Monte Carlo methods to evaluate $\boxed{v_\pi}$

$\rightarrow$ Temporal-Difference methods $(TD)$

## Monte Carlo

Generate an episode wrt $\pi$

Calculate state values for the episode, right to left

Add new values to list $V(s)$ for each $s$

Finally take mean of each $V(s)$ list — Initial estimates of $V(s)$ are discarded

# TD:

"Bootstrapping" — like DP — requires full model

Uses current estimates of $V(s_t)$ to update $V(s_{t+1})$

TD — hybrid of DP & MC

If we were to do MC incrementally

$$V(s_t) \leftarrow V(s_t) + \alpha \left[ G_t - V(s_t) \right]$$

← global value calculated for entire episode

$\uparrow$ learning rate

observed error / deviation

Replace $G_t$ by $\underbrace{R_{t+1}}_{\substack{\text{observed} \\ \text{reward}}} + \underbrace{\gamma V(S_{t+1})}_{\text{old estimate}}$

TD update rule

$$V(S_t) \leftarrow V(S_t) + \underset{\substack{| \\ \text{learning rate}}}{\alpha} \Big[ \underbrace{(R_{t+1} + \gamma V(S_{t+1})) - V(S_t)}_{\text{error}} \Big]$$

TD $[0]$ - zero lookahead

## 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

$$v_\pi(s) = E_\pi[\underline{G_t} \mid S_t = s] \qquad \text{Monte Carlo}$$

$$R_{t+1} + \gamma G_{t+1}$$
$$\overline{\overline{v_\pi(S_{t+1})}}$$

$$E_\pi\left[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s\right] \quad \text{TD[0]}$$

TD-error $\quad \delta_t \overset{def}{=} \left(R_{t+1} + \gamma V(S_{t+1})\right) - V(S_t)$

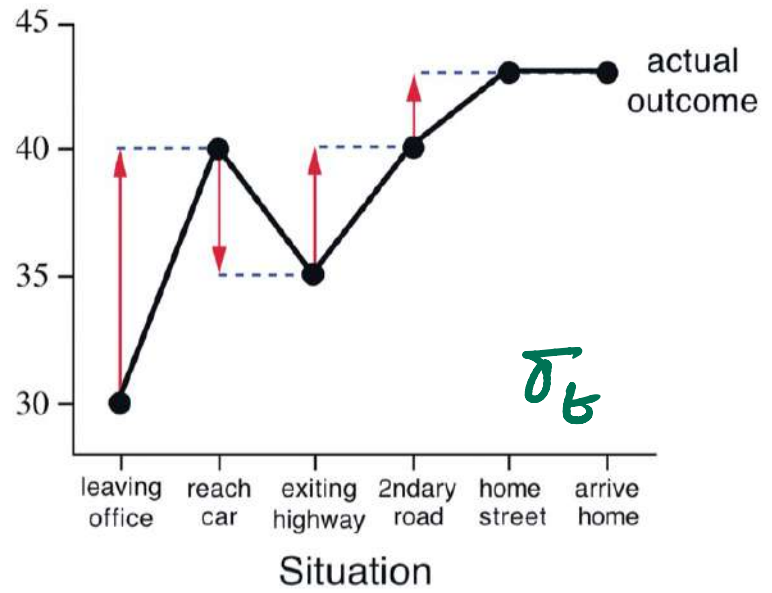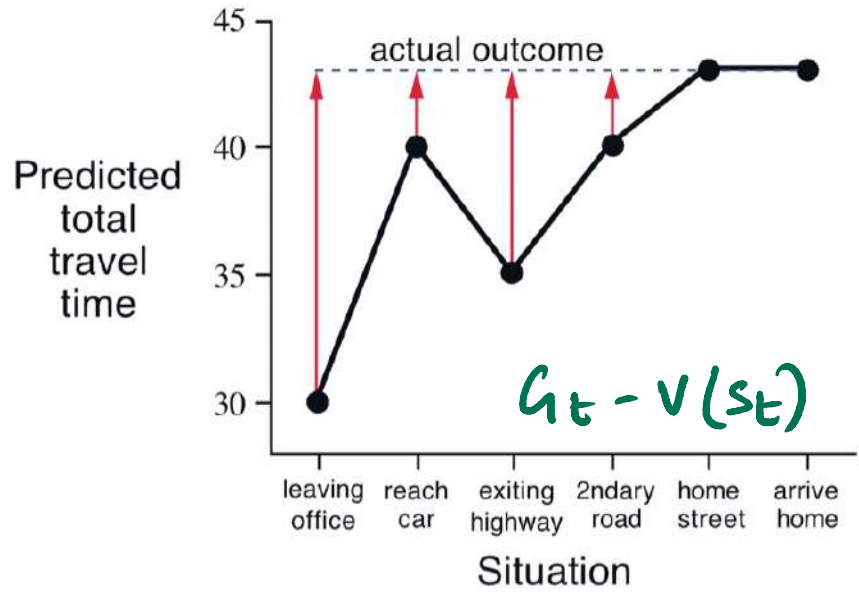MC-error $\quad G_t - V(S_t) = R_{t+1} + \gamma G_{t+1} - V(S_t) + \gamma V(S_{t+1}) - \gamma V(S_{t+1})$

