

→ MDP를 모르는데 일단 해보자

Lecture 4: Model-Free Prediction

MDP를 찾음.

↳ Value를 찾는 것.

David Silver

1강: 강화학습의 개념

2강: MDP에 대한 설명

3강: MDP를 알 때 (Model을 알 때)

무작위로 봄에 대해 설명

Prediction: Value를 찾는 것. (Policy가 정해져 있어야 함)

Control: 더 나은 Policy를 찾는 것.

Outline

1 Introduction

2 Monte-Carlo Learning

3 Temporal-Difference Learning

4 TD(λ)

environment을 모를 때, Model-free Prediction
을 풀기 위해 ① MC 를 사용한다.
② TD

Prediction은 Policy가 정해져 있고
그 Policy를 따를 때의 value fn을
구하는 것. → true value를 아는 것이다.

"optimal value fn"이라는 말이 있다.

↳ 그래서 Prediction에서는

Bellman expectation Eq 사용.

Control은

Bellman optimality Eq 사용.

Model-Free Reinforcement Learning

- Last lecture:
 - Planning by dynamic programming
 - Solve a *known* MDP
- This lecture:
 - Model-free prediction
 - Estimate the value function of an *unknown* MDP
- Next lecture:
 - Model-free control
 - Optimise the value function of an *unknown* MDP

Monte-Carlo Reinforcement Learning

Empirical, 경험에 의해 구하는 방법
실제로 해보면서 값을 추정하는 것이다.

MDP transition / reward 툴라든 드디어
모델-자유, model-free
episode 해보면서 학습하기 때문에

- MC methods learn directly from episodes of experience
- MC is *model-free*: no knowledge of MDP transitions / rewards
- MC learns from *complete episodes*: no bootstrapping
- MC uses the simplest possible idea: *value = mean return*
- Caveat: can only apply MC to *episodic MDPs*
 - All episodes must terminate *또는 episode는 끝날 수 있어야 한다.*

MC는 episode가 끝나야 Return (누적된 reward)을 받아서 학습할 수 있다.
이 return을 평균을 낸 것을 Value라고 하는 것이 MC.

Monte-Carlo Policy Evaluation = Prediction

- Goal: learn v_π from episodes of experience under policy π

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

- Recall that the *return* is the total discounted reward:

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

- Recall that the value function is the expected return:

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

한번에 몇개의 episode를 다루지 않음.

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

First-Visit Monte-Carlo Policy Evaluation

처음 방문한 state만 허락한다.

$$V(s) = \frac{\text{return} \rightarrow S(s)}{\text{probability } N(S)}$$

본 주제를 살피는 실사건의 학습이

동부전 예측에서 오차가 줄어간다.

Every-Visit Monte-Carlo Policy 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 $N(s) \rightarrow \infty$

first
visit
값

Blackjack Example

rule: 그장을 받고 카드를 더받거나 안받을 수 있다.
최대한 고이 가깝게 하는 것이 goal

- States (200 of them):

- Current sum (12-21) → 현재 카드의 합
- Dealer's showing card (ace-10) → 딜러가 보여준 카드
- Do I have a "useable" ace? (yes-no) → ace 있는가
- Action **stick**: Stop receiving cards (and terminate)
- Action **twist**: Take another card (no replacement)

- Reward for **stick**:

- +1 if sum of cards > sum of dealer cards
- 0 if sum of cards = sum of dealer cards
- -1 if sum of cards < sum of dealer cards

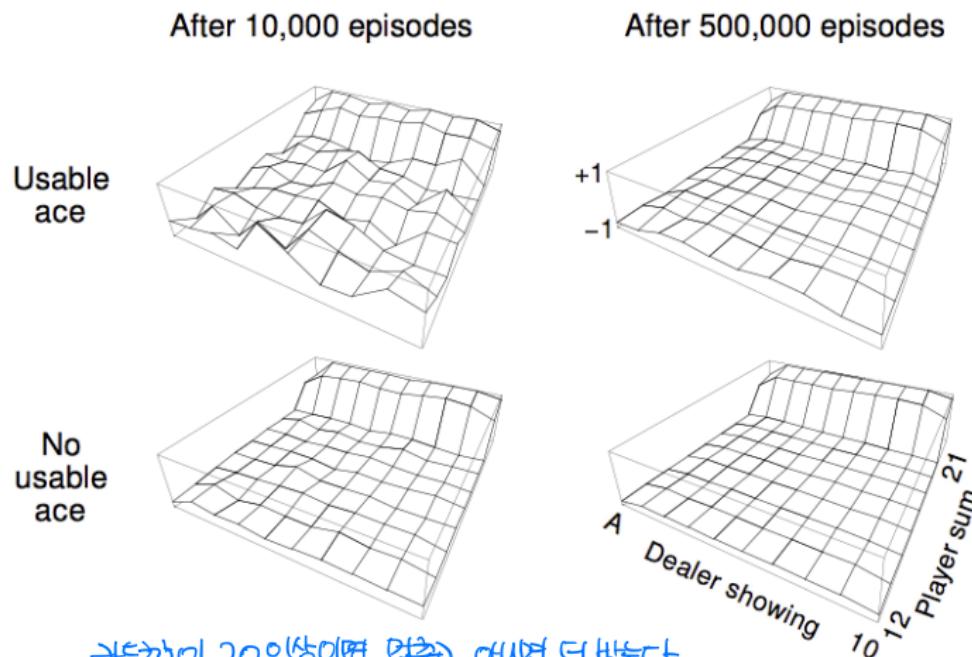
- Reward for **twist**:

- -1 if sum of cards > 21 (and terminate)
- 0 otherwise

- Transitions: automatically **twist** if sum of cards < 12



Blackjack Value Function after Monte-Carlo Learning



Policy: **stick** if sum of cards ≥ 20 , otherwise **twist**

Incremental Mean

The mean μ_1, μ_2, \dots of a sequence x_1, x_2, \dots can be computed incrementally,

M(에서 평균을 구하는 것은)
 DH episode마다 방향한 state가
 만족 값을 가지면서 평균을 찾아온다.
 그러나 Incremental Mean은 사용하면 다 저렴하고 있지 않아도 된다.

