

# Reinforcement Learning

## Lecture 7

TD( $\lambda$ )

Eligibility Traces

Robert Peharz

Institute of Theoretical Computer Science  
Graz University of Technology  
Winter Term 2023/24

# Recap: Model-free Control with Monte Carlo

$$q_{\pi}(s, a)$$

$$= \underbrace{r(s, a) + \gamma \sum_{s'} p(s' | s, a) v_{\pi}(s')}_{\text{MDP, now unknown}}$$

Exploration



Exploitation

On-policy methods

need to find a compromise

Off-policy methods

separate the problem

Behavior policy  $b$

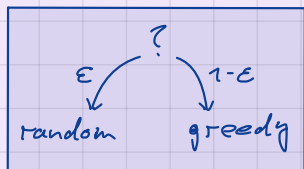
generates behavior, data

Target policy  $\pi$

learns from  $b$ 's data

## Monte Carlo Control with Exploring Starts

$\epsilon$ -greedy Policy: With probability  $\epsilon$ , take random action (uniformly), otherwise take greedy action



$$\pi(a|s) = \begin{cases} 1-\epsilon + \frac{\epsilon}{|A(s)|} & \text{if } a = \underset{a'}{\operatorname{argmax}} q(s, a') \\ \frac{\epsilon}{|A(s)|} & \text{o.w.} \end{cases}$$



Image: [pinterest.uk](https://www.pinterest.uk)

GLIE

Greedy in the Limit with Infinite Exploration

## Off-policy Monte Carlo via Importance Sampling

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

# Recap: TD-based Control

$$q(s,a) \leftarrow q(s,a) + \alpha (\text{goal} - q(s,a))$$

on-policy

off-policy

sample A

Sarsa

$(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$   
SARSA

Goal =  
 $R_{t+1} + \gamma q(S_{t+1}, A_{t+1})$

Importance-weighted Sarsa

Goal =  
 $R_{t+1} + \gamma \frac{\hat{\pi}(A_{t+1} | S_{t+1})}{b(A_{t+1} | S_{t+1})} q(S_{t+1}, A_{t+1})$

$E_A$

Expected Sarsa

Goal =  
 $R_{t+1} + \gamma \sum_{a'} \hat{\pi}(a' | S_{t+1}) q(S_{t+1}, a')$

Q-learning (Sarsa Max)

Goal =  
 $R_{t+1} + \gamma \max_{a'} q(S_{t+1}, a')$

## TD( $\lambda$ ) and Eligibility Traces

## Recall: n-step TD Updates

General update:  $V(s) \leftarrow V(s) + \alpha (\hat{G} - V(s))$

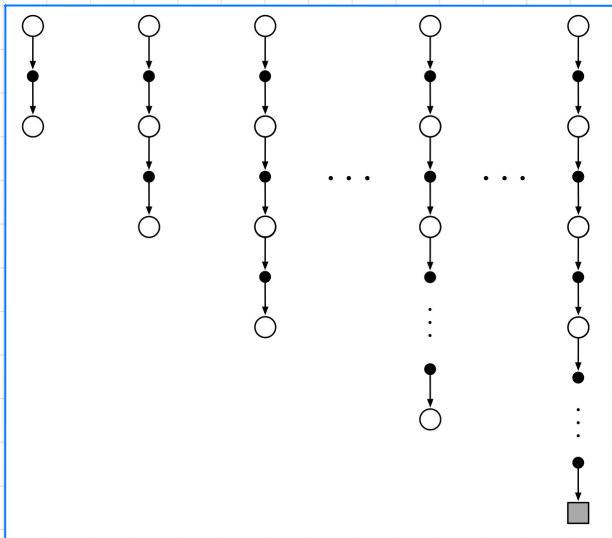
(1-step) TD:  $\hat{g} = R_{t+1} + \gamma v(S_{t+1})$

2-step TD:  $\hat{g} = R_{t+1} + \gamma R_{t+2} + \gamma^2 v(S_{t+2})$

3-step TD:  $\hat{g} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v(S_{t+3})$

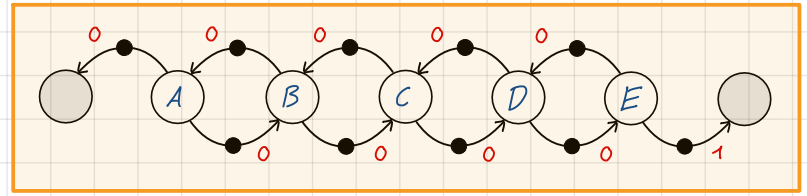
•

$\infty$ -step TD (MC)  $\hat{g} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots = g_t$



Which is the best  $n$ ?  
Which converges fastest?

# Random Walk



Consider the random walk example again, but with 19 states.  
Difference to true  $v_{\pi}$  for various  $n$  and  $\alpha$ :

