



**Iran University  
of Science and  
Technology**

In the Name of God

# Reinforcement Learning in Control

**Dr. Saeed Shamaghdari**

**Electrical Engineering Department  
Control Group**

Fall 2025 | 4041

# Temporal Difference Learning

## Temporal Difference Learning

### Q Box:

Similarity to dynamic programming (DP)?  
Similarity to Monte Carlo methods (MCs)?

### Reminder Box

Value function approximation in DP and MC

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

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)

The Constant-alpha method in MC (requires episode completion)  
Using Return as an estimate of the expected value (Sampling)

## Temporal Difference Prediction

Value function approximation TD

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

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

Value function approximation TD

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

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

## To Recap ...

### TD Learning Approach:

*Temporal Difference Learning*: learning at each time step.

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

New value  
of state t

Former  
estimation of  
value of state  
t

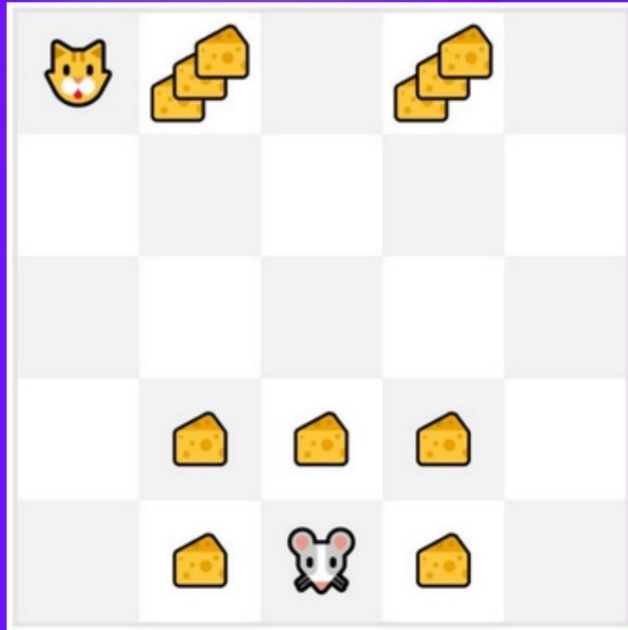
Learning Rate  
Reward

Discounted value of next  
state

TD Target

## To Recap ...

### TD Approach:



At the end of one step (State, Action, Reward, Next State):

- We have  $R_{t+1}$  and  $S_{t+1}$
  - We update  $V(S_t)$ :
    - We estimate  $G_t$  by adding  $R_{t+1}$  and the discounted value of next state.
- TD target :  $R_{t+1} + \gamma V(S_{t+1})$

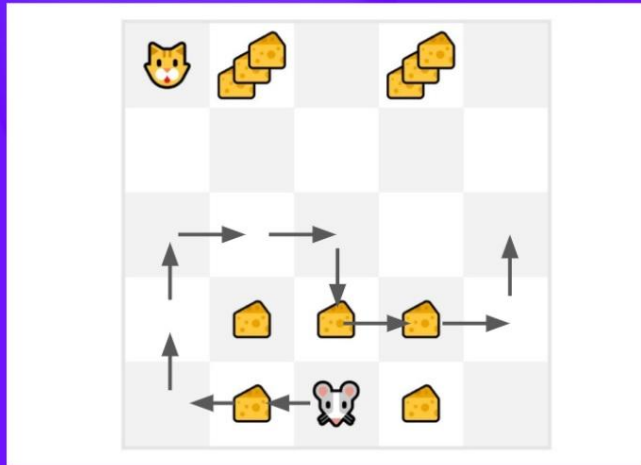
$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

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.**

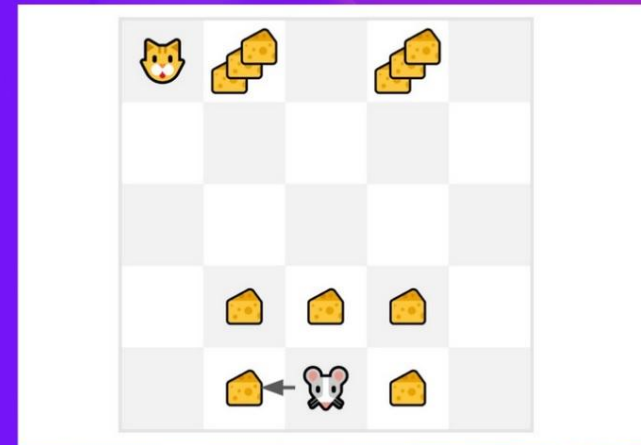
## To Recap ...

### Two Learning Approaches:

*Monte Carlo: learning at the end of the episode*



*Temporal Difference Learning: learning at each step.*





## Temporal Difference Prediction

Value function approximation TD

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

- **DP:** Computing the expected value using the model, utilizing previous value function estimates in new estimates (**bootstrapping**).
- **TD:** Sampling, using previous value function estimates in new estimates (**bootstrapping**).

