# 强化学习及其应用

## Reinforcement Learning and Its Applications

# 第三章 值估计

## Value Evaluation

授课人：周晓飞
zhouxiaofei@iie.ac.cn
2018-6-26

# 第三章 值估计

3.1 随机逼近

3.2 蒙特卡洛值估计

3.3 时序差分值估计

3.4 算法总结

# 第三章 值估计

## 3.1 随机逼近

# 随机逼近

## 累计平均逼近 E(x)

*For time $k$:*

$$S \leftarrow S + x_k \,;$$
$$N \leftarrow N + 1;$$
$$u_k \leftarrow S/N;$$

$$N \rightarrow \infty, \qquad u \rightarrow E(x)$$

# 随机逼近

## 增量均值逼近 E(x)

*For time $k$:*

$N \leftarrow N+1;$

$u_k \leftarrow u_{k-1} + (1/N)(x_k - u_{k-1});$

**相当于均值的更新：**

$$\mu_k = \frac{1}{k}\sum_{j=1}^{k} x_j$$

$$= \frac{1}{k}\left(x_k + \sum_{j=1}^{k-1} x_j\right)$$

$$= \frac{1}{k}(x_k + (k-1)\mu_{k-1})$$

$$= \mu_{k-1} + \frac{1}{k}(x_k - \mu_{k-1})$$

# 随机逼近

## Robbins-Monro 逼近 E(x)

*For time $k$:*

$$N \leftarrow N+1;$$

$$u_k \leftarrow u_{k-1} + a\,(x_k - u_{k-1});$$

**相当于权重比例更新：**

$$u_k \leftarrow (1 - a)u_{k-1} + ax_k$$

**本课程常用 Robbins-Monro 随机逼近公式**

# 第三章 值估计

# 蒙特卡洛值估计

## 问题描述

- Goal: learn $v_\pi$ from episodes of experience under policy $\pi$

$$S_1, A_1, R_2, ..., S_k \sim \pi$$

- Recall that the *return* is the <u>total</u> discounted reward:

$$G_t = R_{t+1} + \gamma R_{t+2} + ... + \gamma^{T-1} R_T$$

- Recall that the value function is the expected return:

$$v_\pi(s) = \mathbb{E}_\pi [G_t \mid S_t = s]$$

- Monte-Carlo policy evaluation uses *empirical mean* return instead of *expected* return

# 蒙特卡洛值估计

## First-Visit MC Evaluation

- To evaluate state $s$
- The first time-step $t$ that state $s$ is visited in an episode,
- Increment counter $N(s) \leftarrow N(s) + 1$
- Increment total return $S(s) \leftarrow S(s) + G_t$
- Value is estimated by mean return $V(s) = S(s)/N(s)$
- By law of large numbers, $V(s) \rightarrow v_\pi(s)$ as $N(s) \rightarrow \infty$

**只对 episode 的起始状态进行统计。**

## Every-Visit MC Evaluation

- To evaluate state $s$
- **Every** time-step $t$ that state $s$ is visited in an episode,
- Increment counter $N(s) \leftarrow N(s) + 1$
- Increment total return $S(s) \leftarrow S(s) + G_t$
- Value is estimated by mean return $V(s) = S(s)/N(s)$
- Again, $V(s) \rightarrow v_\pi(s)$ as $N(s) \rightarrow \infty$

**episode 的每个状态进行统计。**

## Incremental MC Evaluation

- Update $V(s)$ incrementally after episode $S_1, A_1, R_2, ..., S_T$
- For each state $S_t$ with return $G_t$

$$N(S_t) \leftarrow N(S_t) + 1$$

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)}(G_t - V(S_t))$$

## Incremental MC Evaluation

- Update $V(s)$ incrementally after episode $S_1, A_1, R_2, ..., S_T$
- For each state $S_t$ with return $G_t$

$$N(S_t) \leftarrow N(S_t) + 1$$

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)}(G_t - V(S_t))$$

- In non-stationary problems, it can be useful to track a running mean, i.e. forget old episodes.