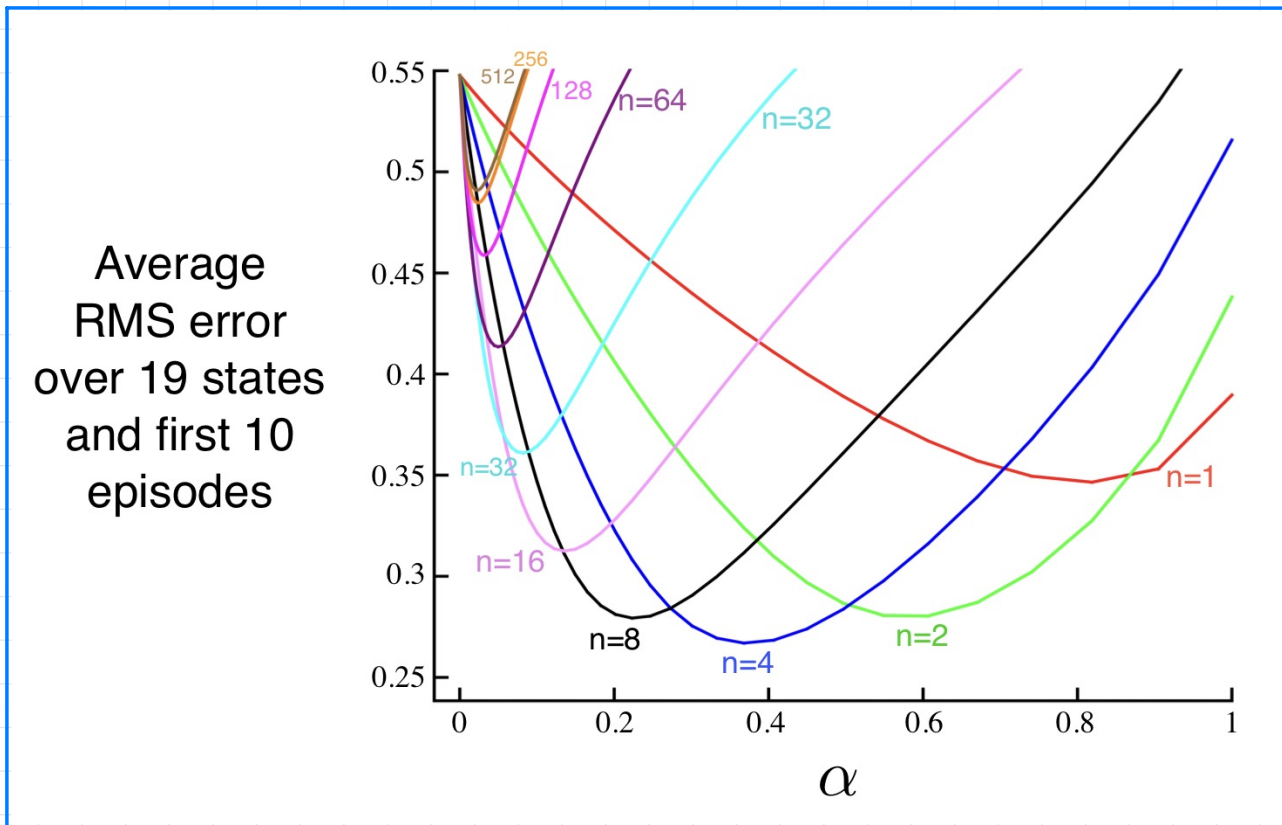


Image: Sutton & Barto

Sweet spot at  $n=4$  and  $\alpha=0.39$ . This is highly problem-dependent.  
Since optimal  $n$  is hard to pick, why not take all of them?

# N-step Return

(1-step) TD:  $\hat{g} = R_{t+1} + \gamma v(S_{t+1})$

2-step TD:  $\hat{g} = R_{t+1} + \gamma R_{t+2} + \gamma^2 v(S_{t+2})$

3-step TD:  $\hat{g} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v(S_{t+3})$

$\vdots$

$\infty$ -step TD (MC)  $\hat{g} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots = g_t$

$=: \hat{g}_t^1$

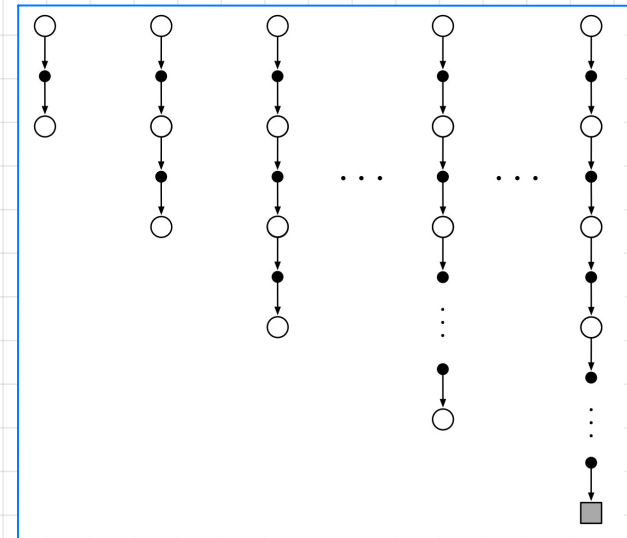
$=: \hat{g}_t^2$

$=: \hat{g}_t^3$

$=: \hat{g}_t^\infty$

## n-step return

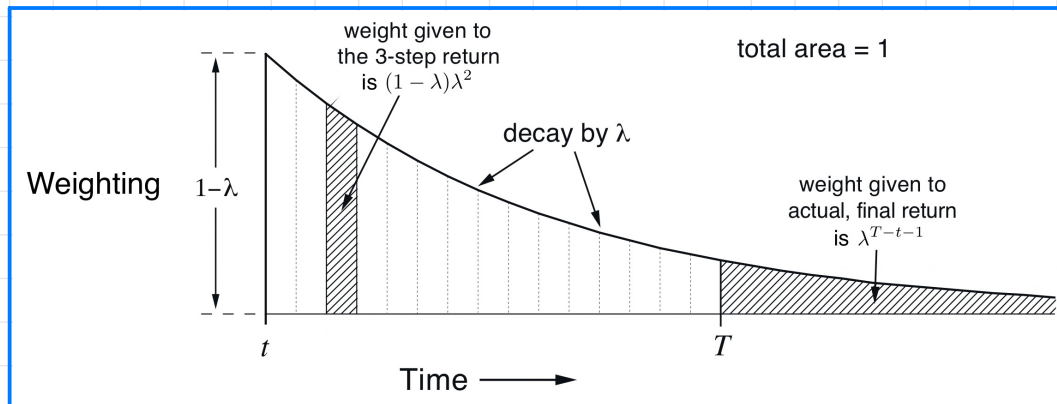
$$\hat{g}_t^n := \sum_{k=0}^{n-1} \gamma^k R_{t+k+1} + \gamma^n v(S_{t+n})$$