## Algorithm

### 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  $V(\text{terminal}) = 0$

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

$A \leftarrow$  action given by  $\pi$  for  $S$

        Take action  $A$ , observe  $R, S'$

$V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

    until  $S$  is terminal

## TD Error

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

## Monte Carlo Error

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

**Example: Driving Home**

Value function approximation TD

<i>State</i>	<i>Elapsed Time (minutes)</i>	<i>Predicted Time to Go</i>	<i>Predicted Total Time</i>
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**

**Example: Driving Home**

Value function approximation TD

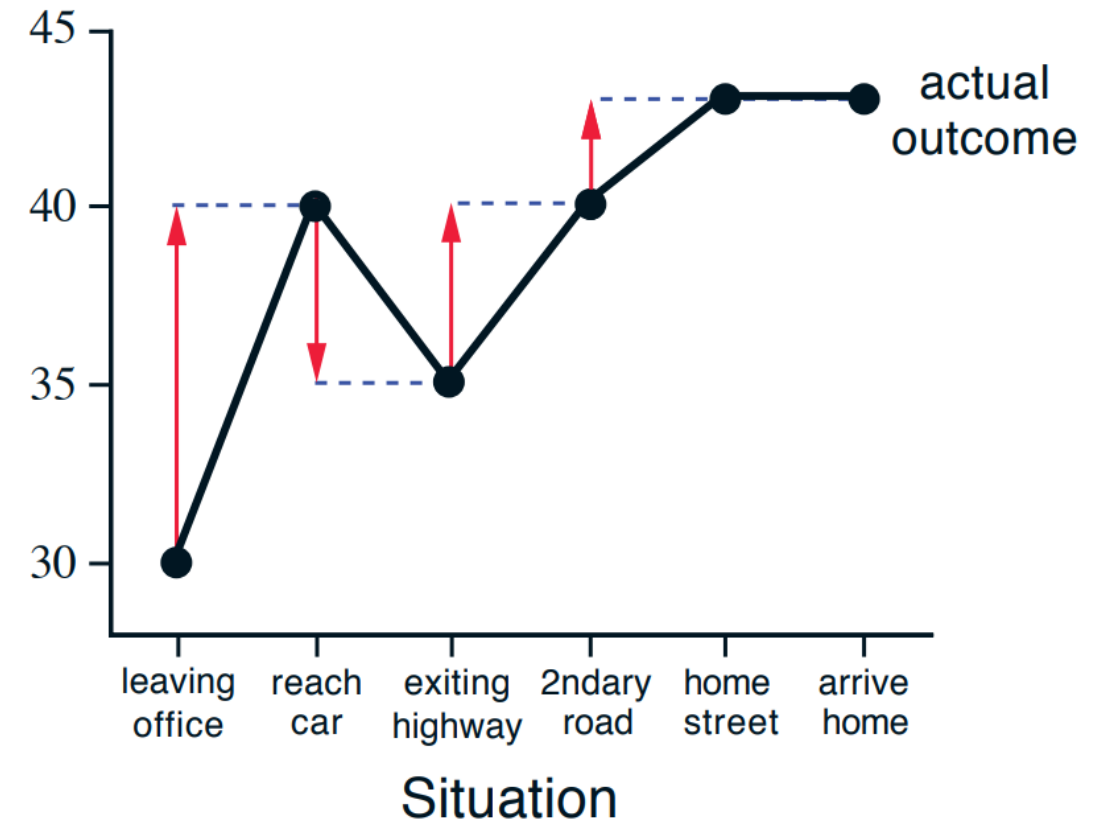
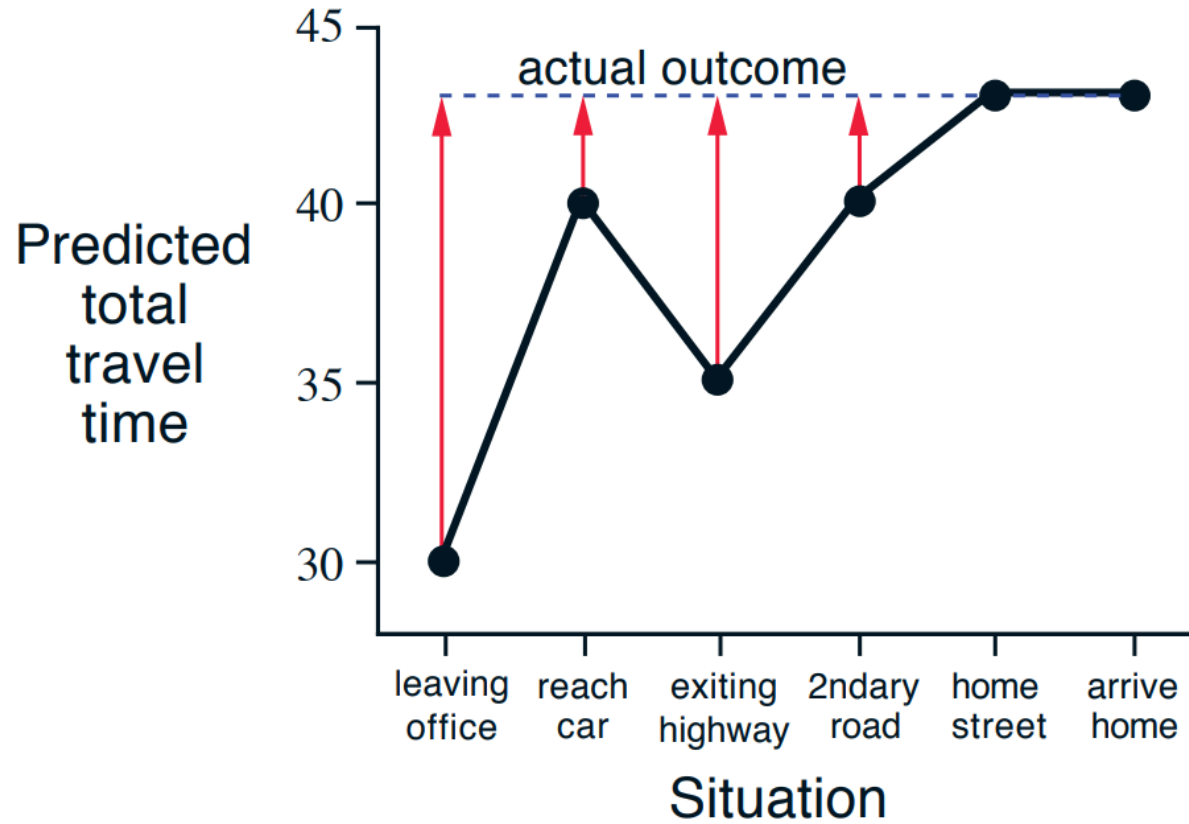
<i>State</i>	<i>Elapsed Time (minutes)</i>	<i>Predicted Time to Go</i>	<i>Predicted Total Time</i>
leaving office, friday at 6	0	30	30
reach car, raining	5	35	40
exiting highway	20	15	35
2ndary road, behind truck	30	10	40
entering home street	40	3	43
arrive home	43	0	43

