**Iran University of Science and Technology**

# یادگیری تقویتی در کنترل

**دکتر سعید شمقدری**

**دانشکده مهندسی برق**
**گروه کنترل**

**نیمسال اول ۱۴۰۵-۱۴۰۴**

# Temporal Difference Learning

# Temporal Difference Learning

**Q Box**

مشابهت با DP؟

مشابهت با MC؟

**Reminder Box**

تخمین تابع Value در DP و MC

$$V(S_t) \leftarrow V(S_t) + \alpha[G_t - V(S_t)]$$

New value of state t — Former estimation of value of state t (= Expected return starting at that state) — Learning Rate — Return at timestep t — Former estimation of value of state t (= Expected return starting at that state)

روش Constant alpha در MC (نیاز به اتمام اپیزود)

استفاده از Return به عنوان تخمین امید ریاضی (Sampling)

## Temporal Difference Prediction

تخمین تابع Value در TD

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

New value of state t

Former estimation of value of state t

Learning Rate

Reward

Discounted value of next state

**Q Box:**
Target in MC and TD?

## Temporal Difference Prediction

تخمین تابع Value در TD

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

New value of state t

Former estimation of value of state t

Learning Rate

Reward

Discounted value of next state

TD Target

**Q Box:**
Target in MC and TD?

TD(0): one step TD

جمع‌بندی ...

جمع‌بندی ...



## TD Approach:

At the end of one step (State, Action, Reward, Next State):
- We have Rt+1 and St+1
- We update V(St):
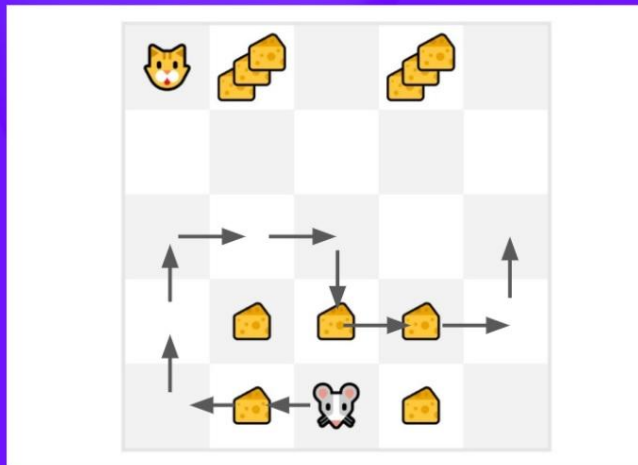  - **We estimate Gt** by adding Rt+1 and the discounted value of next state.
  **TD target** : Rt+1 + gamma * V(St+1)

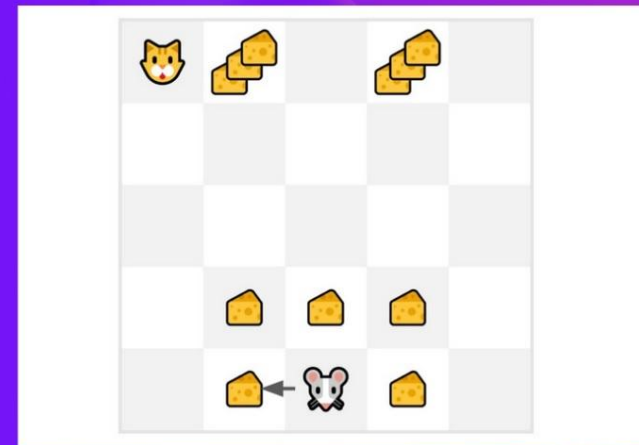$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

Now we **continue to interact with this environment** with our updated value function. By running more and more steps, **the agent will learn to play better and better.**

**جمع‌بندی ...**

## Temporal Difference Prediction

<div dir="rtl">

تخمین تابع Value در TD

</div>

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t \mid S_t = s]$$

$$= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s]$$

$$= \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s]$$

<div dir="rtl">

- DP: محاسبه امید ریاضی با مدل، استفاده از تخمین قبلی تابع Value در تخمین جدید (bootstrapping)
- TD: Sampling، استفاده از تخمین قبلی تابع Value در تخمین جدید (bootstrapping)