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

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

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

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

전체개별연립.

Incremental Monte-Carlo Updates

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

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

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

→ $N(S_t) \rightarrow \infty$
⇒ 점점 작아진다

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

MDP가
급격히 바뀌는 경우.

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

learning rate

고장나면 수 있다.

→ 여전히 한 경험을 암울.

Temporal-Difference Learning

Model free

TD는 예측도 할지 않아도 사용가능함

- 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

MC and TD

MC $\Rightarrow G_t$ 쪽으로 update

TD $\Rightarrow R_{t+1} + \gamma V(S_{t+1})$ 쪽으로 update

- Goal: learn v_π online from experience under policy π
- Incremental every-visit Monte-Carlo
 - Update value $V(S_t)$ toward *actual* return G_t

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

- Simplest temporal-difference learning algorithm: TD(0)

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

$$V(S_t) \leftarrow V(S_t) + \alpha (\underbrace{R_{t+1} + \gamma V(S_{t+1}) - V(S_t)}_{\text{TD error}}) \quad \underbrace{R_{t+1} + \gamma V(S_{t+1})}_{\text{TD target}}$$

한 step 가서 예측한 값으로 $V(S_t)$ 업데이트 한다.

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

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

Driving Home Example

State	Elapsed Time (minutes)	Current Time to Go	Predicted Total Time
leaving office	0	30	30
reach car, raining	5	35	40
exit highway	20	15	35
behind truck	30	10	40
home street	40	3	43
arrive home	43	0	43

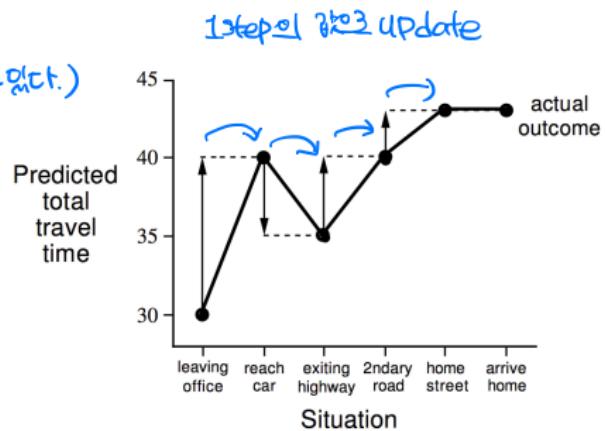
Driving Home Example: MC vs. TD

※ MC와 TD는 같은 알고리즘이다.

Changes recommended by Monte Carlo methods ($\alpha=1$)



Changes recommended by TD methods ($\alpha=1$)



Advantages and Disadvantages of MC vs. TD

MC는 에피소드가 다 끝나면 return을 알 수 있다. (terminate)

TD는 DH Step 후 학습하기 때문에 에피소드가 끝나지 않았을 때도 알 수 있다. (non-terminating)

- 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

Bias/Variance Trade-Off

→ G_t 는 $V_\pi(S_t)$ 을 대체하기 때문에 Bias가 없다. 편향되지 않음.

- Return $G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$ is unbiased estimate of $v_\pi(S_t)$
- True TD target $R_{t+1} + \gamma v_\pi(S_{t+1})$ is unbiased estimate of $v_\pi(S_t)$ $V_\pi(S_{t+1})$ 을 만족하면 Bellman Eq. 대체로 unbiased다.
- 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

→ 실제로 $V_\pi(S_{t+1})$ 을 아주 많이 대체해 $V(S_{t+1})$ 을 사용하는데

이제까지 대체해 몇 번을 해도 $V_\pi(S_t)$ 에 대해서 절대로 더 정확히 안다고 있다.
그래서 biased이다.

Advantages and Disadvantages of MC vs. TD (2)