Estimation error at highway exit (based on Monte Carlo)?

**Q Box**Estimation correction at highway exit (MC-based) for  $\alpha=0.5$ ?

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

## TD: Faster Learning



## Temporal Difference Learning

### TD Convergence

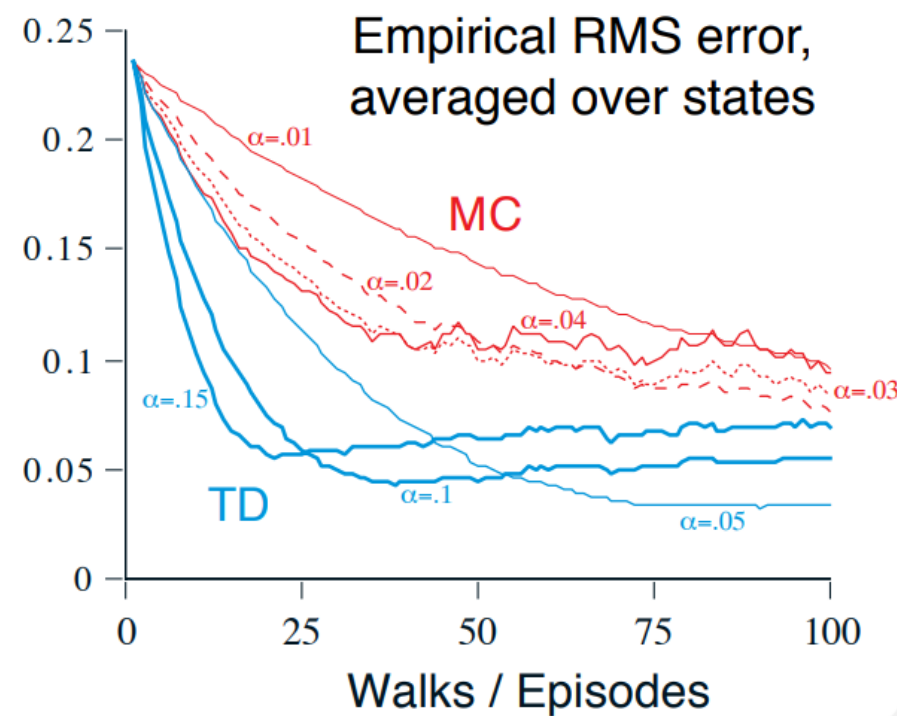
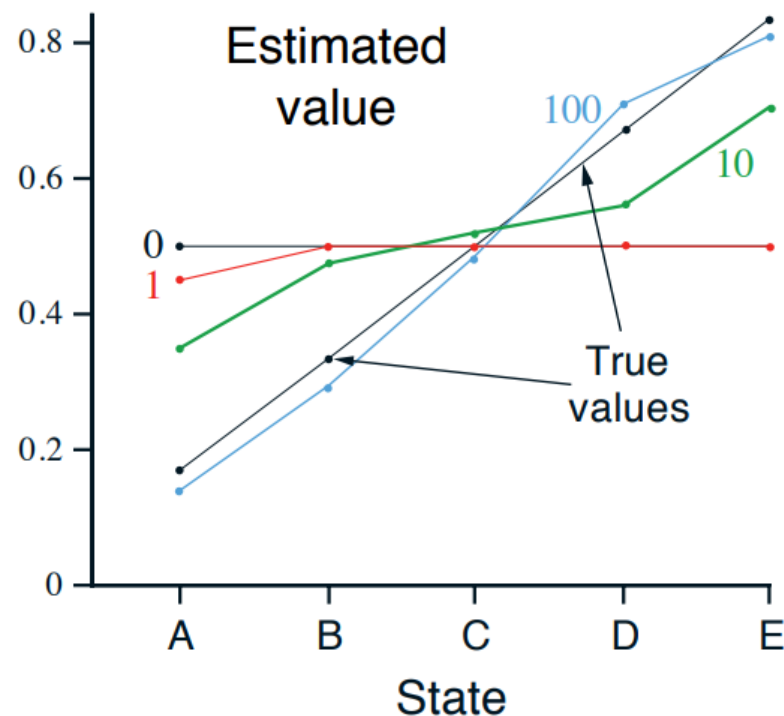
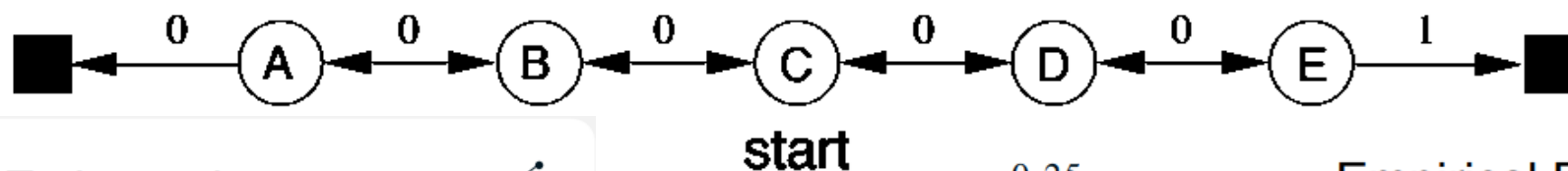
- Very small  $\alpha$
- $\alpha$  with decay conditions