**MC-error** $\quad G_t - V(S_t) = R_{t+1} + \gamma G_{t+1} - V(S_t) + \gamma V(S_{t+1}) - \gamma V(S_{t+1})$

$$= \underbrace{(R_{t+1} + \gamma V(S_{t+1}) - V(S_t)}_{\delta_t} + \underbrace{\gamma G_{t+1} - \gamma V(S_{t+1})}_{\gamma(G_{t+1} - V(S_{t+1}))}$$

$$= \delta_t + \gamma \cdot \delta_{t+1} + \gamma^2 \delta_{t+2} \cdots$$

$$= \sum_{k=t}^{T-1} \gamma^{k-1} \delta_k$$

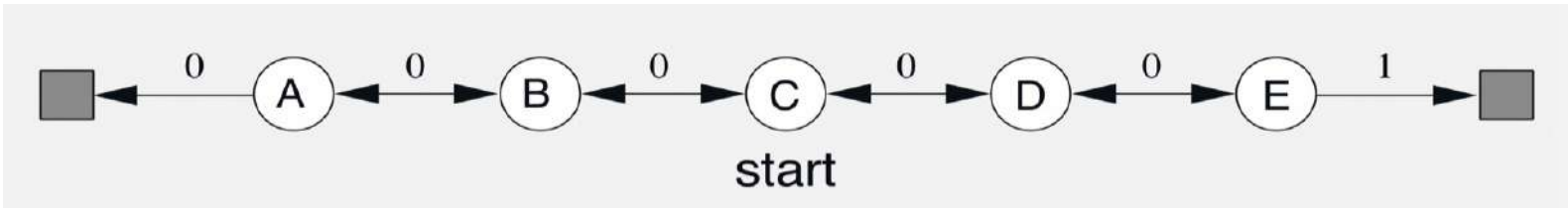| State | Elapsed Time (minutes) | Predicted Time to Go | Predicted Total Time |
|---|---|---|---|
| leaving office, friday at 6 | 0 | 30 | 30 |
| reach car, raining | 5 | 35 | 40 |
| exiting highway | 20 | 15 | 35 |
| 2ndary road, behind truck | 30 | 10 | 40 |
| entering home street | 40 | 3 | 43 |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



$G_t - V(s_t)$

$\delta_t$

# Advantage of TD vs MC

- Incremental update of $V$
- No need to wait till episode ends
    - May not even have finite episodes

