

# **Q-learning and Multi-step TD**

Konpat Preechakul  
Chulalongkorn University  
September 2019

# Recap TD and MC

- TD with moving average

$$v^{TD}(s) \leftarrow v(s) + \alpha [r + v(s') - v(s)]$$

- MC with moving average

$$v^{MC}(s_t) \leftarrow v(s_t) + \alpha \left[ \sum_{\tau=t} r_{\tau+1} - v(s_t) \right]$$

# Recap SARSA

$$q(s_t, a_t) \leftarrow q(s_t, a_t) + \alpha [r(s_t, a_t) + \overset{\text{On-policy algorithm}}{q(s_{t+1}, a_{t+1})} - q(s_t, a_t)]$$

$$q(s_t, a_t) \stackrel{\alpha}{\leftarrow} r(s_t, a_t) + q(s_{t+1}, a_{t+1})$$

- We need (s, a, r, s', a') to update

- Hence the name **SARSA**



# Q-learning

# Off-policy TD control

- SARSA = on-policy TD control
- SARSA = TD prediction + policy iteration
- If we combine **value iteration** with TD prediction
- We get Q-learning
- Q-learning = TD prediction + value iteration

# Q-learning formulation

$$q(s, a) \leftarrow$$

$$q(s, a) + \alpha [r(s, a) + \max_{a'} q(s', a') - q(s, a)]$$

$$q(s, a) \stackrel{\alpha}{\leftarrow} r(s, a) + \max_{a'} q(s', a')$$

- SARSA with value iteration
- **Needs only (s, a, r, s')**
- SARSA needs (s, a, r, s', a')

# Off-policy vs On-policy

## On-policy

- Policy can only use experience generated by itself to improve itself

## Off-policy

- Policy can use “any” experience to improve itself
- Usually more efficient (less interactions)

# Q-learning with epsilon greedy

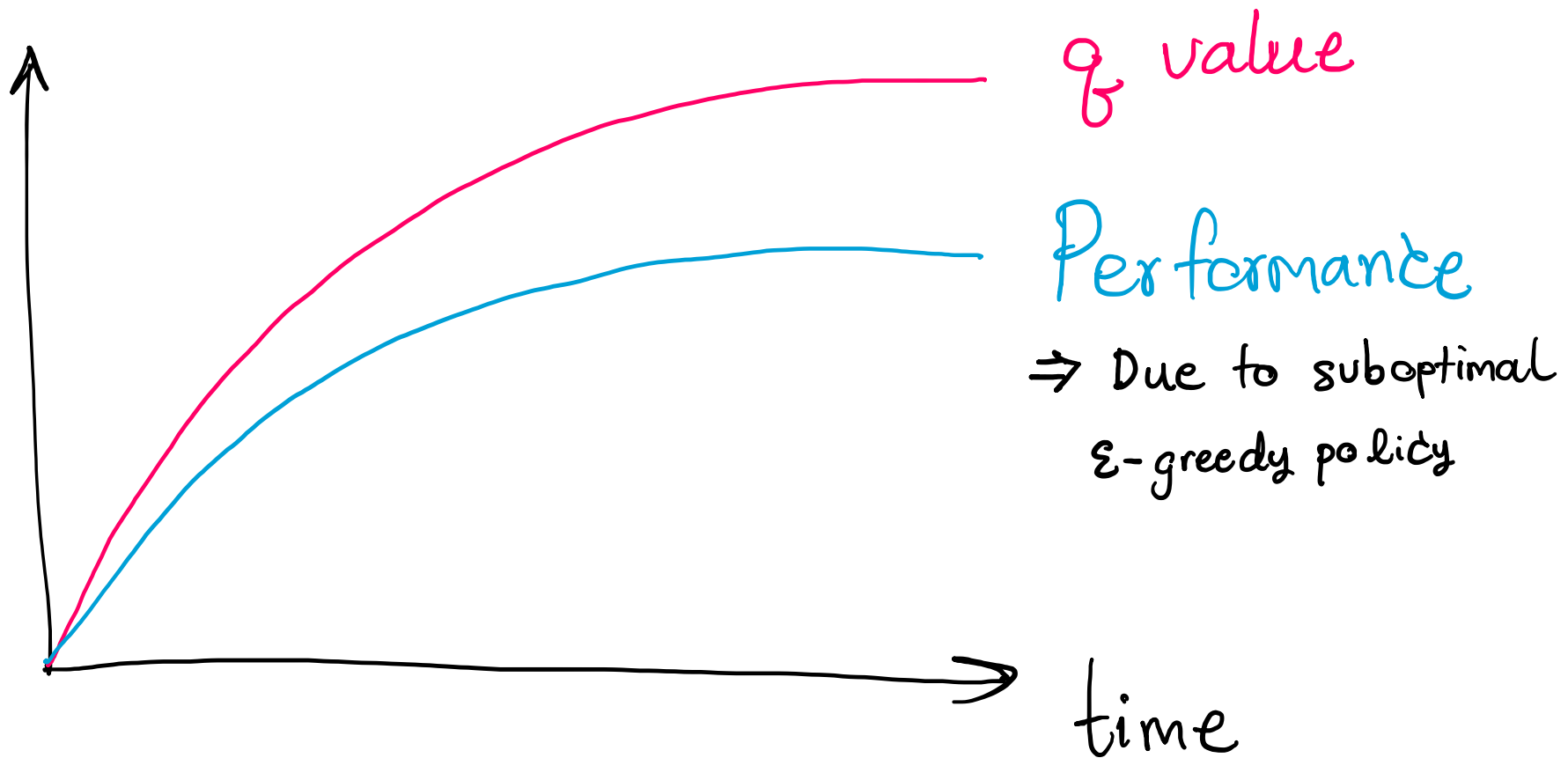
- Q-learning's prediction is not aware of epsilon greedy

$$q(s, a) \stackrel{\alpha}{\leftarrow} r(s, a) + \max_{a'} q(s', a')$$

- Its prediction will **NOT** reflect that
- It could lead to risky behavior
- This underlines annealing epsilon