$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$$

Chapter 3 Value Evaluation          *-12-*          中国科学院大学网络安全学院 2018 年研究生夏季课程

# 第三章 值估计

## 问题描述

**采用<span style="color:red">不完整的 episodes</span>，估计 V 值。**

- TD methods learn directly from episodes of experience
- TD is *model-free*: no knowledge of MDP transitions / rewards
- TD learns from *incomplete* episodes, by *bootstrapping*
- TD updates a guess towards a guess

## TD (0)

- Update value $V(S_t)$ toward *estimated* return $R_{t+1} + \gamma V(S_{t+1})$

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

**Bellman** 迭代的随机形式

- $R_{t+1} + \gamma V(S_{t+1})$ is called the *TD target*
- $\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$ is called the *TD error*

## TD (0)

### V.S. MC

- Incremental every-visit Monte-Carlo
  - Update value $V(S_t)$ toward *actual* return $G_t$

$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t - V(S_t) \right)$$

## TD (0)

- TD can learn *before* knowing the final outcome
  - TD can learn online after every step
  - MC must wait until end of episode before return is known
- TD can learn *without* the final outcome
  - TD can learn from incomplete sequences
  - MC can only learn from complete sequences
  - TD works in continuing (non-terminating) environments
  - MC only works for episodic (terminating) environments

## TD (0)

- Return $G_t = R_{t+1} + \gamma R_{t+2} + ... + \gamma^{T-1} R_T$ is *unbiased* estimate of $v_\pi(S_t)$
- True TD target $R_{t+1} + \gamma v_\pi(S_{t+1})$ is *unbiased* estimate of $v_\pi(S_t)$
- TD target $R_{t+1} + \gamma V(S_{t+1})$ is *biased* estimate of $v_\pi(S_t)$
- TD target is much lower variance than the return:
    - Return depends on *many* random actions, transitions, rewards
    - TD target depends on *one* random action, transition, reward

# 时序差分值估计

## TD (0)

- MC has high variance, zero bias
  - Good convergence properties
  - (even with function approximation)
  - Not very sensitive to initial value
  - Very simple to understand and use
- TD has low variance, some bias
  - Usually more efficient than MC
  - TD(0) converges to $v_\pi(s)$
  - (but not always with function approximation)
  - More sensitive to initial value

## TD (0)

- TD exploits Markov property
  - Usually more efficient in Markov environments
- MC does not exploit Markov property
  - Usually more effective in non-Markov environments

## TD (0)

### AB Example:

Two states $A, B$; no discounting; 8 episodes of experience

$A, 0, B, 0$

$B, 1$

$B, 1$

$B, 1$

$B, 1$

$B, 1$

$B, 1$

$B, 0$

What is $V(A), V(B)$?

## TD (0)

- MC converges to solution with minimum mean-squared error
  - Best fit to the observed returns

$$\sum_{k=1}^{K}\sum_{t=1}^{T_k}\left(G_t^k - V(s_t^k)\right)^2$$

  - In the AB example, $V(A) = 0$
- TD(0) converges to solution of max likelihood Markov model
  - Solution to the MDP $\langle \mathcal{S}, \mathcal{A}, \hat{\mathcal{P}}, \hat{\mathcal{R}}, \gamma \rangle$ that best fits the data

$$\hat{\mathcal{P}}_{s,s'}^a = \frac{1}{N(s,a)}\sum_{k=1}^{K}\sum_{t=1}^{T_k}\mathbf{1}(s_t^k, a_t^k, s_{t+1}^k = s, a, s')$$