</div>

الگوریتم

## Tabular TD(0) for estimating $v_\pi$

Input: the policy $\pi$ to be evaluated
Algorithm parameter: step size $\alpha \in (0, 1]$
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(terminal) = 0$

Loop for each episode:
    Initialize $S$
    Loop for each step of episode:
        $A \leftarrow$ action given by $\pi$ for $S$
        Take action $A$, observe $R$, $S'$
        $V(S) \leftarrow V(S) + \alpha \big[ R + \gamma V(S') - V(S) \big]$
        $S \leftarrow S'$
    until $S$ is terminal

## TD Error

$$\delta_t \doteq R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

## Monte Carlo Error

$$
\begin{aligned}
G_t - V(S_t) &= R_{t+1} + \gamma G_{t+1} - V(S_t) + \gamma V(S_{t+1}) - \gamma V(S_{t+1}) \\
&= \delta_t + \gamma \big( G_{t+1} - V(S_{t+1}) \big) \\
&= \delta_t + \gamma \delta_{t+1} + \gamma^2 \big( G_{t+2} - V(S_{t+2}) \big) \\
&= \delta_t + \gamma \delta_{t+1} + \gamma^2 \delta_{t+2} + \cdots + \gamma^{T-t-1} \delta_{T-1} + \gamma^{T-t} \big( G_T - V(S_T) \big) \\
&= \delta_t + \gamma \delta_{t+1} + \gamma^2 \delta_{t+2} + \cdots + \gamma^{T-t-1} \delta_{T-1} + \gamma^{T-t} (0 - 0) \\
&= \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k.
\end{aligned}
$$

## Driving Home
مثال

Value function approximation TD

| State | Elapsed Time (minutes) | Predicted Time to Go | Predicted Total Time |
|---|---|---|---|
| leaving office, friday at 6 | 0 | 30 | 30 |
| reach car, raining | 5 | 35 | 40 |
| exiting highway | 20 | 15 | 35 |
| 2ndary road, behind truck | 30 | 10 | 40 |
| entering home street | 40 | 3 | 43 |
| arrive home | 43 | 0 | 43 |

Reward?
Return? T_Go
Value? E(T_Go)

**Q Box**

## مثال

**Driving Home**
Value function approximation TD

| State | Elapsed Time (minutes) | Predicted Time to Go | Predicted Total Time |
|---|---|---|---|
| leaving office, friday at 6 | 0 | 30 | 30 |
| reach car, raining | 5 | 35 | 40 |
| exiting highway | 20 | 15 | 35 |
| 2ndary road, behind truck | 30 | 10 | 40 |
| entering home street | 40 | 3 | 43 |
| arrive home | 43 | 0 | 43 |

**Q Box**

خطای تخمین در خروجی بزرگراه (بر اساس MC)؟

اصلاح تخمین در خروجی بزرگراه (بر اساس MC) به ازای alpha=0.5؟

$$\alpha\big(G_t - V(s_t)\big)?$$

**TD: یادگیری سریعتر**

# **Temporal Difference Learning**

<div dir="rtl">

همگرایی TD

- مقدار α بسیار کوچک
- α با شرایط کاهشی

</div>

$$\sum_{n=1}^{\infty} \alpha_n(a) = \infty \qquad \text{and} \qquad \sum_{n=1}^{\infty} \alpha_n^2(a) < \infty$$

<div dir="rtl">

سرعت همگرایی MC و TD

</div>

# Markov Reward Process: Random Walk مثال

In this example we empirically compare the prediction abilities of TD(0) and constant-α MC when applied to the following Markov reward process:

# Random Walk Under Batch Updating



**Figure 6.2:** Performance of TD(0) and constant-$\alpha$ MC under batch training on the random walk task.

**You are the Predictor**

<div dir="rtl">مثال</div>

$A, 0, B, 0$

$B, 1$

$B, 1$

$B, 1$

$B, 1$

$B, 1$

$B, 1$

$B, 0$

## SARSA: On-Policy TD Control



$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t) \right]$$

## SARSA: On-Policy TD Control



$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t) \right]$$

**Target Q Value**

**Updated Q Value**    **Current Q Value**    **Current Q Value**

**Target Policy is always as same as Behaviour Policy**

الگوریتم

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
    Loop for each step of episode:
        Take action $A$, observe $R$, $S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

درک بهتر!