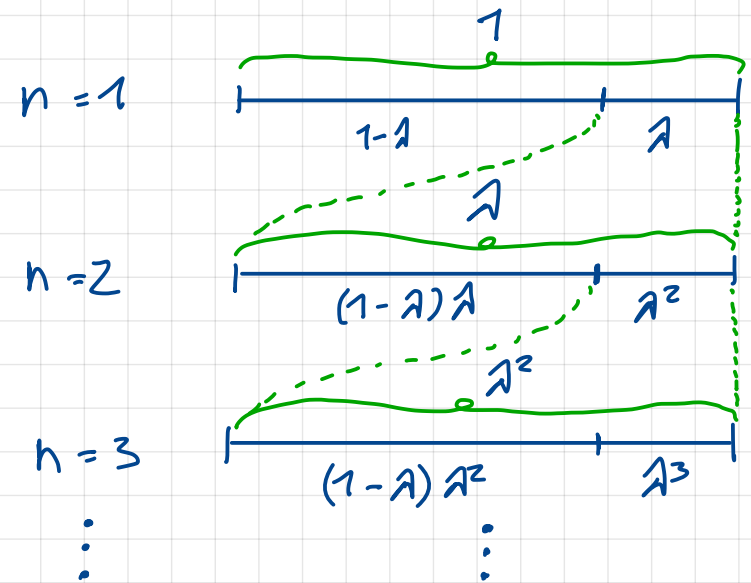
- any  $\hat{g}_t^n$  is a legit goal
- any weighted average of  $\hat{g}_t^n$ 's is also a legit goal!
- e.g. we could use  $\frac{1}{2} \hat{g}_t^2 + \frac{1}{2} \hat{g}_t^4$ , but also average all  $\hat{g}_t^n$ 's

# $\lambda$ -Return

- for  $\lambda \in [0, 1]$  consider weights  $(1-\lambda)\lambda^{n-1}$
- note that  $\sum_{n=1}^{\infty} (1-\lambda)\lambda^{n-1} = 1$



"stick breaking weights"



Average all  $\hat{G}_t^n$  with these weights:  $\lambda$ -Return

$$\hat{G}_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \hat{G}_t^n$$

For finite episode length  $T$ :

$$\hat{G}_t^\lambda = (1-\lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} \hat{G}_t^n + \lambda^{T-t-1} \hat{G}_t^\infty$$

Note:

$$\begin{aligned} \hat{G}_t^0 &= \hat{G}_t^1 \\ \hat{G}_t^1 &= \hat{G}_t^\infty \end{aligned}$$

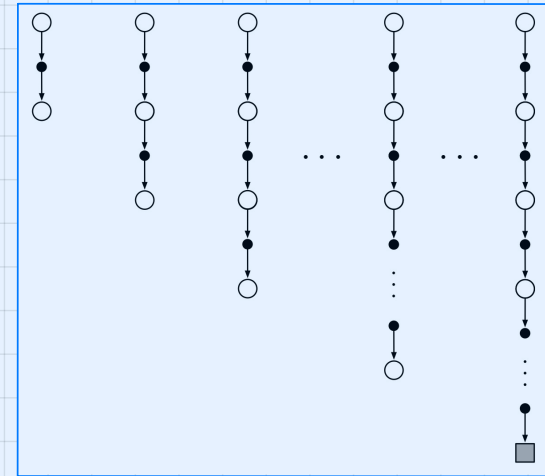


# $\lambda$ -Return for Prediction

## n-step TD

$$\hat{g}_t^n := \sum_{k=0}^{n-1} \gamma^k R_{t+k+1} + \gamma^n v(S_{t+n})$$

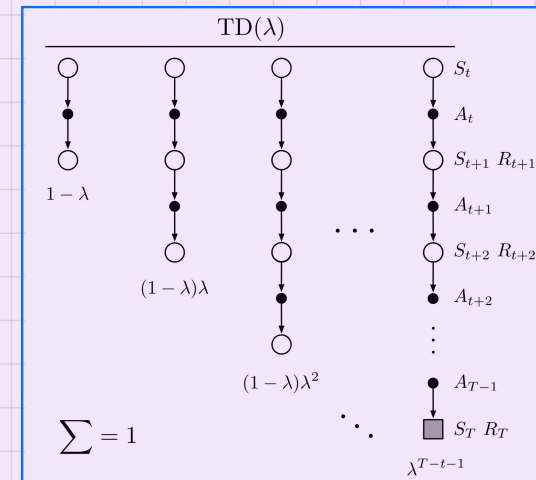
$$v(s) \leftarrow v(s) + \alpha (\hat{g}_t^n - v(s))$$



## $\lambda$ -return algorithm

$$\hat{g}_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \hat{g}_t^n$$

$$v(s) \leftarrow v(s) + \alpha (\hat{g}_t^\lambda - v(s))$$



# Prediction with $\lambda$ -Return

Let  $\alpha \in (0, 1]$ ,  $\lambda \in [0, 1]$

Initialize  $v(s)$  arbitrarily, except  $v(s) = 0$  for terminal  $s$

- repeat

- using  $\tilde{\Pi}$  generate episode

$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, \dots, S_{T-1}, A_{T-1}, R_T, S_T$  ← terminal

- for  $t = T-1 \dots 0$

- for  $n = 1 \dots T-t$   $\hat{g}_t^n = \begin{cases} R_{t+1} + \gamma v(S_{t+1}) & \text{if } n=1 \\ R_{t+1} + \gamma \hat{g}_{t+1}^{n-1} & \text{if } 2 \leq n \leq T-t \end{cases}$  (\*)