Both TD & MC converge asymptotically to correct values

Which is faster?
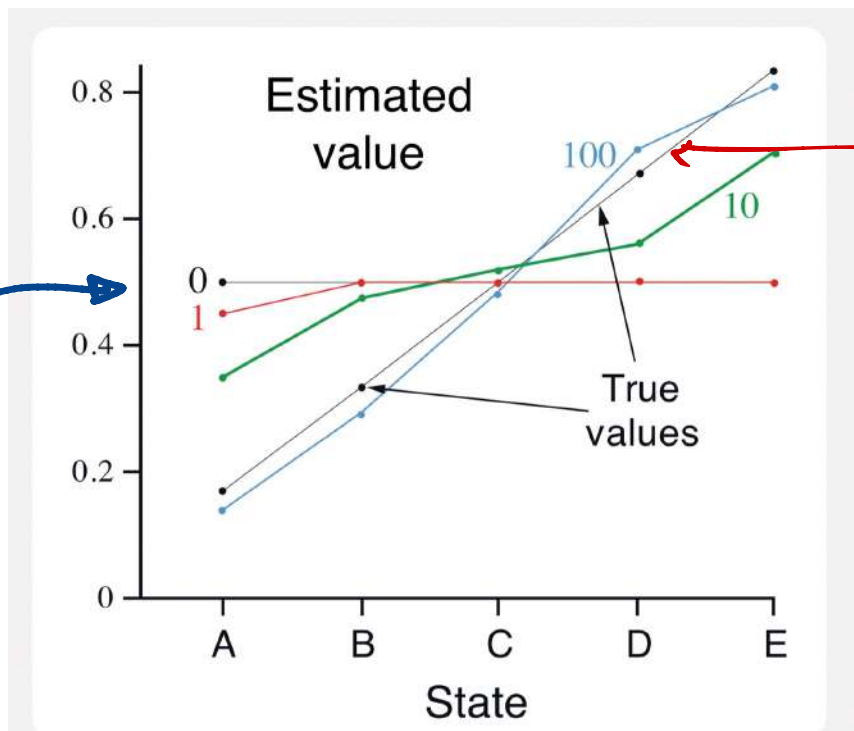
- Not been proven
- In practice TD appears faster

0    A    0    B    0    C    0    D    0    E    1

start

1/6      2/6      V = 1/2      4/6      5/6

TD[0]

starting

with V = 0.5

everywhere

theoretical
value

Estimated value

True values

State

Empirical RMS error, averaged over states

MC

TD

α=.01

α=.02

α=.04

α=.03

α=.15

α=.1

α=.05

learny rates

Walks / Episodes

# Batch learning

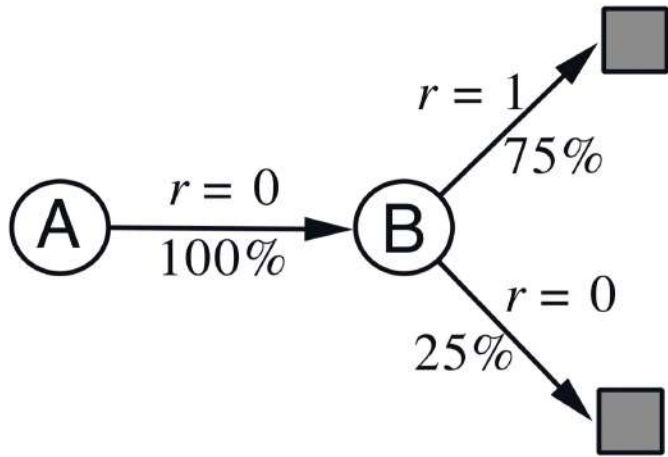Consider entire set of episodes as a single batch update

MC & TD give different answers

## Example 8 observations

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

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

$$V(B) = \frac{6}{8} = \frac{3}{4}$$

What is $V(A)$?

MC sets $V(A) = 0$

TD sets $V(A) = V(B)$

MDP learned
  by TD[0]

r = 0
100%
A → B
r = 1
75%
r = 0
25%

MC — Best estimate wrt Mean Square Error
  of observations

TD[0] - Best estimate wrt. MLE

# From value estimation to policy iteration

MC:    **On policy** — Same $\pi$ generates runs & gets update

       **Off policy** — runs are generated independent of $\pi$ to be updated

## On policy TD

Typically switch from estimating $v_\pi(s)$ to $q_\pi(sa)$

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

Update is a function of $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}) \Rightarrow$ 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 $Q(terminal, \cdot) = 0$

