation methods studied in the last lecture are useful.

## 1 Motivating examples

- 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

## 1 Motivating examples

- 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

## Motivating example: stochastic algorithms

We next consider some stochastic problems and show how to use the RM algorithm to solve them.

First, revisit the mean estimation problem: calculate

$$w = \mathbb{E}[X]$$

based on some iid samples  $\{x\}$  of X. We studied it in the last lecture.

• By writing  $g(w) = w - \mathbb{E}[X]$ , we can reformulate the problem to a root-finding problem

$$g(w) = 0$$

• Since we can only obtain samples  $\{x\}$  of X, the noisy observation is

$$\tilde{g}(w,\eta) = w - x = (w - \mathbb{E}[X]) + (\mathbb{E}[X] - x) \doteq g(w) + \eta$$

• According to the last lecture, we know the RM algorithm for solving g(w) = 0 is

$$w_{k+1} = w_k - \alpha_k \tilde{g}(w_k, \eta_k) = w_k - \alpha_k (w_k - x_k)$$

Second, consider a little more complex problem. That is to estimate the mean of a function v(X),

$$w = \mathbb{E}[v(X)],$$

based on some iid random samples  $\{x\}$  of X.

• To solve this problem, we define

 $g(w) = w - \mathbb{E}[v(X)]$  $\tilde{g}(w,\eta) = w - v(x) = (w - \mathbb{E}[v(X)]) + (\mathbb{E}[v(X)] - v(x)) \doteq g(w) + \eta.$ 

• Then, the problem becomes a root-finding problem: g(w)=0. The corresponding RM algorithm is

$$w_{k+1} = w_k - \alpha_k \tilde{g}(w_k, \eta_k) = w_k - \alpha_k [w_k - v(x_k)]$$

#### Third, consider an even more complex problem: calculate

 $w = \mathbb{E}[R + \gamma v(X)],$ 

where R, X are random variables,  $\gamma$  is a constant, and  $v(\cdot)$  is a function.

• Suppose we can obtain samples  $\{x\}$  and  $\{r\}$  of X and R. We define

$$g(w) = w - \mathbb{E}[R + \gamma v(X)],$$
  

$$\tilde{g}(w, \eta) = w - [r + \gamma v(x)]$$
  

$$= (w - \mathbb{E}[R + \gamma v(X)]) + (\mathbb{E}[R + \gamma v(X)] - [r + \gamma v(x)])$$
  

$$\stackrel{.}{=} g(w) + \eta.$$

• Then, the problem becomes a root-finding problem: g(w) = 0. The corresponding RM algorithm is

 $w_{k+1} = w_k - \alpha_k \tilde{g}(w_k, \eta_k) = w_k - \alpha_k [w_k - (r_k + \gamma v(x_k))]$ 

This algorithm looks like TD algorithms as shown later.

Quick summary:

- The above three examples become more and more complex.
- They can all be solved by the RM algorithm.
- We will see that the TD algorithms have similar expressions.

#### 1 Motivating examples

- 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

Problem statement:

- Given policy  $\pi$ , the aim is to estimate the state values  $\{v_{\pi}(s)\}_{s\in\mathcal{S}}$  under  $\pi$ .
- 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  $\pi$ .

Important notations:

$$\begin{array}{c} v(s) \longrightarrow v_{\pi}(s) \\ & \Downarrow \\ v(s_t) \longrightarrow v_{\pi}(s_t) \\ & \Downarrow \\ v_t(s_t) \longrightarrow v_{\pi}(s_t) \end{array}$$

The TD learning algorithm is

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

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

where t = 0, 1, 2, ...

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 st is updated whereas the values of the unvisited states s ≠ st 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



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,

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

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

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

Therefore,

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

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

Shiyu Zhao

#### 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 two time steps.
- It reflects the difference between  $v_t$  and  $v_{\pi}$ . To see that, denote

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

Note that

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

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

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.

#### Q: What does this TD algorithm do mathematically?

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

• Chapter 2 has introduced the model-based algorithm for solving the Bellman equation: closed-form solution + iterative algorithm.

#### First, a new expression of the Bellman equation.

The definition of state value of  $\pi$  is

$$v_{\pi}(s) = \mathbb{E}[R + \gamma G | S = s], \quad s \in \mathcal{S}$$
(4)

where  $\boldsymbol{G}$  is discounted return. Since

$$\mathbb{E}[G|S=s] = \sum_{a} \pi(a|s) \sum_{s'} p(s'|s, a) v_{\pi}(s') = \mathbb{E}[v_{\pi}(S')|S=s],$$

where S' is the next state, we can rewrite (4) as

$$v_{\pi}(s) = \mathbb{E}[R + \gamma v_{\pi}(S')|S = s], \quad s \in \mathcal{S}.$$
(5)

Equation (5) is another expression of the Bellman equation. It is sometimes called the Bellman expectation equation, an important tool to design and analyze TD algorithms.

