# Reinforcement Learning
## Temporal-Difference Learning (RLbook 6)

Christopher Simpkins

Kennesaw State University

# TD Prediction

$$V(S_t) \leftarrow V(S_t) + \alpha \left[G_t - V(S_t)\right] \qquad (6.1)$$

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

# Tabular 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:
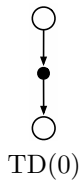        $A \leftarrow$ action given by $\pi$ for $S$
        Take action $A$, observe $R$, $S'$
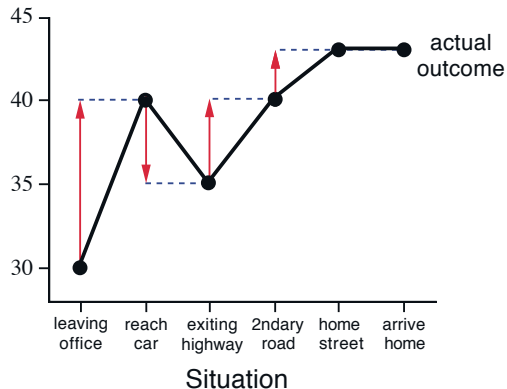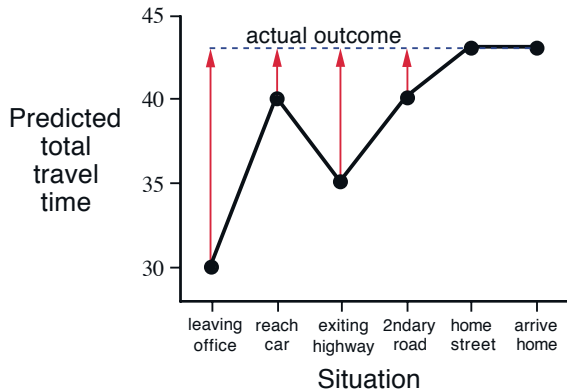        $V(S) \leftarrow V(S) + \alpha\big[R + \gamma V(S') - V(S)\big]$
        $S \leftarrow S'$
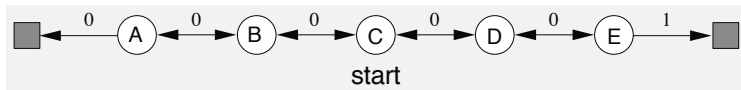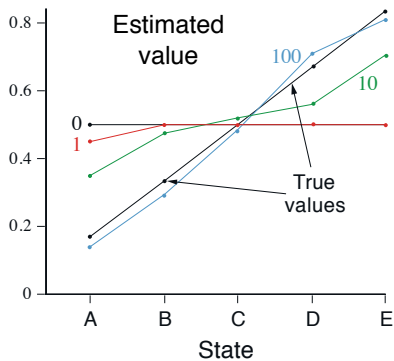    until $S$ is terminal

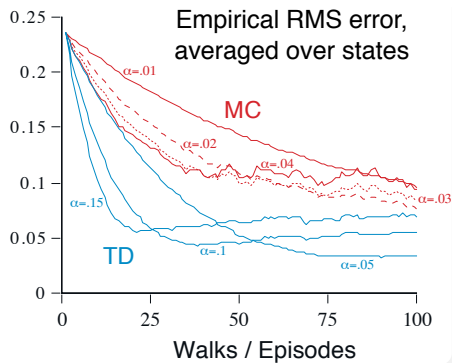# TD(0) Backup



TD(0)

# Example: Driving Home

# Example: Random Walk

# Random Walk State Values

# Random Walk Error Rates



Empirical RMS error, averaged over states
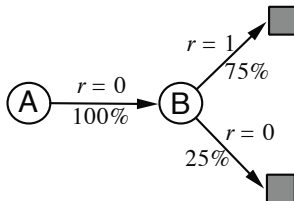
# TD vs MC Performance

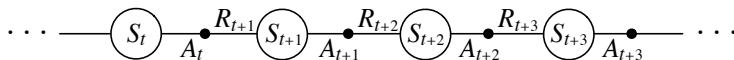# Example: Predicting Returns

Given these 8 episodes:

▶ $A, 0, B, 0; B, 1; B, 1; B, 1; B, 1; B, 1; B, 1; B, 1;$

What are the value estimates for $A$ and $B$?

