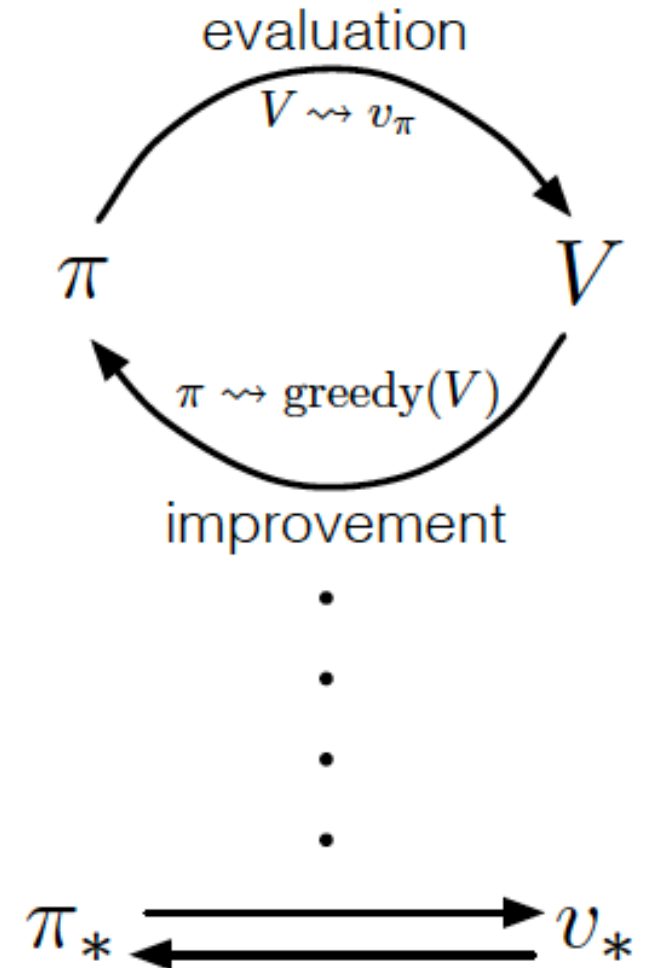


Temporal-Difference Learning

- Combination of learning from experience without a model (like Monte Carlo methods) and bootstrapping (like DP methods)
- Based on GPI and differences between DP, TD and MC are in their approaches to the prediction problem (value function evaluation).



TD prediction

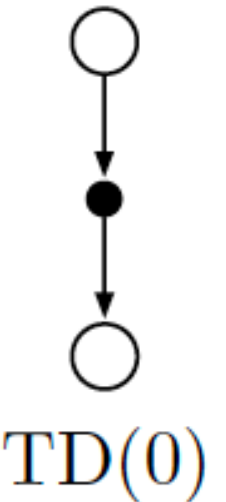
- While MC methods must wait until the end of the episode to calculate return for the value function update...

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

- TD methods need to wait only until the next time step, and the update for TD(0) or one-step TD is

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

- Bootstrapping: updating from an existing estimate



Tabular TD(0) for estimating v_π

Input: the policy π 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 π for S

 Take action A , observe R, S'