# Example



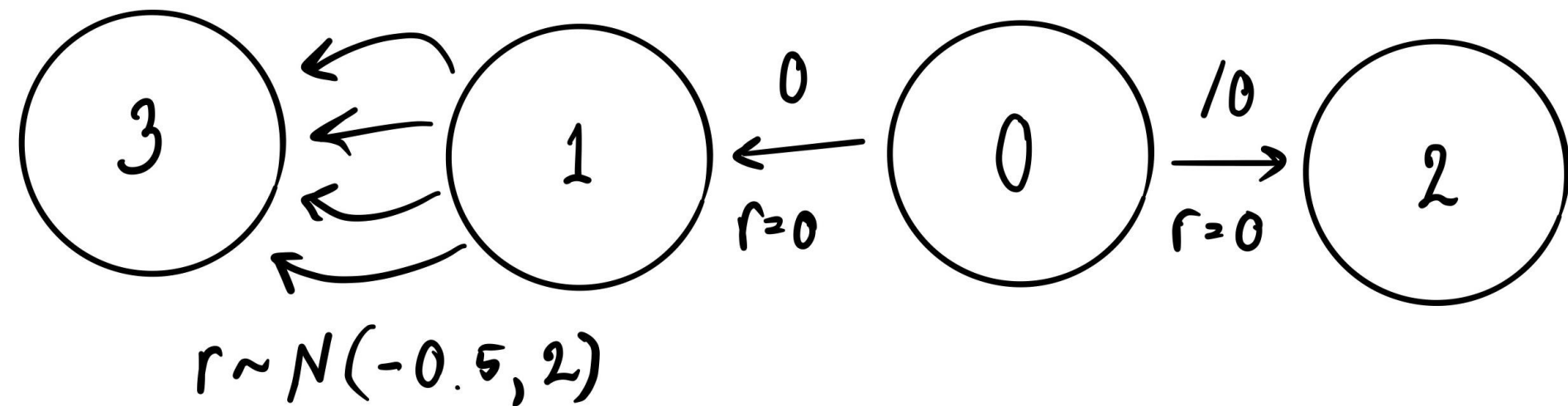
# Maximization bias



# Maximization bias of Q-learning

- Max-operation needs care when used with sampling
- Max of many random with zero mean is usually larger than mean
- This is a maximization bias
- This happens to SARSA and many algos.

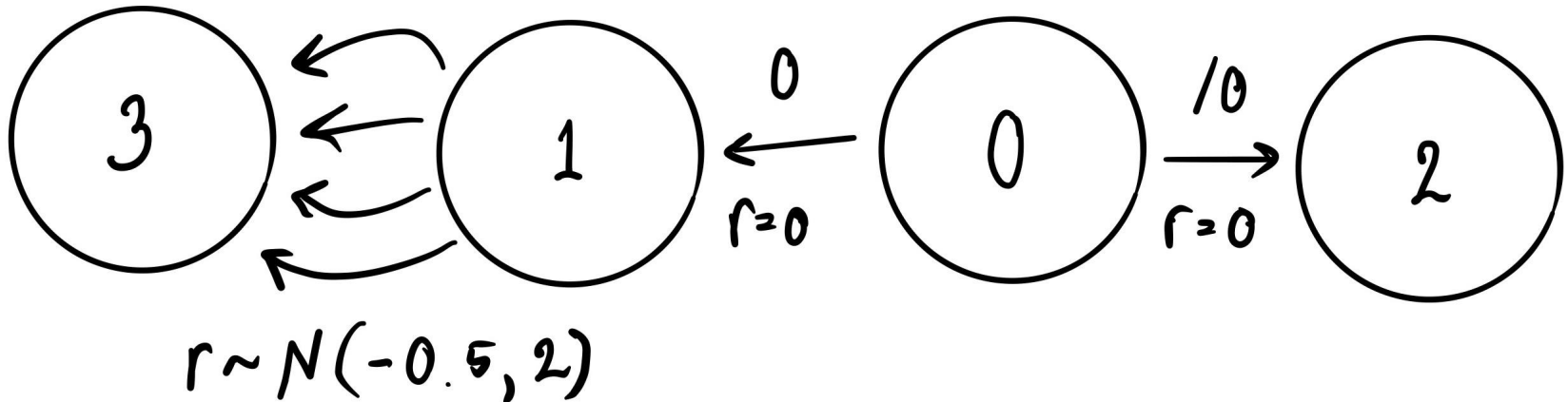
# Example



# Q-learning is overly optimistic

- We want  $\max_a \mathbb{E} [q(s, a)]$
- Q-learning uses  $q(s, a) \approx \mathbb{E} [q(s, a)]$
- This overestimates the value
- Is it possible mitigate this? How?

# Mitigating Maximization Bias



- $\max_a \mathbb{E} [q(s, a)]$  overestimates because we have stochastic  $q$  (due to luck)
- We can reduce this by maintaining 2  $q$  functions (don't share exp.)
- If both agree on the max, it is unlikely due to luck

# Double Q-learning

---

**Algorithm 1** Double Q-learning

---

```
1: Initialize  $Q^A, Q^B, s$ 
2: repeat
3:   Choose  $a$ , based on  $Q^A(s, \cdot)$  and  $Q^B(s, \cdot)$ , observe  $r, s'$ 
4:   Choose (e.g. random) either UPDATE(A) or UPDATE(B)
5:   if UPDATE(A) then
6:     Define  $a^* = \arg \max_a Q^A(s', a)$ 
7:      $Q^A(s, a) \leftarrow Q^A(s, a) + \alpha(s, a) (r + \gamma Q^B(s', a^*) - Q^A(s, a))$ 
8:   else if UPDATE(B) then
9:     Define  $b^* = \arg \max_a Q^B(s', a)$ 
10:     $Q^B(s, a) \leftarrow Q^B(s, a) + \alpha(s, a) (r + \gamma Q^A(s', b^*) - Q^B(s, a))$ 
11:   end if
12:    $s \leftarrow s'$ 
13: until end
```

---

\*Double Q-learning “underestimates”

Hasselt, Hado van. 2010. “Double Q-Learning.” In *NIPS 2010*, 2613–21.

**N-step TD**



# A quest for bias-variance balance

Bias	Variance	Algorithm
Low	High	Monte Carlo
High	Low	Temporal Difference
Med	Med	N-step TD?

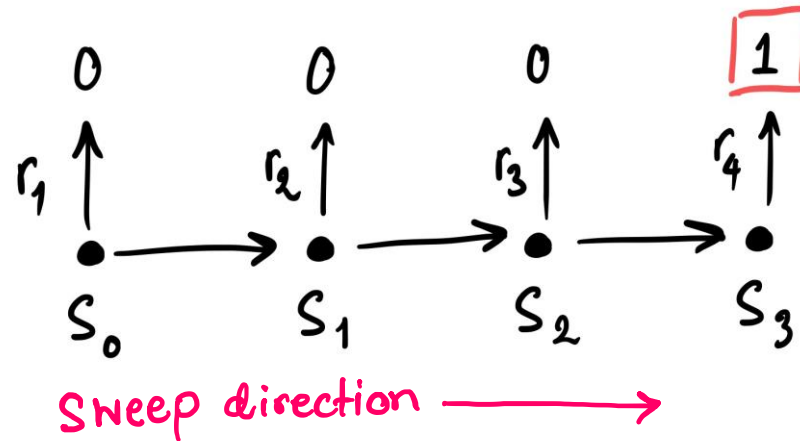
# Why we prefer many-steps?

