# Temporal Difference Learning

CMPUT 366: Intelligent Systems

S&B §6.0-6.2, §6.4-6.5

# Lecture Overview

1. Recap & Logistics

2. TD Prediction

3. On-Policy TD Control (Sarsa)

4. Off-Policy TD Control (Q-Learning)

# Assignment #3

- Assignment #3 is **due tonight** (Mar 25) at **11:59pm**

  - This is a firm deadline

# Recap: Monte Carlo RL

- **Monte Carlo** estimation: Estimate expected returns to a state or action by averaging actual returns over sampled trajectories

  - Estimating action values requires either exploring starts or a soft policy (e.g., $\epsilon$-greedy)

- **Off-policy learning** is the estimation of value functions for a target policy based on episodes generated by a different behaviour policy

- **Off-policy control** is learning the optimal policy (target policy) using episodes from a behaviour policy

# Recap: Off-Policy Monte Carlo Prediction

**Off-policy MC prediction (policy evaluation) for estimating $Q \approx q_\pi$**

Input: an arbitrary target policy $\pi$

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$\quad Q(s,a) \in \mathbb{R}$ (arbitrarily)

$\quad C(s,a) \leftarrow 0$

Loop forever (for each episode):

$\quad b \leftarrow$ any policy with coverage of $\pi$

$\quad$ Generate an episode following $b$: $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$

$\quad G \leftarrow 0$

$\quad W \leftarrow 1$

$\quad$ Loop for each step of episode, $t = T-1, T-2, \ldots, 0$, while $W \neq 0$:

$\quad\quad G \leftarrow \gamma G + R_{t+1}$

$\quad\quad C(S_t, A_t) \leftarrow C(S_t, A_t) + W$

$\quad\quad Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} \left[ G - Q(S_t, A_t) \right]$

$\quad\quad W \leftarrow W \frac{\pi(A_t | S_t)}{b(A_t | S_t)}$

# Recap: Off-Policy Monte Carlo Control

**Off-policy MC control, for estimating $\pi \approx \pi_*$**

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:
    $Q(s, a) \in \mathbb{R}$ (arbitrarily)
    $C(s, a) \leftarrow 0$
    $\pi(s) \leftarrow \arg\max_a Q(s, a)$    (with ties broken consistently)