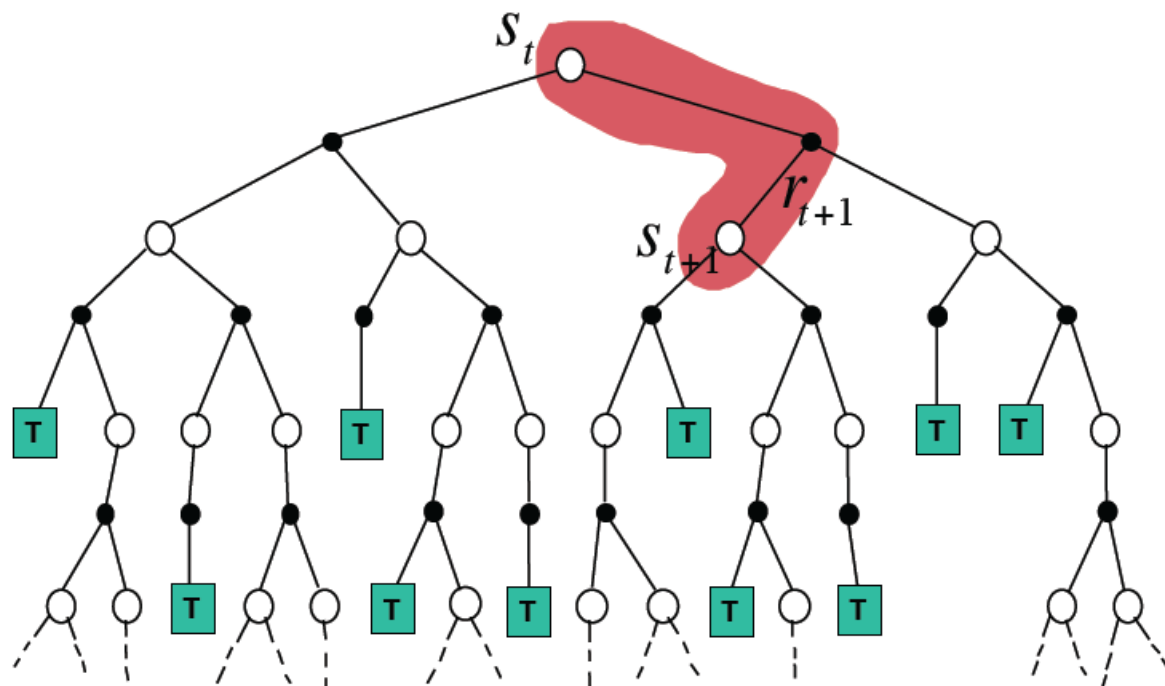$$\hat{\mathcal{R}}_s^a = \frac{1}{N(s,a)}\sum_{k=1}^{K}\sum_{t=1}^{T_k}\mathbf{1}(s_t^k, a_t^k = s, a)r_t^k$$