Second, solve the Bellman equation in (5) using the RM algorithm. In particular, by defining

$$g(v(s)) = v(s) - \mathbb{E}[R + \gamma v_{\pi}(S')|s],$$

we can rewrite (5) as

g(v(s)) = 0.

Since we can only obtain the samples  $r \mbox{ and } s' \mbox{ of } R \mbox{ and } S',$  the noisy observation we have is

$$\tilde{g}(v(s)) = v(s) - [r + \gamma v_{\pi}(s')]$$

$$= \underbrace{\left(v(s) - \mathbb{E}[R + \gamma v_{\pi}(S')|s]\right)}_{g(v(s))} + \underbrace{\left(\mathbb{E}[R + \gamma v_{\pi}(S')|s] - [r + \gamma v_{\pi}(s')]\right)}_{\eta}.$$

## TD learning of state values - The idea of the algorithm

Therefore, the RM algorithm for solving g(v(s)) = 0 is

$$v_{k+1}(s) = v_k(s) - \alpha_k \tilde{g}(v_k(s)) = v_k(s) - \alpha_k \Big( v_k(s) - \big[ \mathbf{r}_k + \gamma v_\pi(s'_k) \big] \Big), \quad k = 1, 2, 3, \dots$$
(6)

where  $v_k(s)$  is the estimate of  $v_{\pi}(s)$  at the kth step;  $r_k, s'_k$  are the samples of R, S' obtained at the kth step.

The RM algorithm in (6) looks very similar to the TD algorithm. However, there are **two differences**.

- Difference 1: The RM algorithm requires  $\{(s, r_k, s'_k)\}$  for k = 1, 2, 3, ...
  - Modification:  $\{(s, r_k, s'_k)\}$  is changed to  $\{(s_t, r_{t+1}, s_{t+1})\}$  so that the algorithm can utilize the sequential samples in an episode.
- Difference 2: The RM algorithm requires  $v_{\pi}(s'_k)$ .
  - Modification:  $v_{\pi}(s'_k)$  is replaced by an estimate  $v_t(s_{t+1})$ .

With the above modifications, the RM algorithms becomes exactly the TD algorithm.

#### Theorem (Convergence of TD Learning)

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

The proof of the theorem can be found in my book. Remarks:

- This theorem says the state value can be found by the TD algorithm for a given a policy  $\pi.$
- $\sum_t \alpha_t(s) = \infty$  and  $\sum_t \alpha_t^2(s) < \infty$  must be valid for all  $s \in S$ .
  - For condition  $\sum_t \alpha_t(s) = \infty$ : At time step t,

 $\diamond \ \, {\rm If} \ s=s_t, \ {\rm then} \ \alpha_t(s)>0;$ 

 $\diamond$  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  $\alpha$  is often selected as a small constant. In this case, the condition that  $\sum_t \alpha_t^2(s) < \infty$  is invalid anymore. When  $\alpha$  is constant, it can still be shown that the algorithm converges in the sense of expectation sense.

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

TD/Sarsa learning	MC learning
<b>Online:</b> TD learning is online. It can up- date the state/action values immediately after receiving a reward.	<b>Offline:</b> MC learning is offline. It has to wait until an episode has been completely collected.
<b>Continuing tasks:</b> Since TD learning is online, it can handle both episodic and continuing tasks.	<b>Episodic tasks:</b> Since MC learning is of- fline, it can only handle episodic tasks that has terminate states.

Table: Comparison between TD learning and MC learning.

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

TD/Sarsa learning	MC learning
<b>Bootstrapping:</b> TD bootstraps because the update of a value relies on the pre- vious estimate of this value. Hence, it requires initial guesses.	<b>Non-bootstrapping:</b> MC is not boot- strapping, because it can directly es- timate state/action values without any initial guess.
<b>Low estimation variance:</b> TD has lower than MC because there are fewer random variables. For instance, Sarsa requires $R_{t+1}, S_{t+1}, A_{t+1}$ .	<b>High estimation variance:</b> 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.

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

#### 1 Motivating examples

- 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

- 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.

First, our aim is to estimate the action values of a given policy  $\pi$ .

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:

$$\begin{aligned} q_{t+1}(s_t, a_t) &= q_t(s_t, a_t) - \alpha_t(s_t, a_t) \Big[ q_t(s_t, a_t) - [r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})] \Big], \\ q_{t+1}(s, a) &= q_t(s, a), \quad \forall (s, a) \neq (s_t, a_t), \end{aligned}$$

where t = 0, 1, 2, ...

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

- 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:

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

This is another expression of the Bellman equation expressed in terms of action values. The proof is given in my book.

