

Model free evaluation — monte carlo, TD  
 control SARSA  $\alpha$ -learning

MC:  $V(s_t) \leftarrow V(s_t) + \alpha [G_t - V(s_t)]$  — (1)

TD:  $V(s_t) \leftarrow V(s_t) + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$

TD(0)

estimate of return

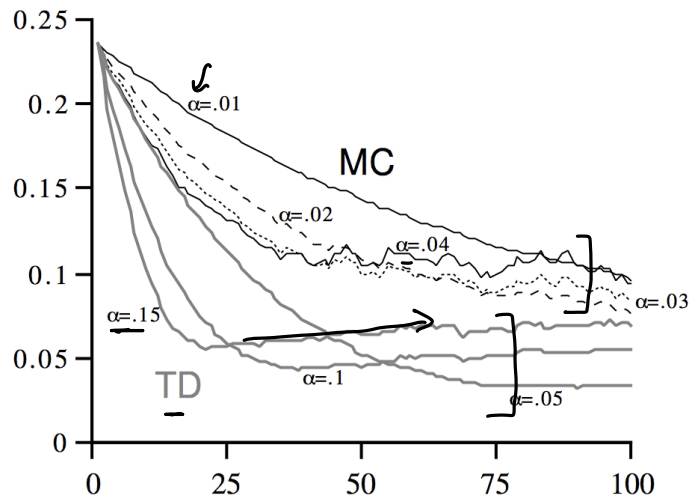
Monte Carlo

TD

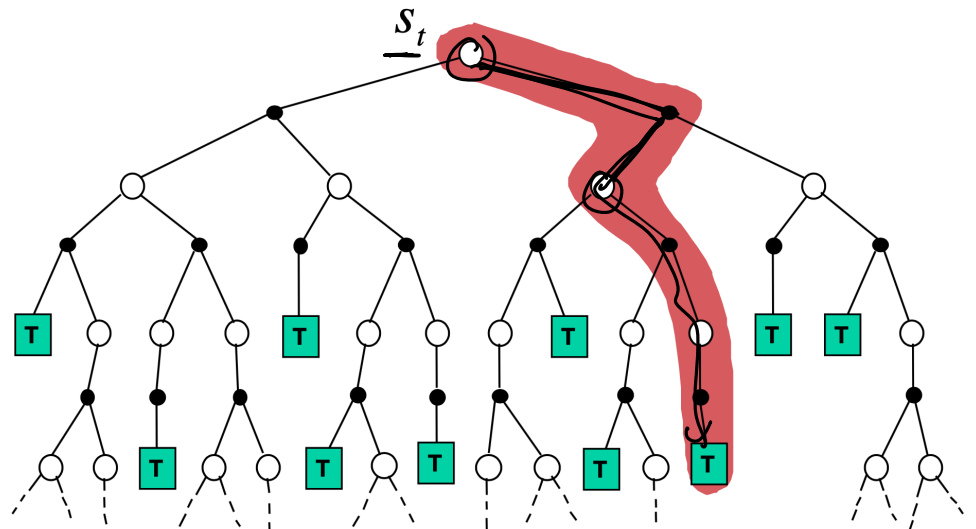
- learn from incomplete episodes
- lower variance but estimates could be biased
- usually more efficient
- exploits markovian property