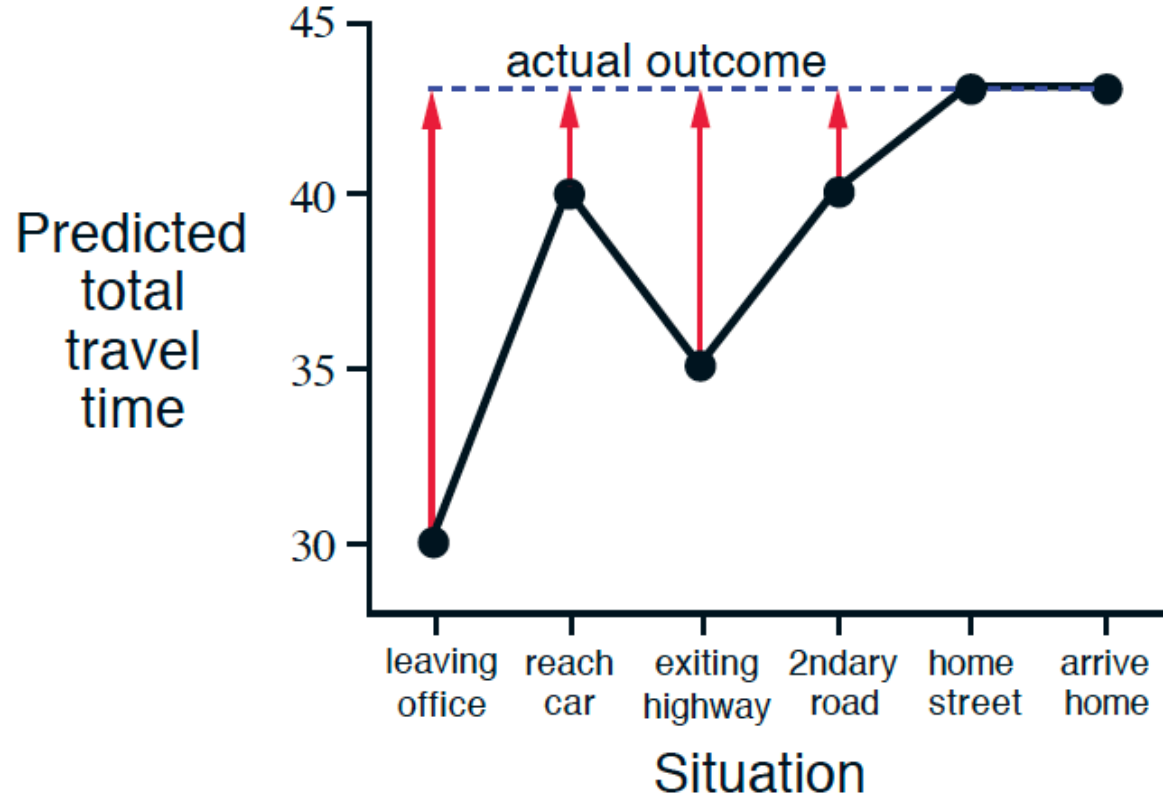
$V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