#### 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  $t \to \infty$  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  $\pi$ .

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, ...) 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) \Big[ q_t(s_t, a_t) - (r_{t+1} + q_t) \Big] \Big] 
                   \gamma q_t(s_{t+1}, a_{t+1}))
             Update policy for s_t:
                   \pi_{t+1}(a|s_t) = 1 - \frac{\epsilon}{|\mathcal{A}(s_t)|}(|\mathcal{A}(s_t)| - 1) if a = \arg \max_a q_{t+1}(s_t, a)
                   \pi_{t+1}(a|s_t) = \frac{\epsilon}{|A(s_t)|} otherwise
             s_t \leftarrow s_{t+1}, a_t \leftarrow a_{t+1}
```

#### 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 *e*-greedy instead of greedy to well balance exploitation and exploration.

#### Be clear about the core idea and complication:

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

#### 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  $\epsilon$  is 0.1.

#### **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.



#### 1 Motivating examples

- 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

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

The definition of action value is

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

The discounted return  $G_t$  can be written in different forms as

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

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  $G_t$ .

Shiyu Zhao

n

• Sarsa aims to solve

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

• MC learning aims to solve

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

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

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

• The algorithm of *n*-step Sarsa is

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

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

- Data: *n*-step Sarsa needs  $(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}, \ldots, r_{t+n}, s_{t+n}, a_{t+n})$ .
- Since (r<sub>t+n</sub>, s<sub>t+n</sub>, a<sub>t+n</sub>) 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<sub>t</sub>, a<sub>t</sub>):

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

- 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.

#### 1 Motivating examples

- 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

- 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.

#### The Q-learning algorithm is

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

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  $r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})$

#### What does Q-learning do mathematically?

It aims to solve

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

This is the Bellman optimality equation expressed in terms of action values. See the proof in my book.

As a result, Q-learning can directly estimate the optimal action values instead of action values of a given 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, such kind of learning is called on-policy.
- When they are different, the learning is called off-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.





#### 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 be confusing to beginners.

# Off-policy vs on-policy

• Sarsa aims to evaluate a given policy  $\boldsymbol{\pi}$  by solving

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

where  $R \sim p(R|s,a)\text{, }S' \sim p(S'|s,a)\text{, }A' \sim \pi(A'|S').$ 

• MC aims to evaluate a given policy  $\pi$  by solving

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

where the samples are generated by  $\pi$ .

#### Both Sarsa and MC are on-policy.

- $\pi$  is the behavior policy because we need the experience samples generated by  $\pi$  to estimate the action values of  $\pi$ .
- π is also the target policy because it is updated continuously so that it approaches the optimal policy.

## Q-learning is off-policy.

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

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

• Second, the algorithm is

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

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



See the book for more detailed pseudocode.

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

 $\begin{array}{l} \textbf{Goal: Learn an optimal target policy } \pi_T \text{ for all states from the experience samples} \\ \textbf{generated by } \pi_b. \\ \textbf{For each episode } \{s_0, a_0, r_1, s_1, a_1, r_2, \ldots\} \text{ generated by } \pi_b, \text{ do} \\ \textbf{For each step } t = 0, 1, 2, \ldots \text{ of the episode, do} \\ \textbf{Update } q\text{-value for } (s_t, a_t): \\ q_{t+1}(s_t, a_t) &= q_t(s_t, a_t) - \alpha_t(s_t, a_t) \Big[ q(s_t, a_t) - (r_{t+1} + \gamma \max_a q_t(s_{t+1}, a)) \Big] \\ \textbf{Update target policy for } s_t: \\ \pi_{T,t+1}(a|s_t) = 1 \text{ if } a = \arg \max_a q_{t+1}(s_t, a) \\ \pi_{T,t+1}(a|s_t) = 0 \text{ otherwise} \end{array}$ 

See the book for more detailed pseudocode.

#### 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.



	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



# 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:



## The importance of exploration: episodes of $10^5$ steps

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



# Q-learning – Examples



#### 1 Motivating examples

- 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

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	$\bar{q}_t = r_{t+1} + \gamma q_t(s_{t+1}, a_{t+1})$
<i>n</i> -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  $q_{t+1}(s_t, a_t) = \bar{q}_t$ .

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

Algorithm	Equation to solve	
Sarsa	$BE: \ q_{\pi}(s,a) = \mathbb{E} \left[ R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1})   S_t = s, A_t = a \right]$	
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: $q(s, a) = \mathbb{E} \left[ R_{t+1} + \gamma \max_{a} q(S_{t+1}, a) \middle  S_t = s, A_t = a \right]$	
Monte Carlo	BE: $q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots   S_t = s, A_t = a]$	

#### 1 Motivating examples

- 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

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