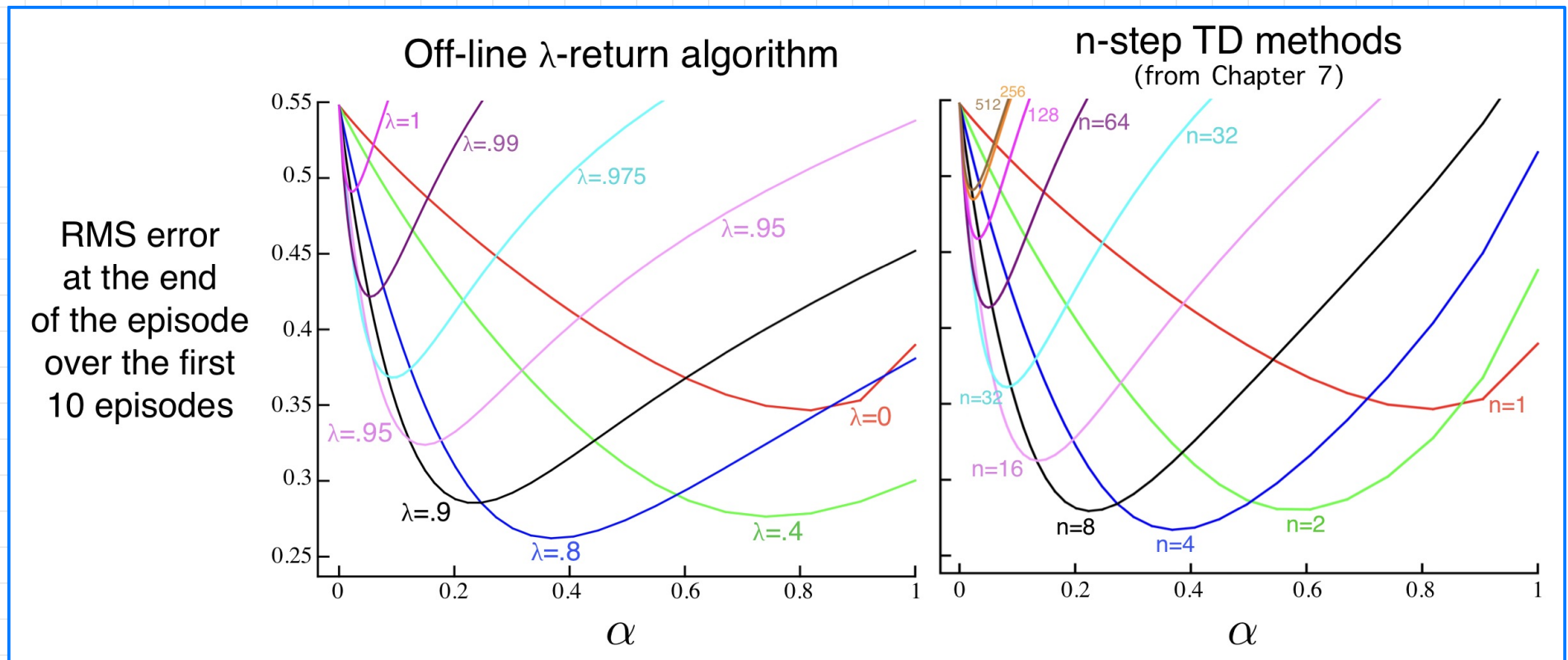
-  $\hat{g}_t^1 = (1-\lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} \hat{g}_t^n + \lambda^{T-t-1} \hat{g}_t^{T-t}$

-  $v(S_t) \leftarrow v(S_t) + \alpha (\hat{g}_t^1 - v(S_t))$

$$\begin{aligned}
 (*) \quad \hat{g}_t^n &:= \sum_{k=0}^{n-1} \gamma^k R_{t+k+1} + \gamma^n v(S_{t+n}) = R_{t+1} + \sum_{k=1}^{n-1} \gamma^k R_{t+k+1} + \gamma^n v(S_{t+n}) \\
 &= R_{t+1} + \gamma \left( \sum_{k=0}^{n-2} \gamma^k R_{t+k+2} + \gamma^{n-1} v(S_{t+1+n-1}) \right) = R_{t+1} + \gamma \hat{g}_{t+1}^{n-1}
 \end{aligned}$$

# N-step TD vs. $\lambda$ -Return

Consider the random walk example with 19 states.  
Difference to true  $v_{\pi}$  for various  $n(\lambda)$  and  $\alpha$



Some intermediate value of  $\lambda=0.5$  "will usually do."

# Control with $\lambda$ -Return

Let  $\epsilon > 0$ ,  $\alpha \in (0, 1]$ ,  $\lambda \in [0, 1]$

Initialize  $q(s, a)$  arbitrarily, except  $q(s, a) = 0$  for terminal  $s$

- repeat

- using  $\epsilon$ -greedy( $q$ ), generate episode

$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, \dots, S_{T-1}, A_{T-1}, R_T, S_T$  ← terminal

- for  $t = T-1 \dots 0$

$$\hat{g}_t^n = \begin{cases} R_{t+1} + \gamma v(S_{t+1}) & \text{if } n=1 \\ R_{t+1} + \gamma \hat{g}_{t+1}^{n-1} & \text{if } 2 \leq n \leq T-t \end{cases}$$

$$\hat{g}_t^\lambda = (1-\lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} \hat{g}_t^n + \lambda^{T-t-1} \hat{g}_t^{T-t}$$

$$q(s^t, A^t) \leftarrow q(s^t, A^t) + \alpha (\hat{g}_t^\lambda - q(s^t, A^t))$$

What is the drawback of this algorithm?