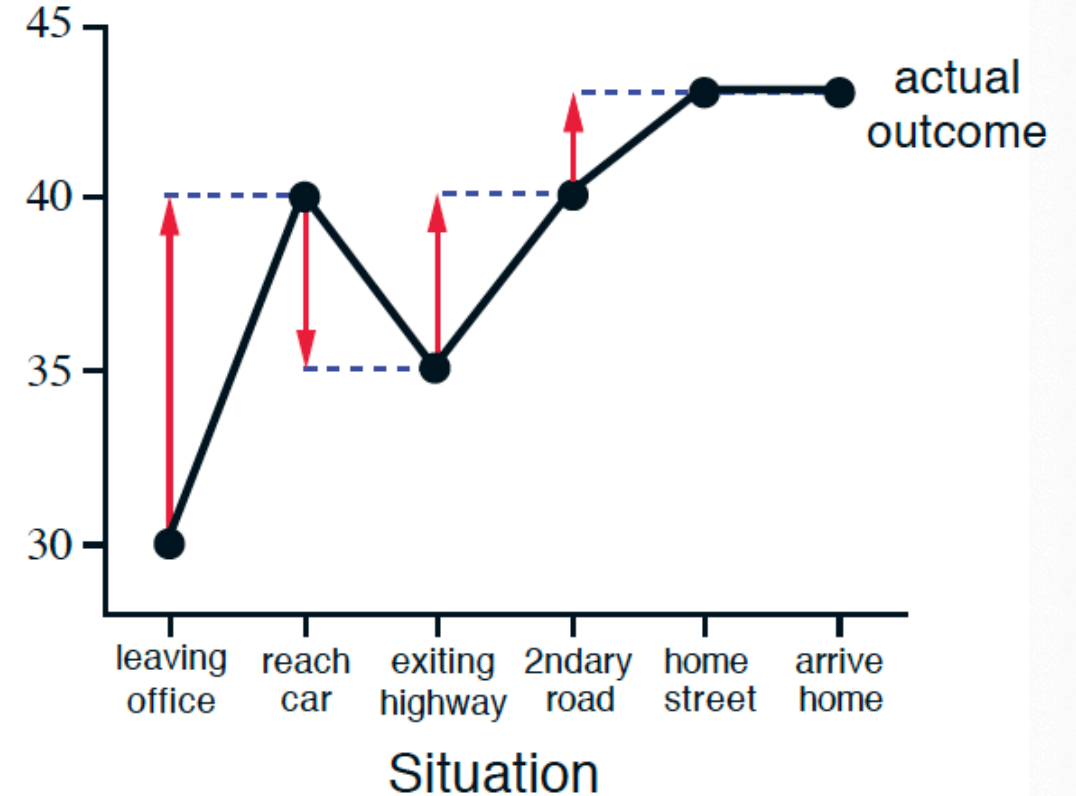
 until S is terminal

Driving Home example

Monte Carlo



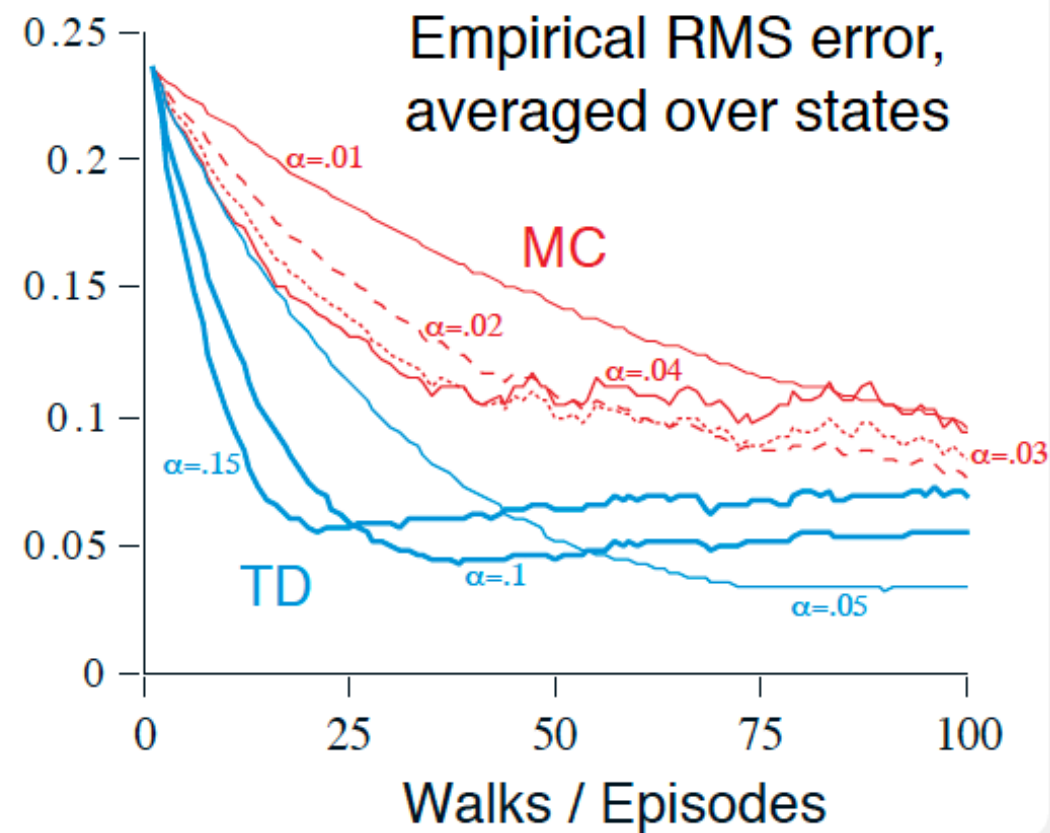
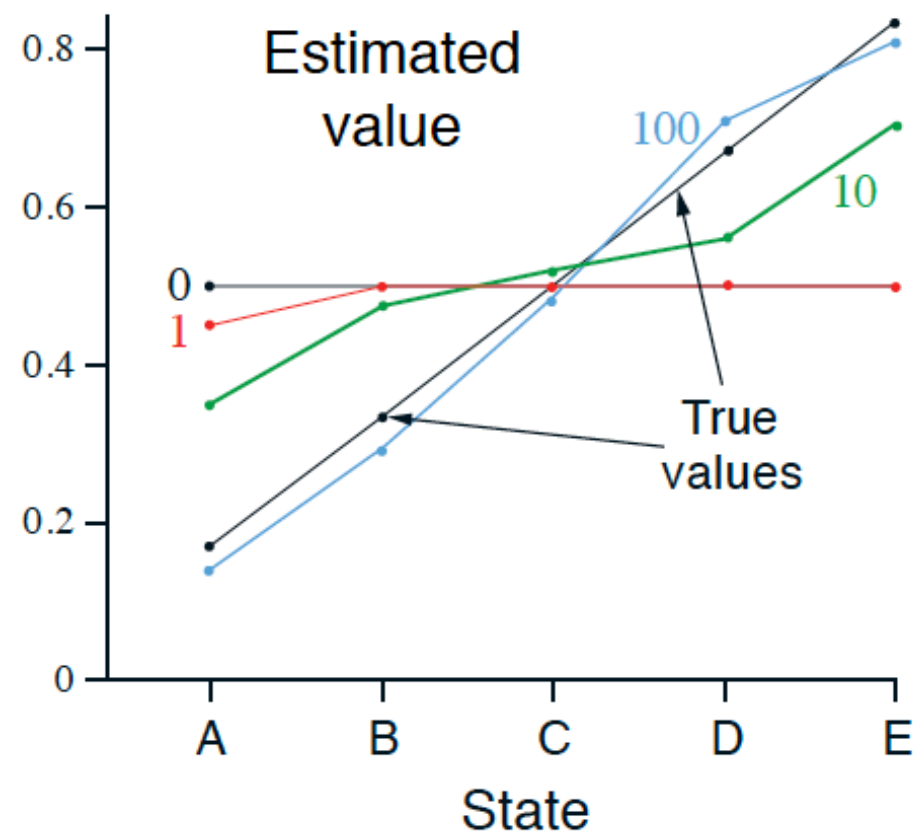
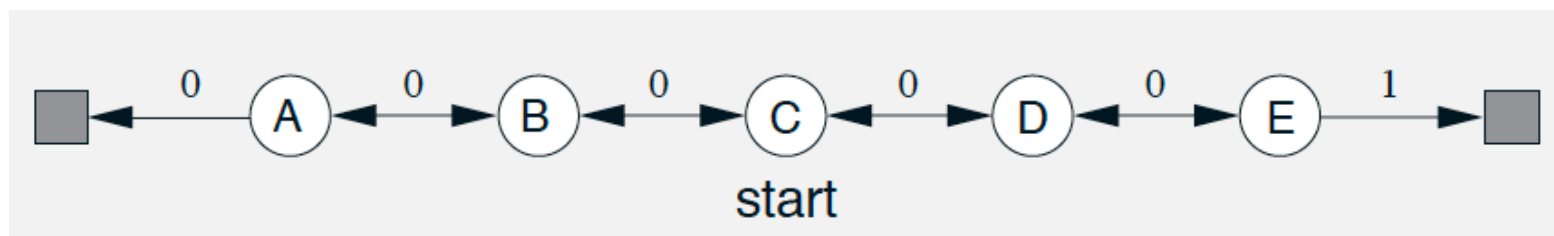
TD methods