- One step is slow to propagate
- One step has high bias

Why we don't prefer TOO many steps?

- Many-step leads to high variance

# Propagating faster example




State	Reward	Next V	$R + V'$	V
0	0			0
1	0			0
2	0			0
3	1	-		0

# From 1 step to Monte Carlo

One-step return

$$g_{t:t+1} = r_{t+1} + \gamma v(s_{t+1})$$

 Bootstrapping

Monte Carlo return

$$g_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \gamma^{T-t} r_T$$

# N-step return

$$g_{t:t+n} = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots + \gamma^n v(s_{t+n})$$



Bootstrapping

- Runs for n steps
- Gets n rewards
- Bootstraps the rest

# N-step TD

$$v(s_t) \leftarrow v(s_t) + \alpha [g_{t:t+n} - v(s_t)]$$

- Using n-step return as a target
- Need to wait for n steps for an update
- **Could we implement it online?**
- Online = do something with **1 reward**

# Online 3-step TD

$$v(s_t) \leftarrow v(s_t) + \alpha [g_{t:t+3} - v(s_t)]$$

$$g_{t:t+n} = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 v(s_{t+n})$$

$$\begin{aligned} \tilde{\delta}_{t+1} &= r_{t+1} + v(s_{t+1}) - v(s_t) \\ \tilde{\delta}_{t+2} &= r_{t+2} + v(s_{t+2}) - v(s_{t+1}) \\ \tilde{\delta}_{t+3} &= r_{t+3} + v(s_{t+3}) - v(s_{t+2}) \end{aligned} \quad \left. \vphantom{\begin{aligned} \tilde{\delta}_{t+1} \\ \tilde{\delta}_{t+2} \\ \tilde{\delta}_{t+3} \end{aligned}} \right\} \text{ignore discount}$$

$$\tilde{\delta}_{t+1} + \tilde{\delta}_{t+2} + \tilde{\delta}_{t+3} = r_{t+1} + r_{t+2} + r_{t+3} + v(s_{t+3})$$

# Online 3-step TD

$$\delta_{t+1} = r_{t+1} + \gamma v(s_{t+1}) - v(s_t)$$

$$\gamma \delta_{t+2} = \gamma r_{t+2} + \gamma^2 v(s_{t+2}) - \gamma v(s_{t+1})$$

$$\gamma^2 \delta_{t+3} = \gamma^2 r_{t+3} + \gamma^3 v(s_{t+3}) - \gamma^2 v(s_{t+2})$$

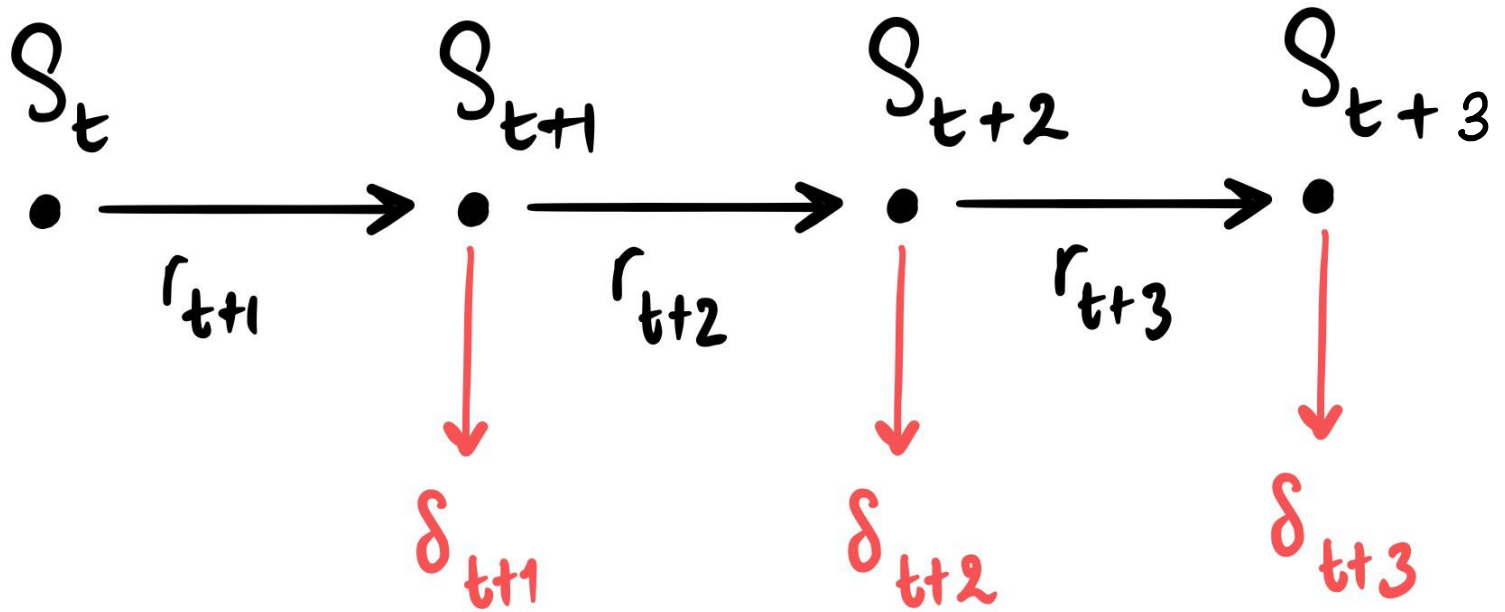
$$\delta_{t+1} + \gamma \delta_{t+2} + \gamma^2 \delta_{t+3}$$