  - In the AB example, $V(A) = 0.75$

# 时序差分值估计

## TD (0)

■ **TD Backup**

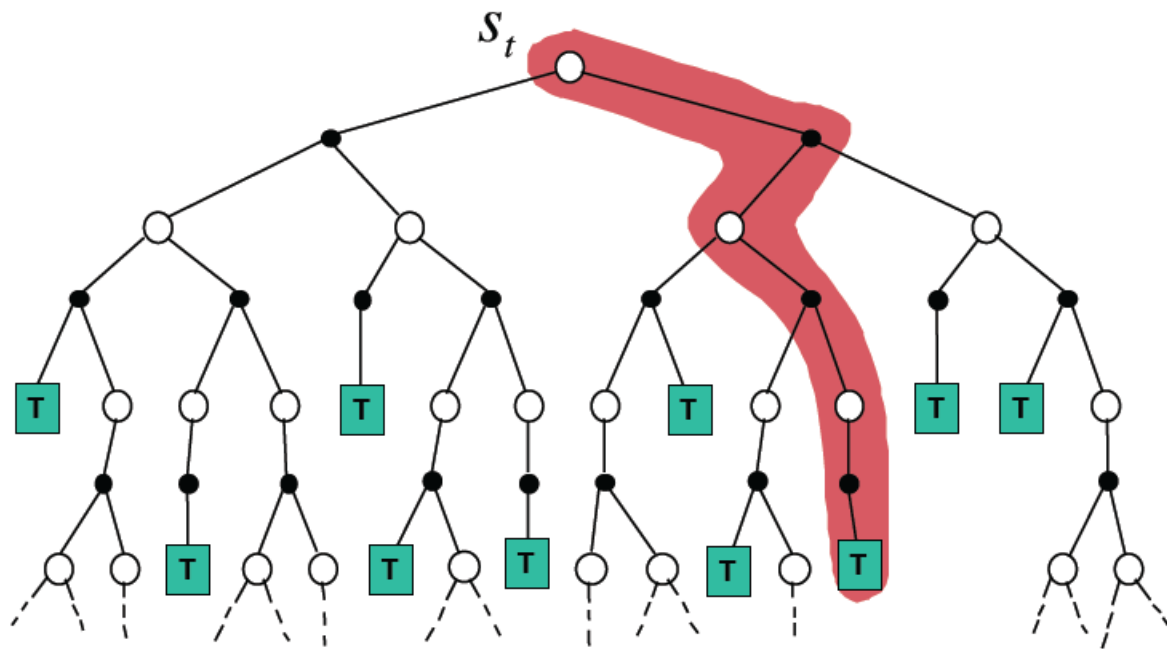$$V(S_t) \leftarrow V(S_t) + \alpha\left(R_{t+1} + \gamma V(S_{t+1}) - V(S_t)\right)$$