- unbiased estimates but high variance

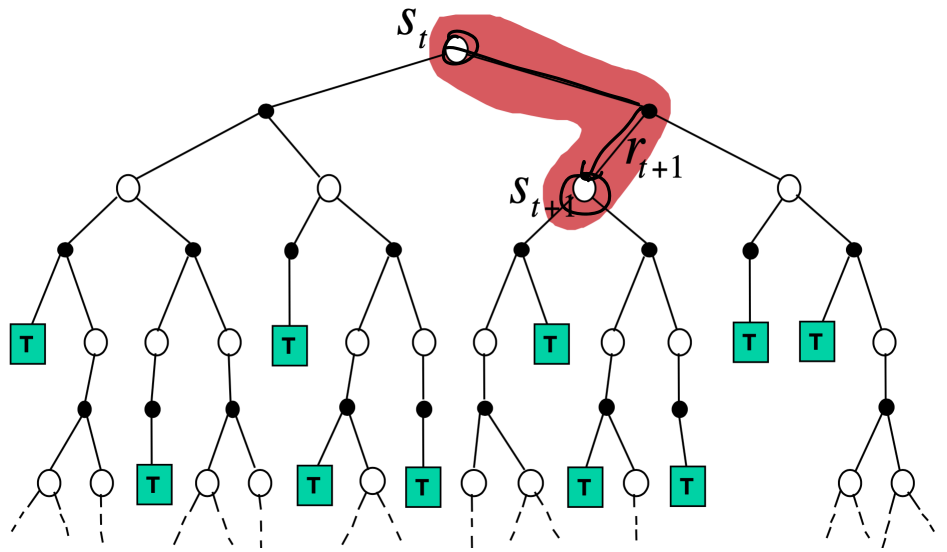
RMS error, averaged over states



monte carlo



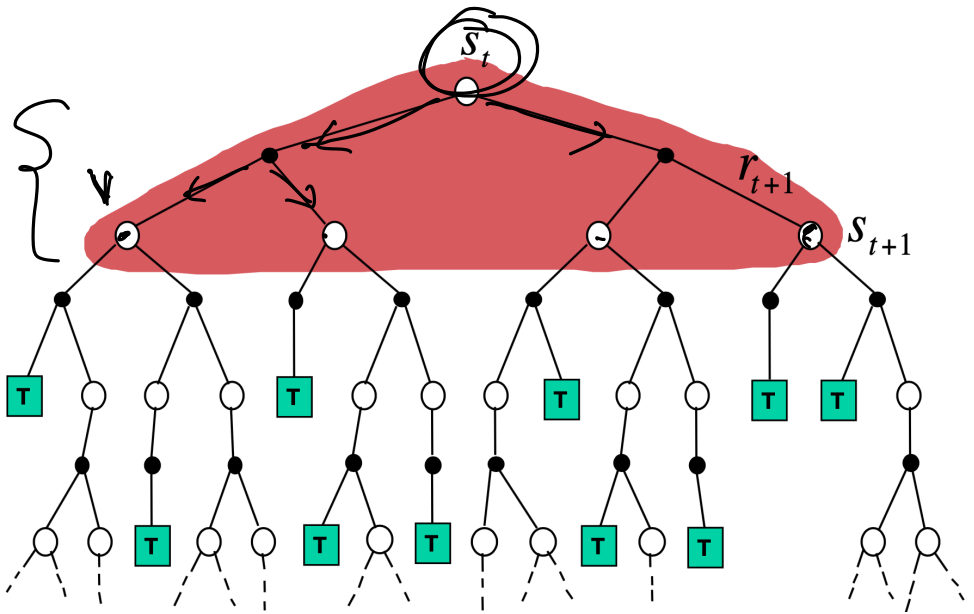
# Temporal Diff. Lec

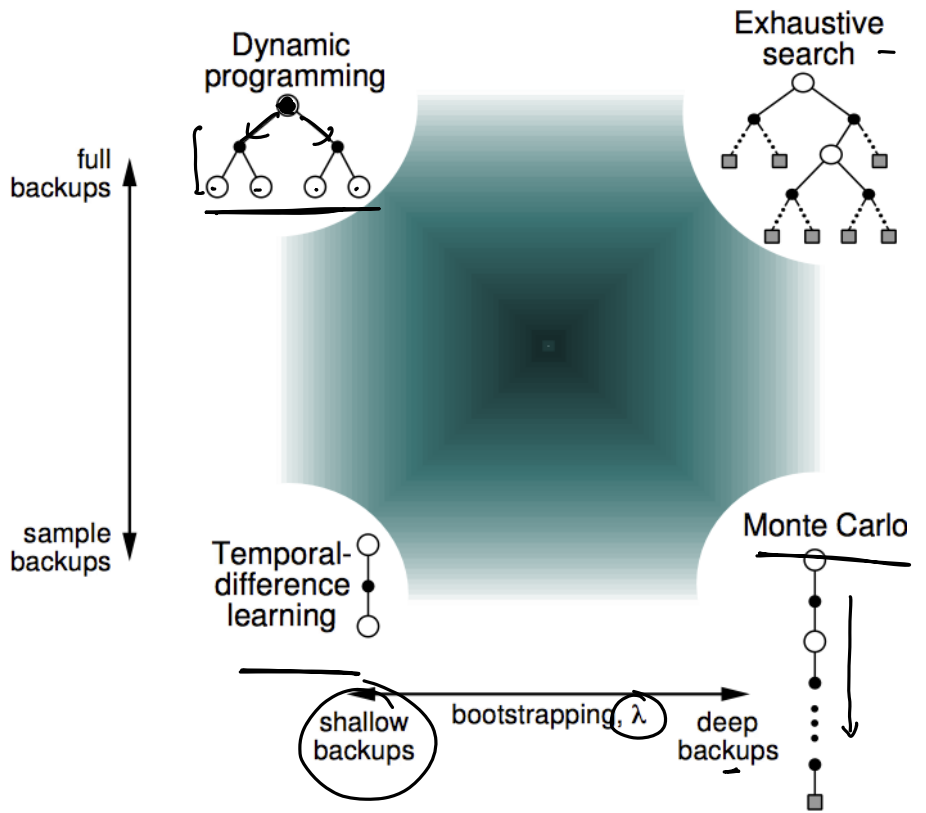


## Dynamic Programming

$$\sum_a \mathcal{R}(a|s) \sum_{s', r} p(s', r | s, a)$$


---





- $A \xrightarrow{0} B \xrightarrow{0} \text{goal}$  ]
- $B \xrightarrow{1} \text{goal}$
- $B \xrightarrow{1} \text{goal}$
- $B \xrightarrow{1} \text{goal}$
- $B \xrightarrow{0} \text{goal}$
- $B \xrightarrow{1} \text{goal}$
- $B \xrightarrow{1} \text{goal}$
- $B \xrightarrow{2} \text{goal}$