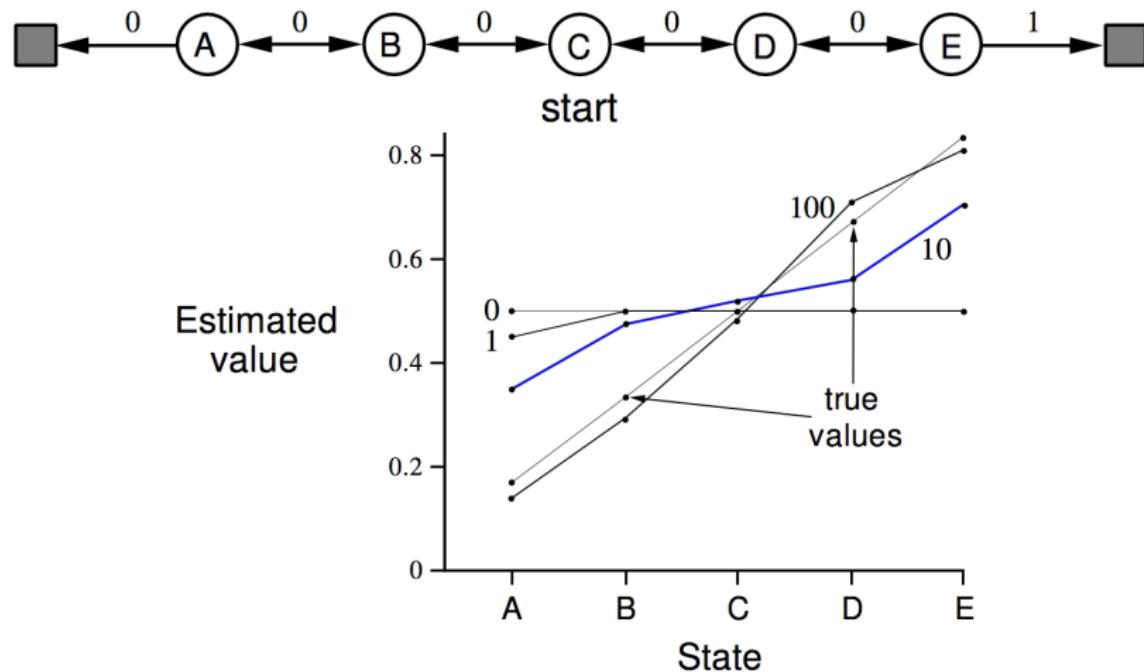
■ MC has high variance, zero bias

- Good convergence properties 속성이 좋다.
- (even with function approximation) $\rightarrow NN$
- Not very sensitive to initial value 초기값에 덜 민감.
- Very simple to understand and use

■ TD has low variance, some bias

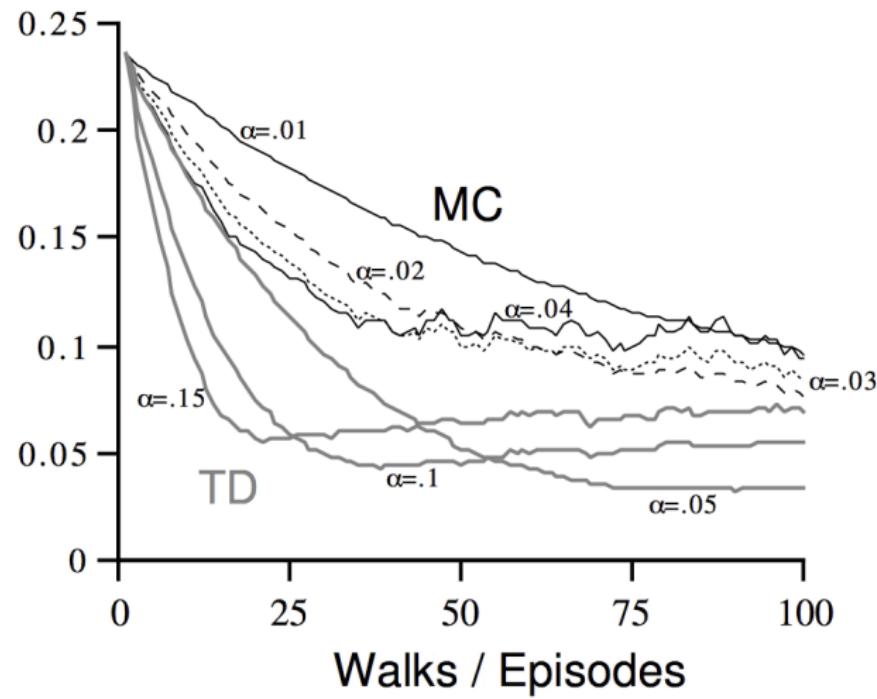
- Usually more efficient than MC
- TD(0) converges to $v_\pi(s)$
- (but not always with function approximation) 수렴 보장이 없다.
- More sensitive to initial value 초기값에 민감.

Random Walk Example



Random Walk: MC vs. TD

RMS error,
averaged
over states



Batch MC and TD → 제한된 업데이트가 있는 경우.

- MC and TD converge: $V(s) \rightarrow v_\pi(s)$ as experience $\rightarrow \infty$
- But what about batch solution for finite experience?

$$s_1^1, a_1^1, r_2^1, \dots, s_{T_1}^1$$

⋮

$$s_1^K, a_1^K, r_2^K, \dots, s_{T_K}^K$$

- e.g. Repeatedly sample episode $k \in [1, K]$
- Apply MC or TD(0) to episode k

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

AB Example

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

$A, 0, B, 0$

$B, 1$

$B, 1$

$B, 1$ $V(B) = 0.75$

$B, 1$

$B, 1$

$B, 1$

$B, 1$

$B, 0$

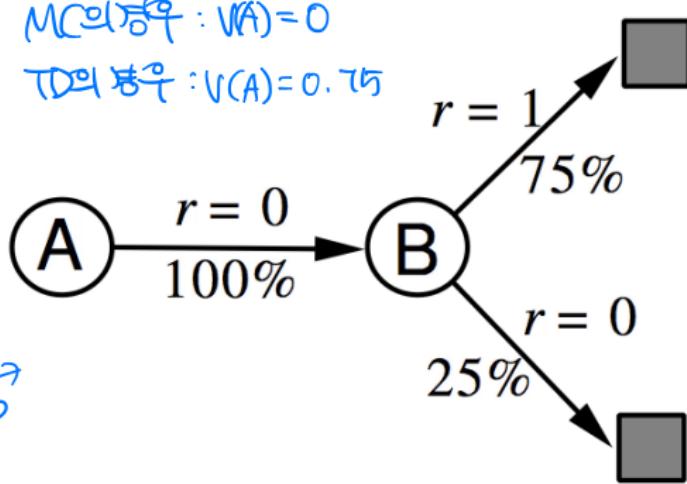
8회중 6번 B 1

$V(B) = 0.75$

단순 MDP

MC의 경우 : $V(A) = 0$

TD의 경우 : $V(A) = 0.75$



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

Certainty Equivalence

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

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

- 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, $V(A) = 0.75$

Advantages and Disadvantages of MC vs. TD (3)