Loop forever (for each episode):
    $b \leftarrow$ any soft policy
    Generate an episode using $b$: $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$
    $G \leftarrow 0$
    $W \leftarrow 1$
    Loop for each step of episode, $t = T-1, T-2, \ldots, 0$:
        $G \leftarrow \gamma G + R_{t+1}$
        $C(S_t, A_t) \leftarrow C(S_t, A_t) + W$
        $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} \left[ G - Q(S_t, A_t) \right]$
        $\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$    (with ties broken consistently)
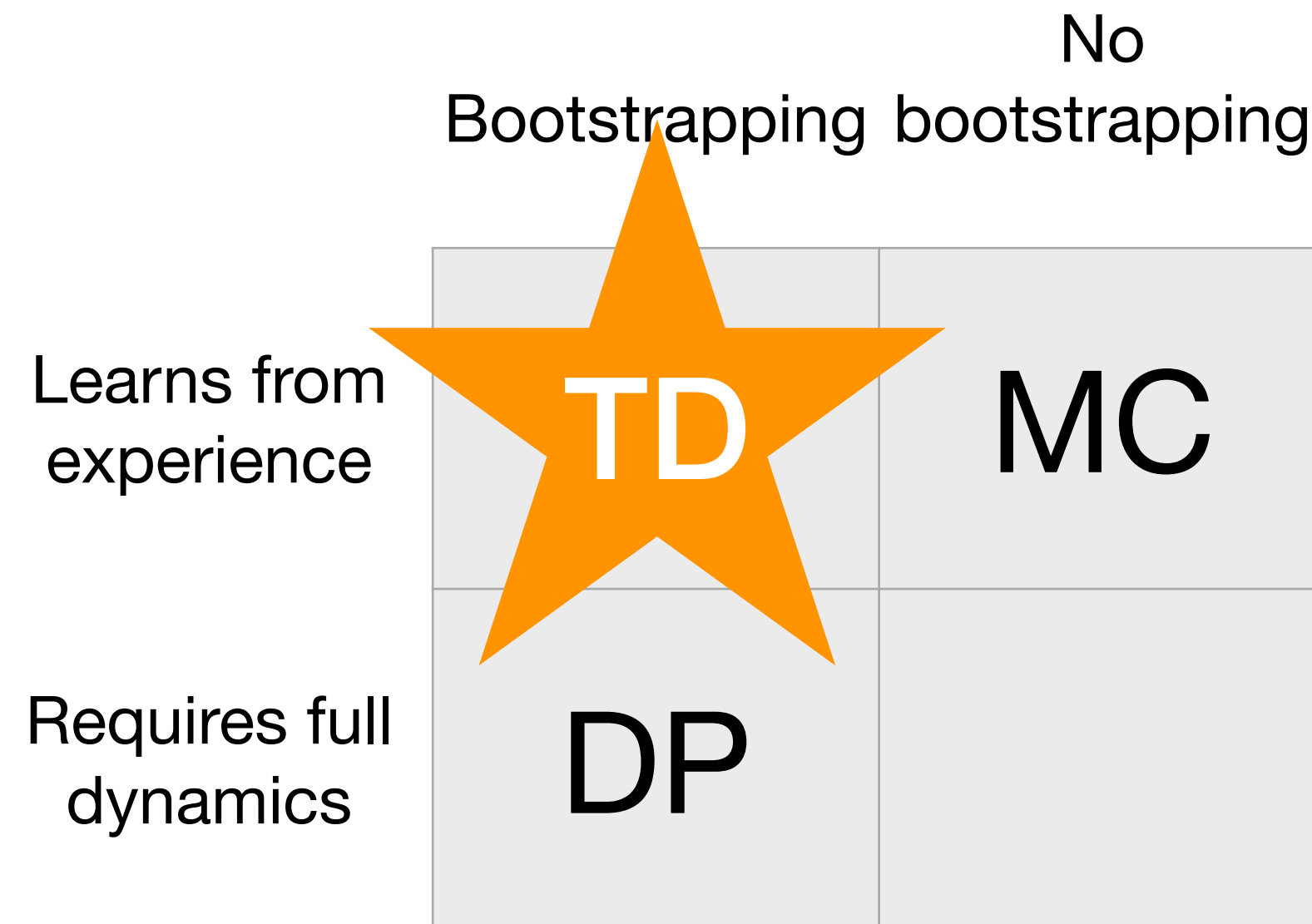        If $A_t \neq \pi(S_t)$ then exit inner Loop (proceed to next episode)
        $W \leftarrow W \frac{1}{b(A_t | S_t)}$

# Learning from Experience

- Suppose we are playing a blackjack-like game **in person**, but we **don't know the rules**.

  - We know the actions we can take, we can see the cards, and we get told when we win or lose

- **Question:** Could we compute an optimal policy using **dynamic programming** in this scenario?

- **Question:** Could we compute an optimal policy using **Monte Carlo**?

  - What would be the **pros and cons** of running Monte Carlo?

# Bootstrapping

|  | Bootstrapping | No bootstrapping |
|---|---|---|
| Learns from experience | TD | MC |
| Requires full dynamics | DP |  |

- Dynamic programming **bootstraps**: Each iteration's estimates are based partly on **estimates from previous iterations**

- Each Monte Carlo estimate is based only on **actual returns**

# Updates

**Dynamic Programming:** $V(S_t) \leftarrow \sum_a \pi(a \,|\, S_t) \sum_{s',r} p(s', r \,|\, S_t, a)\big[r + \gamma V(s')\big]$

**Monte Carlo:** $V(S_t) \leftarrow V(S_t) + \alpha \Big[\boxed{G_t} - V(S_t)\Big]$

**TD(0):** $V(S_t) \leftarrow V(S_t) + \alpha \Big[\boxed{R_{t+1} + \gamma V(S_{t+1})} - V(S_t)\Big]$

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t \mid S_t = s] \qquad \text{Monte Carlo: Approximate because of } \mathbb{E}$$

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

$$= \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s]. \qquad \text{Dynamic programming:}$$