$$= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 v(s_{t+n})$$

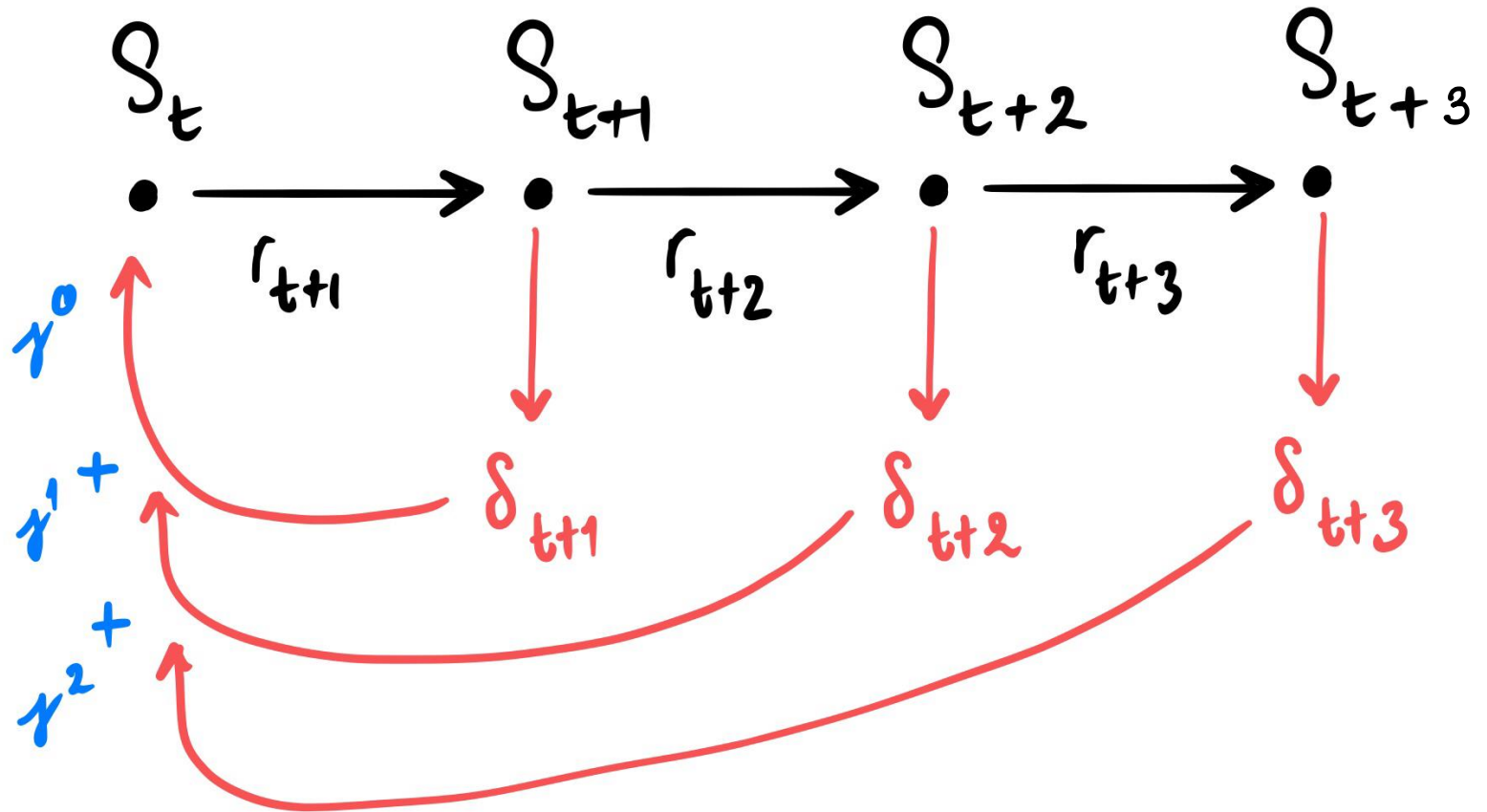
$$\delta_{t+1} + \gamma \delta_{t+2} + \gamma^2 \delta_{t+3} = g_{t:t+3}$$