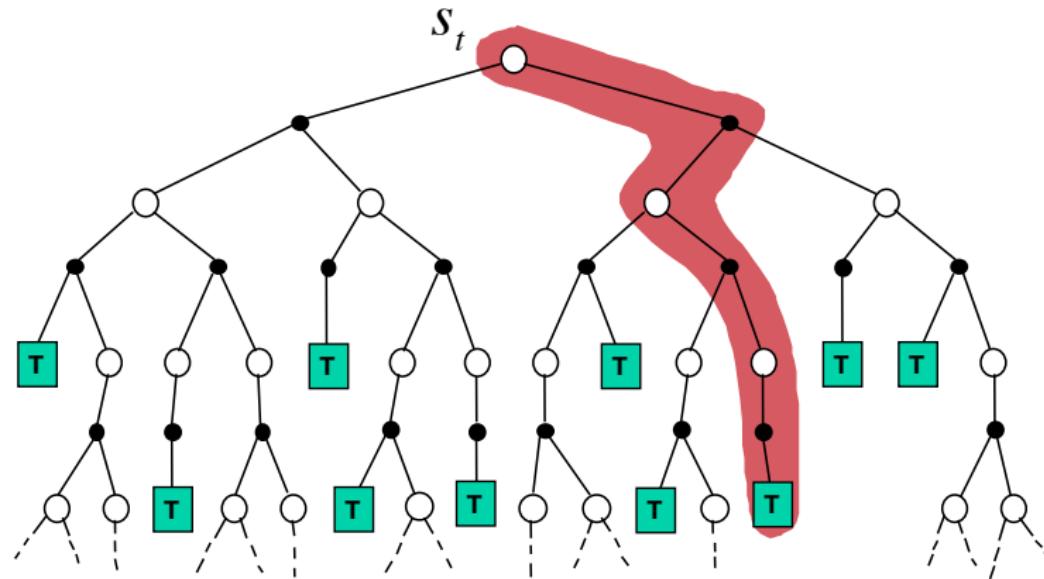
[TD는 MP 사용해서 푸는
MC는 MP 사용 안해도 푸는]

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

TD는 Markov Model의 likelihood를 maximizes.

Monte-Carlo Backup

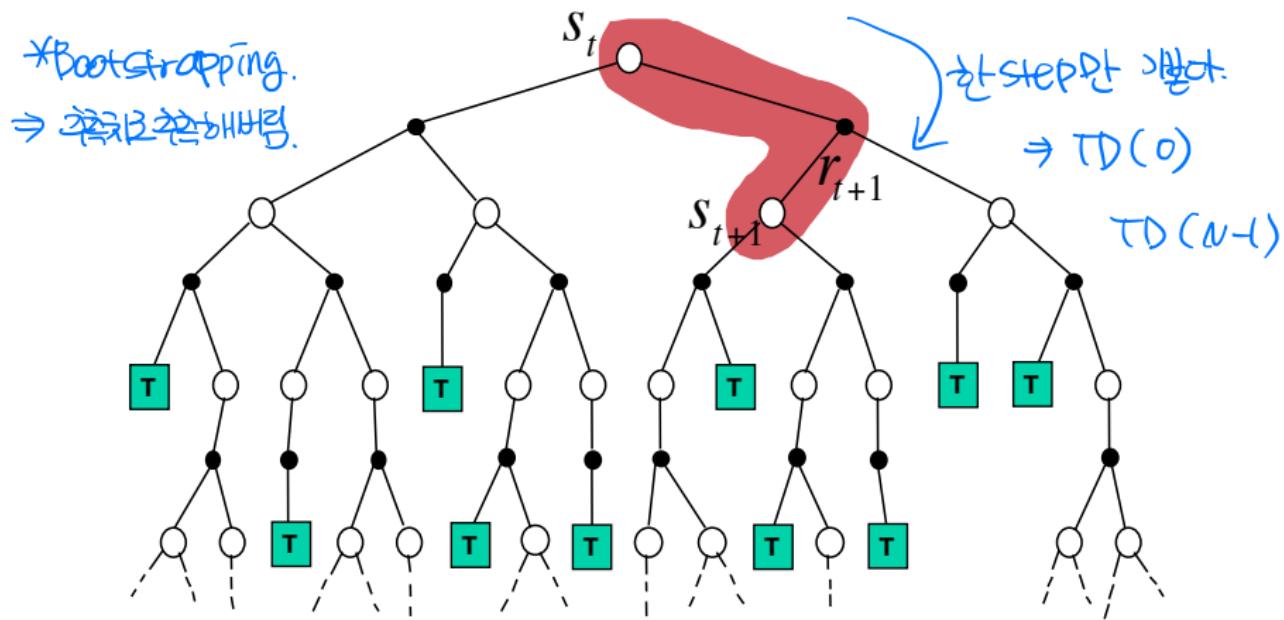
$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



MDP를 알면 hybrid 가능.