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

*Update $Q(S_t, A_t) \longrightarrow$ Update $\pi$*

Fewer steps
per episode

Actions

= vertical wind

Reward ~
−1 every where

− Find shortest path

# Off Policy    Q-learning

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

next action is best
wrt $Q$, not
determined by $\pi$

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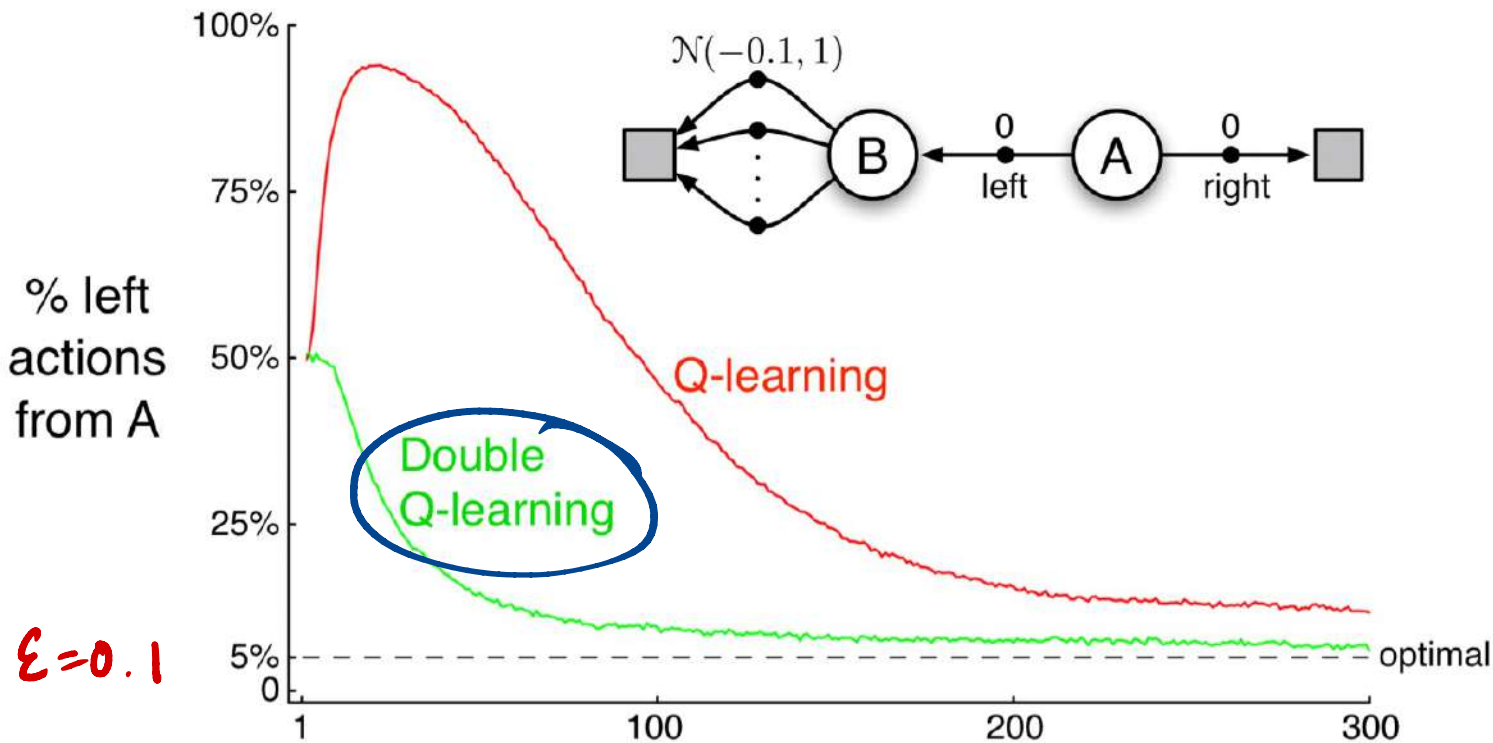    until $S$ is terminal

# Variant of SARSA - Expected SARSA

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

$$\sum_a \pi(a \mid S_{t+1}) \, Q(S_{t+1}, a)$$

# Maximization Bias

— tends to produce inflated estimates



% left actions from A

$\mathcal{N}(-0.1, 1)$

B

0
left

A

0
right

Q-learning

Double Q-learning

100%

75%

50%

25%

5%
0

optimal

$\varepsilon = 0.1$

1          100          200          300

**Double learning** — Simultaneously create 2 estimates, randomly use one to update other.

---

**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