# Sarsa: On-policy TD Control

$$\cdots \underline{\phantom{xx}} \underset{A_t}{\overset{R_{t+1}}{(S_t)}} \bullet \underset{A_{t+1}}{\overset{R_{t+1}}{(S_{t+1})}} \bullet \underset{A_{t+2}}{\overset{R_{t+2}}{(S_{t+2})}} \bullet \underset{A_{t+3}}{\overset{R_{t+3}}{(S_{t+3})}} \bullet \underline{\phantom{xx}} \cdots$$

Sarsa update:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(S_t, A_t) \right]$$

# Sarsa 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 $Q(terminal, \cdot) = 0$

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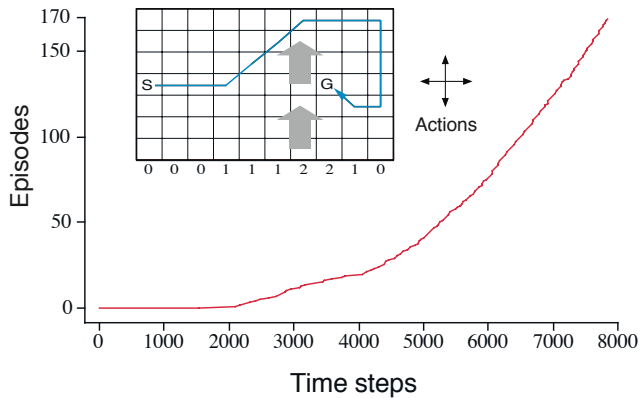    Loop for each step of episode:
        Take action $A$, observe $R$, $S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S,A) \leftarrow Q(S,A) + \alpha\big[R + \gamma Q(S',A') - Q(S,A)\big]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$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

# Example: Windy Grid World

# Q-Learning: Off-policy TD Control

Sarsa update:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(S_t, A_t)\right]$$

Q-learning update:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(S_t, A_t)\right]$$

# Q-Learning 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 $Q(terminal, \cdot) = 0$

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)
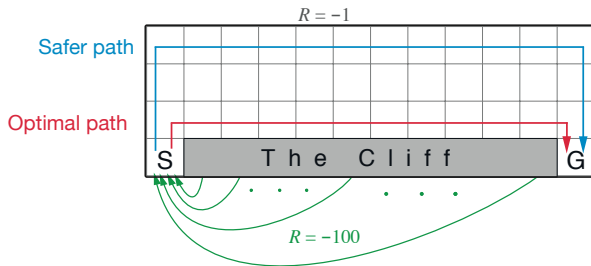        Take action $A$, observe $R, S'$
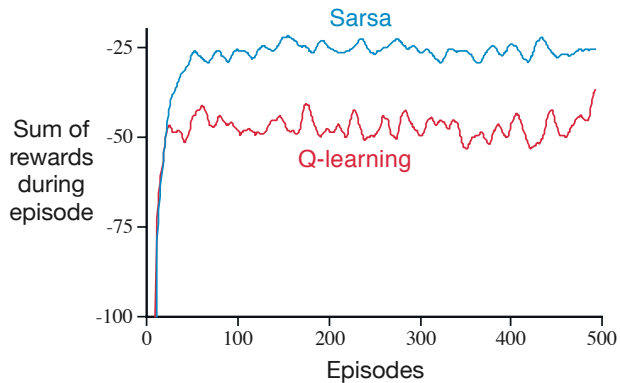        $Q(S, A) \leftarrow Q(S, A) + \alpha \big[ R + \gamma \max_a Q(S', a) - Q(S, A) \big]$
        $S \leftarrow S'$
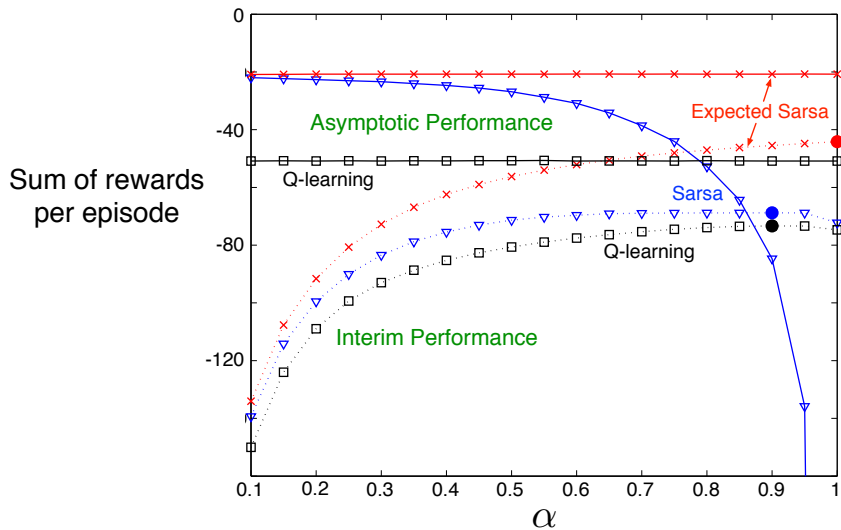    until $S$ is terminal

