# Model-Free Prediction (2)

### Junni Zou

Institute of Media, Information and Network
Dept. of Computer Science and Engineering
Shanghai Jiao Tong University
http://min.sjtu.edu.cn
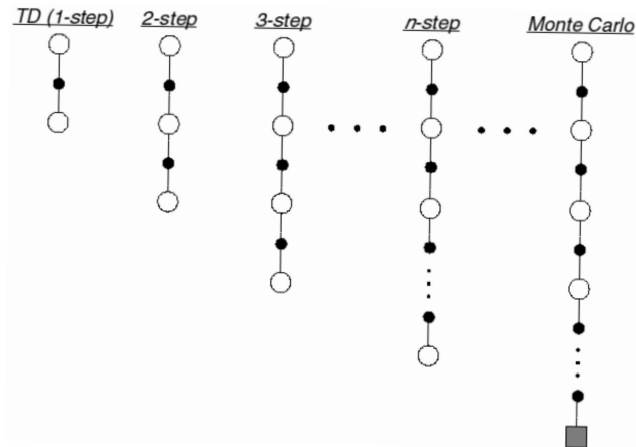
### Spring, 2021

Institute of Media,
Information, and Network

# Outline

Institute of Media,
Information, and Network

# Table of Contents

Institute of Media,
Information, and Network

# n-Step Prediction

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

# *n*-Step Return

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

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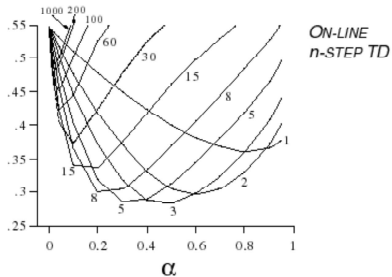\end{array}
$$

- Define the *n*-step return

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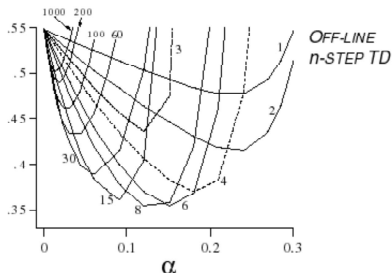$$

- *n*-step temporal-difference learning

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

Institute of Media,
Information, and Network

# Large Random Walk Example

# Averaging *n*-Step Returns

One backup

- We can average *n*-step returns over different *n*
- e.g. average the 2-step and 4-step returns

$$\frac{1}{2} G^{(2)} + \frac{1}{2} G^{(4)}$$

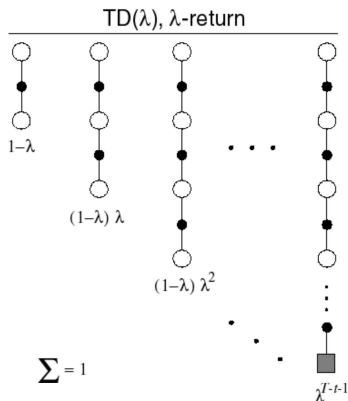- Combines Information from two different time-steps
- Can we efficiently combine information from all time-steps?

$\frac{1}{2}$

$\frac{1}{2}$

Institute of Media,
Information, and Network

# Table of Contents

Institute of Media,
Information, and Network

## $\lambda$-return



TD($\lambda$), $\lambda$-return

$1-\lambda$

$(1-\lambda)\,\lambda$

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

$\Sigma = 1$

$\lambda^{T-t-1}$

- 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 $TD(\lambda)$

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

Institute of Media,
Information, and Network

# $TD(\lambda)$ Weighting Function

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

# Forward-view $TD(\lambda)$



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

Institute of Media, Information, and Network

# Forward-View $TD(\lambda)$ on Large Random Walk

# Table of Contents

Institute of Media,
Information, and Network

# Backward View of $TD(\lambda)$

- Forward view provides theory
- Backward view provides mechanism
- Update online, every step, from incomplete sequences

Institute of Media,
Information, and Network

# Eligibility Traces



- Credit assignment problem: did bell or light cause shock?
- Frequency heuristic: assign credit to most frequent states
- Recency heuristic: assign credit to most recent states
- *Eligibility* traces combine both heuristics

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



accumulating eligibility trace

times of visits to a state

# Backward View of $TD(\lambda)$ (2)

- 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 $E_t(s)$

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

# Table of Contents

Institute of Media,
Information, and Network

# TD($\lambda$) and TD(0)

- When $\lambda = 0$, only current state is updated

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

- This is exactly equivalent to TD(0) update

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

Institute of Media,
Information, and Network

# TD($\lambda$) and MC

- When $\lambda = 1$, credit is deferred until end of episode
- Consider episodic environments with offline updates
- Over the course of an episode, total update for TD(1) is the same as total update for MC

**Theorem**

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

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

Institute of Media, Information, and Network

# Table of Contents

Institute of Media,
Information, and Network

# MC and TD(1)

- Consider an episode <u>where $s$ is visited once at time-step $k$</u>,
- TD(1) eligibility trace discounted time since visit,

$$E_t(s) = \gamma E_{t-1}(s) + \mathbf{1}(S_t = s)$$
$$= \begin{cases} 0 & \text{if } t < k \\ \gamma^{t-k} & \text{if } t \geq k \end{cases}$$

- TD(1) updates accumulate error *online*

$$\sum_{t=1}^{T-1} \alpha \delta_t E_t(s) = \alpha \sum_{t=k}^{T-1} \gamma^{t-k} \delta_t = \alpha \big(G_t^\lambda - V(S_t)\big)$$