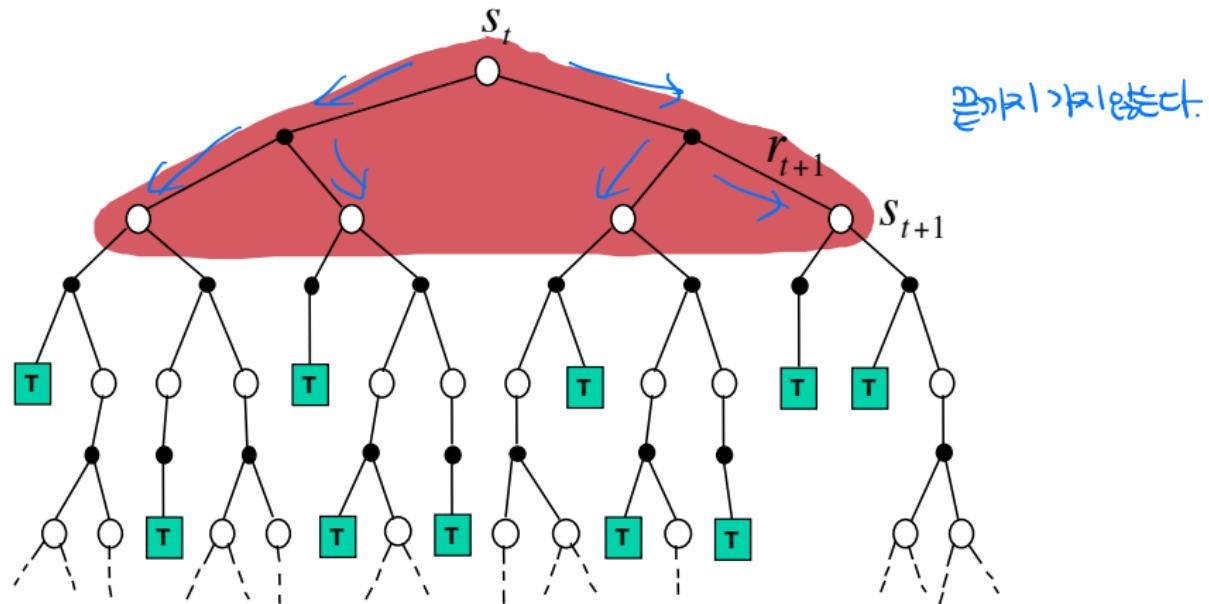
Temporal-Difference Backup

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



Dynamic Programming Backup

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



Bootstrapping and Sampling

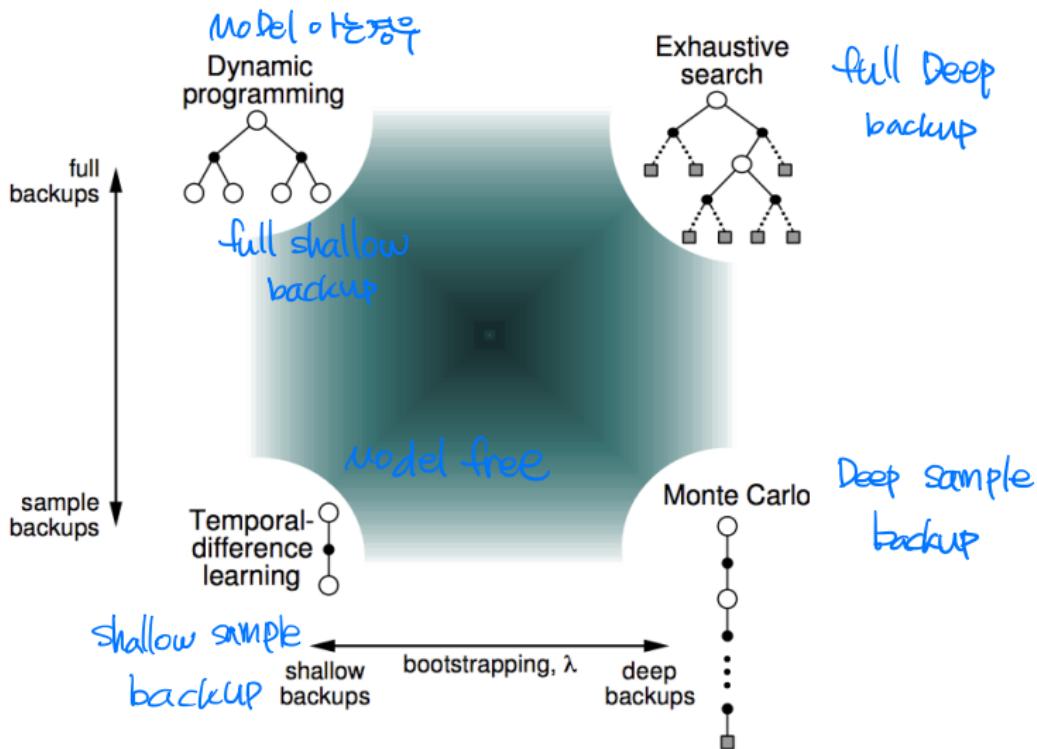
- Bootstrapping: update involves an estimate