Approximate because $v_\pi$ not known

TD(0): Approximate because of $\mathbb{E}$ **and** $v_\pi$ not known

# TD(0) 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 $V(terminal) = 0$

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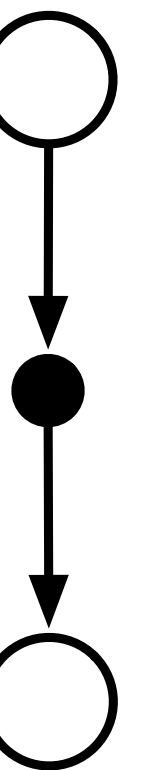    Initialize $S$
    Loop for each step of episode:
        $A \leftarrow$ action given by $\pi$ for $S$
        Take action $A$, observe $R$, $S'$
        $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$
        $S \leftarrow S'$
    until $S$ is terminal

**Question:** What information does this algorithm use?

# TD for Control

- We can plug TD prediction into the **generalized policy iteration** framework

- **Monte Carlo control loop:**

  1. Generate an **episode** using estimated $\pi$

  2. Update estimates of $Q$ and $\pi$

- **On-policy TD control loop:**

  1. Take an **action** according to $\pi$

  2. Update estimates of $Q$ and $\pi$

# On-Policy TD Control

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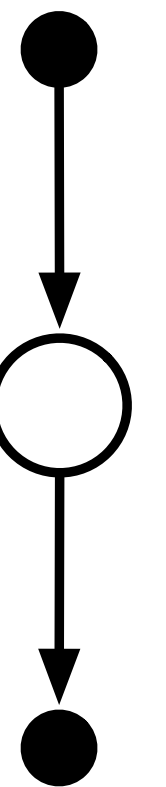    Loop for each step of episode:
        Take action $A$, observe $R, S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha \big[ R + \gamma Q(S', A') - Q(S, A) \big]$
        $S \leftarrow S'; A \leftarrow A';$
    until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

**Question:** What **information** does this algorithm use?

**Question:** Will this estimate the Q-values of the **optimal** policy?

# Actual Q-Values vs. Optimal Q-Values

- Just as with on-policy Monte Carlo control, Sarsa does not converge to the **optimal** policy, because it always chooses an $\epsilon$**-greedy action**

  - And the estimated Q-values are with respect to the **actual actions**, which are $\epsilon$-greedy

- **Question:** Why is it necessary to choose $\epsilon$-greedy actions?

- What if we **acted** $\epsilon$-greedy, but **learned the Q-values** for the optimal policy?