# $\lambda$ -Return as Forward View Algorithm

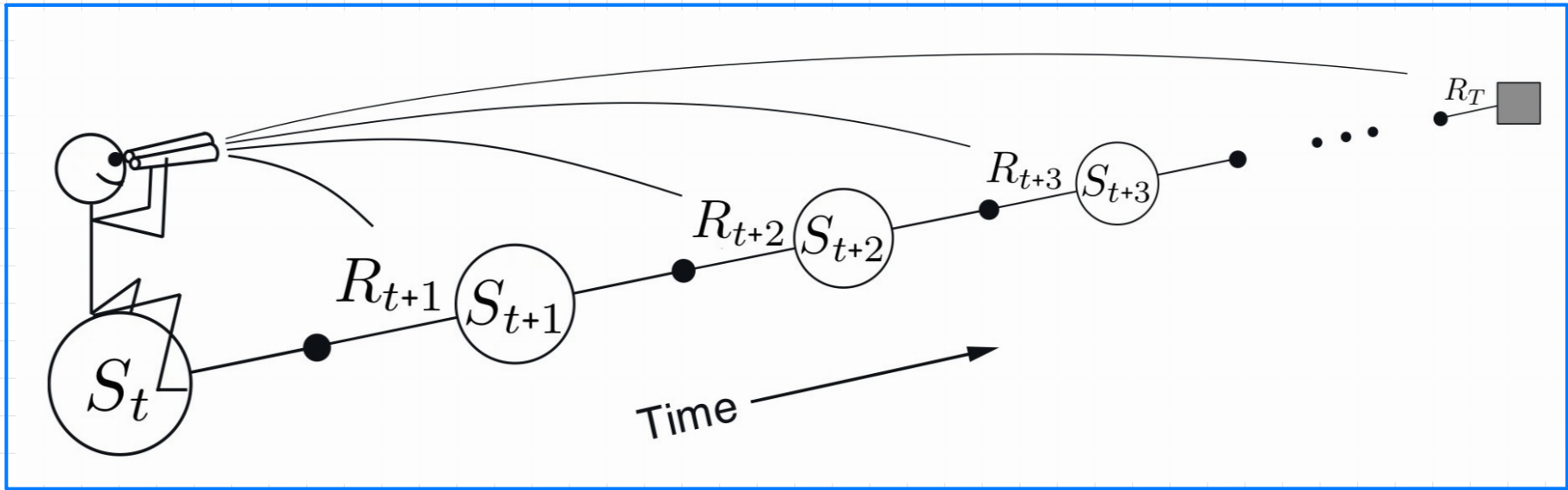


Image: Sutton & Barto

- $\lambda$ -return is a "forward view" algorithm
- we need to wait until the end of the episode
- same disadvantage as MC (since  $\hat{g}_t^\lambda$  contains  $\hat{g}_t^\infty$ )

# Is there a Backward View?

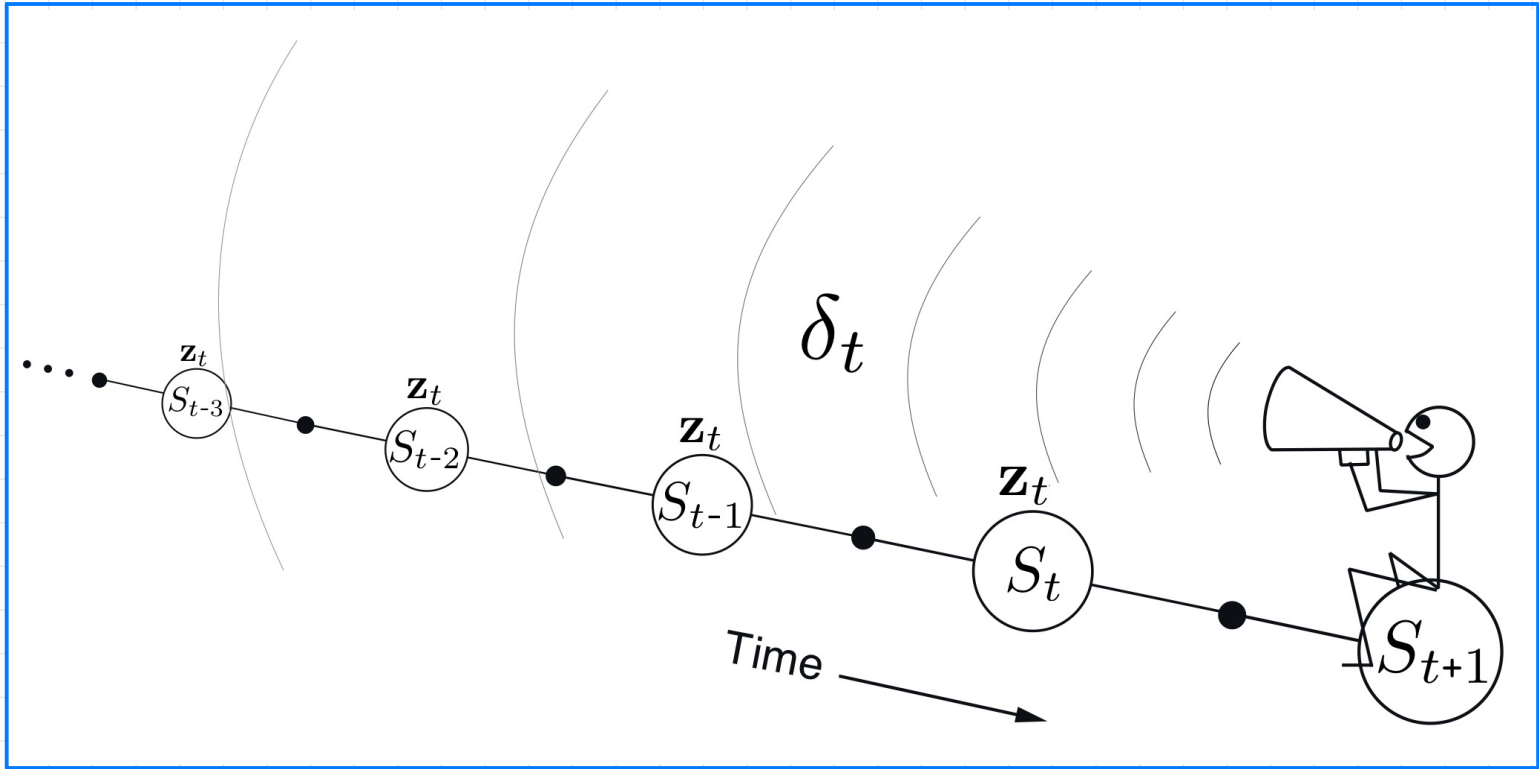
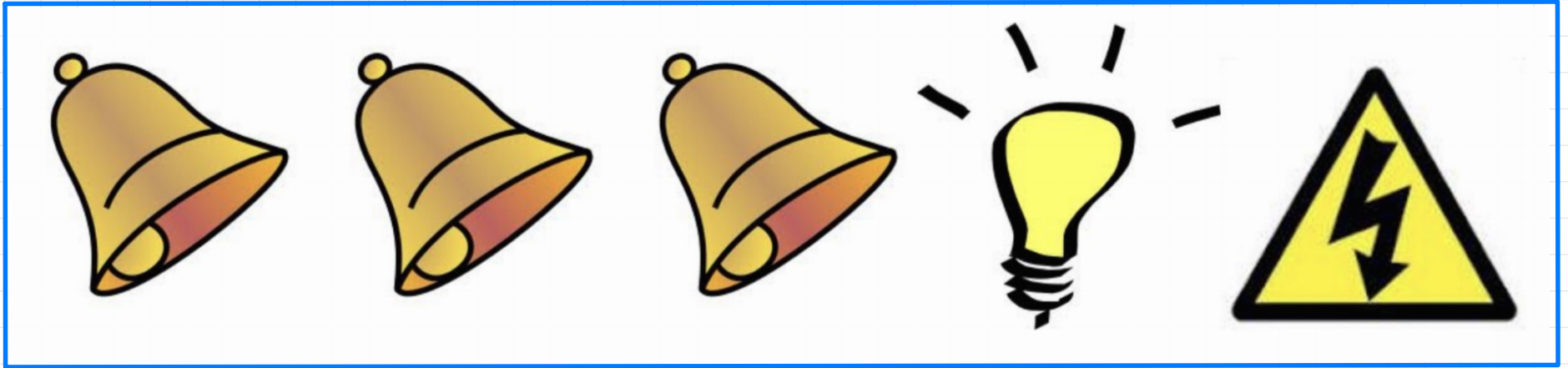