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

Convergence speed of TD and MC

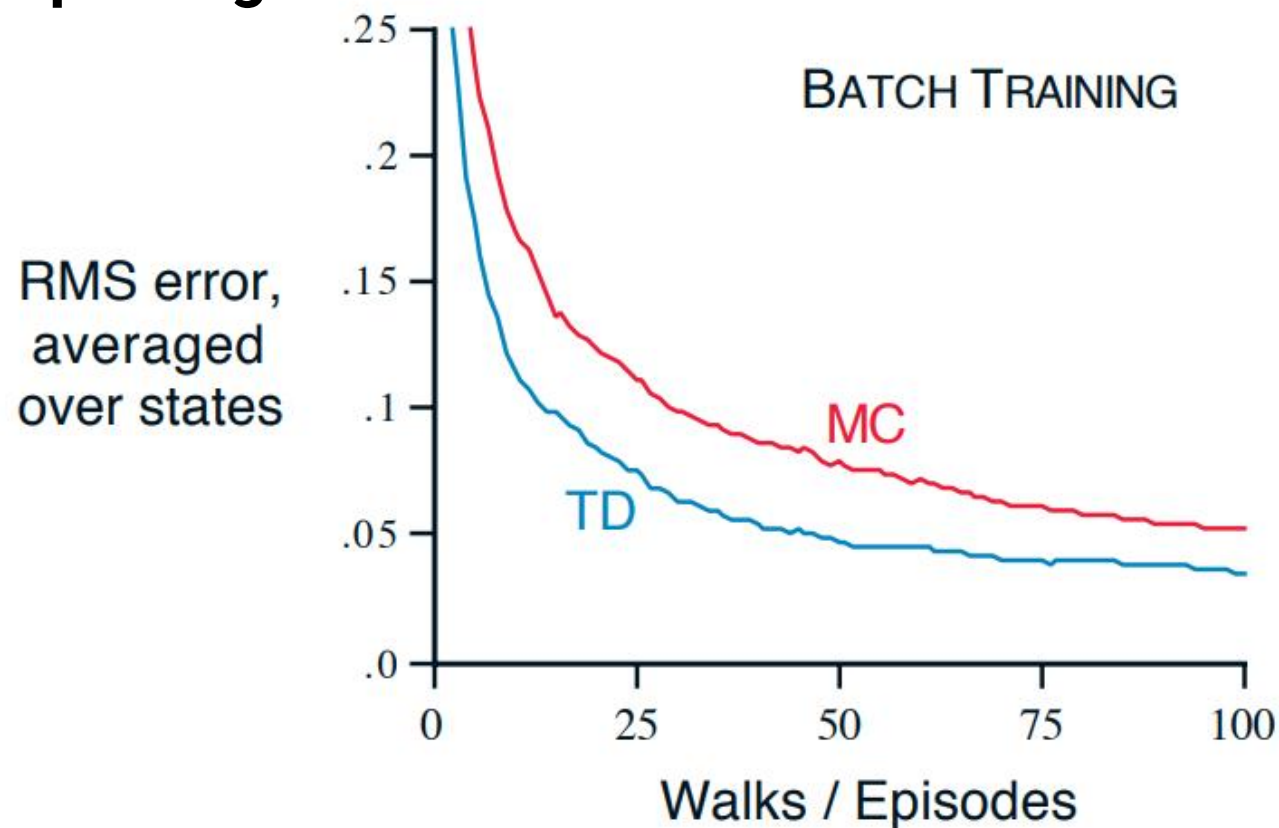
## Markov Reward Process: Random Walk Example

In this example we empirically compare the prediction abilities of TD(0) and constant- $\alpha$  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.