# Off-Policy TD Control



**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:
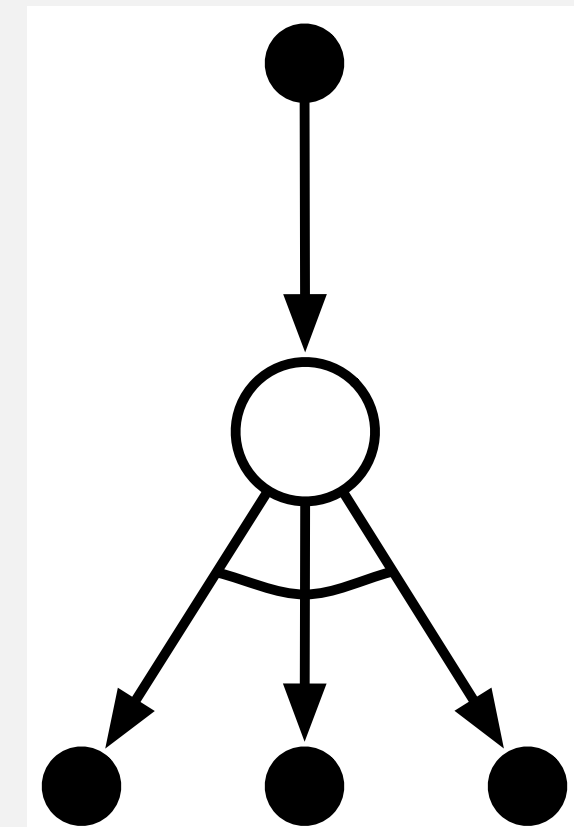        Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        Take action $A$, observe $R$, $S'$
        $Q(S, A) \leftarrow Q(S, A) + \alpha\big[R + \gamma \max_a Q(S', a) - Q(S, A)\big]$
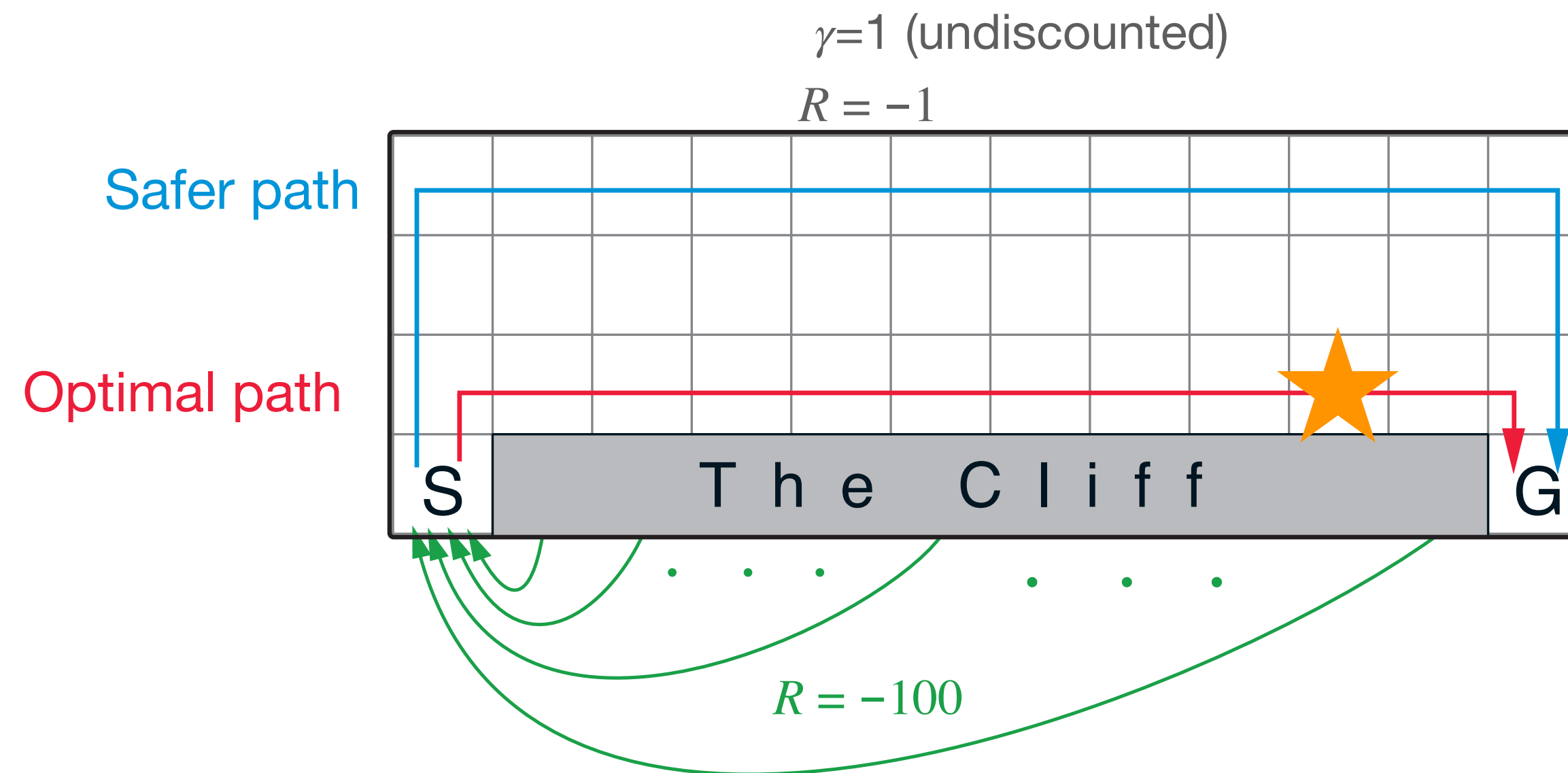        $S \leftarrow S'$
    until $S$ is terminal