 中国科学院大学网络安全学院 2018 年研究生夏季课程

## TD (0)

**V.S. MC**

$$V(S_t) \leftarrow V(S_t) + \alpha\left(G_t - V(S_t)\right)$$

## TD (0)

**V.S. DP**

$$V(S_t) \leftarrow \mathbb{E}_\pi \left[ R_{t+1} + \gamma V(S_{t+1}) \right]$$

## TD (0)

## TD (0)

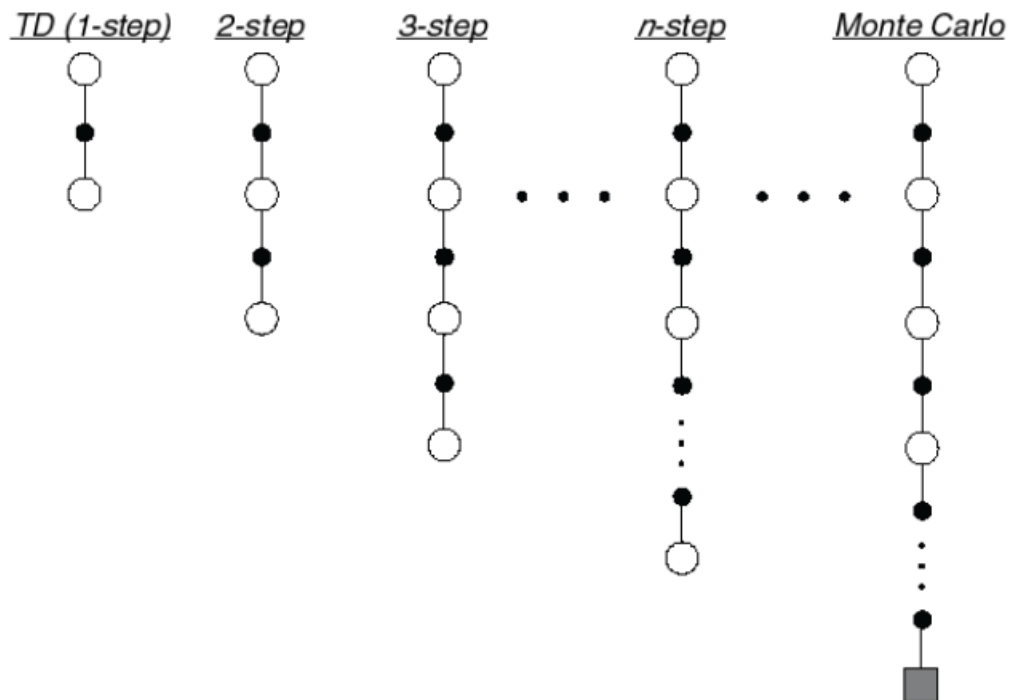**Bootstrapping and Sampling**

- **Bootstrapping**: update involves an estimate
    - MC does not bootstrap
    - DP bootstraps
    - TD bootstraps
- **Sampling**: update samples an expectation
    - MC samples
    - DP does not sample
    - TD samples

## TD(λ)

### n-step TD

■ Let TD target look *n* steps into the future

## TD(λ)

- Consider the following $n$-step returns for $n = 1, 2, \infty$:

$$
\begin{aligned}
n = 1 \quad (TD) \quad & G_t^{(1)} = R_{t+1} + \gamma V(S_{t+1}) \\
n = 2 \quad & G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2}) \\
\vdots \qquad & \qquad \vdots \\
n = \infty \quad (MC) \quad & G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + ... + \gamma^{T-1} R_T
\end{aligned}
$$

- Define the $n$-step return

$$
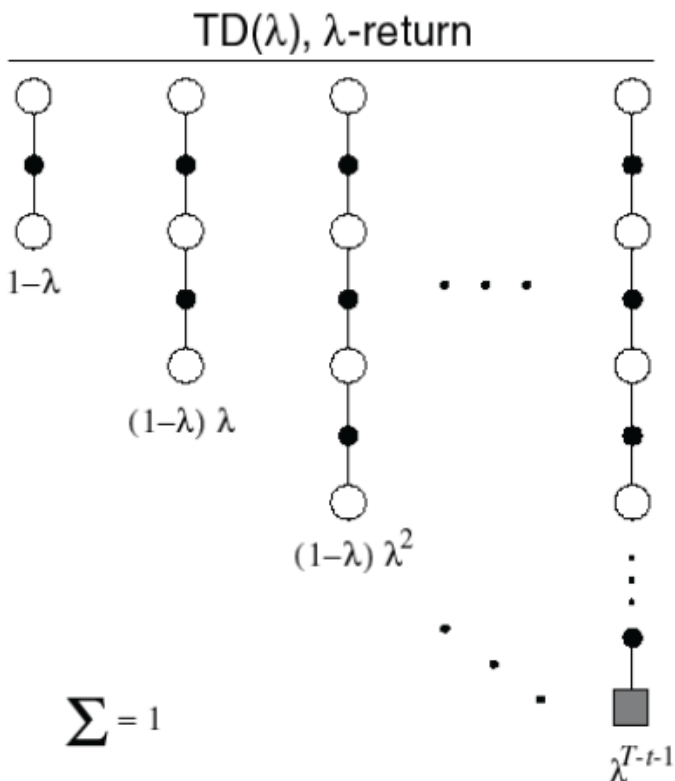G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + ... + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})
$$

- $n$-step temporal-difference learning

$$
V(S_t) \leftarrow V(S_t) + \alpha \left( G_t^{(n)} - V(S_t) \right)
$$

## TD(λ)

■ **Forward Returns and TD(λ)**



TD(λ), λ-return

$1-\lambda$
$(1-\lambda)\lambda$
$(1-\lambda)\lambda^2$
$\sum = 1$
$\lambda^{T-t-1}$

■ The $\lambda$-*return* $G_t^\lambda$ combines all $n$-step returns $G_t^{(n)}$

■ Using weight $(1-\lambda)\lambda^{n-1}$

$$G_t^\lambda = (1-\lambda)\sum_{n=1}^{\infty}\lambda^{n-1}G_t^{(n)}$$

■ Forward-view TD(λ)