- MC does not bootstrap ↳ 끝까지 기다리기
- DP bootstraps
- TD bootstraps

- Sampling: update samples an expectation

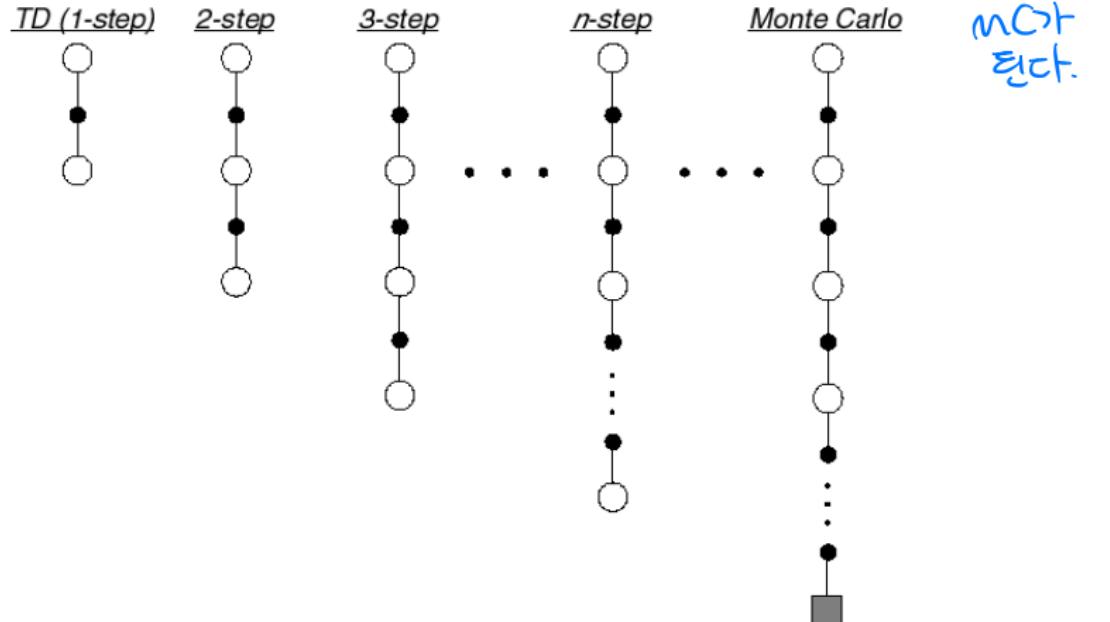
- MC samples
- DP does not sample ↳ 가능한 action 전부 채택(?) non-sample
- TD samples

Unified View of Reinforcement Learning



n-Step Prediction

- Let TD target look n steps into the future



n -Step Return

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

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

- Define the n -step return

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

- n -step temporal-difference learning

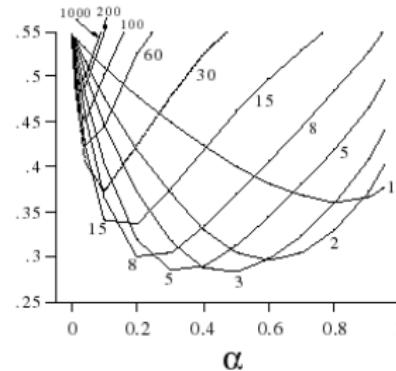
$$V(S_t) \leftarrow V(S_t) + \alpha \underbrace{\left(G_t^{(n)} - V(S_t) \right)}_{\text{n-step TD error}}$$

Large Random Walk Example

$$\text{TD}(1) \sim \text{TD}(1000)$$

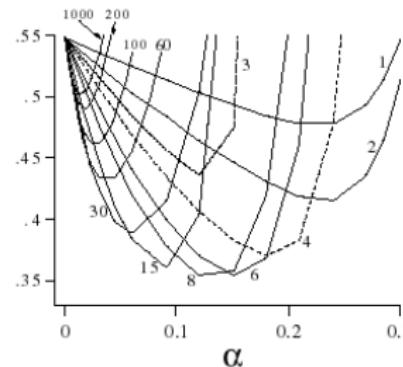
액션 MC가
안 좋은 성과.

RMS error,
averaged over
first 10 episodes



ON-LINE → 실시간 업데이트
 n -STEP TD

RMS error,
averaged over
first 10 episodes



OFF-LINE → 디자인 업데이트
 n -STEP TD

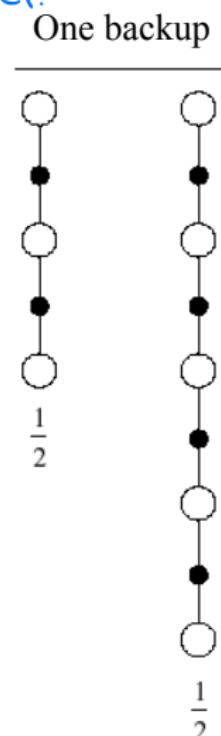
Averaging n -Step Returns

\hookrightarrow TD(0) ≈ MC π[τ] → TD(λ)

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

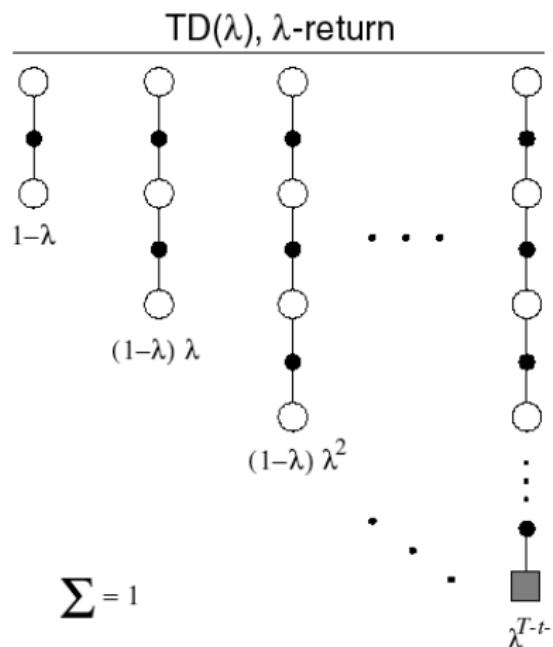
$$\frac{1}{2}G^{(2)} + \frac{1}{2}G^{(4)}$$

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



λ -return

$$\lambda = 0 \sim 1$$



MC로 갈수록 λ 가 작아짐.