Monte Carlo

$\gamma = 1$

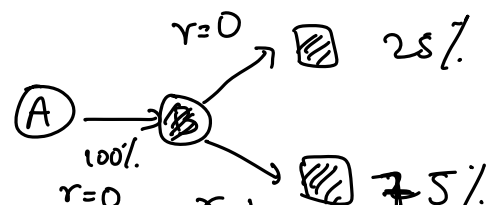
$V(A) = 0$

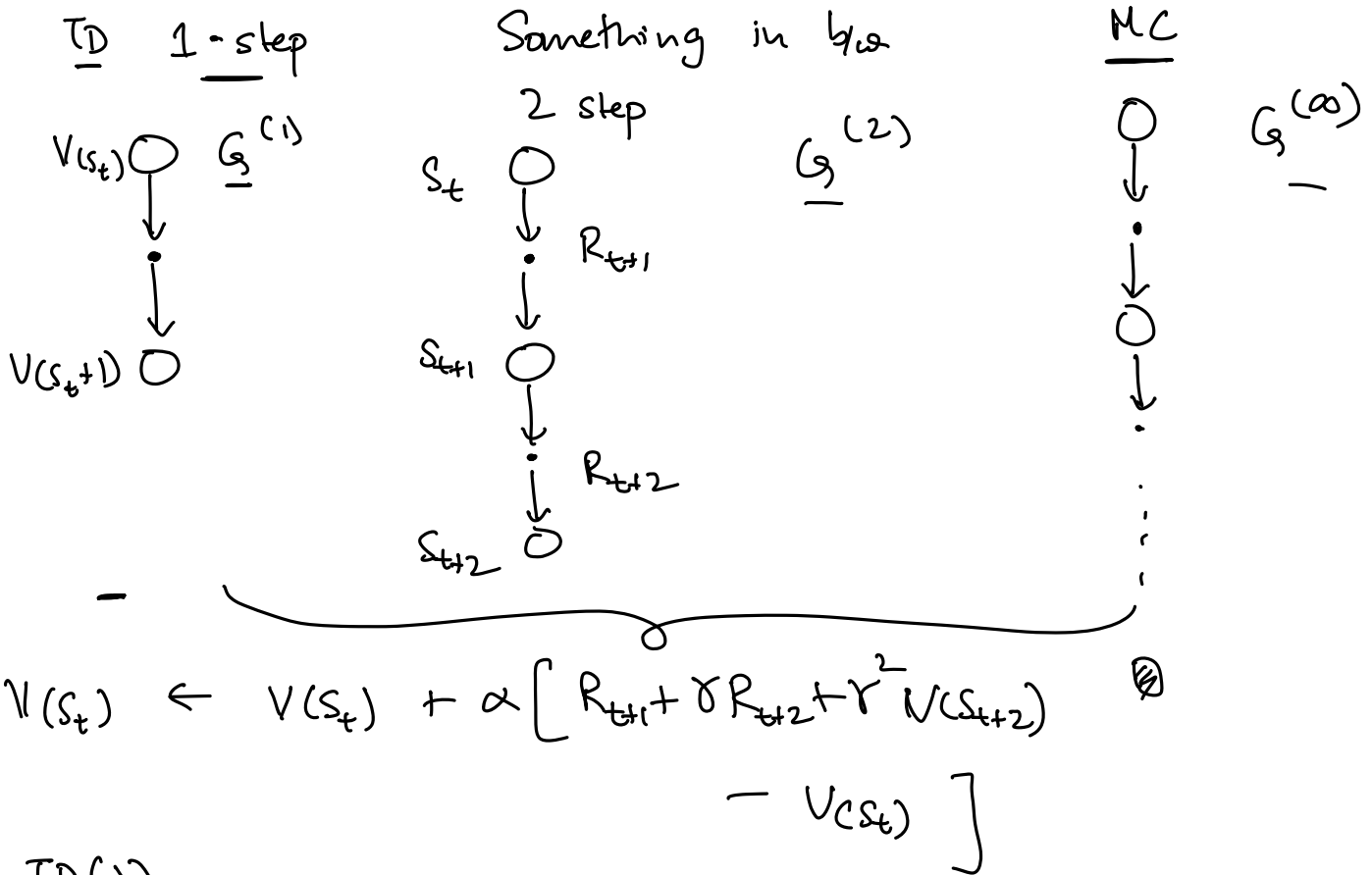
$V(B) = 3/4$

TD

$V(A) = 0.75$  ]