- By end of episode it accumulates total error

$$\delta_k + \gamma\delta_{k+1} + \gamma^2\delta_{k+2} + \cdots + \gamma^{T-1-k}\delta_{T-1}$$

# Telescoping in TD (1)

When $\lambda = 1$, sum of TD errors telescopes into MC error,

$$
\begin{aligned}
G_t^\lambda - V(S_t) = -V(S_t) \quad &+ \quad (1-\lambda)\lambda^0 \left(R_{t+1} + \gamma V(S_{t+1})\right) \\
&+ \quad (1-\lambda)\lambda^1 \left(R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2})\right) \\
&+ \quad (1-\lambda)\lambda^2 \left(R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 V(S_{t+3})\right) \\
&+ \quad ... \\
= -V(S_t) \quad &+ \quad (\gamma\lambda)^0 \left(R_{t+1} + \gamma V(S_{t+1}) - \gamma\lambda V(S_{t+1})\right) \\
&+ \quad (\gamma\lambda)^1 \left(R_{t+2} + \gamma V(S_{t+2}) - \gamma\lambda V(S_{t+2})\right) \\
&+ \quad (\gamma\lambda)^2 \left(R_{t+3} + \gamma V(S_{t+3}) - \gamma\lambda V(S_{t+3})\right) \\
&+ \quad ... \\
= \quad &\quad (\gamma\lambda)^0 \left(R_{t+1} + \gamma V(S_{t+1}) - V(S_t)\right) \\
&+ \quad (\gamma\lambda)^1 \left(R_{t+2} + \gamma V(S_{t+2}) - V(S_{t+1})\right) \\
&+ \quad (\gamma\lambda)^2 \left(R_{t+3} + \gamma V(S_{t+3}) - V(S_{t+2})\right) \\
&+ \quad ... \\
= \delta_t + \gamma\lambda\delta_{t+1} &+ (\gamma\lambda)^2\delta_{t+2} + ...
\end{aligned}
$$

Institute of Media,
Information, and Network

# TD($\lambda$) and TD(1)

- TD(1) is roughly equivalent to every-visit Monte-Carlo
- Error is accumulated online, step-by-step
- If value function is only updated offline at end of episode
- Then total update is exactly the same as MC

# Telescoping in TD ($\lambda$)

For general $\lambda$, TD errors also telescope to $\lambda$-error, $G_t^\lambda - V(S_t)$

$$
\begin{aligned}
G_t^\lambda - V(S_t) = -V(S_t) \quad &+ \quad (1-\lambda)\lambda^0 \left( R_{t+1} + \gamma V(S_{t+1}) \right) \\
&+ \quad (1-\lambda)\lambda^1 \left( R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2}) \right) \\
&+ \quad (1-\lambda)\lambda^2 \left( R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 V(S_{t+3}) \right) \\
&+ \quad ... \\
= -V(S_t) \quad &+ \quad (\gamma\lambda)^0 \left( R_{t+1} + \gamma V(S_{t+1}) - \gamma\lambda V(S_{t+1}) \right) \\
&+ \quad (\gamma\lambda)^1 \left( R_{t+2} + \gamma V(S_{t+2}) - \gamma\lambda V(S_{t+2}) \right) \\
&+ \quad (\gamma\lambda)^2 \left( R_{t+3} + \gamma V(S_{t+3}) - \gamma\lambda V(S_{t+3}) \right) \\
&+ \quad ... \\
= \quad &\phantom{+} \quad (\gamma\lambda)^0 \left( R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right) \\
&+ \quad (\gamma\lambda)^1 \left( R_{t+2} + \gamma V(S_{t+2}) - V(S_{t+1}) \right) \\
&+ \quad (\gamma\lambda)^2 \left( R_{t+3} + \gamma V(S_{t+3}) - V(S_{t+2}) \right) \\
&+ \quad ... \\
= \delta_t + \gamma\lambda\delta_{t+1} &+ (\gamma\lambda)^2\delta_{t+2} + ...
\end{aligned}
$$

Institute of Media,
Information, and Network

# Forwards and Backwards TD($\lambda$)

- Consider an episode where $s$ is visited once at time-step $k$,
- TD($\lambda$) eligibility trace discounts time since visit,

$$E_t(s) = \gamma E_{t-1}(s) + \mathbf{1}(S_t = s)$$
$$= \begin{cases} 0 & \text{if } t < k \\ (\gamma\lambda)^{t-k} & \text{if } t \geq k \end{cases}$$

- Backward TD($\lambda$) updates accumulate error *online*

$$\sum_{t=1}^{T} \alpha\delta_t E_t(s) = \alpha \sum_{t=k}^{T} (\gamma\lambda)^{t-k}\delta_t = \alpha\big(G_t^{\lambda} - V(S_t)\big)$$

- By end of episode it accumulates total error for $\lambda$-return
- For multiple visit to $s$, $E_t(s)$ accumulates many errors

Institute of Media, Information, and Network

# Offline Equivalence of Forward and Backward TD

Offline updates

- Updates are accumulated within episode
- but applied in batch at the end of episode

Institute of Media,
Information, and Network

# Offline Equivalence of Forward and Backward TD (2)

Online updates

- TD($\lambda$) updates are applied online at each step within episode
- Forward and backward-view TD($\lambda$) are slightly different
- NEW: Exact online TD($\lambda$) achieves perfect equivalence
- By using a slightly different form of eligibility trace
- Sutton and von Seijen, ICML 2014

Institute of Media,
Information, and Network

# Summary of Forward and Backward TD($\lambda$)

| 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) |

$=$ here indicates equivalence in total update at end of episode.

Institute of Media,
Information, and Network