## Example: You are the Predictor

A, 0, B, 0

B, 1

B, 1

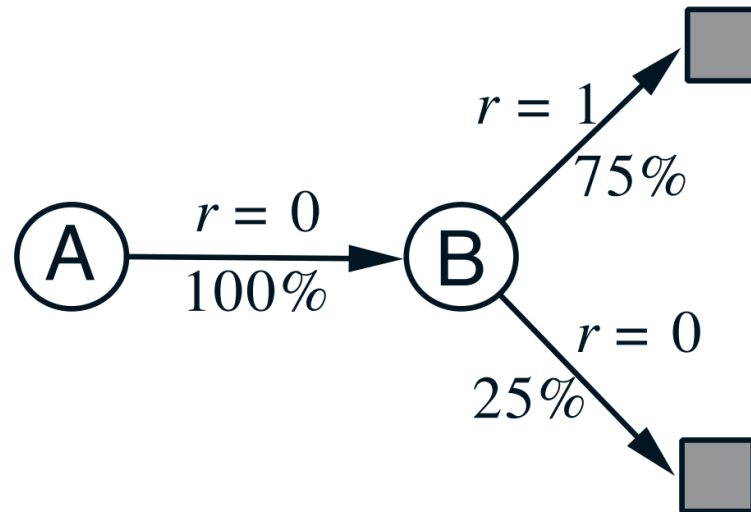
B, 1

B, 1

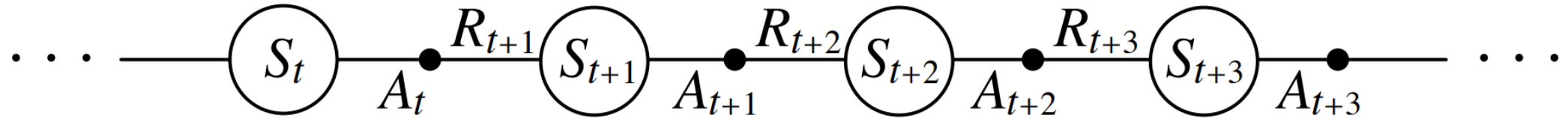
B, 1

B, 1

B, 0

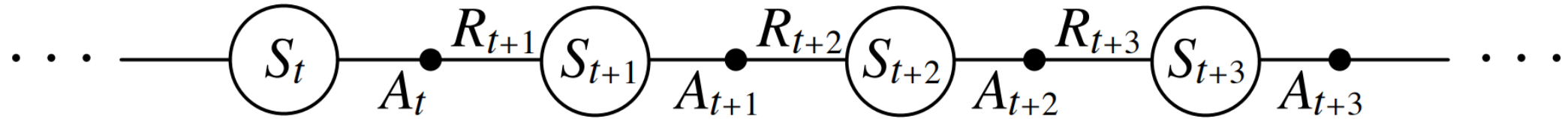


## SARSA: On-Policy TD Control



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

## SARSA: On-Policy TD Control



Target Q Value

$$\underbrace{Q(S_t, A_t)}_{\text{Updated Q Value}} \leftarrow \underbrace{Q(S_t, A_t)}_{\text{Current Q Value}} + \alpha \left[ \overbrace{R_{t+1} + \gamma \underbrace{Q(S_{t+1}, A_{t+1})}_{\text{Target Policy is always as same as Behaviour Policy}}}^{\text{Target Q Value}} - \underbrace{Q(S_t, A_t)}_{\text{Current Q Value}} \right]$$