Image: Sutton & Barto

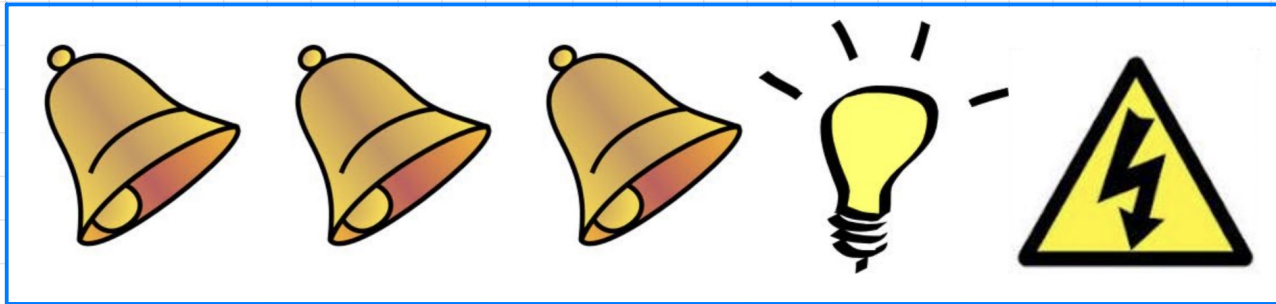
# Eligibility Traces

Do the bells or the light cause the shock?

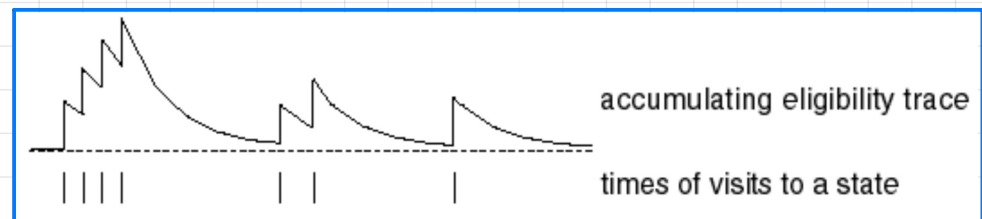




# Eligibility Traces



- did the bells or the light cause the electric shock?
- bells: more frequent  
light: more recent
- heuristic credit assignment
- eligibility traces combine both heuristics



$$t=0: e(s) \leftarrow 0$$

$$t>0: e(s) \leftarrow \lambda \gamma e(s) + \mathbb{1}[S_t = s]$$

exponential decay

"pumps up by 1" if state is visited

## TD( $\lambda$ )

"One the oldest and most widely used algorithms in RL"  
Sutton & Barto, p. 292

Given:  $\hat{\pi}$  // TD( $\lambda$ ) is a prediction algorithm

$\alpha, \lambda, \gamma$

initialize  $v(s)$  arbitrarily  $\forall s$

repeat until  $v$  has converged

$e(s) \leftarrow 0 \quad \forall s$

initialize  $S$

repeat until  $S$  is terminal or  $v$  has converged

$A \sim \hat{\pi}(a|S)$

observe  $S', R$  //  $\sim p(s', a|S, A)$

$e(s) \leftarrow \lambda \gamma e(s) + \mathbb{1}[S=s] \quad \forall s$

$\delta \leftarrow R + \gamma v(S') - v(S)$  // TD error

$v(s) \leftarrow v(s) + \alpha \delta e(s) \quad \forall s$

$S \leftarrow S'$

# Forward $\lambda$ -Return vs TD( $\lambda$ )

Update using  $\lambda$ -return:

$$v(S_t) \leftarrow v(S_t) + \alpha \underbrace{(G_t^\lambda - v(S_t))}_{\text{error term}}$$