Advantages of TD Prediction Methods

- Implemented in an online, fully incremental fashion so suitable for very long episodes or continuous tasks
- For any fixed policy TD(0) prediction converges to the correct value function (with decreasing step-size)
- In practice TD methods converge faster than MC method (not proved)

Random walk



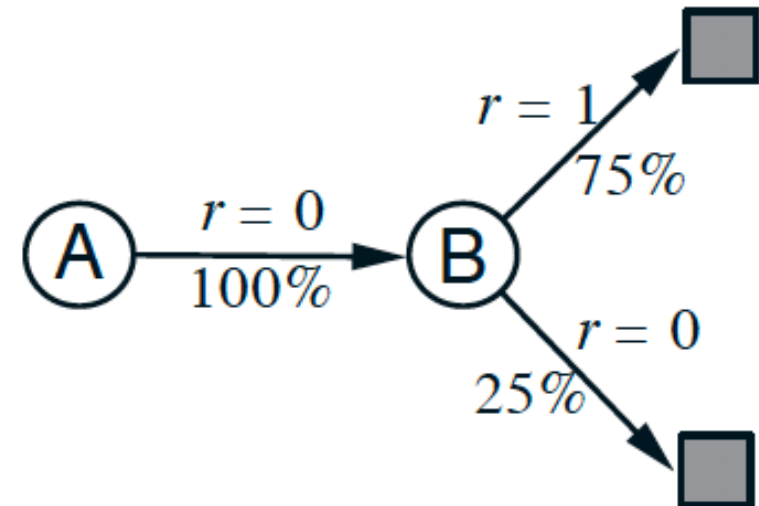
Optimality of TD(0)