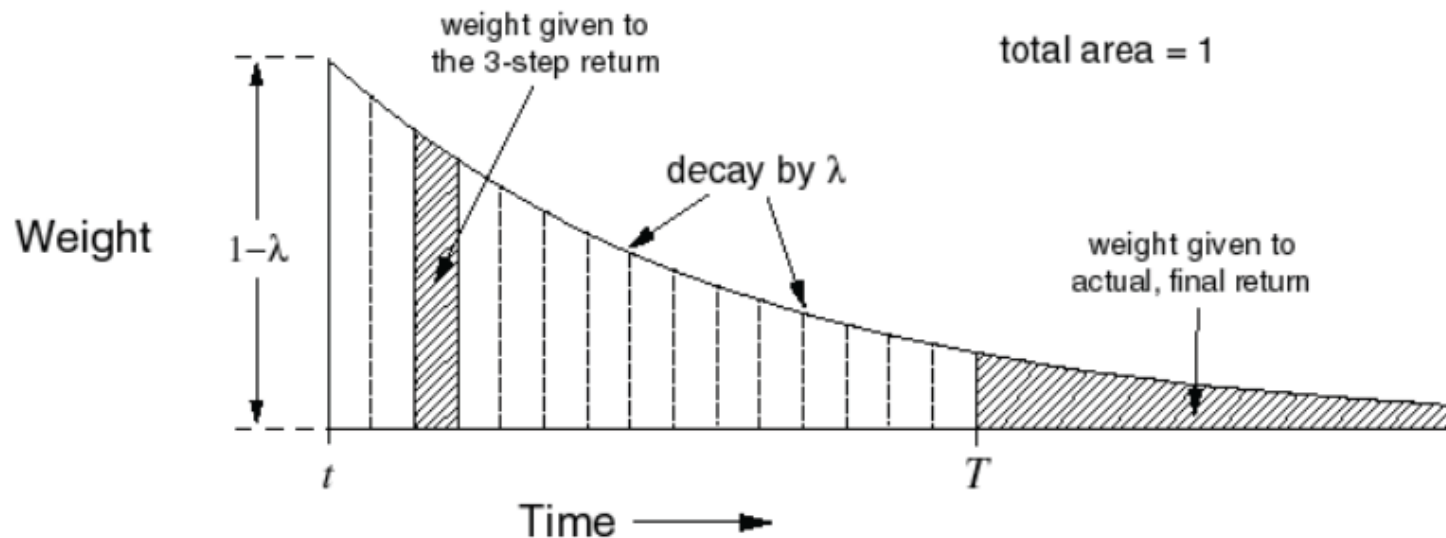$$V(S_t) \leftarrow V(S_t) + \alpha\left(G_t^\lambda - V(S_t)\right)$$

**TD(λ)**

$$\alpha \left( G_k^\lambda - V(S_k) \right) = \alpha \sum_{t=k}^{T} (\gamma\lambda)^{t-k} \delta_t$$