It can be shown that:

$$G_t^\lambda - v(S_t) = \sum_{t=1}^{\infty} (\lambda \gamma)^{t-1} \overset{\text{eligibility trace } e}{\gamma} \overset{\text{TD error } \delta}{(R_{t+1} + \gamma v(S_{t+1}) - v(S_t))}$$

Hence,  $\lambda$ -return and TD( $\lambda$ ) are almost the same algorithm!   
 TD( $\lambda$ ) updates immediately, while  $\lambda$ -return waits until end of episode. Accumulating  $\delta e$  in TD( $\lambda$ ) until end of episode before updating would make the two equivalent

- TD(0) is TD
- TD(1) is akin to MC, but with immediate updates

# Forward $\lambda$ -Return vs TD( $\lambda$ )

RMS error  
at the end  
of the episode  
over the first  
10 episodes

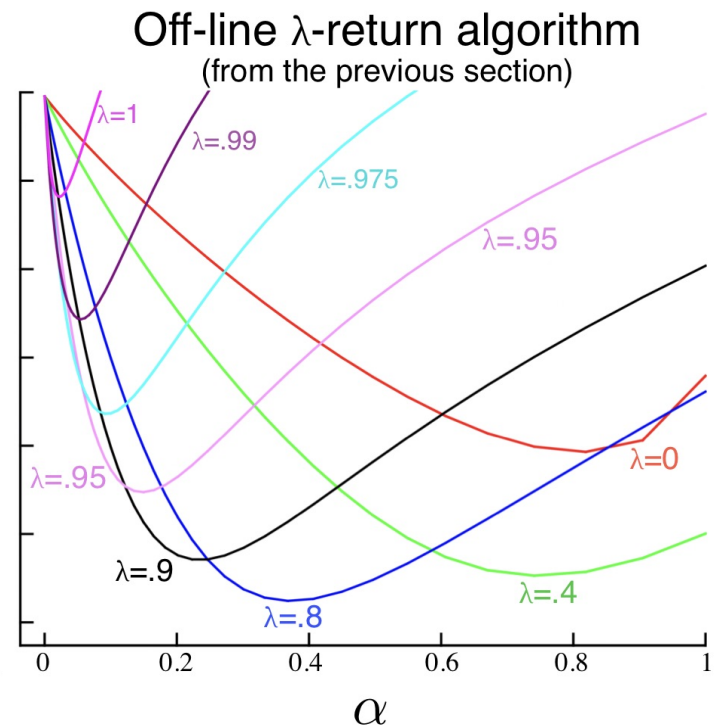
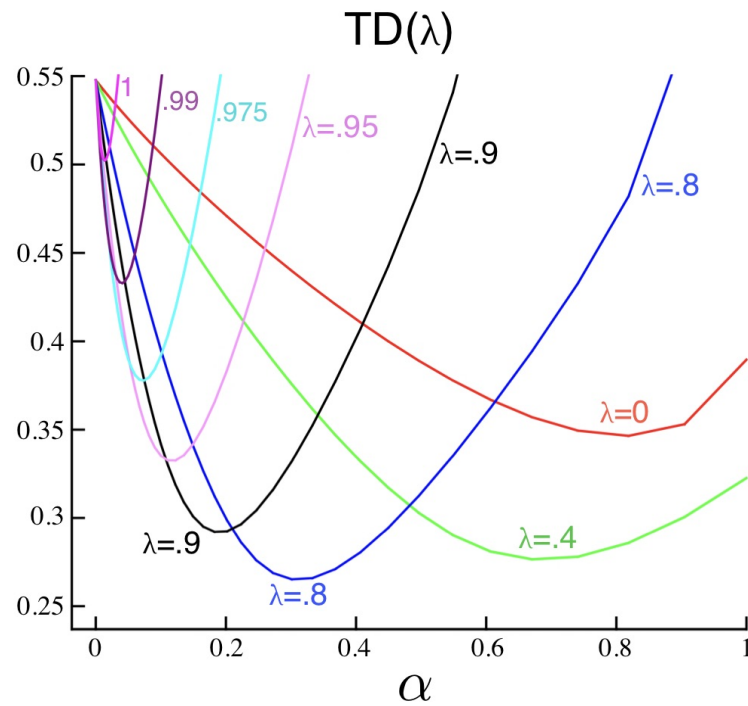


Image: Sutton & Barto

# Sarsa( $\lambda$ )

Given:  $\alpha, \lambda$

initialize  $q(s,a)$  arbitrarily  $\forall s,a$

repeat until  $q$  has converged

$e(s,a) \leftarrow 0 \quad \forall s,a$

initialize  $S$

$A \sim \epsilon\text{-greedy}(q(S, \cdot))$

repeat until  $S$  is terminal or  $q$  has converged

observe  $S', R \quad // \sim p(s', r | S, A)$

$A' \sim \epsilon\text{-greedy}(q(S', \cdot))$

$e(s,a) \leftarrow \lambda \gamma e(s,a) + \mathbb{1}[S=s \wedge A=a] \quad \forall s,a$