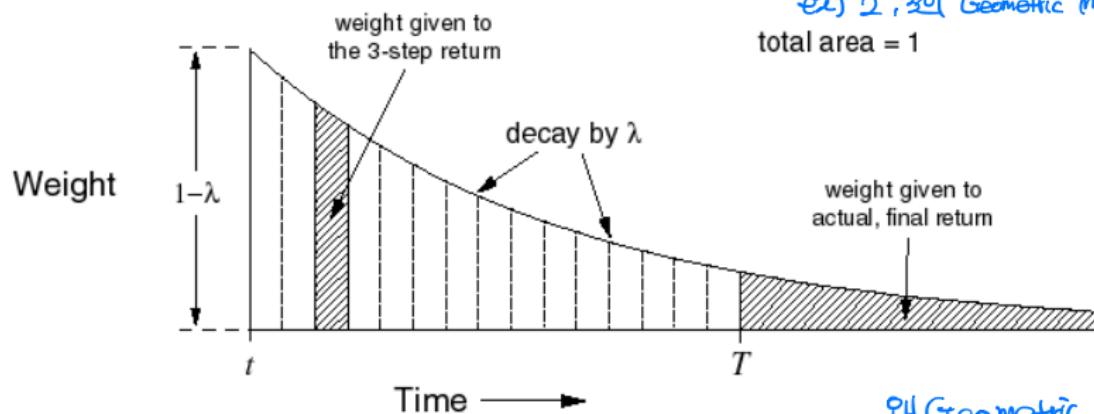
- The λ -return G_t^λ combines all n -step returns $G_t^{(n)}$
- Using weight $(1 - \lambda)\lambda^{n-1}$

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

- Forward-view TD(λ)

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

TD column의 G_t^λ 를 다 더하는데 $(1 - \lambda)\lambda^{n-1}$ 가중치를 사용하는 방법.

TD(λ) Weighting Function

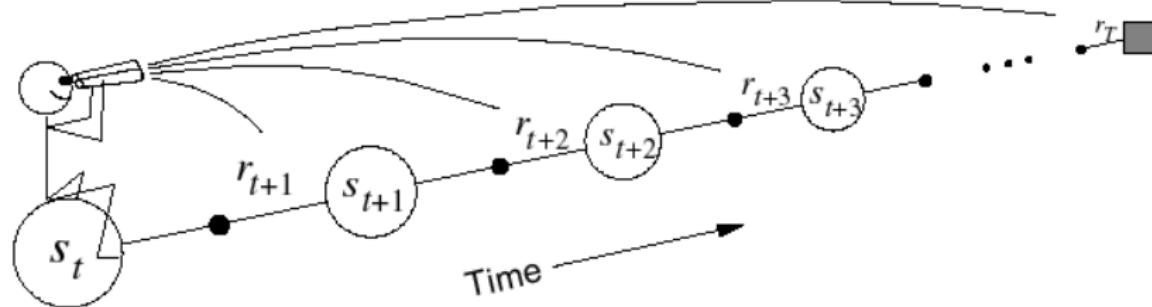
Geometric mean
 $= (a \times b)^{\frac{1}{2}}, (a \cdot b \cdot c)^{\frac{1}{3}}$
 ex) 2, 3의 Geometric mean은 $\sqrt{2 \cdot 3}$

왜 Geometric mean 사용?

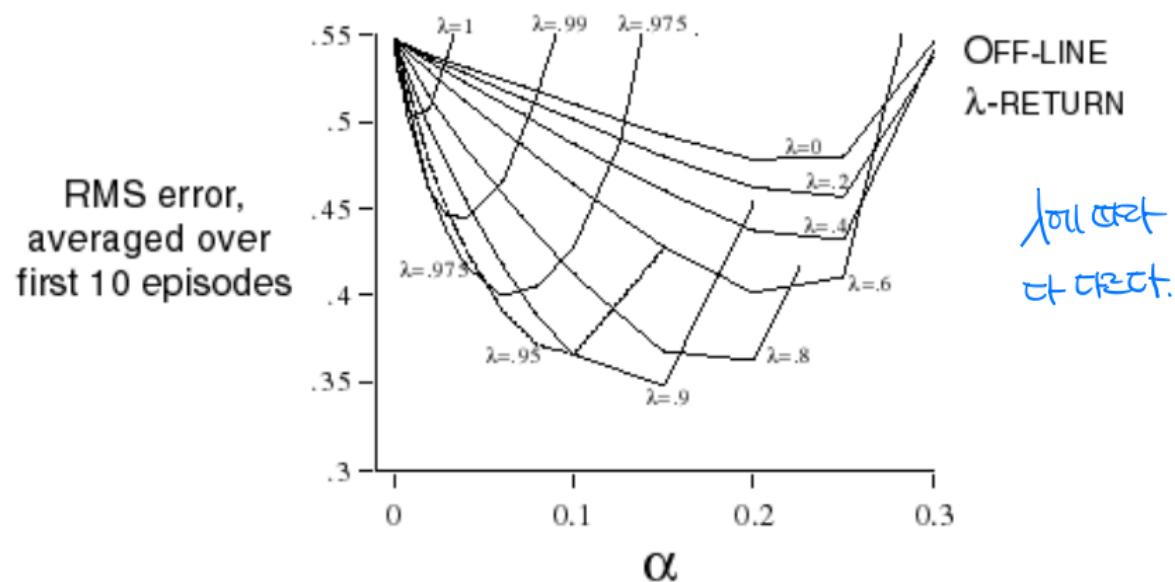
기후 예측 분석 분야.