- Batch updating: the value function is updated only once per batch, by the sum of all the increments in the batch
- TD(0) and MC methods both converge deterministically but to different answers
- Example:

A, 0, B, 0
B, 1
B, 1
B, 1

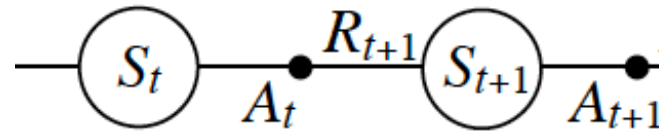
B, 1
B, 1
B, 1
B, 0

	V(A)	V(B)
MC (minimum MSE)	0	3/4
TD (certainty-equivalence - maximum likelihood estimate)	3/4	3/4



Sarsa: On-policy TD Control

- GPI with TD prediction considering transitions from state–action pairs to state–action pairs (with the corresponding reward)



- Learning an action value function with the following update rule

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

- Policy improvement: ε -greedy or ε -soft



Sarsa

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., ε -greedy)

 Loop for each step of episode:

 Take action A , observe R, S'

 Choose A' from S' using policy derived from Q (e.g., ε -greedy)

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

$S \leftarrow S'; A \leftarrow A';$

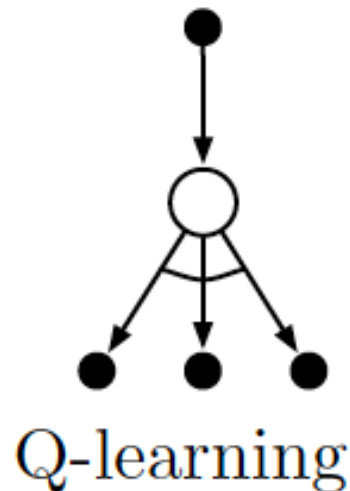
 until S is terminal

Q-learning: Off-policy TD Control

- Learned action-value function directly approximates the optimal action-value function, independent of the policy being followed.

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

- Convergence if all pairs continue to be updated and step-size decay



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., ε -greedy)

 Take action A , observe R, S'

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

$S \leftarrow S'$

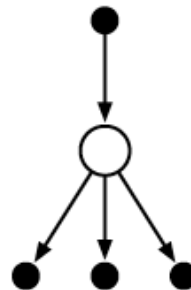
 until S is terminal

Expected Sarsa

- Like Q-learning except that instead of the maximum over next state–action pairs it uses the expected value

$$\begin{aligned} Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \mathbb{E}_{\pi} [Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_t, A_t) \right] \\ &\leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right], \end{aligned}$$