Updated Q Value

Current Q Value

Current Q Value

Target Policy is always as same as Behaviour Policy

## Algorithm

### 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  $Q(\text{terminal}, \cdot) = 0$

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

## Better Intuition

### Select Action

Agent chooses action based on policy

### Observe State

Agent perceives current environment

### Execute Action !

Agent performs selected action

### Observe Reward

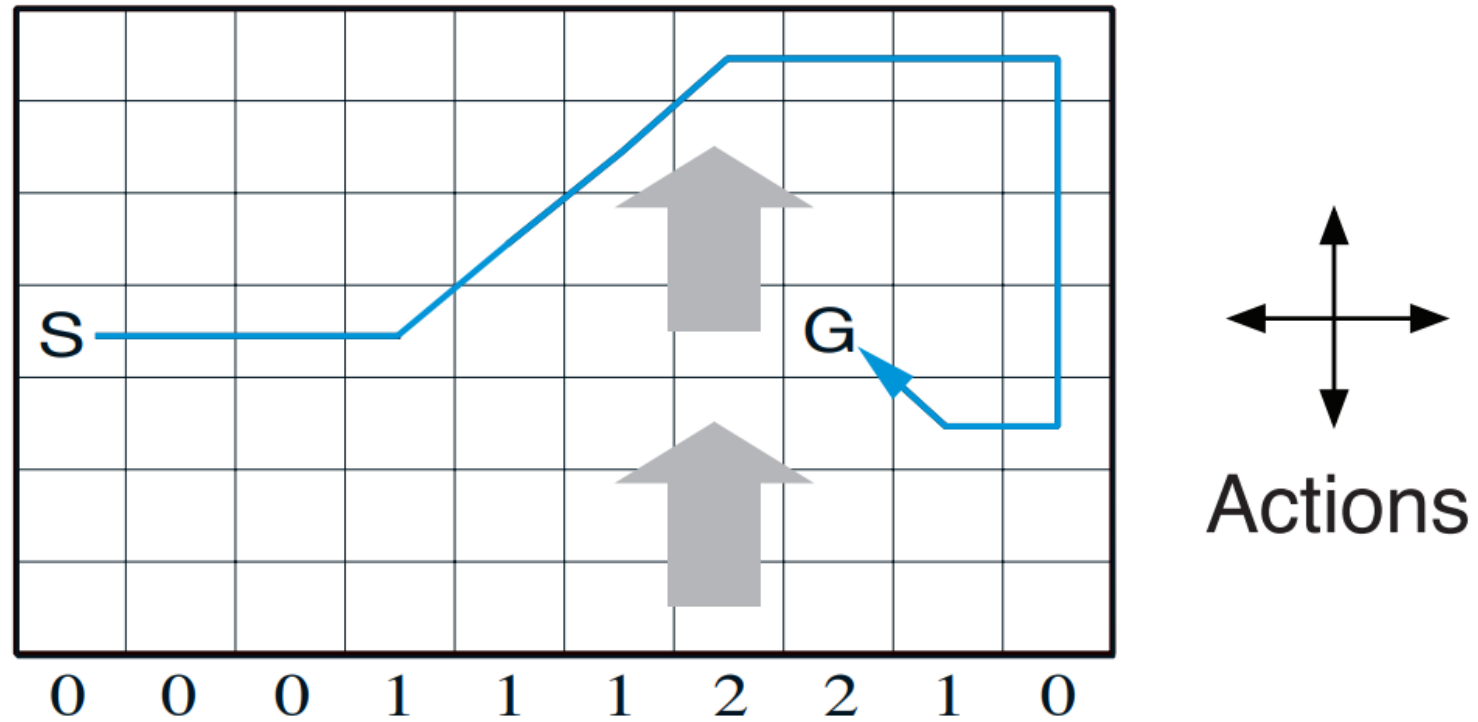
Agent receives feedback from environment

### Update Q-value

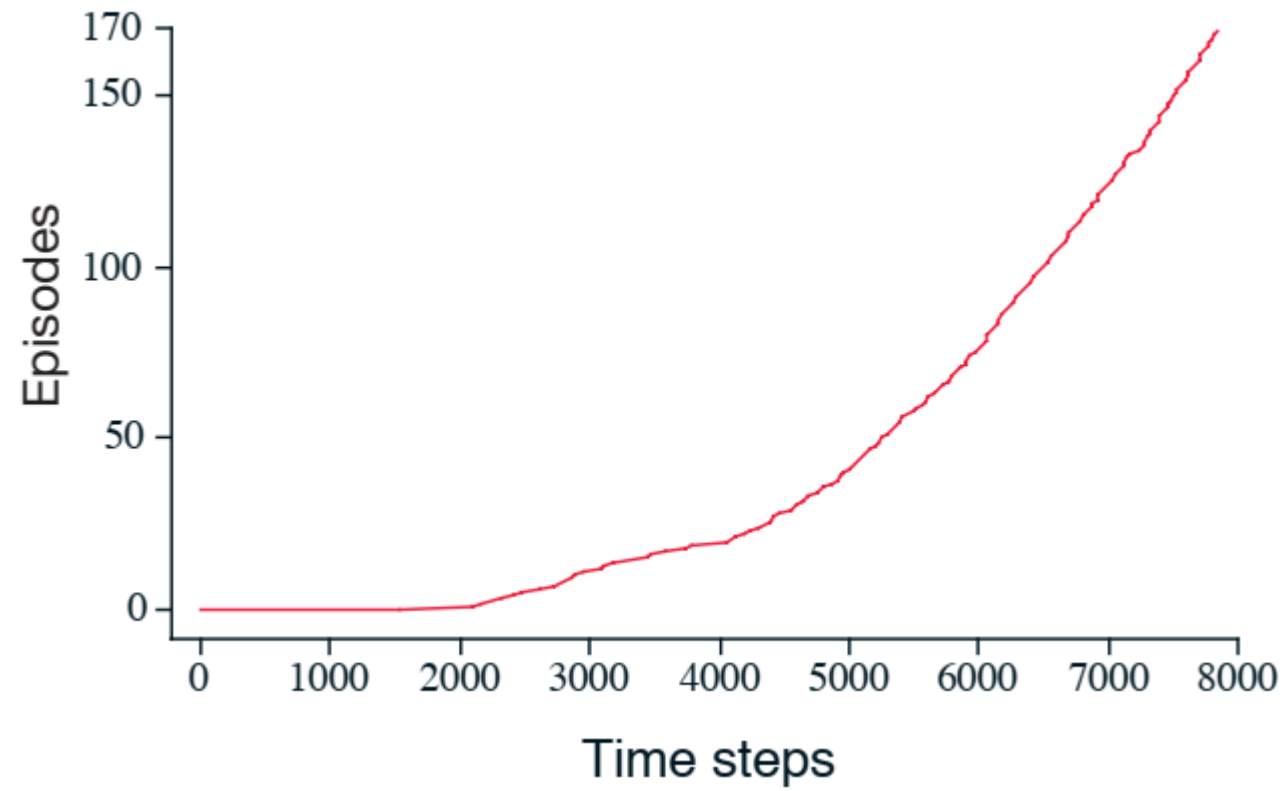
Agent adjusts Q-values based on experience