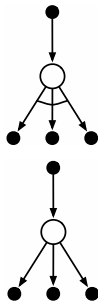# Example: Cliff Walking

# Sarsa vs Q-learning in Cliff Walking

# Expected Sarsa

$$Q(S_t, A_t) \leftarrow Q(s_t, A_t) + \alpha \left[ R_{t+1} + \gamma \mathbb{E}_\pi \left[ Q(S_{t+1}, A_{t+1}) \mid s_{t+1} \right] - Q(S_t, A_t) \right]$$
$$= Q(s_t, A_t) + \alpha \left[ R_{t+1} + \gamma \sum_a \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right] \qquad (6.9)$$
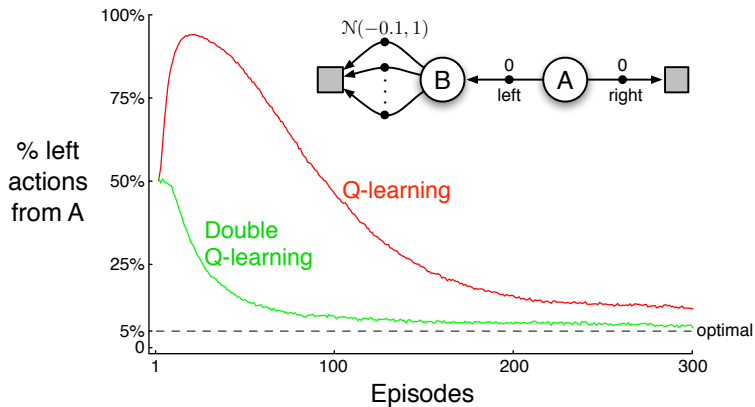
# Asymptotic Performance of TD Control Methods

# Q-learning vs Expected Sarsa Backup

# Double Q-learning Performance

# Double Q-learning Algorithm

**Double Q-learning, for estimating $Q_1 \approx Q_2 \approx q_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q_1(s, a)$ and $Q_2(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, such that $Q(terminal, \cdot) = 0$

Loop for each episode:
   Initialize $S$
   Loop for each step of episode:
      Choose $A$ from $S$ using the policy $\varepsilon$-greedy in $Q_1 + Q_2$
      Take action $A$, observe $R$, $S'$
      With 0.5 probabilility:
$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha\Big(R + \gamma Q_2\big(S', \arg\max_a Q_1(S', a)\big) - Q_1(S, A)\Big)$$
      else:
$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha\Big(R + \gamma Q_1\big(S', \arg\max_a Q_2(S', a)\big) - Q_2(S, A)\Big)$$
      $S \leftarrow S'$
   until $S$ is terminal

# Games and Afterstates