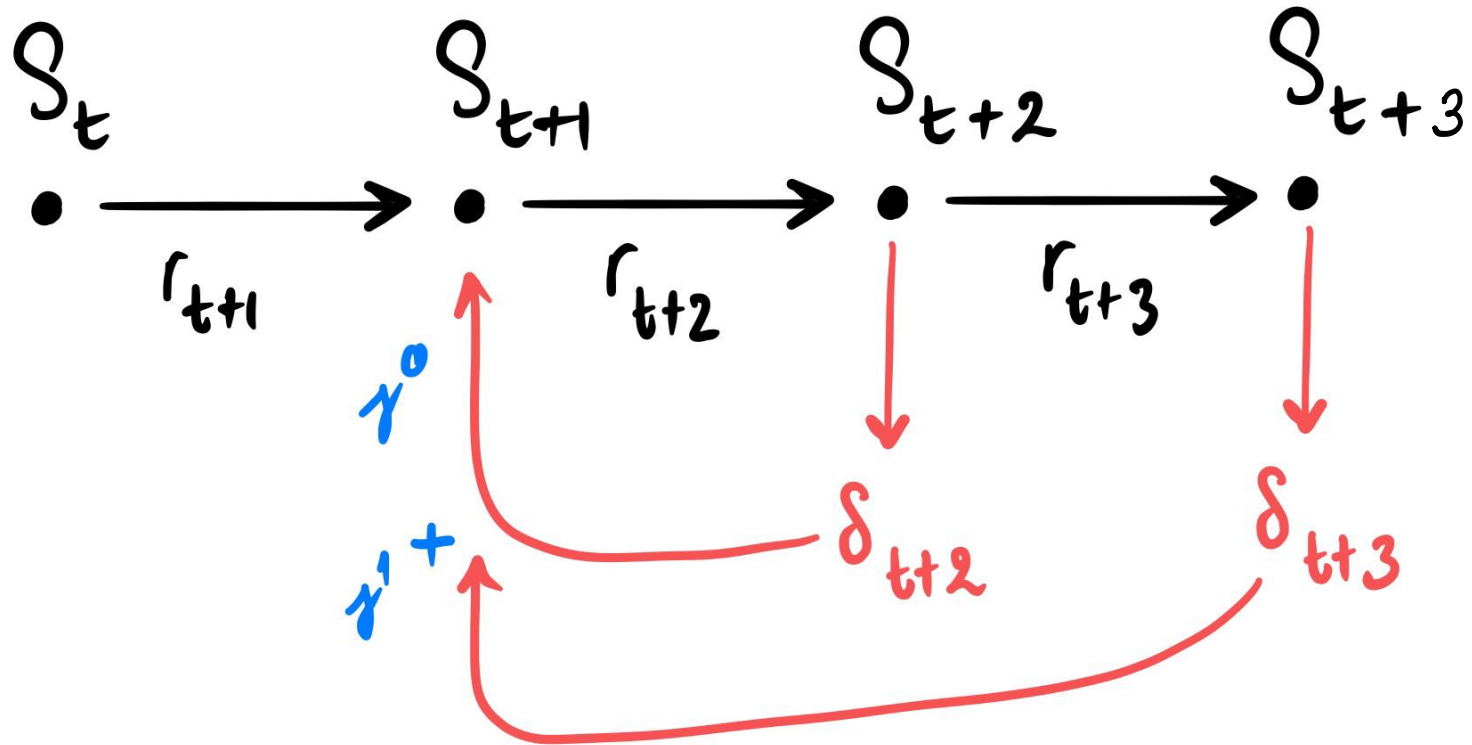
# Example



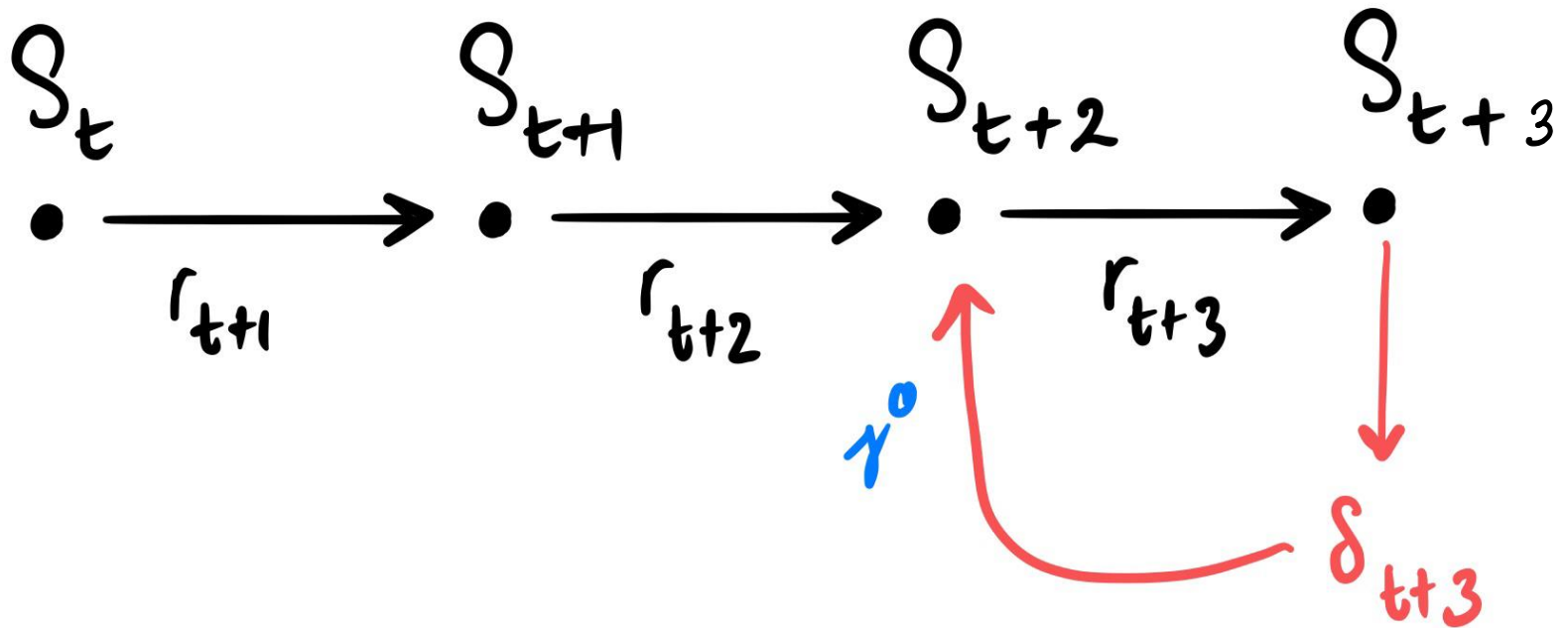
# Properly discounted



# Properly discounted

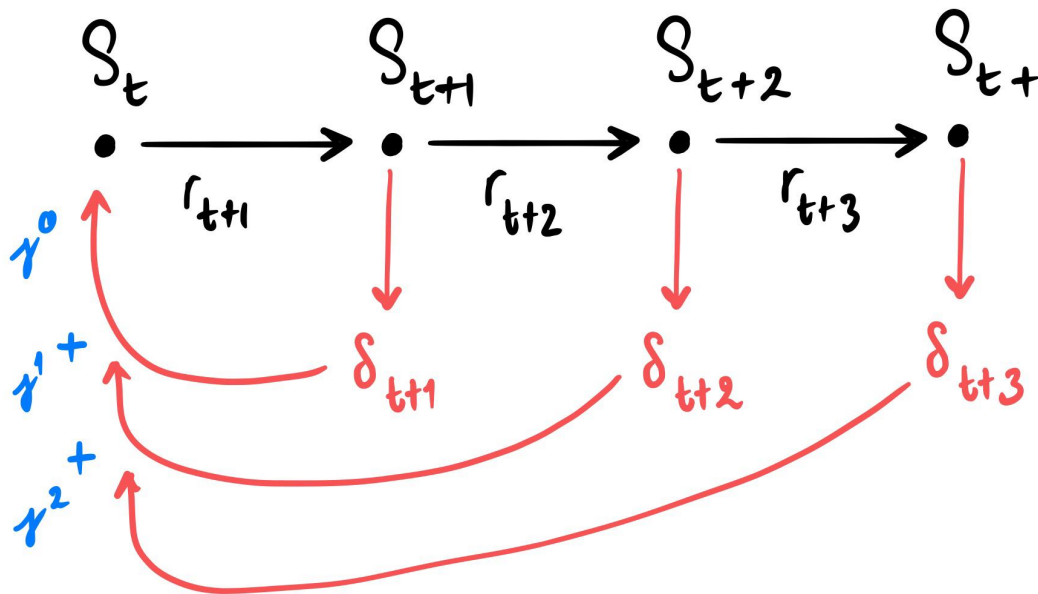


# Properly discounted

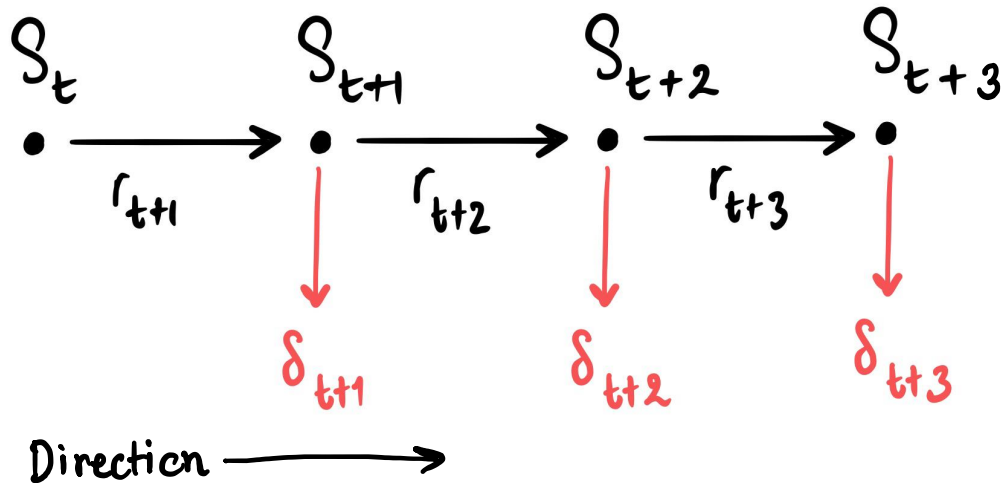


# How to keep track of the discount?

- We need some kind of memory
- Called “Eligibility trace”