## Example: Windy Grid



## Example: Windy Grid





## Q-Learning: Off-Policy TD Control

$$\underbrace{Q(S_t, A_t)}_{\text{New Q-value estimation}} \leftarrow \underbrace{Q(S_t, A_t)}_{\text{Former Q-value estimation}} + \underbrace{\alpha}_{\text{Learning Rate}} [\underbrace{R_{t+1}}_{\text{Immediate Reward}} + \underbrace{\gamma \max_a Q(S_{t+1}, a)}_{\text{Discounted Estimate optimal Q-value of next state}} - \underbrace{Q(S_t, A_t)}_{\text{Former Q-value estimation}}]$$

TD Target

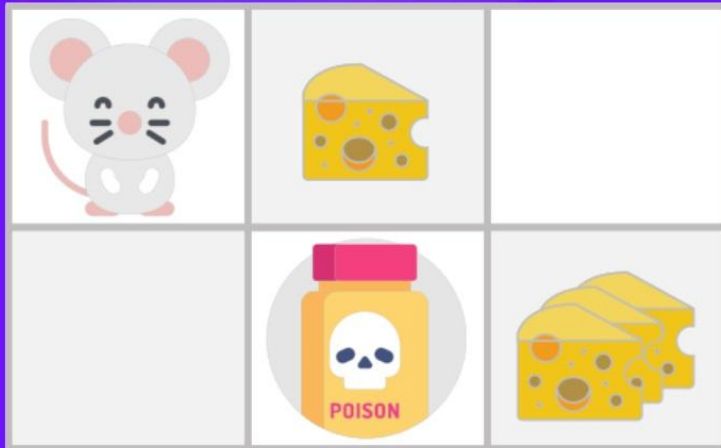
TD Error



## Example

- You always start at the **same starting point**.
- The goal: eat the **big pile of cheese** (at the bottom right-hand corner) and **avoid the poison**.
- The episode ends if we eat the poison, eat the big pile of cheese or if we spent more than 5 steps.
- Learning rate = 0.1
- Gamma = 0.99











## 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

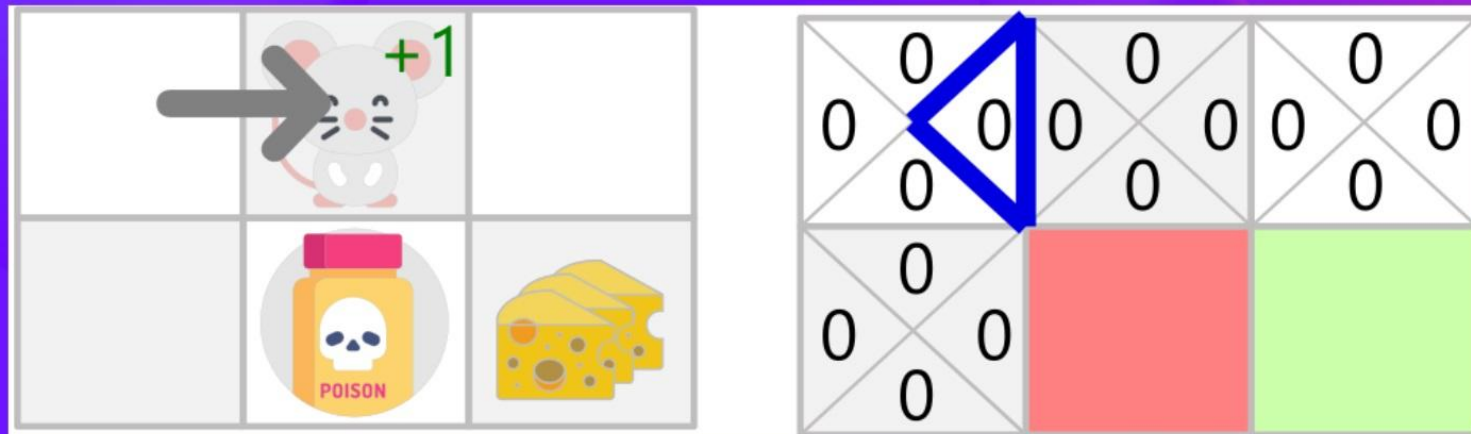
Initialize  $Q$  arbitrarily (e.g.,  $Q(s, a) = 0$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$ , and  $Q(\text{terminal-state}, \cdot) = 0$ )