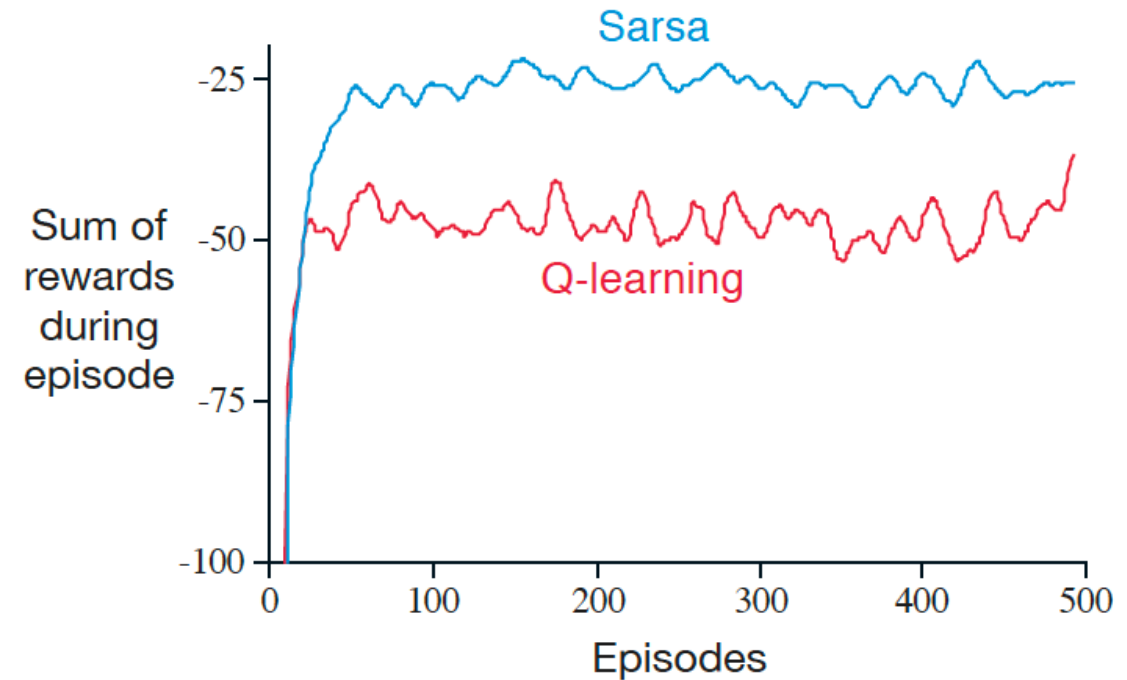
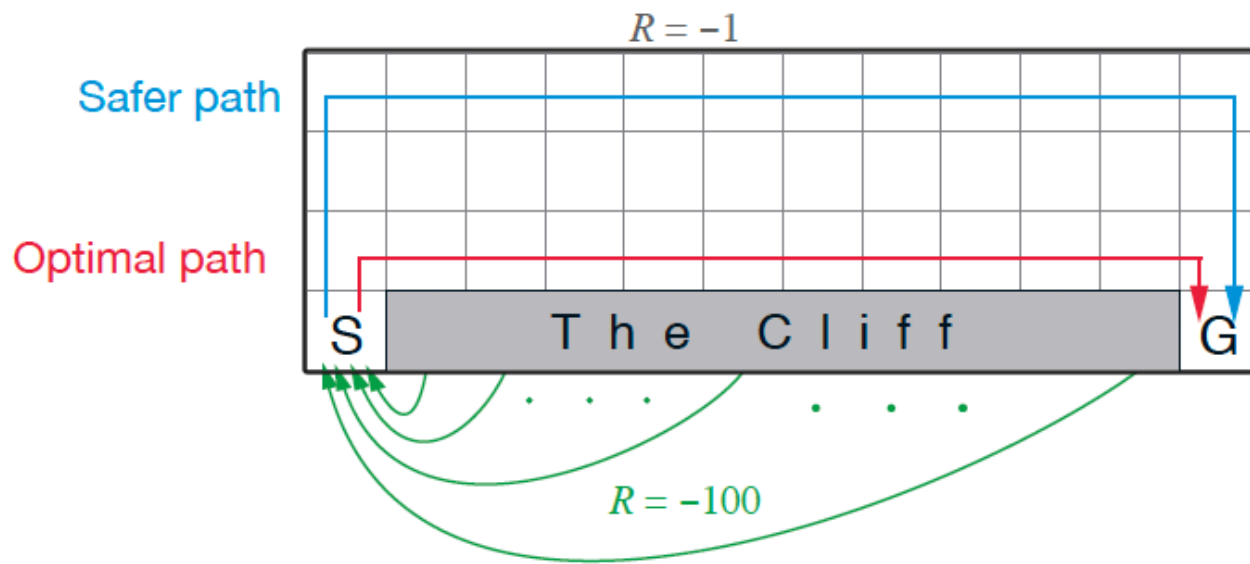
- Eliminates the variance of random action selection A_{t+1} of Sarsa