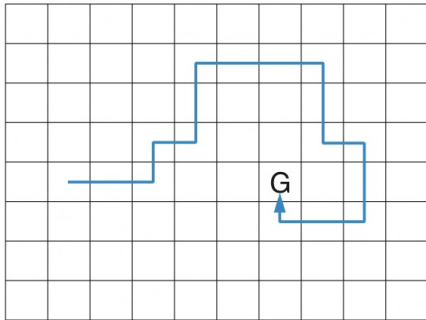
$\delta \leftarrow R + \gamma q(S', A') - q(S, A) \quad // \text{TD error}$

$q(s,a) \leftarrow q(s,a) + \alpha \delta e(s,a) \quad \forall s,a$

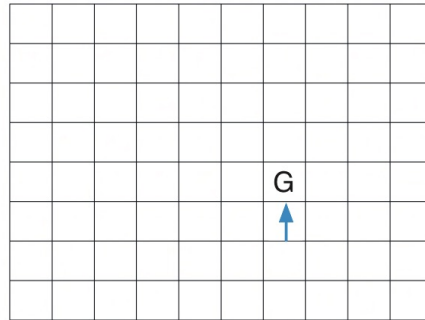
$S, A \leftarrow S', A'$

# Eligibility Trace of Sarsa( $\lambda$ )

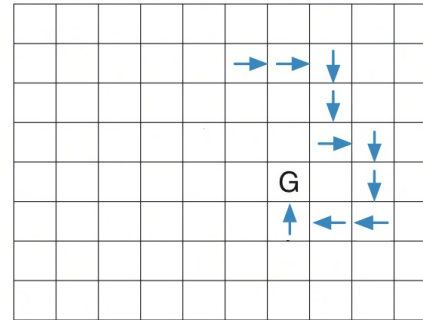
Path taken



Action values increased  
by one-step Sarsa



Action values increased  
by 10-step Sarsa



Action values increased  
by Sarsa( $\lambda$ ) with  $\lambda=0.9$

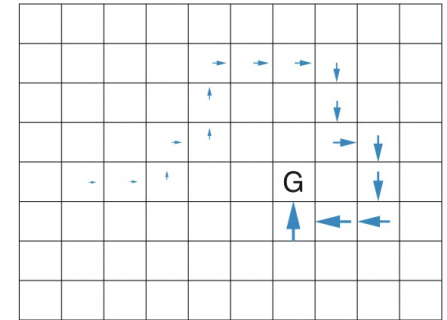
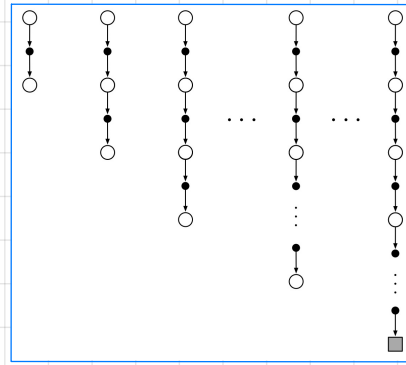


Image: Sutton & Barto

# Summary

## n-step return

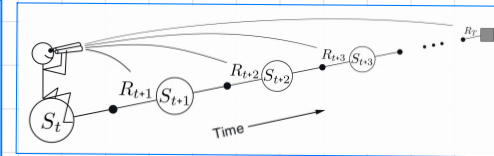
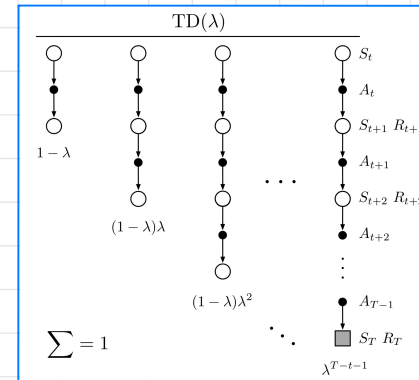
$$\hat{g}_t^n := \sum_{k=0}^{n-1} \gamma^k R_{t+k+1} + \gamma^n v(S_{t+n})$$



## $\lambda$ -return algorithm

$$\hat{g}_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \hat{g}_t^n$$

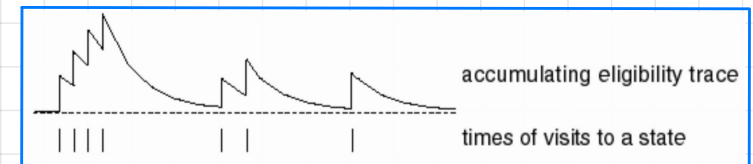
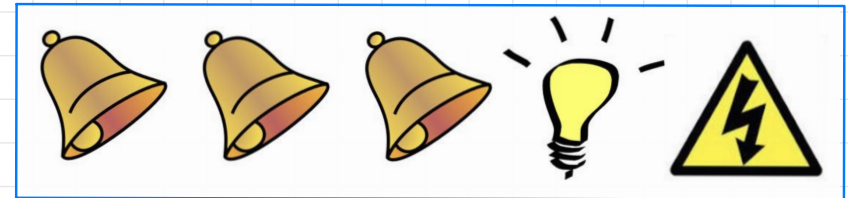
$$v(s) \leftarrow v(s) + \alpha (\hat{g}_t^\lambda - v(s))$$



## eligibility traces

$$t=0: \quad e(s) \leftarrow 0$$

$$t>0: \quad e(s) \leftarrow \lambda \gamma e(s) + \mathbb{1}[S_t = s]$$



$$\begin{aligned} e(s) &\leftarrow \lambda \gamma e(s) + \mathbb{1}[S=s] \quad \forall s \\ \delta &\leftarrow R + \gamma v(S') - v(S) \\ v(s) &\leftarrow v(s) + \alpha \delta e(s) \end{aligned}$$