**Select Action**

Agent chooses
action based on
policy

**Observe State**

Agent perceives
current environment

**Execute Action** !

Agent performs
selected action

**Observe Reward**

Agent receives
feedback from
environment

**Update Q-value**

Agent adjusts Q-
values based on
experience

# Windy Grid

مثال

# Windy Grid

مثال

# Q-Learning: Off-Policy TD Control



$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

New
Q-value
estimation

Former
Q-value
estimation

Learning
Rate

Immediate
Reward

Discounted Estimate
optimal Q-value
of next state

Former
Q-value
estimation

TD Target

TD Error

I TD Learning

# Example

- The reward function:

  - 0: Going to a state **with no cheese in it.**

  - +1: Going to a state with a **small cheese in it.**

  - +10: Going to the state with **the big pile of cheese.**

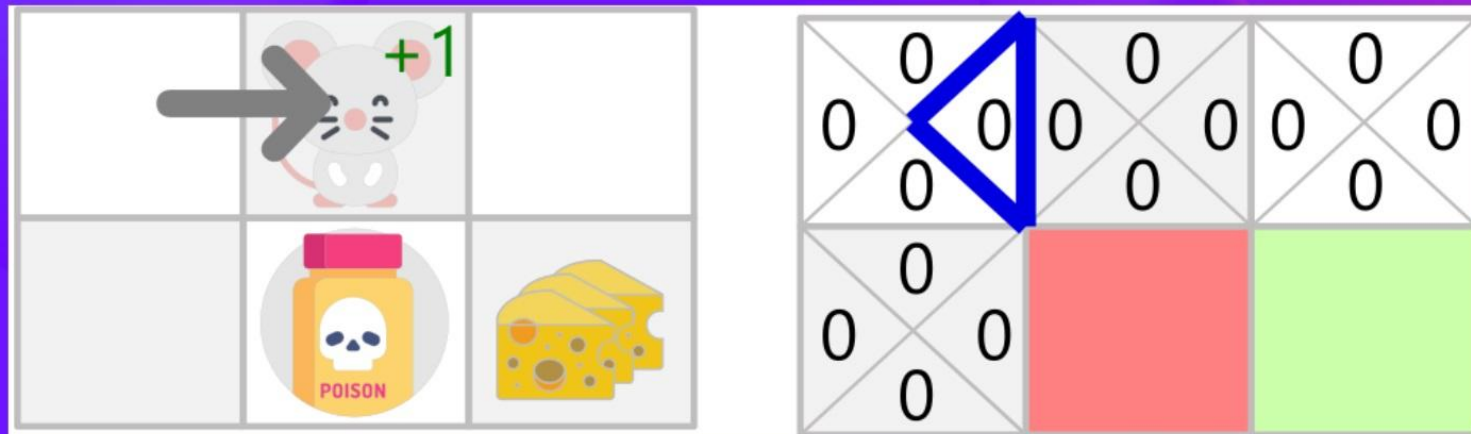  - -10: Going to the state **with the poison and thus die.**

# Example, Step 1

Initialize $Q$ arbitrarily (e.g., $Q(s, a) = 0$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$, and $Q(terminal\text{-}state, \cdot) = 0$)

| | ← | → | ↑ | ↓ |
|---|---|---|---|---|
| 🐭 | 0 | 0 | 0 | 0 |
| 🧀 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 |
| 🧪 | 0 | 0 | 0 | 0 |
| 🧀 | 0 | 0 | 0 | 0 |

## We initialize the Q-Table

# Example, Step 4

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

New
Q-value
estimation

Former
Q-value
estimation

Learning
Rate

Immediate
Reward

Discounted Estimate
optimal Q-value
of next state

Former
Q-value
estimation

TD Target

TD Error

## Update our Q-value estimation

# Example, Step 4

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t))$$

Q(Initial state, Right)  = 0 + 0.1 * [1 + 0.99 * 0 - 0]
Q(Initial state, Right)  = 0.1

| | ← | → | ↑ | ↓ |
|---|---|---|---|---|
| 🐭 | 0 | 0.1 | 0 | 0 |
| 🧀 | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 |
| | 0 | 0 | 0 | 0 |
| 🫙 | 0 | 0 | 0 | 0 |
| 🧀 | 0 | 0 | 0 | 0 |