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

$$V(s) \leftarrow V(s) + \alpha \sum_{t=k}^{T} (\gamma\lambda)^{t-k} \delta_t$$
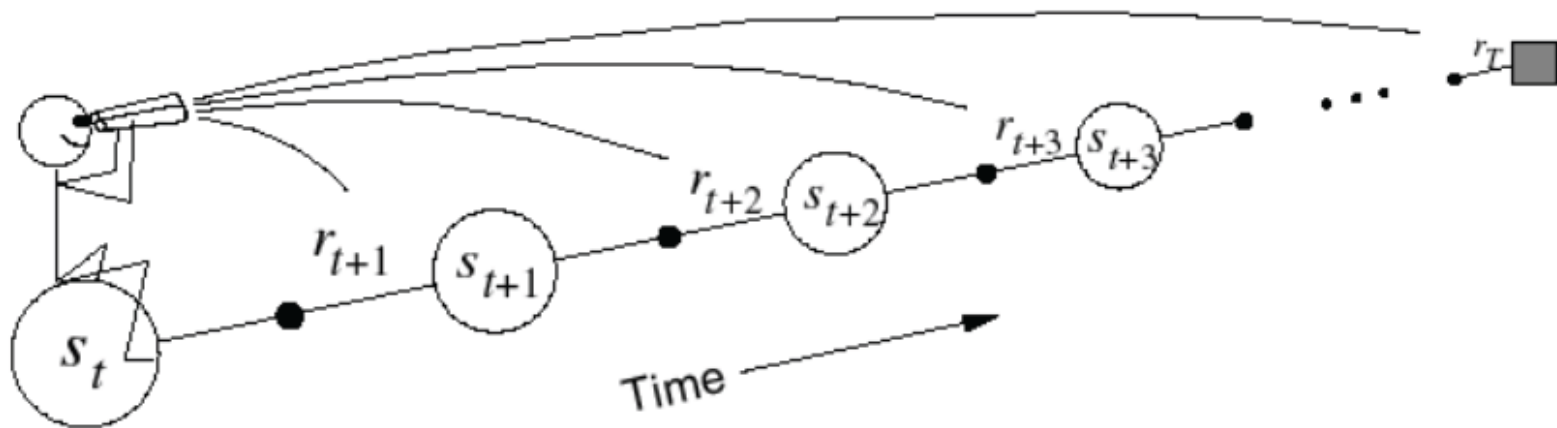
# 时序差分值估计

## TD(λ)



weight given to the 3-step return

decay by λ

total area = 1

weight given to actual, final return

Weight

$1-\lambda$

$t$

$T$

Time

$$G_t^{\lambda} = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

## TD(λ)



- Update value function towards the λ-return
- Forward-view looks into the future to compute $G_t^\lambda$
- Like MC, can only be computed from complete episodes

## TD(λ)

**Backward TD(λ)**

$$E_0(s) = 0$$
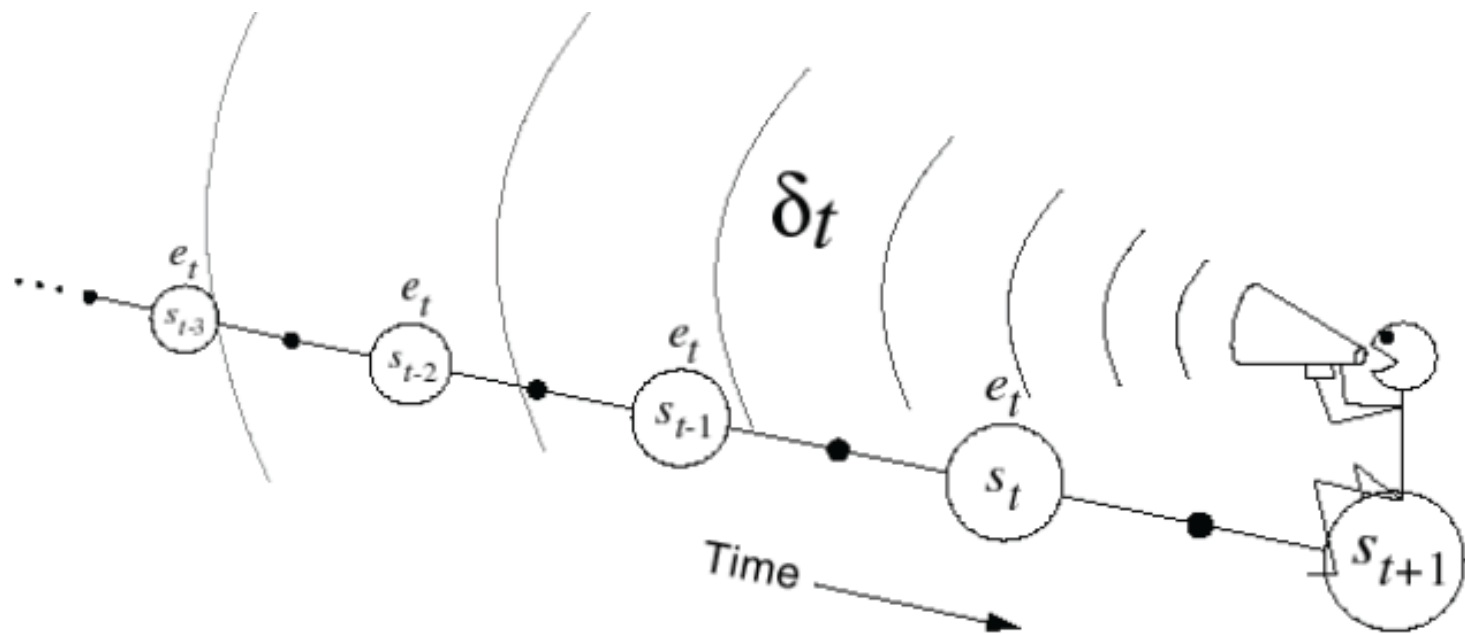$$E_t(s) = \gamma\lambda E_{t-1}(s) + \mathbf{1}(S_t = s)$$