$V(B) = 0.75$  ]





TD( $\lambda$ )

$$G^\lambda = \frac{G^{(1)} + \lambda G^{(2)} + \lambda^2 G^{(3)} + \dots}{1 - \lambda}$$

$$G^0 = G^{(1)} \quad \text{TD}(0)$$

## Model-free Control

- Policy iteration w/ monte carlo estimates.

- policy evaluation [monte carlo  $V = \frac{G_\pi(s)}{\pi}$ ]

- policy improvements [greedy ...]

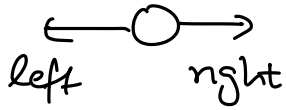
wait quite work.

$$\left\{ \begin{array}{l} \arg \max_a q_\pi(s, a) \\ \arg \max_a \sum p(s', r | s, a) [r + \gamma V_\pi(s')] \end{array} \right.$$

$q_\pi(s, a)$

$V = \frac{G_\pi(s)}{\pi}$

# Stochastic Environment.



$$Q(\text{left}) = 0$$

$$Q(\text{right}) = 1$$

## $\epsilon$ -greedy policies.

m action

$$\Rightarrow \pi(a|s) = \begin{cases} \underline{1-\epsilon} + \underline{\epsilon/m} & \text{if } a = \text{argmax } \underline{Q_{\pi}(s,a)} \\ \underline{\epsilon/m} & \text{OW} \end{cases}$$

## $\epsilon$ greedy policy improvement

for any  $\epsilon$  soft policy  $\pi$ , the  $\epsilon$  greedy policy  $\pi'$  wrt  $q_{\pi}(s,a)$  is an improvement over  $\pi$ .

$$v_{\pi'}(s) \geq v_{\pi}(s) \quad \forall s.$$

## Monte Carlo Control.

- policy evaluation  $\rightarrow$  estimate  $q_{\pi}(s,a)$  & not  $v_{\pi}(s)$
- policy improvement  $\rightarrow$   $\epsilon$  greedy wrt  $Q_{\pi}(s,a)$

**TD Control** Estimate  $Q_{\pi}(s,a)$  with TD.

## SARSA

Repeat (for episodes)

Sample initial state  $s$

Act,  $A$  in state  $s$  [using  $\epsilon$  greedy policy wrt  $Q(s,a)$ ]

Repeat (for steps in episode)

execute  $A$  & observe  $s', R$

sample  $A'$  using  $\epsilon$  greedy wrt  $Q(s, A)$

$$* Q[s, A] \leftarrow Q[s, A] + \alpha [R + \gamma \underline{Q(s, A')} - Q(s, A)]$$

$S \leftarrow S'$

$A \leftarrow A'$

$$Q(s, A) \rightarrow \underline{Q_{\epsilon}^*}(s, A) \quad \epsilon \rightarrow 1/t$$

It turns out SARSA can be converted into an off policy algorithm.

### Q-learning

- act as per some behavior policy  $\mu(\cdot | s_t)$   
continuing exploration: as long as it continues to experience all state action pairs.

- But, when updating  $Q$  function, use the  $Q$ -value of the policy that you are trying to learn.

$$\underline{Q}(s_t, \underline{A}_t) \leftarrow \underline{Q}(s_t, \underline{A}_t) + \alpha (R_{t+1} + \gamma \underline{Q}(s_{t+1}, \underline{A}_{t+1}) - \underline{Q}(s_t, \underline{A}_t))$$

$(\underline{s}_t, \underline{A}_t, R_{t+1}, s_{t+1}) \xrightarrow{\mu} (s_{t+1}, \underline{A}_{t+1}, R_{t+2}, s_{t+2}) \rightarrow \dots$   
 $A' \sim \pi(s_{t+1})$   
 $Q, \pi$

Q learning for control

- Behavioural policy  $\mu \leftarrow \in$  greedy wrt  $Q_{\pi}$  function
- Target policy  $\pi \rightarrow$  greedy policy wrt  $Q$ .

$$\downarrow$$
$$\operatorname{argmax}_a Q_{\pi}(s, a)$$

$$Q_{\pi}(s_t, A_t) \leftarrow Q_{\pi}(s_t, A_t) + \alpha \left( R_{t+1} + \gamma \max_a Q_{\pi}(s_{t+1}, a) - Q_{\pi}(s_t, A_t) \right)$$

Watkins 9.2

Q learning converges to  $Q^*$  (S, A).