이전 것 기억 안하고 계산 가능.

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

Forward-view TD(λ)TD(λ)는 어떠한가 끝낼 수 있어야 가능

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

Forward-View TD(λ) on Large Random Walk

\vdash TD(λ) \vdash Backward View of TD(λ)Backward View TD(λ)TD(0) 짧음을 가진 TD(λ)

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

↓
작업하면서 움직인다.

offline은 epi 끝나면 학습 후 움직임.

Eligibility Traces

어떤 시점에 책임이 있는지에 update 했을 때 해준다.

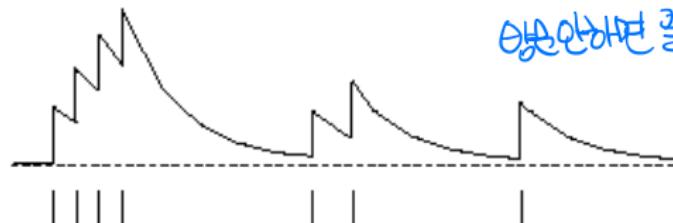


- Credit assignment problem: did bell or light cause shock?
- Frequency heuristic: assign credit to most frequent states
- Recency heuristic: assign credit to most recent states
- Eligibility traces combine both heuristics

$$E_0(s) = 0$$

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

이동하면서 점이면 하면 된다.



accumulating eligibility trace

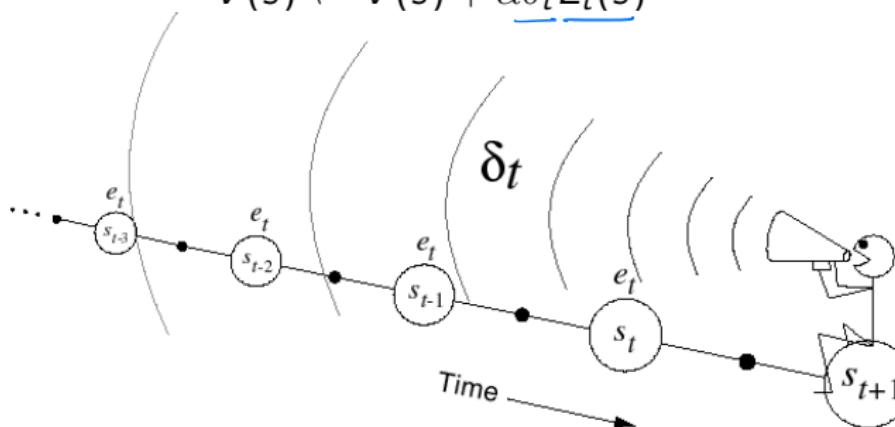
times of visits to a state

Backward View TD(λ)

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

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

$$V(s) \leftarrow V(s) + \alpha \underline{\delta_t} \underline{E_t(s)}$$



TD(λ) 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

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

TD(λ) 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(λ)

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

MC and TD(1)

- Consider an episode where s is visited once at time-step k ,
- TD(1) eligibility trace discounts time since visit,

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

- TD(1) updates accumulate error *online*

$$\sum_{t=1}^{T-1} \alpha \delta_t E_t(s) = \alpha \sum_{t=k}^{T-1} \gamma^{t-k} \delta_t = \alpha (G_k - V(S_k))$$

- By end of episode it accumulates total error

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

Telescoping in TD(1)

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

$$\begin{aligned} & \delta_t + \gamma\delta_{t+1} + \gamma^2\delta_{t+2} + \dots + \gamma^{T-1-t}\delta_{T-1} \\ &= R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \\ &\quad + \gamma R_{t+2} + \gamma^2 V(S_{t+2}) - \gamma V(S_{t+1}) \\ &\quad + \gamma^2 R_{t+3} + \gamma^3 V(S_{t+3}) - \gamma^2 V(S_{t+2}) \\ &\quad \vdots \\ &\quad + \gamma^{T-1-t} R_T + \gamma^{T-t} V(S_T) - \gamma^{T-1-t} V(S_{T-1}) \\ &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots + \gamma^{T-1-t} R_T - V(S_t) \\ &= G_t - V(S_t) \end{aligned}$$

TD(λ) 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(λ)

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

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

Forwards and Backwards TD(λ)

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

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

- Backward TD(λ) updates accumulate error *online*

$$\sum_{t=1}^T \alpha \delta_t E_t(s) = \alpha \sum_{t=k}^T (\gamma\lambda)^{t-k} \delta_t = \alpha \left(G_k^\lambda - V(S_k) \right)$$

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

Offline Equivalence of Forward and Backward TD

Offline updates

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

Online Equivalence of Forward and Backward TD

Online updates

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

Summary of Forward and Backward TD(λ)

Offline updates	$\lambda = 0$	$\lambda \in (0, 1)$	$\lambda = 1$
Backward view	TD(0) 	TD(λ) 	TD(1)
Forward view	TD(0)	Forward TD(λ)	MC
Online updates	$\lambda = 0$	$\lambda \in (0, 1)$	$\lambda = 1$
Backward view	TD(0) 	TD(λ) ※	TD(1) ※
Forward view	TD(0) 	Forward TD(λ) 	MC
Exact Online	TD(0)	Exact Online TD(λ)	Exact Online TD(1)

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