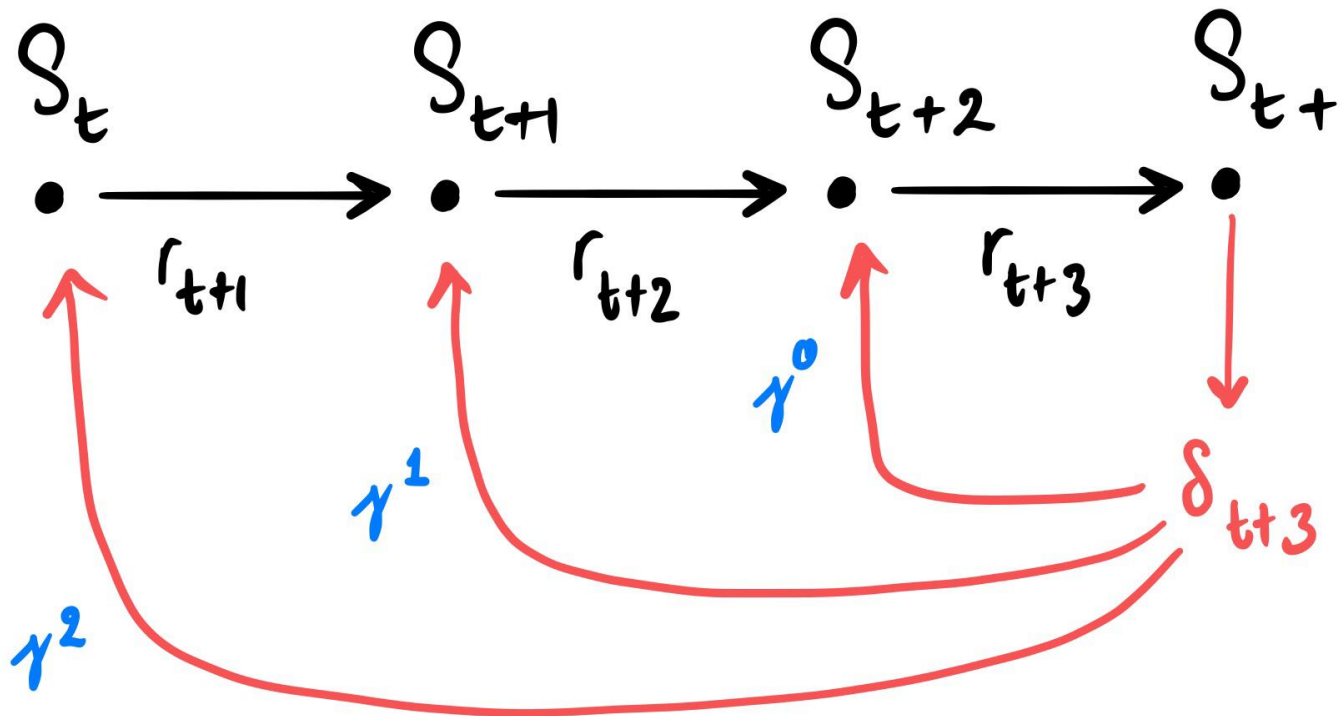


# Eligibility keeps track of discount



State	t	t+1	t+2	t+3	t+4
t	0				
t+1	0				
t+2	0				
t+3	0				

# Delta updates to previous states



# Online N-step TD Algorithm

```
for until  $v$  is stable do
  collect experience  $(s, a, r, s')$  using  $\pi$ 
   $\delta \leftarrow r + v(s') - v(s)$ 
   $e(s) \leftarrow e(s) + 1$ 
  for  $s$  in  $S$  do
     $v(s) \leftarrow v(s) + \alpha e(s) \delta$ 
     $e(s) \leftarrow \gamma e(s)$ 
    if  $e(s) \leq \gamma^n$  then
       $e(s) \leftarrow 0$ 
    end if
  end for
end for
```



# TD( $\lambda$ )

If we want something even more “medium”



# Lambda return

$$g_\lambda = (1 - \lambda) [g_{t:t+1} + \lambda g_{t:t+2} + \lambda^2 g_{t:t+3} + \dots] + \lambda^{T-t} g_t$$

$$\begin{aligned} g_\lambda = & (1 - \lambda) \lambda^0 [r_{t+1} + \gamma v(s_{t+1})] + \\ & (1 - \lambda) \lambda^1 [r_{t+1} + \gamma r_{t+2} + \gamma^2 v(s_{t+2})] + \\ & (1 - \lambda) \lambda^2 [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 v(s_{t+2})] + \\ & \dots + \\ & \underbrace{\lambda^{T-t}}_{\text{Weight}} [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots] \end{aligned}$$

Weight

$\lambda = 1$

# Lambda return

- Lambda is “share decay rate”
- Each n-step return has its own share
- $\text{Lambda} \Rightarrow 0$  favors small-step return
- $\text{Lambda} \Rightarrow 1$  favors large-step return

Bias	Variance	Lambda
High	Low	Smaller
Low	High	Larger

# TD(lambda)


$$v(s_t) \leftarrow v(s_t) + \alpha [g_\lambda - v(s_t)]$$

$$g_\lambda = (1 - \lambda) [g_{t:t+1} + \lambda g_{t:t+2} + \lambda^2 g_{t:t+3} + \dots] + \lambda^{T-t} g_t$$

- Naïve implementation needs to wait for a whole episode
- Do we have an online implementation?