# Q-Learning Recap

الگوریتم

**Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:

    Initialize $S$

    Loop for each step of episode:

        Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
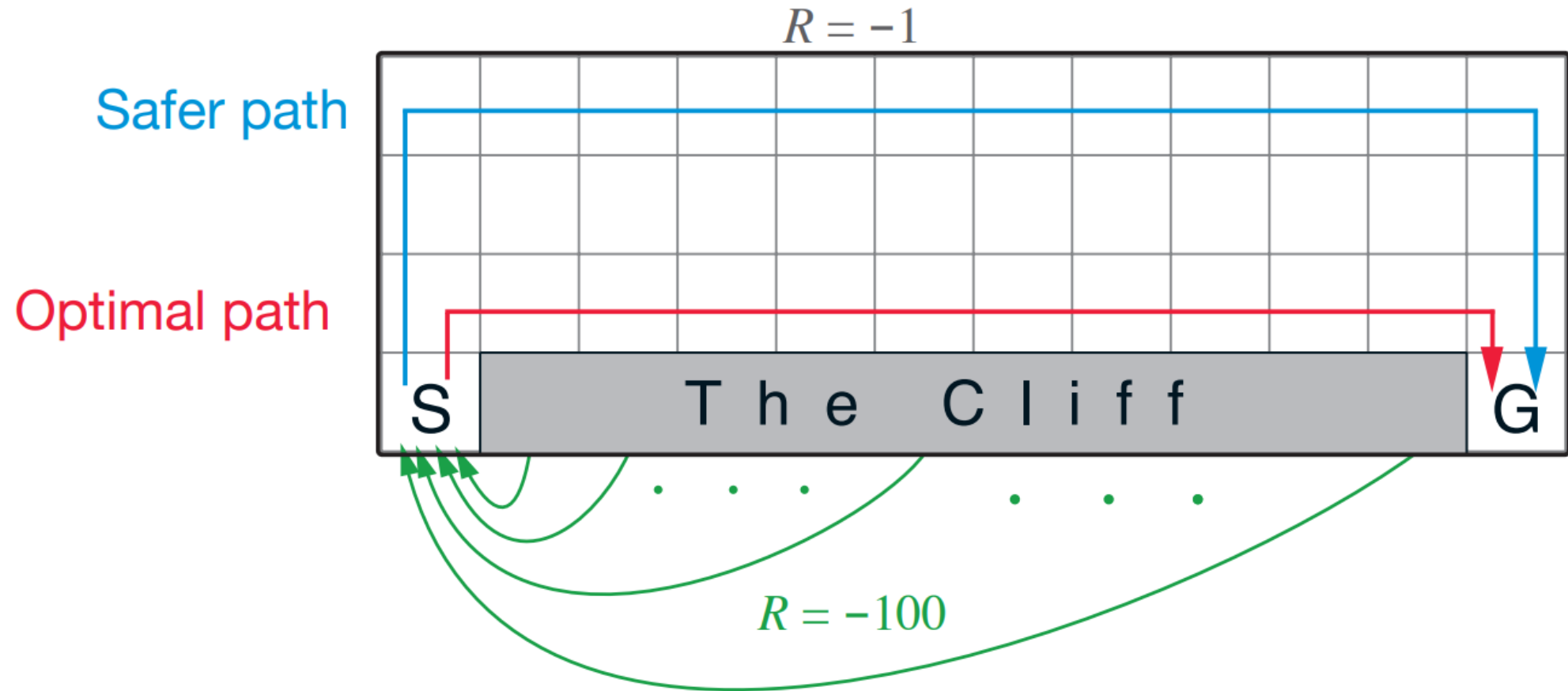
        Take action $A$, observe $R, S'$

        $Q(S, A) \leftarrow Q(S, A) + \alpha \left[ R + \gamma \max_a Q(S', a) - Q(S, A) \right]$

        $S \leftarrow S'$

    until $S$ is terminal

# Cliff Walking

مثال

# Cliff Walking

مثال