**Question:** What **information** does this algorithm use?

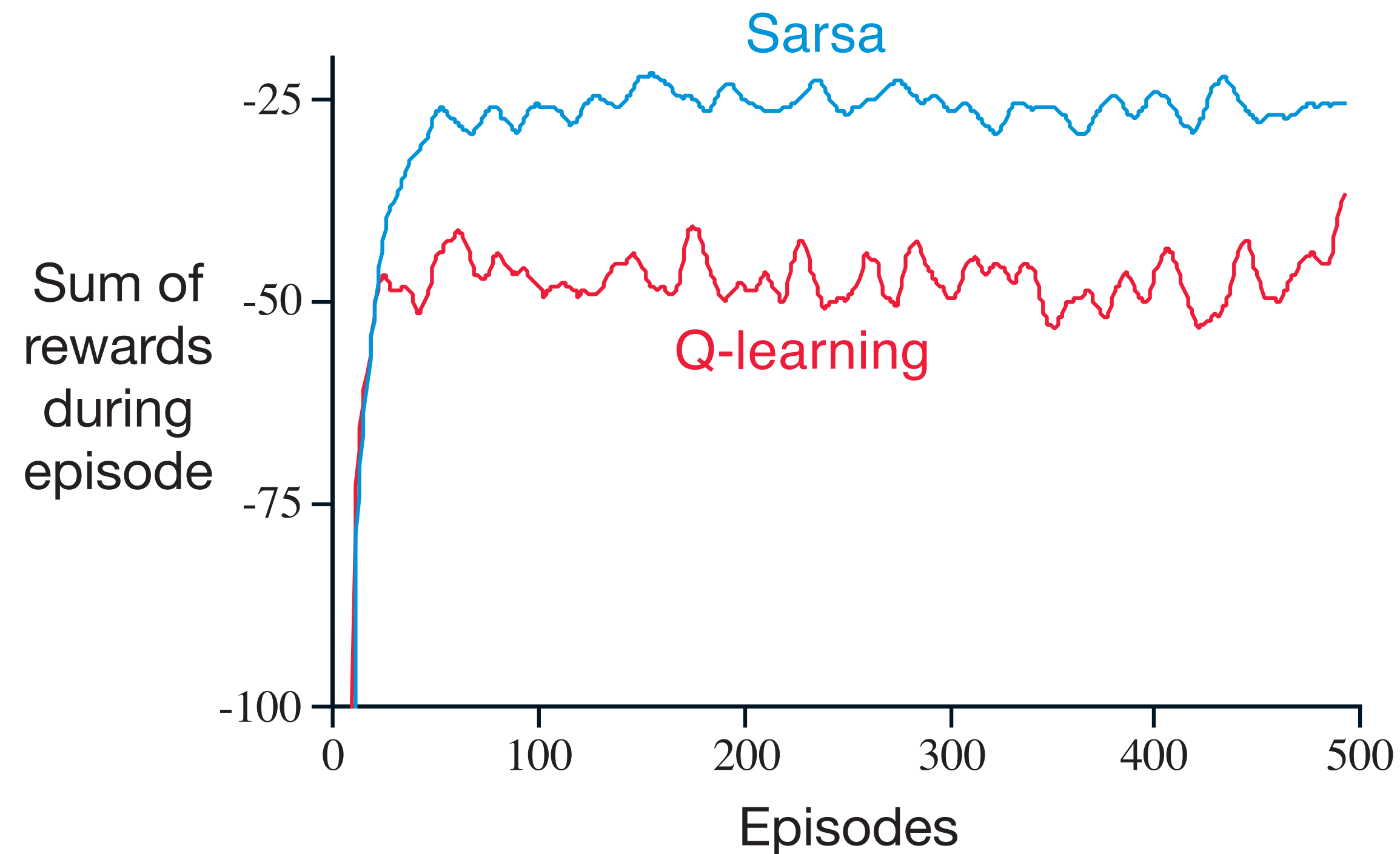**Question:** Why aren't we estimating the **policy** $\pi$ explicitly?

# Example: The Cliff

$\gamma = 1$ (undiscounted)

$R = -1$

Safer path

Optimal path

$R = -1$

S    T h e    C l i f f    G

safe path!

sa

$R = -100$

optimal path

op

S    T h e    C l i f f    G

- Agent gets -1 reward until they reach the goal state

- Step into the Cliff region, get reward -100 and go back to start

- **Question:** How will **Q-Learning** estimate the value of **this** state?

- **Question:** How will **Sarsa** Sarsa estimate the value of **this** state?

# Performance on The Cliff



Q-Learning estimates **optimal policy**, but Sarsa consistently **outperforms** Q-Learning.  (**why?**)

# Summary

- Temporal Difference Learning **bootstraps** *and* learns from **experience**

  - Dynamic programming bootstraps, but doesn't learn from experience (requires full dynamics)

  - Monte Carlo learns from experience, but doesn't bootstrap

- Prediction: **TD(0) algorithm**

- **Sarsa** estimates action-values of **actual $\epsilon$-greedy policy**

- **Q-Learning** estimates action-values of **optimal** policy while **executing** an **$\epsilon$-greedy** policy