# Tabular TD Lambda

$$v(s_t) \leftarrow v(s_t) + \alpha [g_\lambda - v(s_t)]$$


$$\begin{aligned} g_\lambda = & (1 - \lambda) \lambda^0 [r_{t+1} + \gamma v(s_{t+1})] + \\ & (1 - \lambda) \lambda^1 [r_{t+1} + \gamma r_{t+2} + \gamma^2 v(s_{t+2})] + \\ & (1 - \lambda) \lambda^2 [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 v(s_{t+2})] + \\ & \dots + \\ & \lambda^{T-t} [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots] \end{aligned}$$

# Tabular TD Lambda

$$v(s_t) \leftarrow v(s_t) + \alpha \left[ \underline{g_\lambda} - v(s_t) \right]$$

$$\begin{aligned} \underline{g_\lambda - v(s_t)} = & \\ & (1 - \lambda) \lambda^0 [r_{t+1} + \gamma v(s_{t+1}) - v(s_t)] + \\ & (1 - \lambda) \lambda^1 [r_{t+1} + \gamma r_{t+2} + \gamma^2 v(s_{t+2}) - v(s_t)] + \\ & (1 - \lambda) \lambda^2 [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 v(s_{t+2}) - v(s_t)] + \\ & \dots + \\ & \lambda^{T-t} [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots] \end{aligned}$$

# Tabular TD Lambda

$$\delta_{t+1} = r_{t+1} + \gamma v(s_{t+1}) - v(s_t)$$

$$\gamma \delta_{t+2} = \gamma r_{t+2} + \gamma^2 v(s_{t+2}) - \gamma v(s_{t+1})$$

$$\gamma^2 \delta_{t+3} = \gamma^2 r_{t+3} + \gamma^3 v(s_{t+3}) - \gamma^2 v(s_{t+2})$$

$$\delta_{t+1} + \gamma \delta_{t+2} + \gamma^2 \delta_{t+3}$$

$$= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 v(s_{t+n})$$

$$\underline{\delta_{t+1} + \gamma \delta_{t+2} + \gamma^2 \delta_{t+3} = g_{t:t+3}}$$

# Tabular TD Lambda

$$v(s_t) \leftarrow v(s_t) + \alpha [g_\lambda - v(s_t)]$$

$$g_\lambda - v(s_t) =$$

$$(1 - \lambda)\lambda^0 [\delta_{t+1}] +$$

$$(1 - \lambda)\lambda^1 [\delta_{t+1} + \gamma\delta_{t+2}] +$$

$$(1 - \lambda)\lambda^2 [\delta_{t+1} + \gamma\delta_{t+2} + \gamma^2\delta_{t+3}] +$$

$$\dots +$$

$$\lambda^{T-t} [\delta_{t+1} + \gamma\delta_{t+2} + \gamma^2\delta_{t+3} + \dots]$$

$$g_\lambda - v(s_t) = \underbrace{?}_{\text{red}} \delta_{t+1} + \underbrace{?}_{\text{red}} \gamma \delta_{t+2} + \underbrace{?}_{\text{red}} \gamma^2 \delta_{t+3} + \dots$$



# Tabular TD Lambda

$$\begin{aligned} g_\lambda - v(s_t) = & \\ (1 - \lambda)\lambda^0 [\delta_{t+1}] + & \\ (1 - \lambda)\lambda^1 [\delta_{t+1} + \gamma\delta_{t+2}] + & \\ (1 - \lambda)\lambda^2 [\delta_{t+1} + \gamma\delta_{t+2} + \gamma^2\delta_{t+3}] + & \\ \dots + & \\ \lambda^{T-t} [\delta_{t+1} + \gamma\delta_{t+2} + \gamma^2\delta_{t+3} + \dots] & \end{aligned}$$

$$1 - (1 - \lambda) = \lambda$$

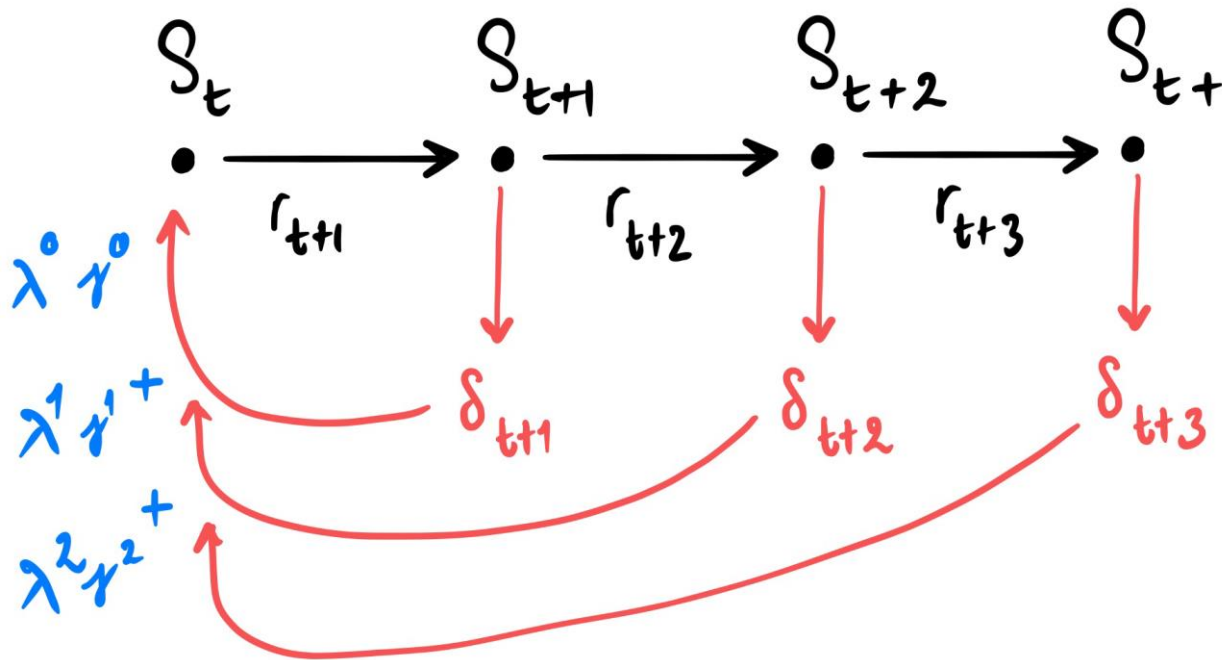
$$1 - (1 - \lambda) - (1 - \lambda)\lambda = \lambda^2$$

$$1 - (1 - \lambda) - (1 - \lambda)\lambda - (1 - \lambda)\lambda^2 = \lambda^3$$