				
	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 2

Choose action  $A_t$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)



We took a random action (exploration)

## Example, Step 3

Take action  $A_t$  and observe  $R_{t+1}, S_{t+1}$





## Example, Step 4

$$\underbrace{Q(S_t, A_t)}_{\text{New Q-value estimation}} \leftarrow \underbrace{Q(S_t, A_t)}_{\text{Former Q-value estimation}} + \underbrace{\alpha}_{\text{Learning Rate}} [\underbrace{R_{t+1}}_{\text{Immediate Reward}} + \underbrace{\gamma \max_a Q(S_{t+1}, a)}_{\text{Discounted Estimate optimal Q-value of next state}} - \underbrace{Q(S_t, A_t)}_{\text{Former Q-value estimation}}]$$

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

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

$Q(\text{Initial state, Right}) = 0 + 0.1 * [1 + 0.99 * 0 - 0]$

$Q(\text{Initial state, Right}) = 0.1$

	←	→	↑	↓
	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: Algorithm

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  $Q(\text{terminal}, \cdot) = 0$

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)

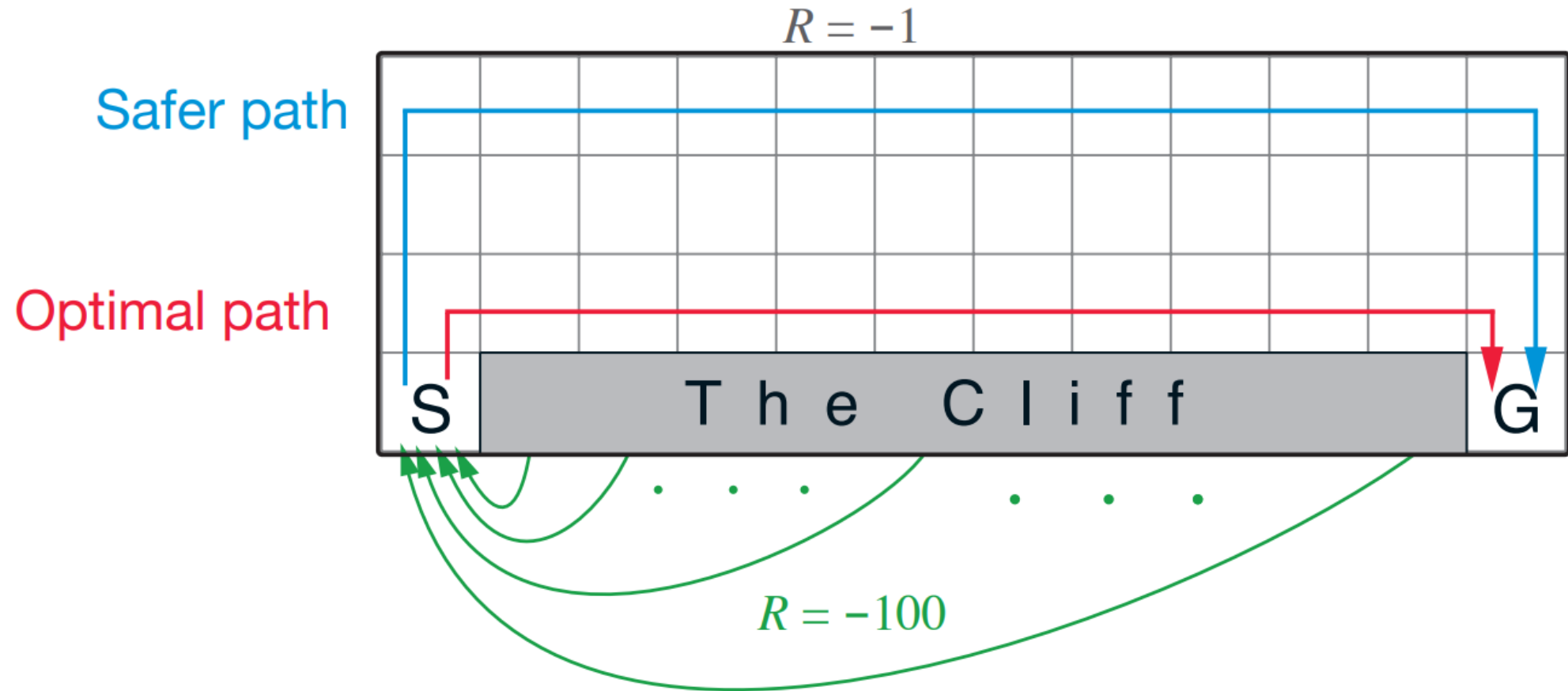
        Take action  $A$ , observe  $R, S'$

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

$S \leftarrow S'$

    until  $S$  is terminal

## Example: Cliff Walking



## Example: Cliff Walking