Expected Sarsa

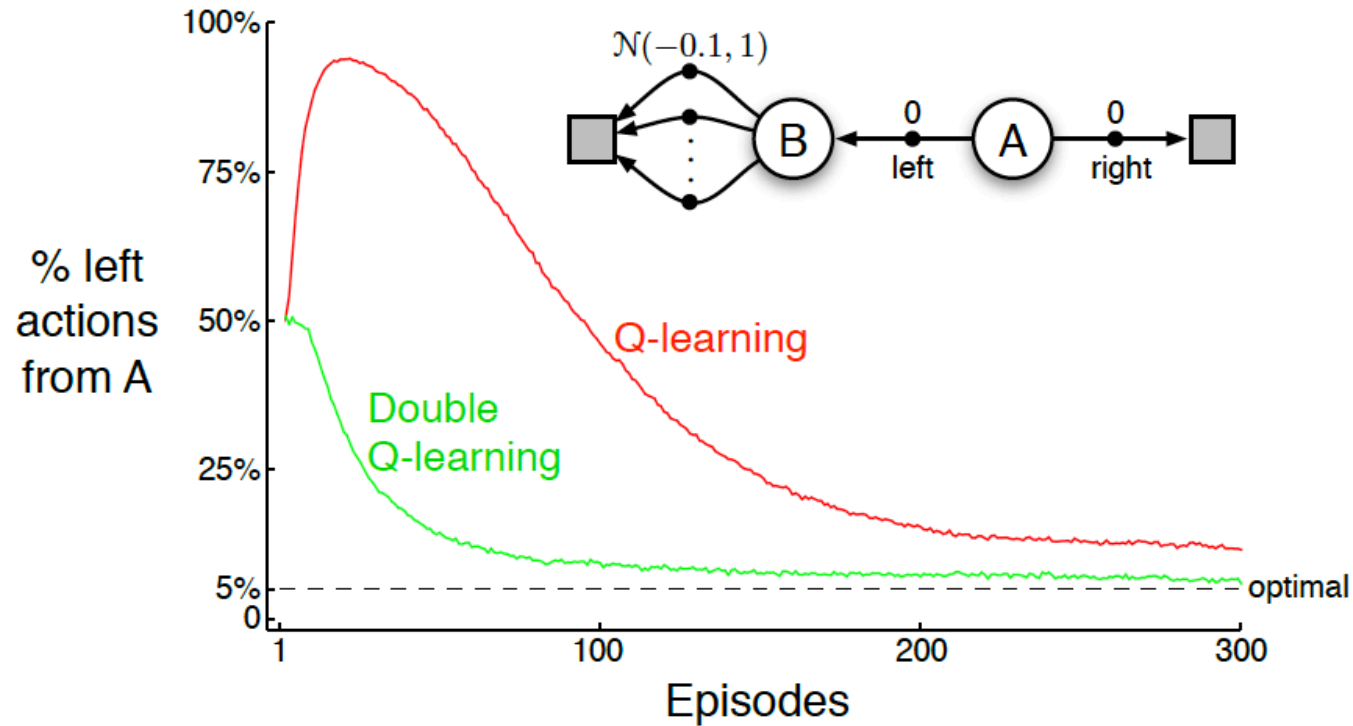
Cliff walking example

- ϵ -greedy with $\epsilon = 0.1$



Maximization Bias and Double Learning

- Maximization (Q-learning target or ϵ -greedy action selection) leads to positive bias
- Double learning: 2 Q-functions where Q_1 determines the maximizing action and Q_2 estimates its value (target)



$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha \left[R_{t+1} + \gamma Q_2(S_{t+1}, \arg \max_a Q_1(S_{t+1}, a)) - Q_1(S_t, A_t) \right].$$

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

Loop for each episode:

 Initialize S

 Loop for each step of episode:

 Choose A from S using the policy ε -greedy in $Q_1 + Q_2$

 Take action A , observe R, S'

 With 0.5 probability:

$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha \left(R + \gamma Q_2(S', \arg\max_a Q_1(S', a)) - Q_1(S, A) \right)$$

 else:

$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha \left(R + \gamma Q_1(S', \arg\max_a Q_2(S', a)) - Q_2(S, A) \right)$$

$S \leftarrow S'$

until S is terminal

Summary

Algorithm	Sarsa	Q-learning	Expected sarsa
Target	$R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$	$R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$	$R_{t+1} + \gamma \sum_a \pi(a S_{t+1}) Q(S_{t+1}, a)$
Backup diagram	