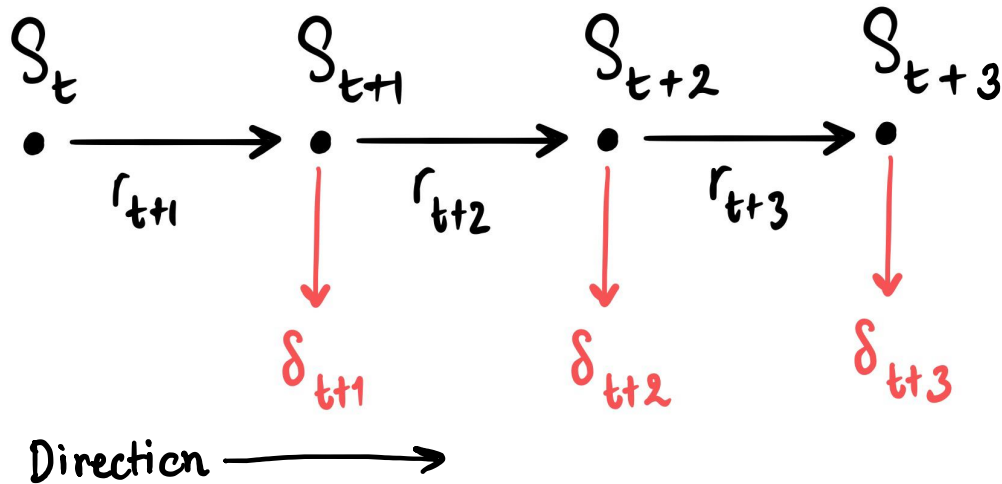
$$g_\lambda - v(s_t) = \delta_{t+1} + \underline{\lambda\gamma}\delta_{t+2} + \underline{\lambda^2\gamma^2}\delta_{t+3} + \dots$$

# Eligibility trace for lambda

- We need to keep track of both “gamma” and “lambda”



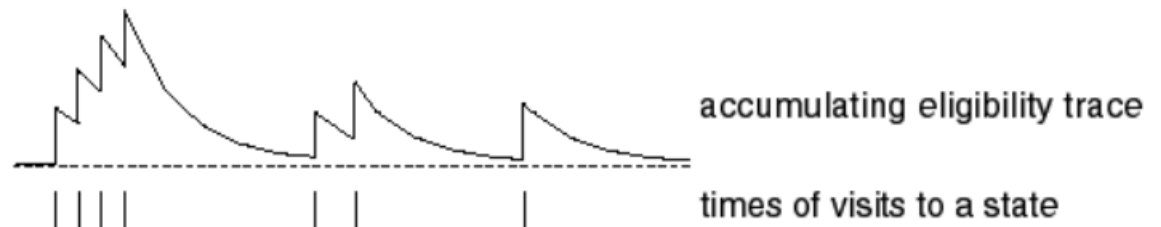
# Eligibility trace for lambda



State	t	t+1	t+2	t+3	t+4
t	0				
t+1	0				
t+2	0				
t+3	0				

# Online TD(lambda) algorithm

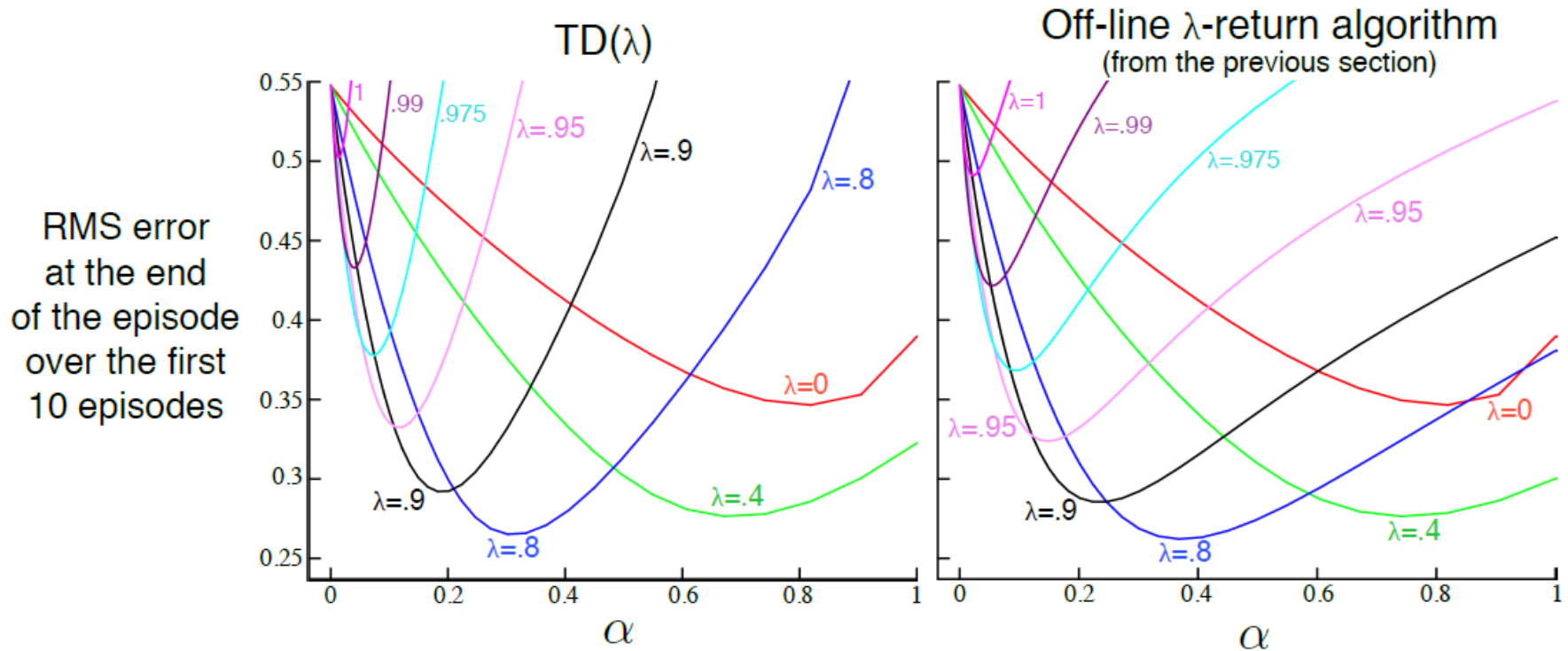
```
for until  $v$  is stable do
  collect experience  $(s, a, r, s')$  using  $\pi$ 
   $\delta \leftarrow r + v(s') - v(s)$ 
   $e(s) \leftarrow e(s) + 1$ 
  for  $s$  in  $S$  do
     $v(s) \leftarrow v(s) + \alpha e(s) \delta$ 
     $e(s) \leftarrow \gamma \lambda e(s)$ 
  end for
end for
```



# Online and Offline equivalency

- They are equivalent on paper
- Only applies to “offline” updates
- If you update online, they are no longer the same
- Offline is more “correct”
- More sophisticated kind of trace is needed for better online performance

# Online vs Offline performance



**Figure 12.6:** 19-state Random walk results (Example 7.1): Performance of TD( $\lambda$ ) alongside that of the offline  $\lambda$ -return algorithm. The two algorithms performed virtually identically at low (less than optimal)  $\alpha$  values, but TD( $\lambda$ ) was worse at high  $\alpha$  values.

# Eligibility trace cons

- We need to update to “all” states
- Inefficient
- In non-tabular case, there is a more efficient eligibility trace

# Assignment tour

- Ex 4.1
- Don't forget to pull