- Keep an eligibility trace for every state $s$
- Update value $V(s)$ for every state $s$
- In proportion to TD-error $\delta_t$ and eligibility trace $E_t(s)$

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$
$$V(s) \leftarrow V(s) + \alpha\delta_t E_t(s)$$

## TD(λ)

## TD(λ)

- **Forward and Backward TD**

### Theorem

The sum of offline updates is identical for forward-view and backward-view TD($\lambda$)

$$\sum_{t=1}^{T} \alpha \delta_t E_t(s) = \sum_{t=1}^{T} \alpha \left( G_t^\lambda - V(S_t) \right) \mathbf{1}(S_t = s)$$

# 时序差分值估计

## TD(λ)

■ **TD(0), TD(λ), TD(1)**

**TD(0)**

$$E_t(s) = \mathbf{1}(S_t = s)$$
$$V(s) \leftarrow V(s) + \alpha \delta_t E_t(s)$$

This is exactly equivalent to TD(0) update

$$V(S_t) \leftarrow V(S_t) + \alpha \delta_t$$

## TD(λ)

### TD(λ)

$$E_t(s) = \gamma\lambda E_{t-1}(s) + \mathbf{1}(S_t = s)$$

$$V(s) \leftarrow V(s) + \alpha\delta_t E_t(s)$$

### TD(1)等价于MC

$$\alpha\left(G_k^\lambda - V(S_k)\right) = \alpha\sum_{t=k}^{T}(\gamma\lambda)^{t-k}\delta_t = \alpha\sum_{t=k}^{T}\gamma^{t-k}\delta_t = \alpha\left(G_k - V(S_k)\right)$$

大家请自行证明

# 第三章 值估计

# 算法总结

| Offline updates | $\lambda = 0$ | $\lambda \in (0, 1)$ | $\lambda = 1$ |
|---|---|---|---|
| Backward view | TD(0) | TD($\lambda$) | TD(1) |
| | $\parallel$ | $\parallel$ | $\parallel$ |
| Forward view | TD(0) | Forward TD($\lambda$) | MC |
| Online updates | $\lambda = 0$ | $\lambda \in (0, 1)$ | $\lambda = 1$ |
| Backward view | TD(0) | TD($\lambda$) | TD(1) |
| | $\parallel$ | $\not\parallel$ | $\not\parallel$ |
| Forward view | TD(0) | Forward TD($\lambda$) | MC |
| | $\parallel$ | $\parallel$ | $\parallel$ |
| Exact Online | TD(0) | Exact Online TD($\lambda$) | Exact Online TD(1) |

# 本讲参考文献

1. Richard S. Sutton and Andrew G. Barto. Reinforcement Learning: An Introduction. （Second edition, in progress，draft.

2. David Silver，Slides@《Reinforcement Learning:An Introduction》,2016.

3. Simon Haykin，申富饶等译，神经网络与学习机器，第三版。